



Community Experience Distilled

Machine Learning in Java

Design, build, and deploy your own machine learning applications
by leveraging key Java machine learning libraries

Boštjan Kaluža

[PACKT] open source*
PUBLISHING community experience distilled

Machine Learning in Java

Design, build, and deploy your own machine learning applications by leveraging key Java machine learning libraries

Boštjan Kaluža



BIRMINGHAM - MUMBAI

Machine Learning in Java

Copyright © 2016 Packt Publishing

All rights reserved. No part of this book may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, without the prior written permission of the publisher, except in the case of brief quotations embedded in critical articles or reviews.

Every effort has been made in the preparation of this book to ensure the accuracy of the information presented. However, the information contained in this book is sold without warranty, either express or implied. Neither the author, nor Packt Publishing, and its dealers and distributors will be held liable for any damages caused or alleged to be caused directly or indirectly by this book.

Packt Publishing has endeavored to provide trademark information about all of the companies and products mentioned in this book by the appropriate use of capitals. However, Packt Publishing cannot guarantee the accuracy of this information.

First published: April 2016

Production reference: 1260416

Published by Packt Publishing Ltd.

Livery Place

35 Livery Street

Birmingham B3 2PB, UK.

ISBN 978-1-78439-658-9

www.packtpub.com

Credits

Author

Boštjan Kaluža

Project Coordinator

Izzat Contractor

Reviewers

Abhik Banerjee

Wei Di

Manjunath Narayana

Ravi Sharma

Proofreader

Safis Editing

Indexer

Mariammal Chettiyar

Commissioning Editor

Amarabha Banerjee

Graphics

Disha Haria

Acquisition Editor

Aaron Lazar

Production Coordinator

Nilesh Mohite

Content Development Editor

Rohit Singh

Cover Work

Nilesh Mohite

Technical Editor

Suvarna Patil

Copy Editor

Vibha Shukla

About the Author

Boštjan Kaluža, PhD, is a researcher in artificial intelligence and machine learning. Boštjan is the chief data scientist at Evolgen, a leading IT operations analytics company, focusing on configuration and change management. He works with machine learning, predictive analytics, pattern mining, and anomaly detection to turn data into understandable relevant information and actionable insight.

Prior to Evolgen, Boštjan served as a senior researcher in the department of intelligent systems at the Jozef Stefan Institute, a leading Slovenian scientific research institution, and led research projects involving pattern and anomaly detection, ubiquitous computing, and multi-agent systems. Boštjan was also a visiting researcher at the University of Southern California, where he studied suspicious and anomalous agent behavior in the context of security applications. Boštjan has extensive experience in Java and Python, and he also lectures on Weka in the classroom.

Focusing on machine learning and data science, Boštjan has published numerous articles in professional journals, delivered conference papers, and authored or contributed to a number of patents. In 2013, Boštjan published his first book on data science, *Instant Weka How-to*, Packt Publishing, exploring how to leverage machine learning using Weka. Learn more about him at <http://bostjankaluza.net>.

About the Reviewers

Abhik Banerjee has been a great data science leader, leading teams comprising of data scientists and engineers. He completed his masters from the University of Cincinnati, where his research focused on various data mining techniques related to itemset mining and biclustering techniques applied to biomedical informatics datasets. He has been working in the areas of machine learning and data mining in the industry for the past 7-8 years, solving various problems related to supervised learning (classification and regression techniques, such as SVM, Bayes net, GBM, GLM, neural networks, deep nets, and so on), unsupervised learning (clustering, blustering, LDA, and so on), and various NLP techniques. He had been working on how these various techniques can be applied to e-mail, biomedical informatics, and retail domains in order to understand the customer better and improve their experience.

Abhik has a strong acumen of problem solving skills, spanning various technological solutions and architectures, such as Hadoop, MapReduce, Spark, Java, Python, machine learning, data mining, NLP, and so on.

Wei Di is a data scientist. She is passionate about creating smart and scalable analytics and data mining solutions that can impact millions of individuals and empower successful business.

Her interests cover wide areas, including artificial intelligence, machine learning, and computer vision. She was previously associated with eBay Human Language Technology team and *eBay Research Labs*, with a focus on image understanding for large-scale applications and joint learning from both visual and text information. Prior to this, she was with *Ancestry.com*, working on large-scale data mining and machine learning models in the areas of record linkage, search relevance, and ranking. She received her PhD from Purdue University in 2011, focusing on data mining and image classification.

Manjunath Narayana received his PhD in computer science from the University of Massachusetts, Amherst, in 2014. He obtained his MS degree in computer engineering from the University of Kansas in 2007 and his BE degree in electronics and communications engineering from B. M. S. College of Engineering, Bangalore, India, in 2004. He is currently a robotics scientist at iRobot Corporation, USA, developing algorithms for consumer robots. Prior to iRobot, he was a research engineer at *metaio, Inc.*, working on computer vision research for augmented reality applications and 3D reconstruction. He has worked in the Computer Vision Systems Toolbox group in The *MathWorks, Inc.*, developing object detection algorithms. His research interests include machine learning, robotics, computer vision, deep learning, and augmented reality. His research has been published at top conferences such as CVPR, ICCV, and BMVC.

Ravi Sharma is a lead data scientist and has expertise in both artificial intelligence and natural language processing. He is currently leading the data science research team at *Msg.ai Inc.*, his commercial applications of data science include developing artificial chat bots for CPG brands, health care industry and entertainment industry. He has designed data collection systems and other strategies that optimize statistical efficiency and data quality. He has implemented a corporate big data-based data warehouse systems and distributed algorithms for high traffic. His areas of interest comprises the big data management platform, feature engineering, model building and tuning, exploratory data analysis, pattern analysis, outlier detection, collaborative filtering algorithms to provide recommendations and text analysis using NLP.

www.PacktPub.com

eBooks, discount offers, and more

Did you know that Packt offers eBook versions of every book published, with PDF and ePub files available? You can upgrade to the eBook version at www.PacktPub.com and as a print book customer, you are entitled to a discount on the eBook copy. Get in touch with us at customer@packtpub.com for more details.

At www.PacktPub.com, you can also read a collection of free technical articles, sign up for a range of free newsletters and receive exclusive discounts and offers on Packt books and eBooks.



<https://www2.packtpub.com/books/subscription/packtlib>

Do you need instant solutions to your IT questions? PacktLib is Packt's online digital book library. Here, you can search, access, and read Packt's entire library of books.

Why subscribe?

- Fully searchable across every book published by Packt
- Copy and paste, print, and bookmark content
- On demand and accessible via a web browser

Table of Contents

Preface	ix
Chapter 1: Applied Machine Learning Quick Start	1
Machine learning and data science	1
What kind of problems can machine learning solve?	2
Applied machine learning workflow	3
Data and problem definition	4
Measurement scales	5
Data collection	6
Find or observe data	7
Generate data	8
Sampling traps	9
Data pre-processing	9
Data cleaning	9
Fill missing values	10
Remove outliers	11
Data transformation	11
Data reduction	12
Unsupervised learning	13
Find similar items	13
Euclidean distances	13
Non-Euclidean distances	14
The curse of dimensionality	15
Clustering	16
Supervised learning	17
Classification	17
Decision tree learning	18
Probabilistic classifiers	18
Kernel methods	18
Artificial neural networks	18
Ensemble learning	19

Evaluating classification	19
Regression	21
Linear regression	22
Evaluating regression	22
Generalization and evaluation	24
Underfitting and overfitting	24
Train and test sets	26
Cross-validation	26
Leave-one-out validation	26
Stratification	27
Summary	27
Chapter 2: Java Libraries and Platforms for Machine Learning	29
The need for Java	30
Machine learning libraries	30
Weka	30
Java machine learning	34
Apache Mahout	35
Apache Spark	36
Deeplearning4j	38
MALLET	39
Comparing libraries	41
Building a machine learning application	42
Traditional machine learning architecture	42
Dealing with big data	43
Big data application architecture	43
Summary	44
Chapter 3: Basic Algorithms – Classification, Regression, and Clustering	45
Before you start	46
Classification	46
Data	47
Loading data	48
Feature selection	49
Learning algorithms	50
Classify new data	53
Evaluation and prediction error metrics	54
Confusion matrix	54
Choosing a classification algorithm	55
Regression	56
Loading the data	56
Analyzing attributes	58

Building and evaluating regression model	59
Linear regression	59
Regression trees	60
Tips to avoid common regression problems	63
Clustering	63
Clustering algorithms	64
Evaluation	66
Summary	66
Chapter 4: Customer Relationship Prediction with Ensembles	67
Customer relationship database	68
Challenge	68
Dataset	69
Evaluation	71
Basic naive Bayes classifier baseline	71
Getting the data	72
Loading the data	73
Basic modeling	75
Evaluating models	75
Implementing naive Bayes baseline	76
Advanced modeling with ensembles	77
Before we start	78
Data pre-processing	79
Attribute selection	80
Model selection	81
Performance evaluation	85
Summary	86
Chapter 5: Affinity Analysis	87
Market basket analysis	87
Affinity analysis	89
Association rule learning	90
Basic concepts	90
Database of transactions	90
Itemset and rule	91
Support	92
Confidence	92
Apriori algorithm	92
FP-growth algorithm	93
The supermarket dataset	94
Discover patterns	94
Apriori	94
FP-growth	96

Other applications in various areas	97
Medical diagnosis	97
Protein sequences	97
Census data	98
Customer relationship management	98
IT Operations Analytics	99
Summary	99
Chapter 6: Recommendation Engine with Apache Mahout	101
Basic concepts	101
Key concepts	102
User-based and item-based analysis	103
Approaches to calculate similarity	103
Collaborative filtering	103
Content-based filtering	104
Hybrid approach	104
Exploitation versus exploration	105
Getting Apache Mahout	105
Configuring Mahout in Eclipse with the Maven plugin	106
Building a recommendation engine	108
Book ratings dataset	108
Loading the data	109
Loading data from file	109
Loading data from database	112
In-memory database	113
Collaborative filtering	114
User-based filtering	115
Item-based filtering	118
Adding custom rules to recommendations	119
Evaluation	120
Online learning engine	121
Content-based filtering	123
Summary	124
Chapter 7: Fraud and Anomaly Detection	125
Suspicious and anomalous behavior detection	126
Unknown-unknowns	126
Suspicious pattern detection	127
Anomalous pattern detection	128
Analysis types	128
Pattern analysis	128
Transaction analysis	129
Plan recognition	129

Fraud detection of insurance claims	129
Dataset	130
Modeling suspicious patterns	132
Vanilla approach	133
Dataset rebalancing	134
Anomaly detection in website traffic	137
Dataset	137
Anomaly detection in time series data	138
Histogram-based anomaly detection	138
Loading the data	140
Creating histograms	141
Density based k-nearest neighbors	142
Summary	144
Chapter 8: Image Recognition with Deeplearning4j	145
<hr/>	
Introducing image recognition	145
Neural networks	147
Perceptron	147
Feedforward neural networks	148
Autoencoder	149
Restricted Boltzmann machine	150
Deep convolutional networks	151
Image classification	153
Deeplearning4j	153
Getting DL4J	153
MNIST dataset	154
Loading the data	154
Building models	155
Building a single-layer regression model	156
Building a deep belief network	158
Build a Multilayer Convolutional Network	160
Summary	163
Chapter 9: Activity Recognition with Mobile Phone Sensors	165
<hr/>	
Introducing activity recognition	166
Mobile phone sensors	166
Activity recognition pipeline	168
The plan	169
Collecting data from a mobile phone	170
Installing Android Studio	170
Loading the data collector	172
Feature extraction	174
Collecting training data	175

Building a classifier	178
Reducing spurious transitions	180
Plugging the classifier into a mobile app	182
Summary	184
Chapter 10: Text Mining with Mallet – Topic Modeling and Spam Detection	185
Introducing text mining	185
Topic modeling	186
Text classification	187
Installing Mallet	188
Working with text data	190
Importing data	192
Importing from directory	192
Importing from file	193
Pre-processing text data	194
Topic modeling for BBC news	196
BBC dataset	196
Modeling	197
Evaluating a model	200
Reusing a model	202
Saving a model	202
Restoring a model	203
E-mail spam detection	203
E-mail spam dataset	204
Feature generation	205
Training and testing	206
Model performance	207
Summary	208
Chapter 11: What is Next?	209
Machine learning in real life	209
Noisy data	210
Class unbalance	210
Feature selection is hard	210
Model chaining	211
Importance of evaluation	211
Getting models into production	212
Model maintenance	212
Standards and markup languages	213
CRISP-DM	213
SEMMA methodology	214
Predictive Model Markup Language	215

Machine learning in the cloud	215
Machine learning as a service	216
Web resources and competitions	217
Datasets	217
Online courses	218
Competitions	219
Websites and blogs	219
Venues and conferences	220
Summary	220
Appendix: References	221
Index	225

Preface

Machine learning is a subfield of artificial intelligence. It helps computers to learn and act like human beings with the help of algorithms and data. With a given set of data, an ML algorithm learns different properties of the data and infers the properties of the data that it may encounter in future.

This book will teach the readers how to create and implement machine learning algorithms in Java by providing fundamental concepts as well as practical examples. In this process, it will also talk about some machine learning libraries that are frequently used, such as Weka, Apache Mahout, Mallet, and so on. This book will help the user to select appropriate approaches for particular problems and compare and evaluate the results of different techniques. This book will also cover performance improvement techniques, including input preprocessing and combining output from different methods.

Without shying away from the technical details, you will explore machine learning with Java libraries using clear and practical examples. You will also explore how to prepare data for analysis, choose a machine learning method, and measure the success of the process.

What this book covers

Chapter 1, Applied Machine Learning Quick Start, introduces the basics of machine learning, laying down the common concepts, machine learning principles, and applied machine learning workflow.

Chapter 2, Java Libraries and Platforms for Machine Learning, reviews the various Java libraries and platforms dedicated to machine learning, what each library brings to the table, and what kind of problems it is able to solve. The review includes Weka, Java-ML, Apache Mahout, Apache Spark, deeplearning4j, and Mallet.

Chapter 3, Basic Algorithms – Classification, Regression, and Clustering, starts with basic machine learning tasks, introducing the key algorithms for classification, regression, and clustering, using small, easy-to-understand datasets.

Chapter 4, Customer Relationship Prediction with Ensembles, dives into a real-world marketing database, where the task is to predict the customer that will churn, upsell, and cross-sell. The problem is attacked with ensemble methods, following the steps of KDD Cup-winning solution.

Chapter 5, Affinity Analysis, discusses how to analyze co-occurrence relationships using association rule mining. We will look into market basket analysis to understand the purchasing behavior of customers and discuss applications of the approach to other domains.

Chapter 6, Recommendation Engine with Apache Mahout, explains the basic concepts required to understand recommendation engine principles, followed by two applications leveraging Apache Mahout to build content-based filtering and collaborative recommender.

Chapter 7, Fraud and Anomaly Detection, introduces the background to anomalous and suspicious pattern detection, followed by two practical applications on detecting frauds in insurance claims and detecting anomalies in website traffic.

Chapter 8, Image Recognition with Deeplearning4j, introduces image recognition and reviews fundamental neural network architectures. We will then discuss how to implement various deep learning architectures with deeplearning4j library to recognize handwritten digits.

Chapter 9, Activity Recognition with Mobile Phone Sensors, tackles the problem of recognizing patterns from sensor data. This chapter introduces the activity recognition process, explains how to collect data with an Android device, and presents a classification model to recognize activities of daily living.

Chapter 10, Text Mining with Mallet – Topic Modeling and Spam Detection, explains the basics of text mining, introduces the text processing pipeline, and shows how to apply this to two real-world problems: topic modeling and document classification.

Chapter 11, What is Next?, concludes the book with practical advice about how to deploy models and gives you further pointers about where to find additional resources, materials, venues, and technologies to dive deeper into machine learning.

What you need for this book

To follow the examples throughout the book, you'll need a personal computer with the JDK installed. All the examples and source code that you can download assume Eclipse IDE with support for Maven, a dependency management and build automation tool; and Git, a version control system. Examples in the chapters rely on various libraries, including Weka, `deeplearning4j`, Mallet, and Apache Mahout. Instructions on how to get and install the libraries are provided in the chapter where the library will be first used.

Who this book is for

The book is intended for those who want to learn how to use Java's machine learning libraries to gain insights from your data. Perhaps you already know a bit about machine learning, but have never used Java; or perhaps you know a little Java, but are new to machine learning. In either case, this book will get you up and running quickly, providing you with the skills that you need to successfully create, customize, and deploy machine learning applications in real life. It would be helpful to be a little familiar with basic programming and data mining concepts, but no prior experience with data mining packages is necessary.

Supporting materials

The book has a dedicated web site, <http://machine-learning-in-java.com>, where you can find all the example code, errata, and additional materials that will help you to get started.

Conventions

In this book, you will find a number of text styles that distinguish between different kinds of information. Here are some examples of these styles and an explanation of their meaning.

Code words in text, database table names, folder names, filenames, file extensions, pathnames, dummy URLs, user input, and Twitter handles are shown as follows: "For example, Bob has attributes named `height`, `eye color`, and `hobbies` with values `185cm`, `blue`, `climbing`, and `sky diving`, respectively."



A block of code is set as follows:



```
Bob = {  
  height: 185cm,  
  eye color: blue,  
  hobbies: climbing, sky diving  
}
```

Any command-line input or output is written as follows:

```
12,3,7,2,0,1,8,9,13,4,11,5,15,10,6,14,16
```

New terms and **important words** are shown in bold. Words that you see on the screen, for example, in menus or dialog boxes, appear in the text like this: "Right-click on the project properties, select **Java Build Path**, click on the **Libraries** tab, and select **Add External JARs**."

 Warnings or important notes appear in a box like this. 

 Tips and tricks appear like this. 

Reader feedback

Feedback from our readers is always welcome. Let us know what you think about this book – what you liked or disliked. Reader feedback is important for us as it helps us develop titles that you will really get the most out of.

To send us general feedback, simply e-mail feedback@packtpub.com, and mention the book's title in the subject of your message.

If there is a topic that you have expertise in and you are interested in either writing or contributing to a book, see our author guide at www.packtpub.com/authors.

Customer support

Now that you are the proud owner of a Packt book, we have a number of things to help you to get the most from your purchase.

Downloading the example code

You can download the example code files from the supplementary web page <http://machine-learning-in-java.com>. Navigate to the Downloads section and follow the link to a Git repository.

Optionally, you can also download the example code from your account at <http://www.packtpub.com>. If you purchased this book elsewhere, you can visit <http://www.packtpub.com/support> and register to have the files e-mailed directly to you.

You can download the code files by following these steps:

1. Log in or register to our website using your e-mail address and password.
2. Hover the mouse pointer on the **SUPPORT** tab at the top.
3. Click on **Code Downloads & Errata**.
4. Enter the name of the book in the **Search** box.
5. Select the book for which you're looking to download the code files.
6. Choose from the drop-down menu where you purchased this book from.
7. Click on **Code Download**.

Once the file is downloaded, please make sure that you unzip or extract the folder using the latest version of:

- WinRAR / 7-Zip for Windows
- Zipeg / iZip / UnRarX for Mac
- 7-Zip / PeaZip for Linux

Errata

Although we have taken every care to ensure the accuracy of our content, mistakes do happen. If you find a mistake in one of our books – maybe a mistake in the text or the code – we would be grateful if you could report this to us. By doing so, you can save other readers from frustration and help us improve subsequent versions of this book. If you find any errata, please report them by visiting <http://www.packtpub.com/submit-errata>, selecting your book, clicking on the **Errata Submission Form** link, and entering the details of your errata. Once your errata are verified, your submission will be accepted and the errata will be uploaded to our website or added to any list of existing errata under the Errata section of that title.

To view the previously submitted errata, go to <https://www.packtpub.com/books/content/support> and enter the name of the book in the search field. The required information will appear under the **Errata** section.

Piracy

Piracy of copyrighted material on the Internet is an ongoing problem across all media. At Packt, we take the protection of our copyright and licenses very seriously. If you come across any illegal copies of our works in any form on the Internet, please provide us with the location address or website name immediately so that we can pursue a remedy.

Please contact us at copyright@packtpub.com with a link to the suspected pirated material.

We appreciate your help in protecting our authors and our ability to bring you valuable content.

Questions

If you have a problem with any aspect of this book, you can contact us at questions@packtpub.com, and we will do our best to address the problem.

1

Applied Machine Learning Quick Start

This chapter introduces the basics of machine learning, laying down the common themes and concepts and making it easy to follow the logic and familiarize yourself with the topic. The goal is to quickly learn the step-by-step process of **applied machine learning** and grasp the main machine learning principles. In this chapter, we will cover the following topics:

- Introducing machine learning and its relation to data science
- Discussing the basic steps in applied machine learning
- Discussing the kind of data we are dealing with and its importance
- Discussing approaches of collecting and preprocessing the data
- Making sense of data using machine learning
- Using machine learning to extract insights from data and build predictors

If you are already familiar with machine learning and are eager to start coding, then quickly jump to the following chapters. However, if you need to refresh your memory or clarify some concepts, then it is strongly recommend to revisit the topics presented in this chapter.

Machine learning and data science

Nowadays, everyone talks about machine learning and data science. So, what exactly is machine learning anyway? How does it relate to data science? These two terms are commonly confused, as they often employ the same methods and overlap significantly. Therefore, let's first clarify what they are.

Josh Wills tweeted:

"Data scientist is a person who is better at statistics than any software engineer and better at software engineering than any statistician".

– (Josh Wills)

More specifically, data science encompasses the entire process of obtaining knowledge from data by integrating methods from statistics, computer science, and other fields to gain insight from data. In practice, data science encompasses an iterative process of data harvesting, cleaning, analysis and visualization, and deployment.

Machine learning, on the other hand, is mainly concerned with fairly generic algorithms and techniques that are used in analysis and modeling phases of data science process. Arthur Samuel proposed the following definition back in 1995:

"Machine Learning relates with the study, design and development of the algorithms that give computers the capability to learn without being explicitly programmed."

– Arthur Samuel

What kind of problems can machine learning solve?

Among the different machine learning approaches, there are three main ways of learning, as shown in the following list:

- Supervised learning
- Unsupervised learning
- Reinforcement learning

Given a set of example inputs, X , and their outcomes, Y , supervised learning aims to learn a general mapping function, f , that transforms inputs to outputs, as $f: X \rightarrow Y$

An example of supervised learning is credit card fraud detection, where the learning algorithm is presented with credit card transactions (matrix X) marked as normal or suspicious. (vector Y). The learning algorithm produces a decision model that marks unseen transactions as normal or suspicious (that is the f function).

In contrast, unsupervised learning algorithms do not assume given outcome labels, Y as they focus on learning the structure of the data, such as grouping similar inputs into clusters. Unsupervised learning can, hence, discover hidden patterns in the data. An example of unsupervised learning is an item-based recommendation system, where the learning algorithm discovers similar items bought together, for example, *people who bought book A also bought book B*.

Reinforcement learning addresses the learning process from a completely different angle. It assumes that an agent, which can be a robot, bot, or computer program, interacts with a dynamic environment to achieve a specific goal. The environment is described with a set of states and the agent can take different actions to move from one state to another. Some states are marked as goal states and if the agent achieves this state, it receives a large reward. In other states, the reward is smaller, non-existing, or even negative. The goal of reinforcement learning is to find an optimal policy, that is, a mapping function that specifies the action to take in each of the states without a teacher explicitly telling whether this leads to the goal state or not. An example of reinforcement learning is a program for driving a vehicle, where the states correspond to the driving conditions – for example, current speed, road segment information, surrounding traffic, speed limits, and obstacles on the road – and the actions can be driving maneuvers such as turn left or right, stop, accelerate, and continue. The learning algorithm produces a policy that specifies the action that is to be taken in specific configuration of driving conditions.

In this book, we will focus on supervised and unsupervised learning only, as they share many concepts. If reinforcement learning sparked your interest, a good book to start with is *Reinforcement Learning: An Introduction* by Richard S. Sutton and Andrew Barto.

Applied machine learning workflow

This book's emphasis is on applied machine learning. We want to provide you with the practical skills needed to get learning algorithms to work in different settings. Instead of math and theory of machine learning, we will spend more time on the practical, hands-on skills (and dirty tricks) to get this stuff to work well on an application. We will focus on supervised and unsupervised machine learning and cover the essential steps from data science to build the applied machine learning workflow.

A typical workflow in applied machine learning applications consists of answering a series of questions that can be summarized in the following five steps:

1. **Data and problem definition:** The first step is to ask interesting questions. What is the problem you are trying solve? Why is it important? Which format of result answers your question? Is this a simple yes/no answer? Do you need to pick one of the available questions?

2. **Data collection:** Once you have a problem to tackle, you will need the data. Ask yourself what kind of data will help you answer the question. Can you get the data from the available sources? Will you have to combine multiple sources? Do you have to generate the data? Are there any sampling biases? How much data will be required?
3. **Data preprocessing:** The first data preprocessing task is data cleaning. For example, filling missing values, smoothing noisy data, removing outliers, and resolving inconsistencies. This is usually followed by integration of multiple data sources and data transformation to a specific range (normalization), to value bins (discretized intervals), and to reduce the number of dimensions.
4. **Data analysis and modeling with unsupervised and supervised learning:** Data analysis and modeling includes unsupervised and supervised machine learning, statistical inference, and prediction. A wide variety of machine learning algorithms are available, including k-nearest neighbors, naïve Bayes, decision trees, support vector machines, logistic regression, k-means, and so on. The choice of method to be deployed depends on the problem definition discussed in the first step and the type of collected data. The final product of this step is a model inferred from the data.
5. **Evaluation:** The last step is devoted to model assessment. The main issue models built with machine learning face is how well they model the underlying data – if a model is too specific, that is, it **overfits** to the data used for training, it is quite possible that it will not perform well on a new data. The model can be too generic, meaning that it **underfits** the training data. For example, when asked how the weather is in California, it always answers sunny, which is indeed correct most of the time. However, such a model is not really useful for making valid predictions. The goal of this step is to correctly evaluate the model and make sure it will work on new data as well. Evaluation methods include separate test and train set, cross-validation, and leave-one-out validation.

In the following sections, we will take a closer look at each of the steps. We will try to understand the type of questions we must answer during the applied machine learning workflow and also look at the accompanying concepts of data analysis and evaluation.

Data and problem definition

Data is simply a collection of measurements in the form of numbers, words, measurements, observations, descriptions of things, images, and so on.

Measurement scales

The most common way to represent the data is using a set of attribute-value pairs. Consider the following example:

```
Bob = {
  height: 185cm,
  eye color: blue,
  hobbies: climbing, sky diving
}
```

For example, Bob has attributes named `height`, `eye color`, and `hobbies` with values 185cm, blue, climbing, sky diving, respectively.

A set of data can be simply presented as a table, where columns correspond to attributes or features and rows correspond to particular data examples or instances. In supervised machine learning, the attribute whose value we want to predict the outcome, Y , from the values of the other attributes, X , is denoted as class or the target variable, as follows:

Name	Height [cm]	Eye color	Hobbies
Bob	185.0	Blue	Climbing, sky diving
Anna	163.0	Brown	Reading
...

The first thing we notice is how varying the attribute values are. For instance, height is a number, eye color is text, and hobbies are a list. To gain a better understanding of the value types, let's take a closer look at the different types of data or measurement scales. Stevens (1946) defined the following four scales with increasingly more expressive properties:

- **Nominal data** are mutually exclusive, but not ordered. Their examples include eye color, marital status, type of car owned, and so on.
- **Ordinal data** correspond to categories where order matters, but not the difference between the values, such as pain level, student letter grade, service quality rating, IMDB movie rating, and so on.
- **Interval data** where the difference between two values is meaningful, but there is no concept of zero. For instance, standardized exam score, temperature in Fahrenheit, and so on.

- **Ratio data** has all the properties of an interval variable and also a clear definition of zero; when the variable equals to zero, there is none of this variable. Variables such as height, age, stock price, and weekly food spending are ratio variables.

Why should we care about measurement scales? Well, machine learning heavily depends on the statistical properties of the data; hence, we should be aware of the limitations each data type possesses. Some machine learning algorithms can only be applied to a subset of measurement scales.

The following table summarizes the main operations and statistics properties for each of the measurement types:

Property	Nominal	Ordinal	Interval	Ratio
Frequency of distribution	✓	✓	✓	✓
Mode and median		✓	✓	✓
Order of values is known		✓	✓	✓
Can quantify difference between each value			✓	✓
Can add or subtract values			✓	✓
Can multiply and divide values				✓
Has true zero				✓

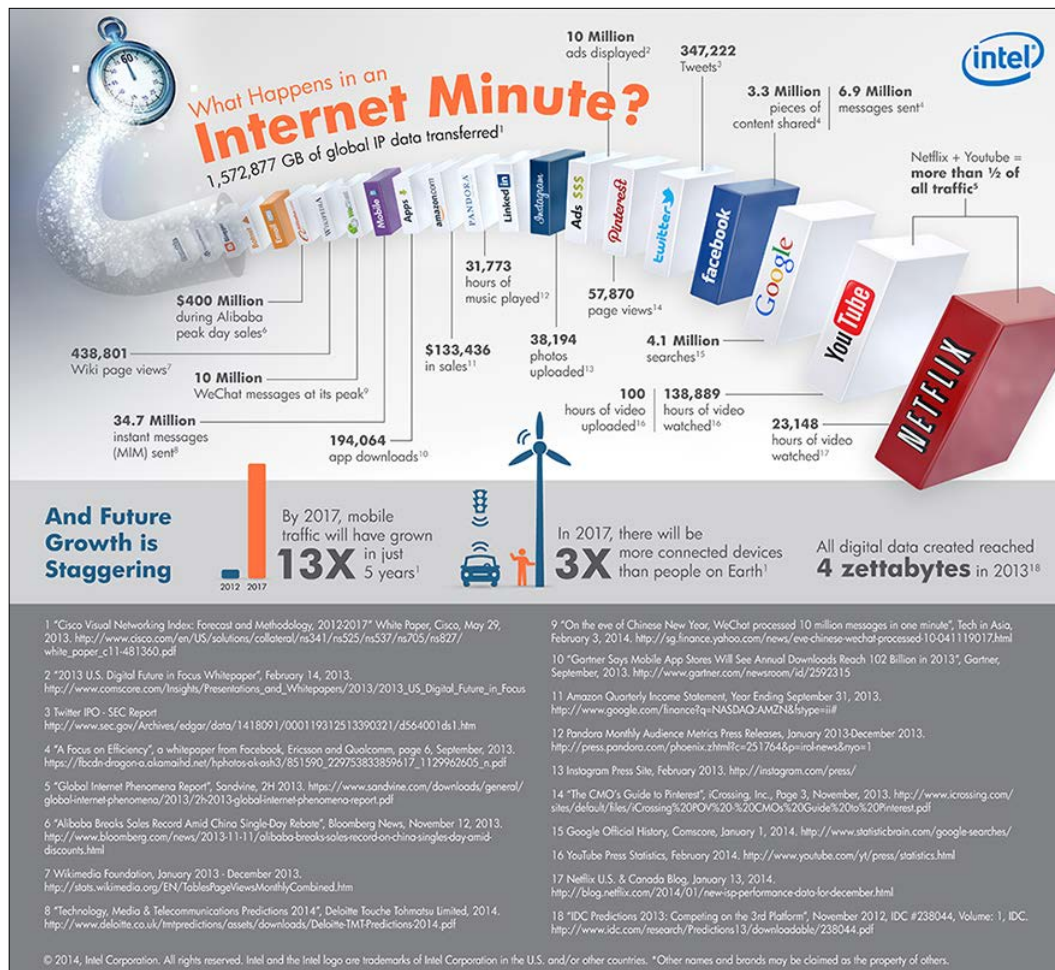
Furthermore, nominal and ordinal data correspond to discrete values, while interval and ratio data can correspond to continuous values as well. In supervised learning, the measurement scale of the attribute values that we want to predict dictates the kind of machine algorithm that can be used. For instance, predicting discrete values from a limited list is called classification and can be achieved using decision trees; while predicting continuous values is called regression, which can be achieved using model trees.

Data collection

So, where does the data come from? We have two choices: observe the data from existing sources or generate the data via surveys, simulations, and experiments. Let's take a closer look at both the approaches.

Find or observe data

Data can be found or observed at many places. An obvious data source is the Internet. Intel (2013) presented the following iconographic, showing the massive amount of data collected by different Internet services. In 2013, digital devices created four zettabytes (10^{21} = billion terabytes) of data. In 2017, it is expected that the number of connected devices will reach three times the number of people on earth; hence, the amount of data generated and collected will increase even further:



To get the data from the Internet, there are multiple options, as shown in the following:

- Bulk downloads from websites such as Wikipedia, IMDb, and Million Song database.
- Accessing the data through API (NY Times, Twitter, Facebook, Foursquare).
- Web scraping – It is OK to scrape public, non-sensitive, and anonymized data. Be sure to check terms of conditions and to fully reference information.

The main drawbacks of found data are that it takes time and space to accumulate the data; they cover only what happened, for instance, intentions, motivations, or internal motivations are not collected. Finally, such data might be noisy, incomplete, inconsistent, and may even change over time.

Another option is to collect measurements from sensors such as inertial and location sensors in mobile devices, environmental sensors, and software agents monitoring key performance indicators.

Generate data

An alternative approach is to generate the data by yourself, for example, with a survey. In survey design, we have to pay attention to data sampling, that is, who are the respondents answering the survey. We only get data from the respondents who are accessible and willing to respond. Also, respondents can provide answers that are in line with their self-image and researcher's expectations.

Next, the data can be collected with simulations, where a domain expert specifies behavior model of users at a micro level. For instance, crowd simulation requires specifying how different types of users will behave in crowd, for example, following the crowd, looking for an escape, and so on. The simulation can be then run under different conditions to see what happens (Tsai et al. 2011). Simulations are appropriate for studying macro phenomena and emergent behavior; however, they are typically hard to validate empirically.

Furthermore, you can design experiments to thoroughly cover all the possible outcomes, where you keep all the variables constant and only manipulate one variable at a time. This is the most costly approach, but usually provides the best quality of data.

Sampling traps

Data collection may involve many traps. To demonstrate one, let me share a story. There is supposed to be a global, unwritten rule for sending regular mail between students for free. If you write *student to student* to the place where the stamp should be, the mail is delivered to the recipient for free. Now, suppose Jacob sends a set of postcards to Emma, and given that Emma indeed receives some of the postcards, she concludes that all the postcards are delivered and that the rule indeed holds true. Emma reasons that as she received the postcards, all the postcards are delivered. However, she does not possess the information about the postcards that were sent by Jacob, but were undelivered; hence, she is unable to account this into her inference. What Emma experienced is **survivorship bias**, that is, she drew the conclusion based on the survived data only. For your information, the postcards that are being sent with *student to student* stamp get a circled black letter *T* stamp on them, which means *postage is due* and that receiver should pay it, including a small fine. However, mail services often have higher costs on applying such fee and hence do not do it (Magalhães, 2010).

Another example is a study, which found that the profession with the lowest average age of death was student. Being a student does not cause you to die at an early age, being a student means you are young. This is what makes the average of those that die so low (Gelman and Nolan, 2002).

Furthermore, a study that found that only 1.5% of drivers in accidents reported they were using a cell phone, whereas 10.9% reported another occupant in the car distracted them. Can we conclude that using a cell phone is safer than speaking with another occupant (Uts, 2003)? To answer this question, we need to know the prevalence of the cell phone use. It is likely that a higher number of people talked to another occupant in the car while driving than talking on the cell during the period when the data was collected.

Data pre-processing

The goal of data pre-processing tasks is to prepare the data for a machine learning algorithm in the best possible way as not all algorithms are capable of addressing issues with missing data, extra attributes, or denormalized values.

Data cleaning

Data cleaning, also known as data cleansing or data scrubbing, is the process of the following:

- Identifying inaccurate, incomplete, irrelevant, or corrupted data to remove it from further processing

- Parsing data, extracting information of interest, or validating whether a string of data is in an acceptable format
- Transforming data into a common encoding format, for example, utf-8 or int32, time scale, or normalized range
- Transforming data into a common data schema, for instance, if we collect temperature measurements from different types of sensors, we might want them to have the same structure

Now, let's look at some more concrete pre-processing steps.

Fill missing values

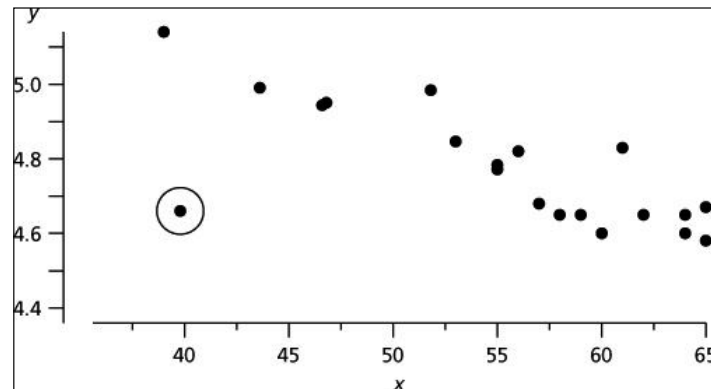
Machine learning algorithms generally do not work well with missing values. Rare exceptions include decision trees, naïve Bayes classifier, and some rule-based learners. It is very important to understand why a value is missing. It can be missing due to many reasons such as random error, systematic error, and sensor noise. Once we identified the reason, there are multiple ways to deal with the missing values, as shown in the following list:

- **Remove the instance:** If there is enough data, and only a couple of non-relevant instances have some missing values, then it is safe to remove these instances.
- **Remove the attribute:** Removing an attribute makes sense when most of the values are missing, values are constant, or attribute is strongly correlated with another attribute.
- **Assign a special value N/A:** Sometimes a value is missing due to valid reasons such as the value is out of scope discrete attribute value is not defined, or it is not possible to obtain or measure the value, which can be an indicator as well. For example, a person never rates a movie, so his rating on this movie is nonexistent.
- **Take the average attribute value:** In case we have a limited number of instances, we might not be able to afford removing instances or attributes. In that case, we can estimate the missing values, for example, by assigning the average attribute value or the average value over similar instances.
- **Predict the value from other attributes:** Predict the value from the previous entries if the attribute possesses time dependencies.

As we have seen, the value can be missing for many reasons, and hence, it is important to understand why the value is missing, absent, or corrupted.

Remove outliers

Outliers in data are values that are unlike any other values in the series and affect all learning methods to various degrees. These can be extreme values, which could be detected with confidence intervals and removed by threshold. The best approach is to visualize the data and inspect the visualization to detect irregularities. An example is shown in the following diagram. Visualization applies to low-dimensional data only:



Data transformation

Data transformation techniques tame the dataset to a format that a machine learning algorithm expects as an input, and may even help the algorithm to learn faster and achieve better performance. Standardization, for instance, assumes that data follows Gaussian distribution and transforms the values in such a way that the mean value is zero and the deviation is 1, as follows:

$$X = \frac{X - \text{mean}(X)}{\text{st.dev}}$$

Normalization, on the other hand, scales the values of attributes to a small, specified range, usually between 0 and 1:

$$X = \frac{X - \text{min}}{\text{max} - \text{min}}$$

Many machine learning toolboxes automatically normalize and standardize the data for you.

The last transformation technique is discretization, which divides the range of a continuous attribute into intervals. Why should we care? Some algorithms, such as decision trees and naïve Bayes prefer discrete attributes. The most common ways to select the intervals are shown in the following:

- **Equal width:** The interval of continuous variable is divided into k equal-width intervals
- **Equal frequency:** Suppose there are N instances, each of the k intervals contains approximately N/k instances
- **Min entropy:** The approach recursively splits the intervals until the entropy, which measures disorder, decreases more than the entropy increase, introduced by the interval split (Fayyad and Irani, 1993)

The first two methods require us to specify the number of intervals, while the last method sets the number of intervals automatically; however, it requires the class variable, which means, it won't work for unsupervised machine learning tasks.

Data reduction

Data reduction deals with abundant attributes and instances. The number of attributes corresponds to the number of dimensions in our dataset. Dimensions with low prediction power do not only contribute very little to the overall model, but also cause a lot of harm. For instance, an attribute with random values can introduce some random patterns that will be picked up by a machine learning algorithm.

To deal with this problem, the first set of techniques removes such attributes, or in other words, selects the most promising ones. This process is known as feature selection or attribute selection and includes methods such as ReliefF, information gain, and Gini index. These methods are mainly focused on discrete attributes.

Another set of tools, focused on continuous attributes, transforms the dataset from the original dimensions into a lower-dimensional space. For example, if we have a set of points in three-dimensional space, we can make a projection into a two-dimensional space. Some information is lost, but in case the third dimension is irrelevant, we don't lose much as the data structure and relationships are almost perfectly preserved. This can be performed by the following methods:

- **Singular value decomposition (SVD)**
- **Principal Component Analysis (PCA)**
- Neural nets auto encoders

The second problem in data reduction is related to too many instances; for example, they can be duplicates or coming from a very frequent data stream. The main idea is to select a subset of instances in such a way that distribution of the selected data still resembles the original data distribution, and more importantly, the observed process. Techniques to reduce the number of instances involve random data sampling, stratification, and others. Once the data is prepared, we can start with the data analysis and modeling.

Unsupervised learning

Unsupervised learning is about analyzing the data and discovering hidden structures in unlabeled data. As no notion of the right labels is given, there is also no error measure to evaluate a learned model; however, unsupervised learning is an extremely powerful tool. Have you ever wondered how Amazon can predict what books you'll like? How Netflix knows what you want to watch before you do? The answer can be found in unsupervised learning. The following is one such example.

Find similar items

Many problems can be formulated as finding similar sets of elements, for example, customers who purchased similar products, web pages with similar content, images with similar objects, users who visited similar websites, and so on.

Two items are considered similar if they are a *small distance* apart. The main questions are how each item is represented and how is the distance between the items defined. There are two main classes of distance measures: Euclidean distances and non-Euclidean distances.

Euclidean distances

In the Euclidean space, with the n dimension, the distance between two elements is based on the locations of the elements in such a space, which is expressed as **p-norm distance**. Two commonly used distance measures are L_2 and L_1 norm distances.

L_2 norm, also known as Euclidean distance, is the most frequently applied distance measure that measures how far apart two items in a two-dimensional space are. It is calculated as a square root of the sum of the squares of the differences between elements a and b in each dimension, as follows:

$$L_2 \text{ norm } d(a, b) = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}$$

L1 norm, also known as Manhattan distance, city block distance, and taxicab norm, simply sums the absolute differences in each dimension, as follows:

$$L_1 \text{ norm } d(a, b) = \sum_{i=1}^n |a_i - b_i|$$

Non-Euclidean distances

A non-Euclidean distance is based on the properties of the elements, but not on their location in space. Some well-known distances are Jaccard distance, cosine distance, edit distance, and Hamming distance.

Jaccard distance is used to compute the distance between two sets. First, we compute the Jaccard similarity of two sets as the size of their intersection divided by the size of their union, as follows:

$$\text{sim}(A, B) = \frac{A \cap B}{A \cup B}$$

The Jaccard distance is then defined as 1 minus Jaccard similarity, as shown in the following:

$$d(A, B) = 1 - \text{sim}(A, B) = 1 - \frac{A \cap B}{A \cup B}$$

Cosine distance between two vectors focuses on the orientation and not magnitude, therefore, two vectors with the same orientation have cosine similarity 1, while two perpendicular vectors have cosine similarity 0. Suppose we have two multidimensional points, think of a point as a vector from origin $(0, 0, \dots, 0)$ to its location. Two vectors make an angle, whose cosine distance is a normalized dot-product of the vectors, as follows:

$$d(A, B) = \arccos \frac{A \cdot B}{\|A\| \|B\|}$$

Cosine distance is commonly used in a high-dimensional feature space; for instance, in text mining, where a text document represents an instance, features that correspond to different words, and their values corresponds to the number of times the word appears in the document. By computing cosine similarity, we can measure how likely two documents match in describing similar content.

Edit distance makes sense when we compare two strings. The distance between the $a=a_1a_2a_3\dots a_n$ and $b=b_1b_2b_3\dots b_n$ strings is the smallest number of the insert/delete operation of single characters required to convert the string from a to b . For example, $a = abcd$ and $b = abbd$. To convert a to b , we have to delete the second b and insert c in its place. No smallest number of operations would convert a to b , thus the distance is $d(a, b)=2$.

Hamming distance compares two vectors of the same size and counts the number of dimensions in which they differ. In other words, it measures the number of substitutions required to convert one vector into another.

There are many distance measures focusing on various properties, for instance, correlation measures the linear relationship between two elements: **Mahalanobis distance** that measures the distance between a point and distribution of other points and **SimRank**, which is based on graph theory, measures similarity of the structure in which elements occur, and so on. As you can already imagine selecting and designing the right similarity measure for your problem is more than half of the battle. An impressive overview and evaluation of similarity measures is collected in *Chapter 2, Similarity and Dissimilarity Measures* in the book *Image Registration: Principles, Tools and Methods* by A. A. Goshtasby (2012).

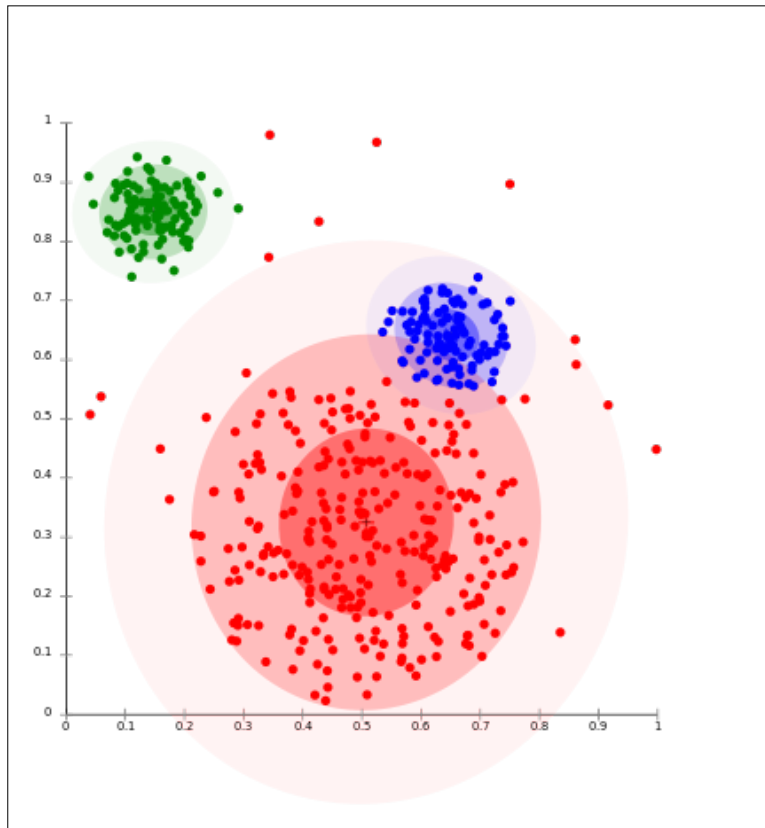
The curse of dimensionality

The curse of dimensionality refers to a situation where we have a large number of features, often hundreds or thousands, which lead to an extremely large space with sparse data and, consequently, to distance anomalies. For instance, in high dimensions, almost all pairs of points are equally distant from each other; in fact, almost all the pairs have distance close to the average distance. Another manifestation of the curse is that any two vectors are almost orthogonal, which means all the angles are close to 90 degrees. This practically makes any distance measure useless.

A cure for the curse of dimensionality might be found in one of the data reduction techniques, where we want to reduce the number of features; for instance, we can run a feature selection algorithm such as ReliefF or feature extraction/reduction algorithm such as PCA.

Clustering

Clustering is a technique for grouping similar instances into clusters according to some distance measure. The main idea is to put instances that are similar (that is, close to each other) into the same cluster, while keeping the dissimilar points (that is, the ones further apart from each other) in different clusters. An example of how clusters might look is shown in the following diagram:



The clustering algorithms follow two fundamentally different approaches. The first is a **hierarchical** or **agglomerative** approach that first considers each point as its own cluster, and then iteratively merges the most similar clusters together. It stops when further merging reaches a predefined number of clusters or if the clusters to be merged are spread over a large region.

The other approach is based on point assignment. First, initial cluster centers (that is, centroids) are estimated – for instance, randomly – and then, each point is assigned to the closest cluster, until all the points are assigned. The most well-known algorithm in this group is **k-means clustering**.

The k-means clustering picks initial cluster centers either as points that are as far as possible from one another or (hierarchically) clusters a sample of data and picks a point that is the closest to the center of each of the k clusters.

Supervised learning

Supervised learning is the key concept behind amazing things such as voice recognition, e-mail spam filtering, face recognition in photos, and detecting credit card frauds. More formally, given a set D of learning examples described with features, X , the goal of supervised learning is to find a function that predicts a target variable, Y . The function f that describes the relation between features X and class Y is called a model:

$$f(X) \rightarrow Y$$

The general structure of supervised learning algorithms is defined by the following decisions (Hand et al., 2001):

1. Define the task.
2. Decide on the machine learning algorithm, which introduces specific inductive bias, that is, apriori assumptions that it makes regarding the target concept.
3. Decide on the **score** or **cost** function, for instance, information gain, root mean square error, and so on.
4. Decide on the optimization/search method to optimize the score function.
5. Find a function that describes the relation between X and Y .

Many decisions are already made for us by the type of the task and dataset that we have. In the following sections, we will take a closer look at the classification and regression methods and the corresponding score functions.

Classification

Classification can be applied when we deal with a discrete class, and the goal is to predict one of the mutually-exclusive values in the target variable. An example would be credit scoring, where the final prediction is whether the person is credit liable or not. The most popular algorithms include decision tree, naïve Bayes classifier, support vector machines, neural networks, and ensemble methods.

Decision tree learning

Decision tree learning builds a classification tree, where each node corresponds to one of the attributes, edges correspond to a possible value (or intervals) of the attribute from which the node originates, and each leaf corresponds to a class label. A decision tree can be used to visually and explicitly represent the prediction model, which makes it a very transparent (white box) classifier. Notable algorithms are ID3 and C4.5, although many alternative implementations and improvements (for example, J48 in Weka) exist.

Probabilistic classifiers

Given a set of attribute values, a probabilistic classifier is able to predict a distribution over a set of classes, rather than an exact class. This can be used as a degree of certainty, that is, how sure the classifier is in its prediction. The most basic classifier is naïve Bayes, which happens to be the optimal classifier if, and only if, the attributes are conditionally independent. Unfortunately, this is extremely rare in practice.

There is really an enormous subfield denoted as probabilistic graphical models, comprising of hundreds of algorithms; for example, Bayesian network, dynamic Bayesian networks, hidden Markov models, and conditional random fields that can handle not only specific relationships between attributes, but also temporal dependencies. Karkera (2014) wrote an excellent introductory book on this topic, *Building Probabilistic Graphical Models with Python*, while Koller and Friedman (2009) published a comprehensive theory bible, *Probabilistic Graphical Models*.

Kernel methods

Any linear model can be turned into a non-linear model by applying the kernel trick to the model – replacing its features (predictors) by a kernel function. In other words, the kernel implicitly transforms our dataset into higher dimensions. The kernel trick leverages the fact that it is often easier to separate the instances in more dimensions. Algorithms capable of operating with kernels include the kernel perceptron, **Support Vector Machines (SVM)**, Gaussian processes, PCA, canonical correlation analysis, ridge regression, spectral clustering, linear adaptive filters, and many others.

Artificial neural networks

Artificial neural networks are inspired by the structure of biological neural networks and are capable of machine learning, as well as pattern recognition. They are commonly used for both regression and classification problems, comprising a wide variety of algorithms and variations for all manner of problem types. Some popular classification methods are **perceptron**, **restricted Boltzmann machine (RBM)**, and **deep belief networks**.

Ensemble learning

Ensemble methods compose of a set of diverse weaker models to obtain better predictive performance. The individual models are trained separately and their predictions are then combined in some way to make the overall prediction. Ensembles, hence, contain multiple ways of modeling the data, which hopefully leads to better results. This is a very powerful class of techniques, and as such, it is very popular; for instance, boosting, bagging, AdaBoost, and Random Forest. The main differences among them are the type of weak learners that are to be combined and the ways in which to combine them.

Evaluating classification

Is our classifier doing well? Is this better than the other one? In classification, we count how many times we classify something right and wrong. Suppose there are two possible classification labels—yes and no—then there are four possible outcomes, as shown in the next figure:

- True positive—hit: This indicates a *yes* instance correctly predicted as *yes*
- True negative—correct rejection: This indicates a *no* instance correctly predicted as *no*
- False positive—false alarm: This indicates a *no* instance predicted as *yes*
- False negative—miss: This indicates a *yes* instance predicted as *no*

		Predicted as positive?	
		Yes	No
Really positive?	Yes	TP—true positive	FN—false negative
	No	FP—false positive	TN—true negative

The basic two performance measures of a classifier are classification error and accuracy, as shown in the following image:

$$\text{Classification error} = \frac{\text{errors}}{\text{totals}} = \frac{FP + FN}{FP + FN + TP + TN}$$

$$\text{Classification accuracy} = 1 - \text{error} = \frac{\text{correct}}{\text{totals}} = \frac{TP + TN}{FP + FN + TP + TN}$$

The main problem with these two measures is that they cannot handle unbalanced classes. Classifying whether a credit card transaction is an abuse or not is an example of a problem with unbalanced classes, there are 99.99% normal transactions and just a tiny percentage of abuses. Classifier that says that every transaction is a normal one is 99.99% accurate, but we are mainly interested in those few classifications that occur very rarely.

Precision and recall

The solution is to use measures that don't involve TN (correct rejections). Two such measures are as follows:

- *Precision*: This is the proportion of positive examples correctly predicted as positive (TP) out of all examples predicted as positive ($TP + FP$):

$$Precision = \frac{TP}{TP + FP}$$

- *Recall*: This is the proportion of positives examples correctly predicted as positive (TP) out of all positive examples ($TP + FN$):

$$Recall = \frac{TP}{TP + FN}$$

It is common to combine the two and report the *F-measure*, which considers both precision and recall to calculate the score as a weighted average, where the score reaches its best value at 1 and worst at 0, as follows:

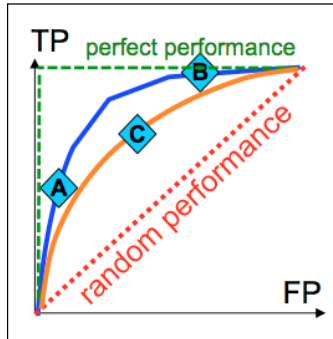
$$F - measure = \frac{2 * Precision * Recall}{Precision + Recall}$$

Roc curves

Most classification algorithms return a classification confidence denoted as $f(X)$, which is, in turn, used to calculate the prediction. Following the credit card abuse example, a rule might look similar to the following:

$$F(X) = \begin{cases} abuse, & \text{if } f(X) > threshold \\ not\ abuse, & \text{else} \end{cases}$$

The threshold determines the error rate and the true positive rate. The outcomes for all the possible threshold values can be plotted as a **Receiver Operating Characteristics (ROC)** as shown in the following diagram:



A random predictor is plotted with a red dashed line and a perfect predictor is plotted with a green dashed line. To compare whether the **A** classifier is better than **C**, we compare the area under the curve.

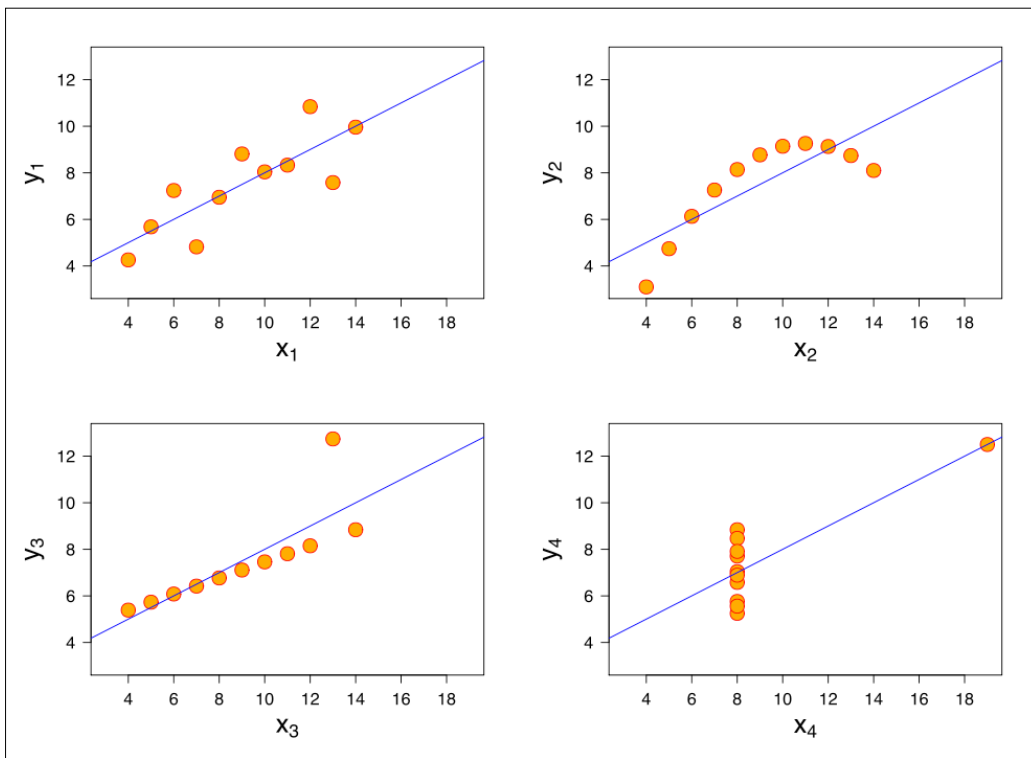
Most of the toolboxes provide all of the previous measures out-of-the-box.

Regression

Regression deals with continuous target variable, unlike classification, which works with a discrete target variable. For example, in order to forecast the outside temperature of the following few days, we will use regression; while classification will be used to predict whether it will rain or not. Generally speaking, regression is a process that estimates the relationship among features, that is, how varying a feature changes the target variable.

Linear regression

The most basic regression model assumes linear dependency between features and target variable. The model is often fitted using least squares approach, that is, the best model minimizes the squares of the errors. In many cases, linear regression is not able to model complex relations, for example, the next figure shows four different sets of points having the same linear regression line: the upper-left model captures the general trend and can be considered as a proper model, the bottom-left model fits points much better, except an outlier – this should be carefully checked – and the upper and lower-right side linear models completely miss the underlying structure of the data and cannot be considered as proper models.



Evaluating regression

In regression, we predict numbers Y from inputs X and the predictions are usually wrong and not exact. The main question that we ask is by how much? In other words, we want to measure the distance between the predicted and true values.

Mean squared error

Mean squared error is an average of the squared difference between the predicted and true values, as follows:

$$MSE(X, Y) = \sqrt{\frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2}$$

The measure is very sensitive to the outliers, for example, 99 exact predictions and *one prediction* off by 10 is scored the same as all predictions wrong by 1. Moreover, the measure is sensitive to the mean. Therefore, relative squared error, which compares the MSE of our predictor to the MSE of the mean predictor (which always predicts the mean value) is often used instead.

Mean absolute error

Mean absolute error is an average of the absolute difference between the predicted and the true values, as follows:

$$MAS(X, Y) = \frac{1}{n} \sum_{i=1}^n |f(X_i) - Y_i|$$

The MAS is less sensitive to the outliers, but it is also sensitive to the mean and scale.

Correlation coefficient

Correlation coefficient compares the average of prediction relative to the mean multiplied by training values relative to the mean. If the number is negative, it means weak correlation, positive number means strong correlation, and zero means no correlation. The correlation between true values X and predictions Y is defined as follows:

$$CC_{XY} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

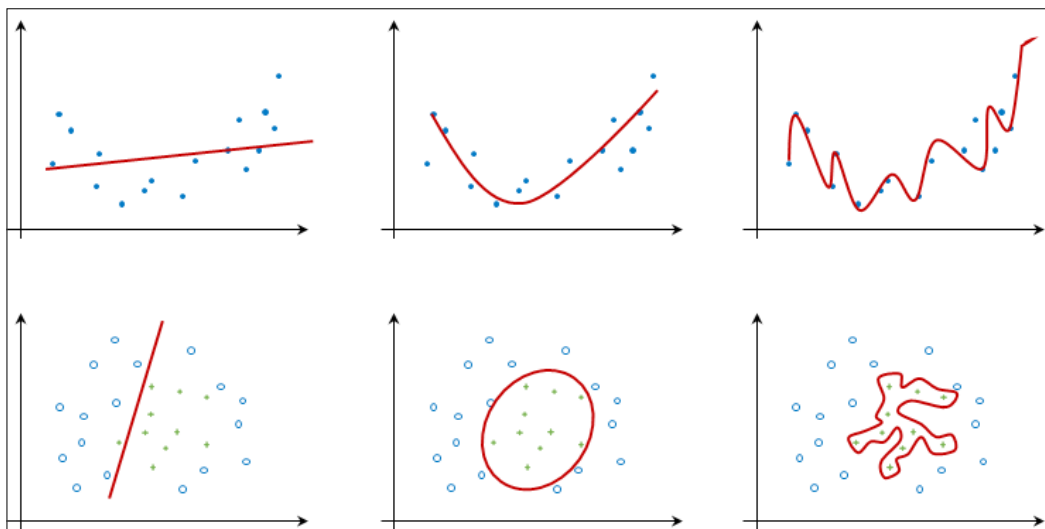
The CC measure is completely insensitive to the mean and scale, and less sensitive to the outliers. It is able to capture the relative ordering, which makes it useful to rank the tasks such as document relevance and gene expression.

Generalization and evaluation

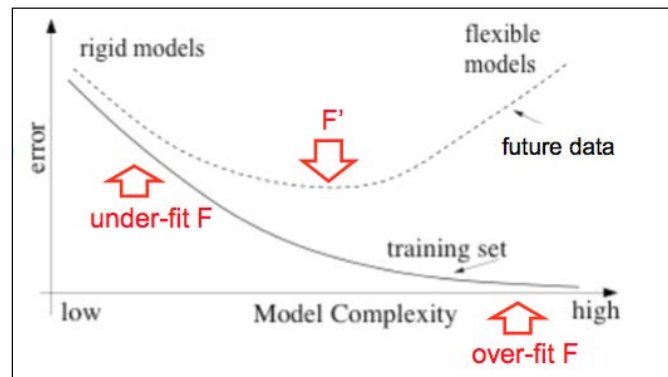
Once the model is built, how do we know it will perform on new data? Is this model any good? To answer these questions, we'll first look into the model generalization and then, see how to get an estimate of the model performance on new data.

Underfitting and overfitting

Predictor training can lead to models that are too complex or too simple. The model with low complexity (the leftmost models) can be as simple as predicting the most frequent or mean class value, while the model with high complexity (the rightmost models) can represent the training instances. Too rigid models, which are shown on the left-hand side, cannot capture complex patterns; while too flexible models, shown on the right-hand side, fit to the noise in the training data. The main challenge is to select the appropriate learning algorithm and its parameters, so that the learned model will perform well on the new data (for example, the middle column):



The following figure shows how the error in the training set decreases with the model complexity. Simple rigid models underfit the data and have large errors. As model complexity increases, it describes the underlying structure of the training data better, and consequentially, the error decreases. If the model is too complex, it overfits the training data and its prediction error increases again:



Depending on the task complexity and data availability, we want to tune our classifiers towards less or more complex structures. Most learning algorithms allow such tuning, as follows:

- **Regression:** This is the order of the polynomial
- **Naive Bayes:** This is the number of the attributes
- **Decision trees:** This is the number of nodes in the tree, pruning confidence
- **k-nearest neighbors:** This is the number of neighbors, distance-based neighbor weights
- **SVM:** This is the kernel type, cost parameter
- **Neural network:** This is the number of neurons and hidden layers

With tuning, we want to minimize the generalization error, that is, how well the classifier performs on future data. Unfortunately, we can never compute the true generalization error; however, we can estimate it. Nevertheless, if the model performs well on the training data, but performance is much worse on the test data, the model most likely overfits.

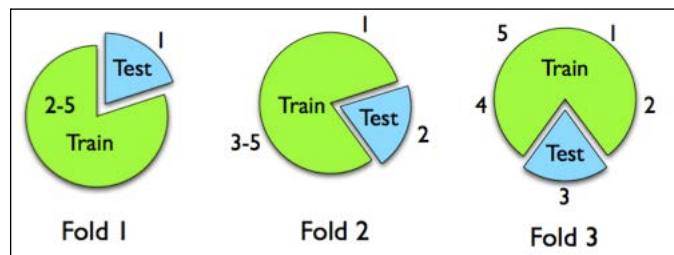
Train and test sets

To estimate the generalization error, we split our data into two parts: training data and testing data. A general rule of thumb is to split them in the *training:testing* ratio, that is, 70:30. We first train the predictor on the training data, then predict the values for the test data, and finally, compute the error – the difference between the predicted and the true values. This gives us an estimate of the true generalization error.

The estimation is based on the following two assumptions: first, we assume that the test set is an unbiased sample from our dataset; and second, we assume that the actual new data will reassemble the distribution as our training and testing examples. The first assumption can be mitigated by cross-validation and stratification. Also, if it is scarce one can't afford to leave out a considerable amount of data for separate test set as learning algorithms do not perform well if they don't receive enough data. In such cases, cross-validation is used instead.

Cross-validation

Cross-validation splits the dataset into k sets of approximately the same size, for example, to five sets as shown in the following figure. First, we use the 2-5 sets for learning and set 1 for training. We then repeat the procedure five times, leaving out one set at a time for testing, and average the error over the five repetitions.



This way, we used all the data for learning and testing as well, while we avoided using the same data to train and test a model.

Leave-one-out validation

An extreme example of cross-validation is the leave-one-out validation. In this case, the number of folds is equal to the number of instances; we learn on all but one instance, and then test the model on the omitted instance. We repeat this for all instances, so that each instance is used exactly once for the validation. This approach is recommended when we have a limited set of learning examples, for example, less than 50.

Stratification

Stratification is a procedure to select a subset of instances in such a way that each fold roughly contains the same proportion of class values. When a class is continuous, the folds are selected so that the mean response value is approximately equal in all the folds. Stratification can be applied along with cross-validation or separate training and test sets.

Summary

In this chapter, we refreshed the machine learning basics. We revisited the workflow of applied machine learning and clarified the main tasks, methods, and algorithms. In the next chapter, we will review the kind of Java libraries that are available and the kind of tasks they can perform.

2

Java Libraries and Platforms for Machine Learning

Implementing machine learning algorithms by yourself is probably the best way to learn machine learning, but you can progress much faster if you step on the shoulders of the giants and leverage one of the existing open source libraries.

This chapter reviews various libraries and platforms for machine learning in Java. The goal is to understand what each library brings to the table and what kind of problems it is able to solve?

In this chapter, we will cover the following topics:

- The requirement of Java to implement a machine learning app
- Weka, a general purpose machine learning platform
- Java machine learning library, a collection of machine learning algorithms
- Apache Mahout, a scalable machine learning platform
- Apache Spark, a distributed machine learning library
- Deeplearning4j, a deep learning library
- MALLET, a text mining library

We'll also discuss how to architect the complete machine learning app stack for both single-machine and big data apps using these libraries with other components.

The need for Java

New machine learning algorithms are often first scripted at university labs, gluing together several languages such as shell scripting, Python, R, MATLAB Java, Scala, or C++ to prove a new concept and theoretically analyze its properties. An algorithm might then take a long path of refactoring before it lands in a library with standardized input/output and interfaces. While Python, R, and MATLAB are quite popular, they are mainly used for scripting, research, and experimenting. Java, on the other hand, is the de-facto enterprise language, which could be attributed to static typing, robust IDE support, good maintainability, as well as decent threading model, and high-performance concurrent data structure libraries. Moreover, there are already many Java libraries available for machine learning, which make it really convenient to interface them in existing Java applications and leverage powerful machine learning capabilities.

Machine learning libraries

There are over 70 Java-based open source machine learning projects listed on the [MLOSS.org](http://mloss.org) website and probably many more unlisted projects live at university servers, GitHub, or Bitbucket. In this section, we will review the major libraries and platforms, the kind of problems they can solve, the algorithms they support, and the kind of data they can work with.

Weka

Weka, which is short for **Waikato Environment for Knowledge Analysis**, is a machine learning library developed at the University of Waikato, New Zealand, and is probably the most well-known Java library. It is a general-purpose library that is able to solve a wide variety of machine learning tasks, such as classification, regression, and clustering. It features a rich graphical user interface, command-line interface, and Java API. You can check out Weka at <http://www.cs.waikato.ac.nz/ml/weka/>.

At the time of writing this book, Weka contains 267 algorithms in total: data pre-processing (82), attribute selection (33), classification and regression (133), clustering (12), and association rules mining (7). Graphical interfaces are well-suited for exploring your data, while Java API allows you to develop new machine learning schemes and use the algorithms in your applications.



Weka is distributed under **GNU General Public License (GNU GPL)**, which means that you can copy, distribute, and modify it as long as you track changes in source files and keep it under GNU GPL. You can even distribute it commercially, but you must disclose the source code or obtain a commercial license.

In addition to several supported file formats, Weka features its own default data format, ARFF, to describe data by attribute-data pairs. It consists of two parts. The first part contains header, which specifies all the attributes (that is, features) and their type; for instance, nominal, numeric, date, and string. The second part contains data, where each line corresponds to an instance. The last attribute in the header is implicitly considered as the target variable, missing data are marked with a question mark. For example, returning to the example from *Chapter 1, Applied Machine Learning Quick Start*, the Bob instance written in an ARFF file format would be as follows:

```
@RELATION person_dataset

@ATTRIBUTE `Name`  STRING
@ATTRIBUTE `Height`  NUMERIC
@ATTRIBUTE `Eye color` {blue, brown, green}
@ATTRIBUTE `Hobbies`  STRING

@DATA
'Bob', 185.0, blue, 'climbing, sky diving'
'Anna', 163.0, brown, 'reading'
'Jane', 168.0, ?, ?
```

The file consists of three sections. The first section starts with the `@RELATION <String>` keyword, specifying the dataset name. The next section starts with the `@ATTRIBUTE` keyword, followed by the attribute name and type. The available types are `STRING`, `NUMERIC`, `DATE`, and a set of categorical values. The last attribute is implicitly assumed to be the target variable that we want to predict. The last section starts with the `@DATA` keyword, followed by one instance per line. Instance values are separated by comma and must follow the same order as attributes in the second section.

More Weka examples will be demonstrated in *Chapter 3, Basic Algorithms – Classification, Regression, and Clustering*, and *Chapter 4, Customer Relationship Prediction with Ensembles*.



To learn more about Weka, pick up a quick-start book, *Weka How-to* by Kaluza (Packt Publishing) to start coding or look into *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations* by Witten and Frank (Morgan Kaufmann Publishers) for theoretical background and in-depth explanations.

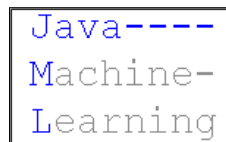
Weka's Java API is organized in the following top-level packages:

- `weka.associations`: These are data structures and algorithms for association rules learning, including **Apriori**, **predictive apriori**, `FilteredAssociator`, **FP-Growth**, **Generalized Sequential Patterns (GSP)**, **Hotspot**, and **Tertius**.
- `weka.classifiers`: These are supervised learning algorithms, evaluators, and data structures. The package is further split into the following components:
 - `weka.classifiers.bayes`: This implements Bayesian methods, including naive Bayes, Bayes net, Bayesian logistic regression, and so on
 - `weka.classifiers.evaluation`: These are supervised evaluation algorithms for nominal and numerical prediction, such as evaluation statistics, confusion matrix, ROC curve, and so on
 - `weka.classifiers.functions`: These are regression algorithms, including linear regression, isotonic regression, Gaussian processes, support vector machine, multilayer perceptron, voted perceptron, and others
 - `weka.classifiers.lazy`: These are instance-based algorithms such as k-nearest neighbors, K^* , and lazy Bayesian rules

-
- `weka.classifiers.meta`: These are supervised learning meta-algorithms, including AdaBoost, bagging, additive regression, random committee, and so on
 - `weka.classifiers.mi`: These are multiple-instance learning algorithms, such as citation k-nn, diverse density, MI AdaBoost, and others
 - `weka.classifiers.rules`: These are decision tables and decision rules based on the separate-and-conquer approach, Ripper, Part, Prism, and so on
 - `weka.classifiers.trees`: These are various decision trees algorithms, including ID3, C4.5, M5, functional tree, logistic tree, random forest, and so on
- `weka.clusterers`: These are clustering algorithms, including k-means, Clope, Cobweb, DBSCAN hierarchical clustering, and farthest.
 - `weka.core`: These are various utility classes, data presentations, configuration files, and so on.
 - `weka.datagenerators`: These are data generators for classification, regression, and clustering algorithms.
 - `weka.estimators`: These are various data distribution estimators for discrete/nominal domains, conditional probability estimations, and so on.
 - `weka.experiment`: These are a set of classes supporting necessary configuration, datasets, model setups, and statistics to run experiments.
 - `weka.filters`: These are attribute-based and instance-based selection algorithms for both supervised and unsupervised data preprocessing.
 - `weka.gui`: These are graphical interface implementing **explorer**, **experimenter**, and **knowledge flow** applications. Explorer allows you to investigate dataset, algorithms, as well as their parameters, and visualize dataset with scatter plots and other visualizations. Experimenter is used to design batches of experiment, but it can only be used for classification and regression problems. Knowledge flows implements a visual drag-and-drop user interface to build data flows, for example, load data, apply filter, build classifier, and evaluate.

Java machine learning

Java machine learning library, or Java-ML, is a collection of machine learning algorithms with a common interface for algorithms of the same type. It only features Java API, therefore, it is primarily aimed at software engineers and programmers. Java-ML contains algorithms for data preprocessing, feature selection, classification, and clustering. In addition, it features several Weka bridges to access Weka's algorithms directly through the Java-ML API. It can be downloaded from <http://java-ml.sourceforge.net>; where, the latest release was in 2012 (at the time of writing this book).



Java-ML is also a general-purpose machine learning library. Compared to Weka, it offers more consistent interfaces and implementations of recent algorithms that are not present in other packages, such as an extensive set of state-of-the-art similarity measures and feature-selection techniques, for example, dynamic time warping, random forest attribute evaluation, and so on. Java-ML is also available under the GNU GPL license.

Java-ML supports any type of file as long as it contains one data sample per line and the features are separated by a symbol such as comma, semi-colon, and tab.

The library is organized around the following top-level packages:

- `net.sf.javaml.classification`: These are classification algorithms, including naive Bayes, random forests, bagging, self-organizing maps, k-nearest neighbors, and so on
- `net.sf.javaml.clustering`: These are clustering algorithms such as k-means, self-organizing maps, spatial clustering, Cobweb, AQBC, and others
- `net.sf.javaml.core`: These are classes representing instances and datasets
- `net.sf.javaml.distance`: These are algorithms that measure instance distance and similarity, for example, **Chebyshev distance**, cosine distance/similarity, Euclidian distance, Jaccard distance/similarity, **Mahalanobis distance**, **Manhattan distance**, **Minkowski distance**, **Pearson correlation coefficient**, **Spearman's footrule distance**, **dynamic time wrapping (DTW)**, and so on

- `net.sf.javaml.featureselection`: These are algorithms for feature evaluation, scoring, selection, and ranking, for instance, gain ratio, ReliefF, Kullback-Liebler divergence, symmetrical uncertainty, and so on
- `net.sf.javaml.filter`: These are methods for manipulating instances by filtering, removing attributes, setting classes or attribute values, and so on
- `net.sf.javaml.matrix`: This implements in-memory or file-based array
- `net.sf.javaml.sampling`: This implements sampling algorithms to select a subset of dataset
- `net.sf.javaml.tools`: These are utility methods on dataset, instance manipulation, serialization, Weka API interface, and so on
- `net.sf.javaml.utils`: These are utility methods for algorithms, for example, statistics, math methods, contingency tables, and others

Apache Mahout

The Apache Mahout project aims to build a scalable machine learning library. It is built atop scalable, distributed architectures, such as Hadoop, using the MapReduce paradigm, which is an approach for processing and generating large datasets with a parallel, distributed algorithm using a cluster of servers.



Mahout features console interface and Java API to scalable algorithms for clustering, classification, and collaborative filtering. It is able to solve three business problems: item recommendation, for example, recommending items such as *people who liked this movie also liked...*; clustering, for example, of text documents into groups of topically-related documents; and classification, for example, learning which topic to assign to an unlabeled document.

Mahout is distributed under a commercially-friendly Apache License, which means that you can use it as long as you keep the Apache license included and display it in your program's copyright notice.

Mahout features the following libraries:

- `org.apache.mahout.cf.taste`: These are collaborative filtering algorithms based on user-based and item-based collaborative filtering and matrix factorization with ALS
- `org.apache.mahout.classifier`: These are in-memory and distributed implementations, including logistic regression, naive Bayes, random forest, **hidden Markov models (HMM)**, and multilayer perceptron
- `org.apache.mahout.clustering`: These are clustering algorithms such as canopy clustering, k-means, fuzzy k-means, streaming k-means, and spectral clustering
- `org.apache.mahout.common`: These are utility methods for algorithms, including distances, MapReduce operations, iterators, and so on
- `org.apache.mahout.driver`: This implements a general-purpose driver to run main methods of other classes
- `org.apache.mahout.ep`: This is the evolutionary optimization using the recorded-step mutation
- `org.apache.mahout.math`: These are various math utility methods and implementations in Hadoop
- `org.apache.mahout.vectorizer`: These are classes for data presentation, manipulation, and MapReduce jobs

Apache Spark

Apache Spark, or simply Spark, is a platform for large-scale data processing builds atop Hadoop, but, in contrast to Mahout, it is not tied to the MapReduce paradigm. Instead, it uses in-memory caches to extract a working set of data, process it, and repeat the query. This is reported to be up to ten times as fast as a Mahout implementation that works directly with disk-stored data. It can be grabbed from <https://spark.apache.org>.



There are many modules built atop Spark, for instance, **GraphX** for graph processing, **Spark Streaming** for processing real-time data streams, and MLlib for machine learning library featuring classification, regression, collaborative filtering, clustering, dimensionality reduction, and optimization.

Spark's MLlib can use a Hadoop-based data source, for example, **Hadoop Distributed File System (HDFS)** or HBase, as well as local files. The supported data types include the following:

- **Local vector** is stored on a single machine. Dense vectors are presented as an array of double-typed values, for example, $(2.0, 0.0, 1.0, 0.0)$; while sparse vector is presented by the size of the vector, an array of indices, and an array of values, for example, $[4, (0, 2), (2.0, 1.0)]$.
- **Labeled point** is used for supervised learning algorithms and consists of a local vector labeled with a double-typed class values. Label can be class index, binary outcome, or a list of multiple class indices (multiclass classification). For example, a labeled dense vector is presented as $[1.0, (2.0, 0.0, 1.0, 0.0)]$.
- **Local matrix** stores a dense matrix on a single machine. It is defined by matrix dimensions and a single double-array arranged in a column-major order.
- **Distributed matrix** operates on data stored in Spark's **Resilient Distributed Dataset (RDD)**, which represents a collection of elements that can be operated on in parallel. There are three presentations: row matrix, where each row is a local vector that can be stored on a single machine, row indices are meaningless; and indexed row matrix, which is similar to row matrix, but the row indices are meaningful, that is, rows can be identified and joins can be executed; and coordinate matrix, which is used when a row cannot be stored on a single machine and the matrix is very sparse.

Spark's MLlib API library provides interfaces to various learning algorithms and utilities as outlined in the following list:

- `org.apache.spark.mllib.classification`: These are binary and multiclass classification algorithms, including linear SVMs, logistic regression, decision trees, and naive Bayes
- `org.apache.spark.mllib.clustering`: These are k-means clustering
- `org.apache.spark.mllib.linalg`: These are data presentations, including dense vectors, sparse vectors, and matrices
- `org.apache.spark.mllib.optimization`: These are the various optimization algorithms used as low-level primitives in MLlib, including gradient descent, stochastic gradient descent, update schemes for distributed SGD, and limited-memory BFGS
- `org.apache.spark.mllib.recommendation`: These are model-based collaborative filtering implemented with alternating least squares matrix factorization

- `org.apache.spark.mllib.regression`: These are regression learning algorithms, such as linear least squares, decision trees, Lasso, and Ridge regression
- `org.apache.spark.mllib.stat`: These are statistical functions for samples in sparse or dense vector format to compute the mean, variance, minimum, maximum, counts, and nonzero counts
- `org.apache.spark.mllib.tree`: This implements classification and regression decision tree-learning algorithms
- `org.apache.spark.mllib.util`: These are a collection of methods to load, save, preprocess, generate, and validate the data

Deeplearning4j

Deeplearning4j, or DL4J, is a deep-learning library written in Java. It features a distributed as well as a single-machine deep-learning framework that includes and supports various neural network structures such as feedforward neural networks, **RBM**, convolutional neural nets, deep belief networks, autoencoders, and others. DL4J can solve distinct problems, such as identifying faces, voices, spam or e-commerce fraud.

Deeplearning4j is also distributed under Apache 2.0 license and can be downloaded from <http://deeplearning4j.org>. The library is organized as follows:

- `org.deeplearning4j.base`: These are loading classes
- `org.deeplearning4j.berkeley`: These are math utility methods
- `org.deeplearning4j.clustering`: This is the implementation of k-means clustering
- `org.deeplearning4j.datasets`: This is dataset manipulation, including import, creation, iterating, and so on
- `org.deeplearning4j.distributions`: These are utility methods for distributions
- `org.deeplearning4j.eval`: These are evaluation classes, including the confusion matrix
- `org.deeplearning4j.exceptions`: This implements exception handlers
- `org.deeplearning4j.models`: These are supervised learning algorithms, including deep belief network, stacked autoencoder, stacked denoising autoencoder, and RBM

- `org.deeplearning4j.nn`: These are the implementation of components and algorithms based on neural networks, such as neural network, multi-layer network, convolutional multi-layer network, and so on
- `org.deeplearning4j.optimize`: These are neural net optimization algorithms, including back propagation, multi-layer optimization, output layer optimization, and so on
- `org.deeplearning4j.plot`: These are various methods for rendering data
- `org.deeplearning4j.rng`: This is a random data generator
- `org.deeplearning4j.util`: These are helper and utility methods

MALLET

Machine Learning for Language Toolkit (MALLET), is a large library of natural language processing algorithms and utilities. It can be used in a variety of tasks such as document classification, document clustering, information extraction, and topic modeling. It features command-line interface as well as Java API for several algorithms such as naive Bayes, HMM, **Latent Dirichlet** topic models, logistic regression, and conditional random fields.



MALLET is available under Common Public License 1.0, which means that you can even use it in commercial applications. It can be downloaded from <http://mallet.cs.umass.edu>. MALLET instance is represented by name, label, data, and source. However, there are two methods to import data into the MALLET format, as shown in the following list:

- Instance per file: Each file, that is, document, corresponds to an instance and MALLET accepts the directory name for the input.
- Instance per line: Each line corresponds to an instance, where the following format is assumed: the `instance_name` label token. Data will be a feature vector, consisting of distinct words that appear as tokens and their occurrence count.

The library comprises the following packages:

- `cc.mallet.classify`: These are algorithms for training and classifying instances, including AdaBoost, bagging, C4.5, as well as other decision tree models, multivariate logistic regression, naive Bayes, and Winnow2.
- `cc.mallet.cluster`: These are unsupervised clustering algorithms, including greedy agglomerative, hill climbing, k-best, and k-means clustering.
- `cc.mallet.extract`: This implements tokenizers, document extractors, document viewers, cleaners, and so on.
- `cc.mallet.fst`: This implements sequence models, including conditional random fields, HMM, maximum entropy Markov models, and corresponding algorithms and evaluators.
- `cc.mallet.grmm`: This implements graphical models and factor graphs such as inference algorithms, learning, and testing. For example, loopy belief propagation, Gibbs sampling, and so on.
- `cc.mallet.optimize`: These are optimization algorithms for finding the maximum of a function, such as gradient ascent, limited-memory BFGS, stochastic meta ascent, and so on.
- `cc.mallet.pipe`: These are methods as pipelines to process data into MALLET instances.
- `cc.mallet.topics`: These are topics modeling algorithms, such as Latent Dirichlet allocation, four-level pachinko allocation, hierarchical PAM, DMRT, and so on.
- `cc.mallet.types`: This implements fundamental data types such as dataset, feature vector, instance, and label.
- `cc.mallet.util`: These are miscellaneous utility functions such as command-line processing, search, math, test, and so on.

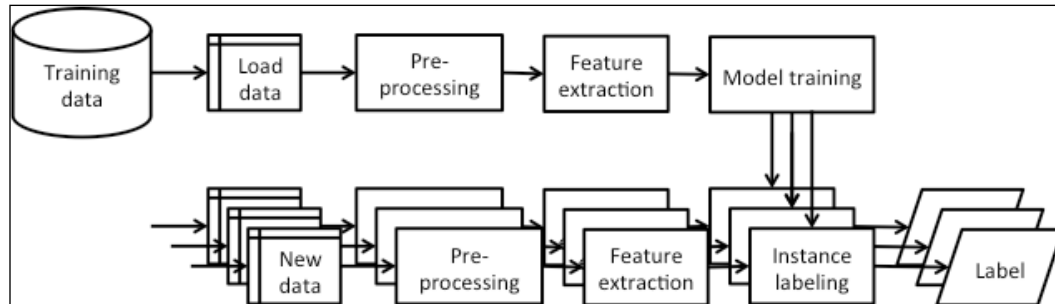
Comparing libraries

The following table summarizes all the presented libraries. The table is, by no means, exhaustive – there are many more libraries covering the specific-problem domains. This review should serve as an overview of the big names in the Java machine learning world:

	Problem domains	License	Architecture	Algorithms
Weka	General purpose	GNU GPL	Single machine	Decision trees, naive Bayes, neural network, random forest, AdaBoost, hierarchical clustering, and so on
Java-ML	General purpose	GNU GPL	Single machine	k-means clustering, self-organizing maps, Markov chain clustering, Cobweb, random forest, decision trees, bagging, distance measures, and so on
Mahout	Classification, recommendation, and clustering	Apache 2.0 License	Distributed, single machine	Logistic regression, naive Bayes, random forest, HMM, multilayer perceptron, k-means clustering, and so on
Spark	General purpose	Apache 2.0 License	Distributed	SVM, logistic regression, decision trees, naive Bayes, k-means clustering, linear least squares, LASSO, ridge regression, and so on
DL4J	Deep learning	Apache 2.0 License	Distributed, single machine	RBM, deep belief networks, deep autoencoders, recursive neural tensor networks, convolutional neural network, and stacked denoising autoencoders
MALLET	Text mining	Common Public License 1.0	Single machine	Naive Bayes, decision trees, maximum entropy, hidden Markov models, and conditional random fields

Building a machine learning application

Machine learning applications, especially those focused on classification, usually follow the same high-level workflow as shown in the following diagram. The workflow comprises two phases: training the classifier and classification of new instances. Both phases share common steps as you can see in the following diagram:



First, we use a set of **Training data**, select a representative subset as the training set, preprocess missing data, and extract features. A selected supervised learning algorithm is used to train a model, which is deployed in the second phase. The second phase puts a new data instance through the same **Pre-processing** and **Feature extraction** procedure and applies the learned model to obtain the instance label. If you are able to collect new labeled data, periodically rerun the learning phase to retrain the model, and replace the old one with the retrained one in the classification phase.

Traditional machine learning architecture

Structured data, such as transactional, customer, analytical, and market data, usually resides within a local relational database. Given a query language, such as SQL, we can query the data used for processing, as shown in the workflow in the previous diagram. Usually, all the data can be stored in the memory and further processed with a machine learning library such as Weka, Java-ML, or MALLET.

A common practice in the architecture design is to create data pipelines, where different steps in the workflow are split. For instance, in order to create a client data record, we might have to scrap the data from different data sources. The record can be then saved in an intermediate database for further processing.

To understand how the high-level aspects of big data architecture differ, let's first clarify when is the data considered big?

Dealing with big data

Big data existed long before the phrase was invented, for instance, banks and stock exchanges have been processing billions of transactions daily for years, and airline companies have worldwide real-time infrastructure for operational management of passenger booking, and so on. So what is big data really? Doug Laney (2001) suggested that big data is defined by three Vs: volume, velocity, and variety. Therefore, to answer the question whether your data is big, we can translate this into the following three subquestions:

- **Volume:** Can you store your data in memory?
- **Velocity:** Can you process new incoming data with a single machine?
- **Variety:** Is your data from a single source?

If you answered all the questions with yes, then your data is probably not big, do not worry, you have just simplified your application architecture.

If your answer to all the questions was no, then your data is big! However, if you have mixed answers, then it's complicated. Some may argue that a *V* is important, other may say the other *V*s. From the machine learning point of view, there is a fundamental difference in algorithm implementation to process the data in memory or from distributed storage. Therefore, a rule of thumb is as follows: if you cannot store your data in the memory, then you should look into a big data machine learning library.

The exact answer depends on the problem that you are trying to solve. If you're starting a new project, I'd suggest you start off with a single-machine library and prototype your algorithm, possibly with a subset of your data if the entire data does not fit into the memory. Once you've established good initial results, consider moving to something more heavy duty such as Mahout or Spark.

Big data application architecture

Big data, such as documents, weblogs, social networks, sensor data, and others, are stored in a NoSQL database, such as MongoDB, or a distributed filesystem, such as HDFS. In case we deal with structured data, we can deploy database capabilities using systems such as Cassandra or HBase built atop Hadoop. Data processing follows the MapReduce paradigm, which breaks data processing problems into smaller subproblems and distributes tasks across processing nodes. Machine learning models are finally trained with machine learning libraries such as Mahout and Spark.



MongoDB is a NoSQL database, which stores documents in a JSON-like format. Read more about it at <https://www.mongodb.org>.

Hadoop is a framework for distributed processing of large datasets across a cluster of computers. It includes its own filesystem format HDFS, job scheduling framework YARN, and implements the MapReduce approach for parallel data processing. More about Hadoop is available at <http://hadoop.apache.org/>.

Cassandra is a distributed database management system build to provide fault-tolerant, scalable, and decentralized storage. More information is available at <http://cassandra.apache.org/>.

HBase is another database that focuses on random read/write access to distributed storage. More information is available at <https://hbase.apache.org/>.

Summary

Selecting a machine learning library has an important impact on your application architecture. The key is to consider your project requirements. What kind of data do you have? What kind of problem are you trying to solve? Is your data big? Do you need distributed storage? What kind of algorithm are you planning to use? Once you figure out what you need to solve your problem, pick a library that best fits your needs.

In the next chapter, we will cover how to complete basic machine learning tasks such as classification, regression, and clustering using some of the presented libraries.

3

Basic Algorithms – Classification, Regression, and Clustering

In the previous chapter, we reviewed the key Java libraries for machine learning and what they bring to the table. In this chapter, we will finally get our hands dirty. We will take a closer look at the basic machine learning tasks such as classification, regression, and clustering. Each of the topics will introduce basic algorithms for classification, regression, and clustering. The example datasets will be small, simple, and easy to understand.

The following is the list of topics that will be covered in this chapter:

- Loading data
- Filtering attributes
- Building classification, regression, and clustering models
- Evaluating models

Before you start

Download the latest version of Weka 3.6 from <http://www.cs.waikato.ac.nz/ml/weka/downloading.html>.

There are multiple download options available. We'll want to use Weka as a library in our source code, so make sure you skip the self-extracting executables and pick the ZIP archive as shown at the following image. Unzip the archive and locate `weka.jar` within the extracted archive:

- **Other platforms (Linux, etc.)**
 - Click [here](#) to download a zip archive containing Weka (weka-3-7-11.zip; 33.2 MB)
 - First unzip the zip file. This will create a new directory called weka-3-7-11. To run Weka, change into that directory and type

```
java -Xmx1000M -jar weka.jar
```

 - Note that Java needs to be installed on your system for this to work. Also note, that using `-jar` will override your current CLASSPATH variable and only use the `weka.jar`.

We'll use the Eclipse IDE to show examples, as follows:

1. Start a new Java project.
2. Right-click on the project properties, select **Java Build Path**, click on the **Libraries** tab, and select **Add External JARs**.
3. Navigate to extract the Weka archive and select the `weka.jar` file.

That's it, we are ready to implement the basic machine learning techniques!

Classification

We will start with the most commonly used machine learning technique, that is, classification. As we reviewed in the first chapter, the main idea is to automatically build a mapping between the input variables and the outcome. In the following sections, we will look at how to load the data, select features, implement a basic classifier in Weka, and evaluate the classifier performance.

Data

For this task, we will have a look at the zoo database [ref]. The database contains 101 data entries of the animals described with 18 attributes as shown in the following table:

animal	aquatic	fins
hair	predator	legs
feathers	toothed	tail
eggs	backbone	domestic
milk	breathes	cat size
airborne	venomous	type

An example entry in the dataset set is a lion with the following attributes:

- animal: lion
- hair: true
- feathers: false
- eggs: false
- milk: true
- airborne: false
- aquatic: false
- predator: true
- toothed: true
- backbone: true
- breathes: true
- venomous: false
- fins: false
- legs: 4
- tail: true
- domestic: false
- catsize: true
- type: mammal

Our task will be to build a model to predict the outcome variable, `animal`, given all the other attributes as input.

Loading data

Before we start with the analysis, we will load the data in Weka's ARRF format and print the total number of loaded instances. Each data sample is held within an `Instances` object, while the complete dataset accompanied with meta-information is handled by the `Instances` object.

To load the input data, we will use the `DataSource` object that accepts a variety of file formats and converts them to `Instances`:

```
DataSource source = new DataSource(args[0]);
Instances data = source.getDataSet();
System.out.println(data.numInstances() + " instances loaded.");
// System.out.println(data.toString());
```

This outputs the number of loaded instances, as follows:

```
101 instances loaded.
```

We can also print the complete dataset by calling the `data.toString()` method.

Our task is to learn a model that is able to predict the `animal` attribute in the future examples for which we know the other attributes but do not know the `animal` label. Hence, we remove the `animal` attribute from the training set. We accomplish this by filtering out the `animal` attribute using the `Remove` filter.

First, we set a string table of parameters, specifying that the first attribute must be removed. The remaining attributes are used as our dataset for training a classifier:

```
Remove remove = new Remove();
String[] opts = new String[]{"-R", "1"};
```

Finally, we call the `Filter.useFilter(Instances, Filter)` static method to apply the filter on the selected dataset:

```
remove.setOptions(opts);
remove.setInputFormat(data);
data = Filter.useFilter(data, remove);
```

Feature selection

As introduced in *Chapter 1, Applied Machine Learning Quick Start*, one of the pre-processing steps is focused on feature selection, also known as attribute selection. The goal is to select a subset of relevant attributes that will be used in a learned model. Why is feature selection important? A smaller set of attributes simplifies the models and makes them easier to interpret by users, this usually requires shorter training and reduces overfitting.

Attribute selection can take into account the class value or not. In the first case, an attribute selection algorithm evaluates the different subsets of features and calculates a score that indicates the quality of selected attributes. We can use different searching algorithms such as exhaustive search, best first search, and different quality scores such as information gain, Gini index, and so on.

Weka supports this process by an `AttributeSelection` object, which requires two additional parameters: evaluator, which computes how informative an attribute is and a ranker, which sorts the attributes according to the score assigned by the evaluator.

In this example, we will use information gain as an evaluator and rank the features by their information gain score:

```
InfoGainAttributeEval eval = new InfoGainAttributeEval();
Ranker search = new Ranker();
```

Next, we initialize an `AttributeSelection` object and set the evaluator, ranker, and data:

```
AttributeSelection attSelect = new AttributeSelection();
attSelect.setEvaluator(eval);
attSelect.setSearch(search);
attSelect.SelectAttributes(data);
```

Finally, we can print an order list of attribute indices, as follows:

```
int[] indices = attSelect.selectedAttributes();
System.out.println(Utils.arrayToString(indices));
```

The method outputs the following result:

```
12,3,7,2,0,1,8,9,13,4,11,5,15,10,6,14,16
```

The top three most informative attributes are 12 (fins), 3 (eggs), 7 (aquatic), 2 (hair), and so on. Based on this list, we can remove additional, non-informative features in order to help learning algorithms achieve more accurate and faster learning models.

What would make the final decision about the number of attributes to keep? There's no rule of thumb related to an exact number – the number of attributes depends on the data and problem. The purpose of attribute selection is choosing attributes that serve your model better, so it is better to focus whether the attributes are improving the model.

Learning algorithms

We have loaded our data, selected the best features, and are ready to learn some classification models. Let's begin with the basic decision trees.

Decision tree in Weka is implemented within the `J48` class, which is a re-implementation of Quinlan's famous C4.5 decision tree learner [Quinlan, 1993].

First, we initialize a new `J48` decision tree learner. We can pass additional parameters with a string table, for instance, tree pruning that controls the model complexity (refer to *Chapter 1, Applied Machine Learning Quick Start*). In our case, we will build an un-pruned tree, hence we will pass a single `U` parameter:

```
J48 tree = new J48();
String[] options = new String[1];
options[0] = "-U";

tree.setOptions(options);
```

Next, we call the `buildClassifier(Instances)` method to initialize the learning process:

```
tree.buildClassifier(data);
```

The built model is now stored in a `tree` object. We can output the entire `J48` unpruned tree calling the `toString()` method:

```
System.out.println(tree);
```

The output is as follows:

```
J48 unpruned tree
-----

feathers = false
|   milk = false
|   |   backbone = false
|   |   |   airborne = false
|   |   |   |   predator = false
```

```

| | | | | legs <= 2: invertebrate (2.0)
| | | | | legs > 2: insect (2.0)
| | | | | predator = true: invertebrate (8.0)
| | | | | airborne = true: insect (6.0)
| | | | | backbone = true
| | | | | fins = false
| | | | | tail = false: amphibian (3.0)
| | | | | tail = true: reptile (6.0/1.0)
| | | | | fins = true: fish (13.0)
| | | | | milk = true: mammal (41.0)
| | | | | feathers = true: bird (20.0)

```

Number of Leaves : .9

Size of the tree : ..17

The outputted tree has 17 nodes in total, 9 of these are terminal (Leaves).

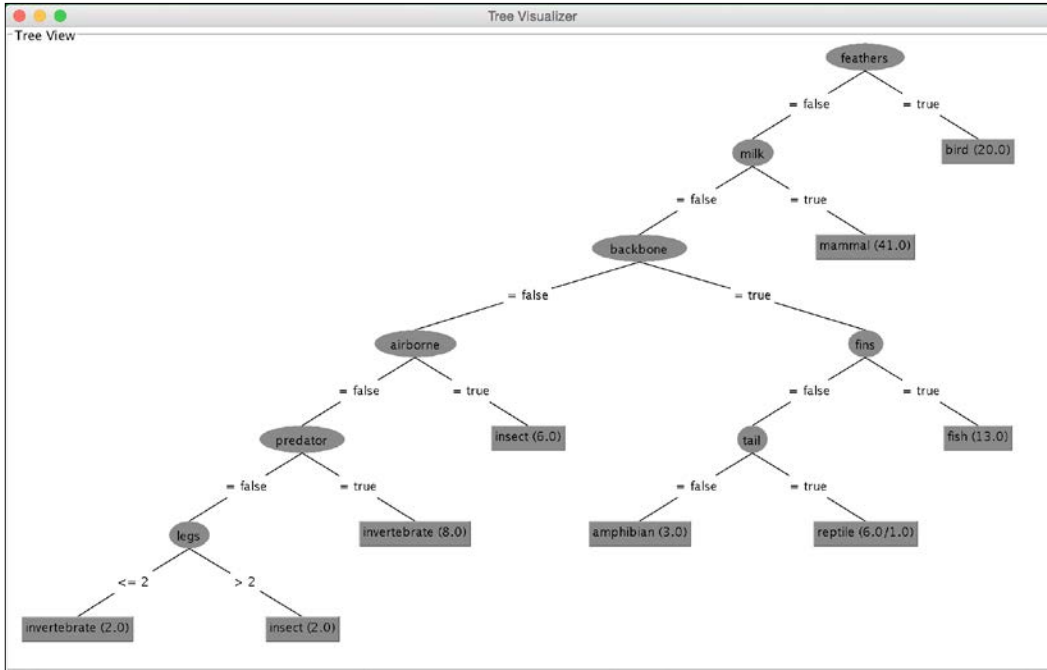
Another way to present the tree is to leverage the built-in `TreeVisualizer` tree viewer, as follows:

```

TreeVisualizer tv = new TreeVisualizer(null, tree.graph(), new
PlaceNode2());
JFrame frame = new javax.swing.JFrame("Tree Visualizer");
frame.setSize(800, 500);
frame.setDefaultCloseOperation(JFrame.EXIT_ON_CLOSE);
frame.getContentPane().add(tv);
frame.setVisible(true);
tv.fitToScreen();

```

The code results in the following frame:



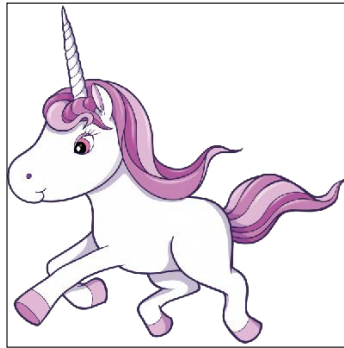
The decision process starts at the top node, also known as the root node. The node label specifies the attribute value that will be checked. In our example, we first check the value of the `feathers` attribute. If the feather is present, we follow the right-hand branch, which leads us to the leaf labeled `bird`, indicating there are 20 examples supporting this outcome. If the feather is not present, we follow the left-hand branch, which leads us to the next `milk` attribute. We check the value of the attribute again and follow the branch that matches the attribute value. We repeat the process until we reach a leaf node.

We can build other classifiers by following the same steps: initialize a classifier, pass the parameters controlling the model complexity, and call the `buildClassifier(Instances)` method.

In the next section, we will learn how to use a trained model to assign a class label to a new example whose label is unknown.

Classify new data

Suppose we record attributes for an animal whose label we do not know, we can predict its label from the learned classification model:



We first construct a feature vector describing the new specimen, as follows:

```
double[] vals = new double[data.numAttributes()];
vals[0] = 1.0; //hair {false, true}
vals[1] = 0.0; //feathers {false, true}
vals[2] = 0.0; //eggs {false, true}
vals[3] = 1.0; //milk {false, true}
vals[4] = 0.0; //airborne {false, true}
vals[5] = 0.0; //aquatic {false, true}
vals[6] = 0.0; //predator {false, true}
vals[7] = 1.0; //toothed {false, true}
vals[8] = 1.0; //backbone {false, true}
vals[9] = 1.0; //breathes {false, true}
vals[10] = 1.0; //venomous {false, true}
vals[11] = 0.0; //fins {false, true}
vals[12] = 4.0; //legs INTEGER [0,9]
vals[13] = 1.0; //tail {false, true}
vals[14] = 1.0; //domestic {false, true}
vals[15] = 0.0; //catsize {false, true}
Instance myUnicorn = new Instance(1.0, vals);
```

Finally, we call the `classify(Instance)` method on the model to obtain the class value. The method returns label index, as follows:

```
double result = tree.classifyInstance(myUnicorn);
System.out.println(data.classAttribute().value((int) result));
```

This outputs the mammal class label.

Evaluation and prediction error metrics

We built a model, but we do not know if it can be trusted. To estimate its performance, we can apply a cross-validation technique explained in *Chapter 1, Applied Machine Learning Quick Start*.

Weka offers an `Evaluation` class implementing cross validation. We pass the model, data, number of folds, and an initial random seed, as follows:

```
Classifier c1 = new J48();
Evaluation eval_roc = new Evaluation(data);
eval_roc.crossValidateModel(c1, data, 10, new Random(1), new
    Object[] {});
System.out.println(eval_roc.toSummaryString());
```

The evaluation results are stored in the `Evaluation` object.

A mix of the most common metrics can be invoked by calling the `toString()` method. Note that the output does not differentiate between regression and classification, so pay attention to the metrics that make sense, as follows:

Correctly Classified Instances	93	92.0792 %
Incorrectly Classified Instances	8	7.9208 %
Kappa statistic	0.8955	
Mean absolute error	0.0225	
Root mean squared error	0.14	
Relative absolute error	10.2478 %	
Root relative squared error	42.4398 %	
Coverage of cases (0.95 level)	96.0396 %	
Mean rel. region size (0.95 level)	15.4173 %	
Total Number of Instances	101	

In the classification, we are interested in the number of correctly/incorrectly classified instances.

Confusion matrix

Furthermore, we can inspect where a particular misclassification has been made by examining the confusion matrix. Confusion matrix shows how a specific class value was predicted:

```
double[][] confusionMatrix = eval_roc.confusionMatrix();
System.out.println(eval_roc.toMatrixString());
```

The resulting confusion matrix is as follows:

```
=== Confusion Matrix ===

  a  b  c  d  e  f  g  <-- classified as
41  0  0  0  0  0  0 | a = mammal
 0 20  0  0  0  0  0 | b = bird
 0  0  3  1  0  1  0 | c = reptile
 0  0  0 13  0  0  0 | d = fish
 0  0  1  0  3  0  0 | e = amphibian
 0  0  0  0  0  5  3 | f = insect
 0  0  0  0  0  2  8 | g = invertebrate
```

The first column names in the first row correspond to labels assigned by the classification mode. Each additional row then corresponds to an actual true class value. For instance, the second row corresponds instances with the `mammal` true class label. In the column line, we read that all mammals were correctly classified as mammals. In the fourth row, reptiles, we notice that three were correctly classified as reptiles, while one was classified as `fish` and one as an `insect`. Confusion matrix hence, gives us an insight into the kind of errors that our classification model makes.

Choosing a classification algorithm

Naive Bayes is one of the most simple, efficient, and effective inductive algorithms in machine learning. When features are independent, which is rarely true in real world, it is theoretically optimal, and even with dependent features, its performance is amazingly competitive (Zhang, 2004). The main disadvantage is that it cannot learn how features interact with each other, for example, despite the fact that you like your tea with lemon or milk, you hate a tea having both of them at the same time.

Decision tree's main advantage is a model, that is, a tree, which is easy to interpret and explain as we studied in our example. It can handle both nominal and numeric features and you don't have to worry about whether the data is linearly separable.

Some other examples of classification algorithms are as follows:

- `weka.classifiers.rules.ZeroR`: This predicts the majority class and is considered as a baseline, that is, if your classifier's performance is worse than the average value predictor, it is not worth considering it.
- `weka.classifiers.trees.RandomTree`: This constructs a tree that considers K randomly chosen attributes at each node.

- `weka.classifiers.trees.RandomForest`: This constructs a set (that is, forest) of random trees and uses majority voting to classify a new instance.
- `weka.classifiers.lazy.IBk`: This is the k-nearest neighbor's classifier that is able to select an appropriate value of neighbors based on cross-validation.
- `weka.classifiers.functions.MultilayerPerceptron`: This is a classifier based on neural networks that use back-propagation to classify instances. The network can be built by hand, or created by an algorithm, or both.
- `weka.classifiers.bayes.NaiveBayes`: This is a naive Bayes classifier that uses estimator classes, where numeric estimator precision values are chosen based on the analysis of the training data.
- `weka.classifiers.meta.AdaBoostM1`: This is the class for boosting a nominal class classifier using the **AdaBoost M1** method. Only nominal class problems can be tackled. This often dramatically improves the performance, but sometimes it overfits.
- `weka.classifiers.meta.Bagging`: This is the class for bagging a classifier to reduce the variance. This can perform classification and regression, depending on the base learner.

Regression

We will explore basic regression algorithms through analysis of energy efficiency dataset (Tsanas and Xifara, 2012). We will investigate the heating and cooling load requirements of the buildings based on their construction characteristics such as surface, wall and roof area, height, hazing area, and compactness. The researchers used a simulator to design 12 different house configurations while varying 18 building characteristics. In total, 768 different buildings were simulated.

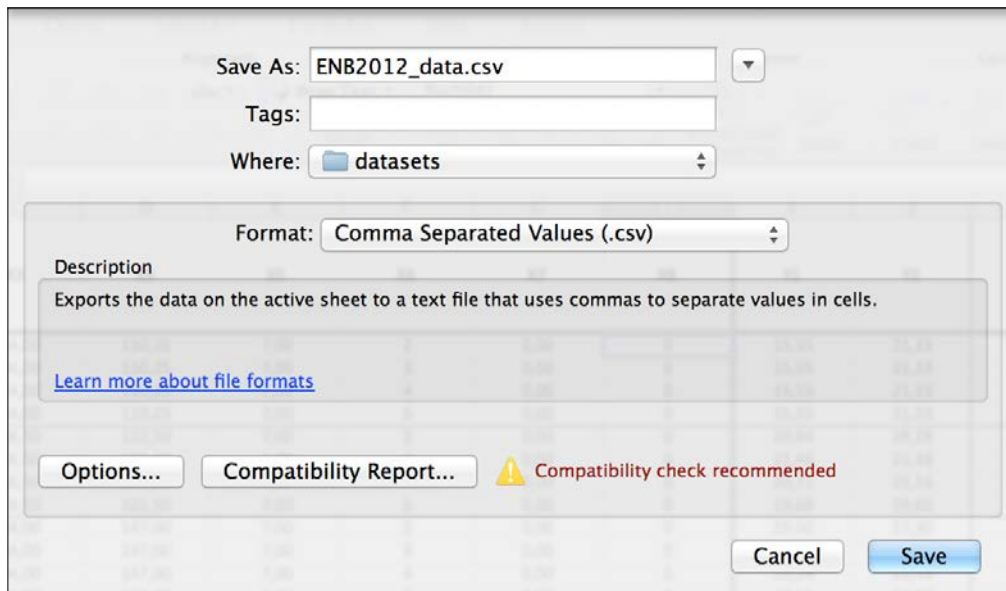
Our first goal is to systematically analyze the impact each building characterizes has on the target variable, that is, heating or cooling load. The second goal is to compare the performance of a classical linear regression model against other methods, such as SVM regression, random forests, and neural networks. For this task, we will use the Weka library.

Loading the data

Download the energy efficiency dataset from

<https://archive.ics.uci.edu/ml/datasets/Energy+efficiency>.

The dataset is in Excel's XLSX format, which cannot be read by Weka. We can convert it to a **Comma Separated Value (CSV)** format by clicking **File | Save As...** and picking CSV in the saving dialog as shown in the following screenshot. Confirm to save only the active sheet (since all others are empty) and confirm to continue to lose some formatting features. Now, the file is ready to be loaded by Weka:



Open the file in a text editor and inspect if the file was indeed correctly transformed. There might be some minor issues that may be potentially causing problems. For instance, in my export, each line ended with a double semicolon, as follows:

```
X1;X2;X3;X4;X5;X6;X7;X8;Y1;Y2;;
0,98;514,50;294,00;110,25;7,00;2;0,00;0;15,55;21,33;;
0,98;514,50;294,00;110,25;7,00;3;0,00;0;15,55;21,33;;
```

To remove the doubled semicolon, we can use the **Find and Replace** function: find ";;" and replace it with ";".

The second problem was that my file had a long list of empty lines at the end of the document, which can be simply deleted:

```
0,62;808,50;367,50;220,50;3,50;5;0,40;5;16,64;16,03;;
;;;;;;;;;;
;;;;;;;;;;
```


Now, we are ready to load the data. Let's open a new file and write a simple data import function using Weka's converter for reading files in CSV format:

```
import weka.core.Instances;
import weka.core.converters.CSVLoader;
import java.io.File;
import java.io.IOException;

public class EnergyLoad {

    public static void main(String[] args) throws IOException {

        // load CSV
        CSVLoader loader = new CSVLoader();
        loader.setSource(new File(args[0]));
        Instances data = loader.getDataSet();

        System.out.println(data);
    }
}
```

The data is loaded. Let's move on.

Analyzing attributes

Before we analyze attributes, let's first try to understand what we are dealing with. In total, there are eight attributes describing building characteristic and two target variables, heating and cooling load, as shown in the following table:

Attribute	Attribute name
X1	Relative compactness
X2	Surface area
X3	Wall area
X4	Roof area
X5	Overall height
X6	Orientation
X7	Glazing area
X8	Glazing area distribution
Y1	Heating load
Y2	Cooling load

Building and evaluating regression model

We will start with learning a model for heating load by setting the class attribute at the feature position:

```
data.setClassIndex(data.numAttributes() - 2);
```

The second target variable—cooling load—can be now removed:

```
//remove last attribute Y2
Remove remove = new Remove();
remove.setOptions(new String[]{"-R", data.numAttributes()+""});
remove.setInputFormat(data);
data = Filter.useFilter(data, remove);
```

Linear regression

We will start with a basic linear regression model implemented with the `LinearRegression` class. Similarly as in the classification example, we will initialize a new model instance, pass parameters and data, and invoke the `buildClassifier(Instances)` method, as follows:

```
import weka.classifiers.functions.LinearRegression;
...

data.setClassIndex(data.numAttributes() - 2);
LinearRegression model = new LinearRegression();
model.buildClassifier(data);
System.out.println(model);
```

The learned model, which is stored in the object, can be outputted by calling the `toString()` method, as follows:

Y1 =

```
-64.774 * X1 +
-0.0428 * X2 +
 0.0163 * X3 +
-0.089 * X4 +
 4.1699 * X5 +
19.9327 * X7 +
 0.2038 * X8 +
83.9329
```

Linear regression model constructed a function that linearly combines the input variables to estimate the heating load. The number in front of the feature explains the feature's impact on the target variable: sign corresponds to positive/negative impact, while magnitude corresponds to its significance. For instance, feature x_1 – relative compactness is negatively correlated with heating load, while glazing area is positively correlated. These two features also significantly impact the final heating load estimation.

The model performance can be similarly evaluated with cross-validation technique.

The 10-fold cross-validation is as follows:

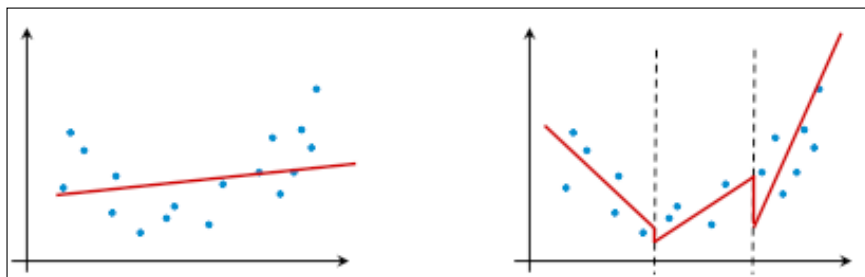
```
Evaluation eval = new Evaluation(data);
eval.crossValidateModel(
    model, data, 10, new Random(1), new String[]{});
System.out.println(eval.toSummaryString());
```

We can output the common evaluation metrics including correlation, mean absolute error, relative absolute error, and so on, as follows:

Correlation coefficient	0.956
Mean absolute error	2.0923
Root mean squared error	2.9569
Relative absolute error	22.8555 %
Root relative squared error	29.282 %
Total Number of Instances	768

Regression trees

Another approach is to construct a set of regression models, each on its own part of the data. The following diagram shows the main difference between a regression model and a regression tree. Regression model constructs a single model that best fits all the data. Regression tree, on the other hand, constructs a set of regression models, each modeling a part of the data as shown on the right-hand side. Compared to the regression model, the regression tree can better fit the data, but the function is a piece-wise linear with jumps between modeled regions:



Regression tree in Weka is implemented within the `M5` class. Model construction follows the same paradigm: initialize model, pass parameters and data, and invoke the `buildClassifier(Instances)` method.

```
import weka.classifiers.trees.M5P;
...
M5P md5 = new M5P();
md5.setOptions(new String[]{" "});
md5.buildClassifier(data);
System.out.println(md5);
```

The induced model is a tree with equations in the leaf nodes, as follows:

M5 pruned model tree:

(using smoothed linear models)

```
X1 <= 0.75 :
|   X7 <= 0.175 :
|   |   X1 <= 0.65 : LM1 (48/12.841%)
|   |   X1 > 0.65 : LM2 (96/3.201%)
|   |   X7 > 0.175 :
|   |   |   X1 <= 0.65 : LM3 (80/3.652%)
|   |   |   X1 > 0.65 : LM4 (160/3.502%)
X1 > 0.75 :
|   X1 <= 0.805 : LM5 (128/13.302%)
|   X1 > 0.805 :
|   |   X7 <= 0.175 :
|   |   |   X8 <= 1.5 : LM6 (32/20.992%)
|   |   |   X8 > 1.5 :
|   |   |   |   X1 <= 0.94 : LM7 (48/5.693%)
|   |   |   |   X1 > 0.94 : LM8 (16/1.119%)
|   |   |   X7 > 0.175 :
|   |   |   |   X1 <= 0.84 :
|   |   |   |   |   X7 <= 0.325 : LM9 (20/5.451%)
|   |   |   |   |   X7 > 0.325 : LM10 (20/5.632%)
|   |   |   |   X1 > 0.84 :
|   |   |   |   |   X7 <= 0.325 : LM11 (60/4.548%)
|   |   |   |   |   X7 > 0.325 :
|   |   |   |   |   |   X3 <= 306.25 : LM12 (40/4.504%)
|   |   |   |   |   |   X3 > 306.25 : LM13 (20/6.934%)
```

LM num: 1

Y1 =

$$\begin{aligned} &72.2602 * X1 \\ &+ 0.0053 * X3 \\ &+ 11.1924 * X7 \\ &+ 0.429 * X8 \\ &- 36.2224 \end{aligned}$$

...

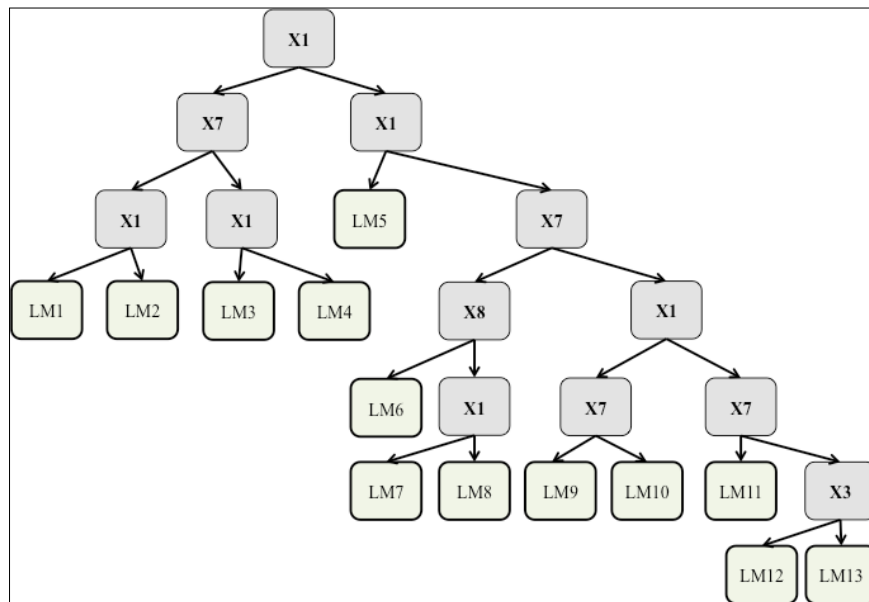
LM num: 13

Y1 =

$$\begin{aligned} &5.8829 * X1 \\ &+ 0.0761 * X3 \\ &+ 9.5464 * X7 \\ &- 0.0805 * X8 \\ &+ 2.1492 \end{aligned}$$

Number of Rules : 13

The tree has 13 leaves, each corresponding to a linear equation. The preceding output is visualized in the following diagram:



The tree can be read similar to a classification tree. The most important features are at the top of the tree. The terminal node, leaf, contains a linear regression model explaining the data that reach this part of the tree.

Evaluation outputs the following results:

Correlation coefficient	0.9943
Mean absolute error	0.7446
Root mean squared error	1.0804
Relative absolute error	8.1342 %
Root relative squared error	10.6995 %
Total Number of Instances	768

Tips to avoid common regression problems

First, use prior studies and domain knowledge to figure out which features to include in regression. Check literature, reports, and previous studies on what kind of features work and reasonable variables for modeling your problem. Suppose you have a large set of features with random data, it is highly likely that several features will be correlated to the target variable (even though the data is random).

Keep the model simple to avoid overfitting. The Occam's razor principle states that you should select a model that best explains your data with the fewest assumptions. In practice, the model can be as simple as 2-4 predictor features.

Clustering

Compared to a supervised classifier, the goal of clustering is to identify intrinsic groups in a set of unlabeled data. It could be applied in identifying representative examples of homogeneous groups, finding useful and suitable groupings, or finding unusual examples, such as outliers.

We'll demonstrate how to implement clustering by analyzing the Bank dataset. The dataset consist of 11 attributes, describing 600 instances with age, sex, region, income, marriage status, children, car ownership status, saving activity, current activity, mortgage status, and PEP. In our analysis, we will try to identify the common groups of clients by applying the **Expectation Maximization (EM)** clustering.

EM works as follows: given a set of clusters, EM first assigns each instance with a probability distribution of belonging to a particular cluster. For example, if we start with three clusters – *A*, *B*, and *C* – an instance might get the probability distribution of *0.70*, *0.10*, and *0.20*, belonging to the *A*, *B*, and *C* clusters, respectively. In the second step, EM re-estimates the parameter vector of the probability distribution of each class. The algorithm iterates these two steps until the parameters converge or the maximum number of iterations is reached.

The number of clusters to be used in EM can be set either manually or automatically by cross validation. Another approach to determining the number of clusters in a dataset includes the **elbow** method. The method looks at the percentage of variance that is explained with a specific number of clusters. The method suggests increasing the number of clusters until the additional cluster does not add much information, that is, explains little additional variance.

Clustering algorithms

The process of building a cluster model is quite similar to the process of building a classification model, that is, load the data and build a model. Clustering algorithms are implemented in the `weka.clusterers` package, as follows:

```
import java.io.BufferedReader;
import java.io.FileReader;

import weka.core.Instances;
import weka.clusterers.EM;

public class Clustering {

    public static void main(String args[]) throws Exception{

        //load data
        Instances data = new Instances(new BufferedReader
            (new FileReader(args[0])));

        // new instance of clusterer
        EM model = new EM();
        // build the clusterer
        model.buildClusterer(data);
        System.out.println(model);

    }
}
```

The model identified the following six clusters:

EM

==

Number of clusters selected by cross validation: 6

Attribute	Cluster					
	0 (0.1)	1 (0.13)	2 (0.26)	3 (0.25)	4 (0.12)	5 (0.14)
=====						
age						
0_34	10.0535	51.8472	122.2815	12.6207	3.1023	1.0948
35_51	38.6282	24.4056	29.6252	89.4447	34.5208	3.3755
52_max	13.4293	6.693	6.3459	50.8984	37.861	81.7724
[total]	62.1111	82.9457	158.2526	152.9638	75.4841	86.2428
sex						
FEMALE	27.1812	32.2338	77.9304	83.5129	40.3199	44.8218
MALE	33.9299	49.7119	79.3222	68.4509	34.1642	40.421
[total]	61.1111	81.9457	157.2526	151.9638	74.4841	85.2428
region						
INNER_CITY	26.1651	46.7431	73.874	60.1973	33.3759	34.6445
TOWN	24.6991	13.0716	48.4446	53.1731	21.617	
17.9946						
...						

The table can be read as follows: the first line indicates six clusters, while the first column shows attributes and their ranges. For example, the attribute `age` is split into three ranges: 0-34, 35-51, and 52-max. The columns on the left indicate how many instances fall into the specific range in each cluster, for example, clients in the 0-34 years age group are mostly in cluster #2 (122 instances).

Evaluation

A clustering algorithm's quality can be estimated using the log likelihood measure, which measures how consistent the identified clusters are. The dataset is split into multiple folds and clustering is run with each fold. The motivation here is that if the clustering algorithm assigns high probability to similar data that wasn't used to fit parameters, then it has probably done a good job of capturing the data structure. Weka offers the `ClusterEvaluation` class to estimate it, as follows:

```
double logLikelihood = ClusterEvaluation.crossValidateModel(model, data,
10, new Random(1));
System.out.println(logLikelihood);
```

It has the following output:

```
-8.773410259774291
```

Summary

In this chapter, you learned how to implement basic machine learning tasks with Weka: classification, regression, and clustering. We briefly discussed attribute selection process, trained models, and evaluated their performance.

The next chapter will focus on how to apply these techniques to solve real-life problems, such as customer retention.

4

Customer Relationship Prediction with Ensembles

Any type of company offering a service, product, or experience needs a solid understanding of relationship with their customers; therefore, **Customer Relationship Management (CRM)** is a key element of modern marketing strategies. One of the biggest challenges that businesses face is the need to understand exactly what causes a customer to buy new products.

In this chapter, we will work on a real-world marketing database provided by the French telecom company, Orange. The task will be to estimate the following likelihoods for customer actions:

- Switch provider (churn)
- Buy new products or services (appetency)
- Buy upgrades or add-ons proposed to them to make the sale more profitable (upselling)

We will tackle the **Knowledge Discovery and Data Mining (KDD) Cup 2009** challenge (KDD Cup, 2009) and show the steps to process the data using Weka. First, we will parse and load the data and implement the basic baseline models. Later, we will address advanced modeling techniques, including data pre-processing, attribute selection, model selection, and evaluation.



KDD Cup is the leading data mining competition in the world. It is organized annually by ACM **Special Interest Group on Knowledge Discovery and Data Mining**. The winners are announced at the **Conference on Knowledge Discovery and Data Mining**, which is usually held in August.

Yearly archives, including all the corresponding datasets, are available here: <http://www.kdd.org/kdd-cup>.

Customer relationship database

The most practical way to build knowledge on customer behavior is to produce scores that explain a target variable such as churn, appetency, or upselling. The score is computed by a model using input variables describing customers, for example, current subscription, purchased devices, consumed minutes, and so on. The scores are then used by the information system, for example, to provide relevant personalized marketing actions.

In 2009, the conference on KDD organized a machine learning challenge on customer-relationship prediction (KDD Cup, 2009).

Challenge

Given a large set of customer attributes, the task was to estimate the following three target variables (KDD Cup, 2009):

- **Churn probability**, in our context, is the likelihood a customer will switch providers:

Churn rate is also sometimes called attrition rate. It is one of two primary factors that determine the steady-state level of customers a business will support. In its broadest sense, churn rate is a measure of the number of individuals or items moving into or out of a collection over a specific period of time. The term is used in many contexts, but is most widely applied in business with respect to a contractual customer base. For instance, it is an important factor for any business with a subscriber-based service model, including mobile telephone networks and pay TV operators. The term is also used to refer to participant turnover in peer-to-peer networks.

- **Appetency probability**, in our context, is the propensity to buy a service or product
- **Upselling probability** is the likelihood that a customer will buy an add-on or upgrade:

Upselling is a sales technique whereby a salesman attempts to have the customer purchase more expensive items, upgrades, or other add-ons in an attempt to make a more profitable sale. Upselling usually involves marketing more profitable services or products, but upselling can also be simply exposing the customer to other options he or she may not have considered previously. Upselling can imply selling something additional, or selling something that is more profitable or otherwise preferable for the seller instead of the original sale.

The challenge was to beat the in-house system developed by Orange Labs. This was an opportunity for the participants to prove that they could handle a large database, including heterogeneous noisy data and unbalanced class distributions.

Dataset

For the challenge, the company Orange released a large dataset of customer data, containing about one million customers, described in ten tables with hundreds of fields. In the first step, they resampled the data to select a less unbalanced subset containing 100,000 customers. In the second step, they used an automatic feature construction tool that generated 20,000 features describing customers, which was then narrowed down to 15,000 features. In the third step, the dataset was anonymized by randomizing the order of features, discarding attribute names, replacing nominal variables with randomly generated strings, and multiplying continuous attributes by a random factor. Finally, all the instances were split randomly into a train and test dataset.

The KDD Cup provided two sets of data: large set and small set, corresponding to fast and slow challenge, respectively. They are described at the KDD Cup site as follows:

Both training and test sets contain 50,000 examples. The data are split similarly for the small and large versions, but the samples are ordered differently within the training and within the test sets. Both small and large datasets have numerical and categorical variables. For the large dataset, the first 14,740 variables are numerical and the last 260 are categorical. For the small dataset, the first 190 variables are numerical and the last 40 are categorical.

In this chapter, we will work with the small dataset consisting of 50,000 instances described with 230 variables each. Each of the 50,000 rows of data correspond to a client and are associated with three binary outcomes – one for each of the three challenges (upsell, churn, and appetency).

To make this clearer, the following image illustrates the dataset in a table format:

230 numeric and nominal attributes										Three binary classes		
Var85	Var123	Var125	Var126	Var132	Var133	Var134	Var225	Var229	Var230	Label Churn	Label Appetency	Label Upselling
12	6	720	8	0	1212385	69134				-1	-1	-1
2	72	0		8	4136430	357038				1	-1	-1
58	114	5967	-28	0	3478905	248932	kG3k	am7c		-1	-1	-1
0	0	0	-14	0	0	0				-1	-1	-1
0	0	15111	58	0	150650	66046	kG3k	mj86		-1	-1	-1
10	0	1935		8	641020	43684		am7c		-1	-1	-1
16	24	13194	-24	0	1664450	104978	kG3k	am7c		-1	-1	-1
2	12	0	-8	8	3839825	1284128				-1	-1	-1
2	90	2754		0	3830510	203586	kG3k	am7c		-1	-1	-1
24	66	6561		32	2577245	210014	kG3k			-1	-1	-1
6	12	5823	58	0	0	7134	kG3k	mj86		-1	-1	-1
28	24	66825	52	8	134105	15166	kG3k			-1	-1	-1
0	0	44154	10	0	0	0		mj86		-1	-1	-1
22	54	5202		0	2772010	1095062	xG3x			-1	-1	-1
0	102	31104	8	0	2170355	57596				-1	-1	1
0	0	2574		0	0	0	ELof	oJmt		-1	-1	-1
14	186	8019		48	3571845	587392	kG3k	am7c		-1	-1	-1
0	30	5319		8	500295	31436		am7c		-1	-1	-1
2	0	13788	4	0	918350	0	kG3k			-1	-1	-1
14	0	7110		0	2055150	392138				1	-1	-1
8	66	0	-8	0	3258940	1121306				-1	-1	-1
0	18	0	-10	0	0	0				-1	-1	-1
12	0	531	36	0	491345	56742	ELof	mj86		-1	-1	-1
0	12	16803	12	0	201110	1693090				1	-1	-1
14	0	25740		0	2932660	313200	xG3x			-1	-1	1

The table depicts the first 25 instances, that is, customers, each described with 250 attributes. For this example, only a selected subset of 10 attributes is shown. The dataset contains many missing values and even empty or constant attributes. The last three columns of the table correspond to the three distinct class labels corresponding to the ground truth, that is, if the customer indeed switched the provider (churn), bought a service (appetency), or bought an upgrade (upsell). However, note that the labels are provided separately from the data in three distinct files, hence, it is essential to retain the order of the instances and corresponding class labels to ensure proper correspondence.

Evaluation

The submissions were evaluated according to the arithmetic mean of the area under the ROC curve (AUC) for the three tasks, that is, churn, appetency, and upselling. ROC curve shows the performance of model as a curve obtained by plotting sensitivity against specificity for various threshold values used to determine the classification result (refer to *Chapter 1, Applied Machine Learning Quick Start*, section *ROC curves*). Now, the AUC is related to the area under this curve, meaning larger the area, better the classifier. Most toolboxes, including Weka, provide an API to calculate AUC score.

Basic naive Bayes classifier baseline

As per the rules of the challenge, the participants had to outperform the basic naive Bayes classifier to qualify for prizes, which makes an assumption that features are independent (refer to *Chapter 1, Applied Machine Learning Quick Start*).

The KDD Cup organizers run the vanilla naive Bayes classifier, without any feature selection or hyperparameter adjustments. For the large dataset, the overall scores of the naive Bayes on the test set were as follows:

- **Churn problem:** $AUC = 0.6468$
- **Appetency problem:** $AUC = 0.6453$
- **Upselling problem:** $AUC=0.7211$

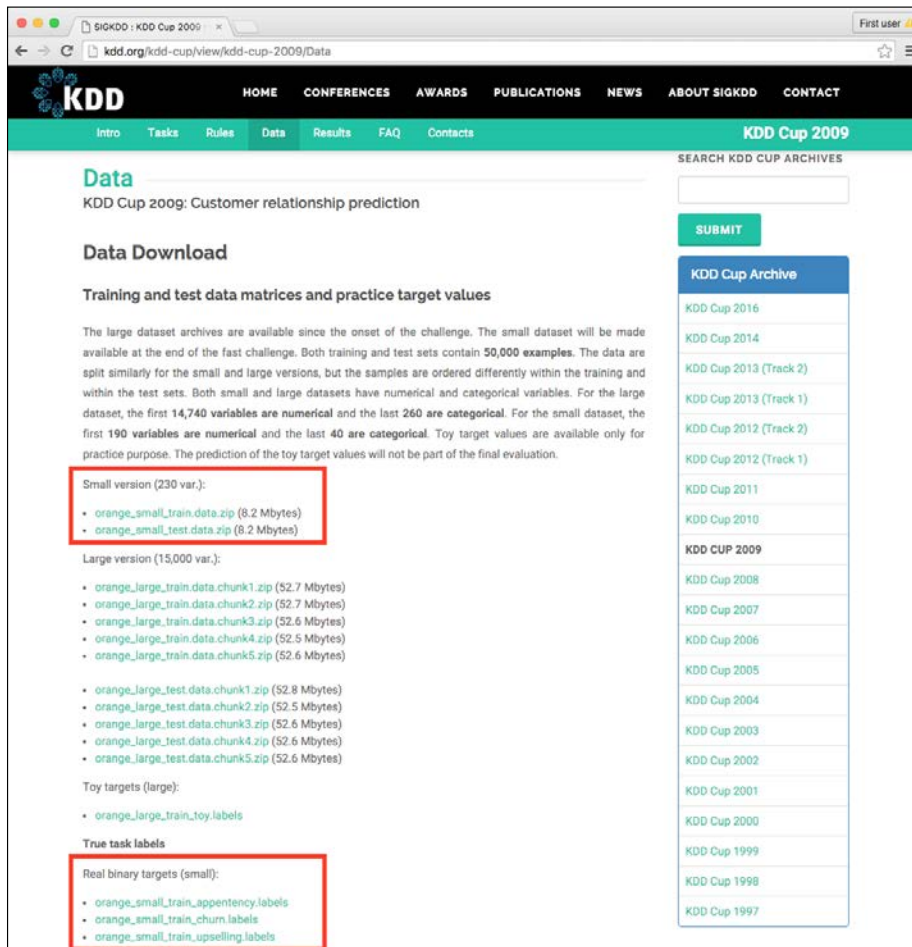
Note that the baseline results are reported for large dataset only. Moreover, while both training and test datasets are provided at the KDD Cup site, the actual true labels for the test set are not provided. Therefore, when we process the data with our models, there is no way to know how well the models will perform on the test set. What we will do is use only the training data and evaluate our models with cross validation. The results will not be directly comparable, but, nevertheless, we have an idea for what a reasonable magnitude of the AUC score is.

Getting the data

At the KDD Cup web page (<http://kdd.org/kdd-cup/view/kdd-cup-2009/Data>), you should see a page that looks like the following screenshot. First, under the **Small version (230 var.)** header, download `orange_small_train.data.zip`. Next, download the three sets of true labels associated with this training data. The following files are found under the **Real binary targets (small)** header:

- `orange_small_train_appentency.labels`
- `orange_small_train_churn.labels`
- `orange_small_train_upselling.labels`

Save and unzip all the files marked in the red boxes, as shown in the following screenshot:



In the following sections, we will first load the data into Weka and apply basic modeling with the naive Bayes to obtain our own baseline AUC scores. Later, we will look into more advanced modeling techniques and tricks.


Loading the data

We will load the data to Weka directly from the `.csvs` format. For this purpose, we will write a function that accepts the path to the data file and the true labels file. The function will load and merge both datasets and remove empty attributes:

```
public static Instances loadData(String pathData, String
    pathLabels) throws Exception {
```

First, we load the data using the `CSVLoader()` class. Additionally, we specify the `\t` tab as a field separator and force the last 40 attributes to be parsed as nominal:

```
// Load data
CSVLoader loader = new CSVLoader();
loader.setFieldSeparator("\t");
loader.setNominalAttributes("191-last");
loader.setSource(new File(pathData));
Instances data = loader.getDataSet();
```

 The `CSVLoader` class accepts many additional parameters specifying column separator, string enclosures, whether a header row is present or not, and so on. Complete documentation is available here:
<http://weka.sourceforge.net/doc.dev/weka/core/converters/CSVLoader.html>

Next, some of the attributes do not contain a single value and Weka automatically recognizes them as the `String` attributes. We actually do not need them, so we can safely remove them using the `RemoveType` filter. Additionally, we specify the `-T` parameters, which means remove attribute of specific type and the attribute type that we want to remove:

```
// remove empty attributes identified as String attribute
RemoveType removeString = new RemoveType();
removeString.setOptions(new String[]{"-T", "string"});
removeString.setInputFormat(data);
Instances filteredData = Filter.useFilter(data, removeString);
```

Alternatively, we could use the `void deleteStringAttributes()` method implemented within the `Instances` class, which has the same effect, for example, `data.removeStringAttributes()`.

Now, we will load and assign class labels to the data. We will again utilize `CSVLoader`, where we specify that the file does not have any header line, that is, `setNoHeaderRowPresent(true)`:

```
// Load labels
loader = new CSVLoader();
loader.setFieldSeparator("\t");
loader.setNoHeaderRowPresent(true);
loader.setNominalAttributes("first-last");
loader.setSource(new File(pathLabels));
Instances labels = loader.getDataSet();
```

Once we have loaded both files, we can merge them together by calling the `Instances.mergeInstances(Instances, Instances)` static method. The method returns a new dataset that has all the attributes from the first dataset plus the attributes from the second set. Note that the number of instances in both datasets must be the same:

```
// Append label as class value
Instances labeledData = Instances.mergeInstances(filteredData,
    labels);
```

Finally, we set the last attribute, that is, the label attribute that we have just added, as a target variable and return the resulting dataset:

```
// set the label attribute as class
labeledData.setClassIndex(labeledData.numAttributes() - 1);

System.out.println(labeledData.toSummaryString());
return labeledData;
}
```

The function outputs a summary as shown in the following and returns the labeled dataset:

```
Relation Name: orange_small_train.data-weka.filters.unsupervised.
attribute.RemoveType-Tstring_orange_small_train_churn.labels.txt
Num Instances: 50000
Num Attributes: 215
```

Name	Type	Nom	Int	Real	Missing	Unique	Dist
1 Var1	Num	0%	1%	0%	49298 / 99%	8 / 0%	18
2 Var2	Num	0%	2%	0%	48759 / 98%	1 / 0%	2
3 Var3	Num	0%	2%	0%	48760 / 98%	104 / 0%	146
4 Var4	Num	0%	3%	0%	48421 / 97%	1 / 0%	4
...							

Basic modeling

In this section, we will implement our own baseline model by following the approach that the KDD Cup organizers took. However, before we go to the model, let's first implement the evaluation engine that will return AUC on all three problems.

Evaluating models

Now, let's take a closer look at the evaluation function. The evaluation function accepts an initialized model, cross-validates the model on all three problems, and reports the results as an area under the ROC curve (AUC), as follows:

```
public static double[] evaluate(Classifier model)
    throws Exception {

    double results[] = new double[4];

    String[] labelFiles = new String[]{
        "churn", "appetency", "upselling"};

    double overallScore = 0.0;
    for (int i = 0; i < labelFiles.length; i++) {
```

First, we call the Instance `loadData(String, String)` function that we implemented earlier to load the train data and merge it with the selected labels:

```
// Load data
Instances train_data = loadData(
    path + "orange_small_train.data",
    path+"orange_small_train_"+labelFiles[i]+".labels.txt");
```

Next, we initialize the `weka.classifiers.Evaluation` class and pass our dataset (the dataset is used only to extract data properties, the actual data are not considered). We call the void `crossValidateModel(Classifier, Instances, int, Random)` method to begin cross validation and select to create five folds. As validation is done on random subsets of the data, we need to pass a random seed as well:

```
// cross-validate the data
Evaluation eval = new Evaluation(train_data);
eval.crossValidateModel(model, train_data, 5,
    new Random(1));
```

After the evaluation completes, we read the results by calling the `double areUnderROC(int)` method. As the metric depends on the target value that we are interested in, the method expects a class value index, which can be extracted by searching the index of the "1" value in the class attribute:

```
// Save results
results[i] = eval.areaUnderROC(
    train_data.classAttribute().indexOfValue("1"));
overallScore += results[i];
}
```

Finally, the results are averaged and returned:

```
// Get average results over all three problems
results[3] = overallScore / 3;
return results;
}
```

Implementing naive Bayes baseline

Now, when we have all the ingredients, we can replicate the naive Bayes approach that we are expected to outperform. This approach will not include any additional data pre-processing, attribute selection, and model selection. As we do not have true labels for test data, we will apply the five-fold cross validation to evaluate the model on a small dataset.

First, we initialize a naive Bayes classifier, as follows:

```
Classifier baselineNB = new NaiveBayes();
```

Next, we pass the classifier to our evaluation function, which loads the data and applies cross validation. The function returns an area under the ROC curve score for all three problems and overall results:

```
double resNB[] = evaluate(baselineNB);
System.out.println("Naive Bayes\n" +
    "\tchurn:      " + resNB[0] + "\n" +
    "\tappetency:  " + resNB[1] + "\n" +
    "\tup-sell:    " + resNB[2] + "\n" +
    "\toverall:   " + resNB[3] + "\n");
```

In our case, the model achieves the following results:

Naive Bayes

```
churn:      0.5897891153549814
appetency:  0.630778394752436
```

```
up-sell:    0.6686116692438094
overall:    0.6297263931170756
```

These results will serve as a baseline when we tackle the challenge with more advanced modeling. If we process the data with significantly more sophisticated, time-consuming, and complex techniques, we expect the results to be much better. Otherwise, we are simply wasting the resources. In general, when solving machine learning problems, it is always a good idea to create a simple baseline classifier that serves us as an orientation point.

Advanced modeling with ensembles

In the previous section, we implemented an orientation baseline, so let's focus on heavy machinery. We will follow the approach taken by the KDD Cup 2009 winning solution developed by the **IBM Research** team (Niculescu-Mizil and others, 2009).

Their strategy to address the challenge was using the **Ensemble Selection** algorithm (Caruana and Niculescu-Mizil, 2004). This is an ensemble method, which means it constructs a series of models and combines their output in a specific way to provide the final classification. It has several desirable properties as shown in the following list that make it a good fit for this challenge:

- It was proven to be robust, yielding excellent performance
- It can be optimized for a specific performance metric, including AUC
- It allows different classifiers to be added to the library
- It is an anytime method, meaning that, if we run out of time, we have a solution available

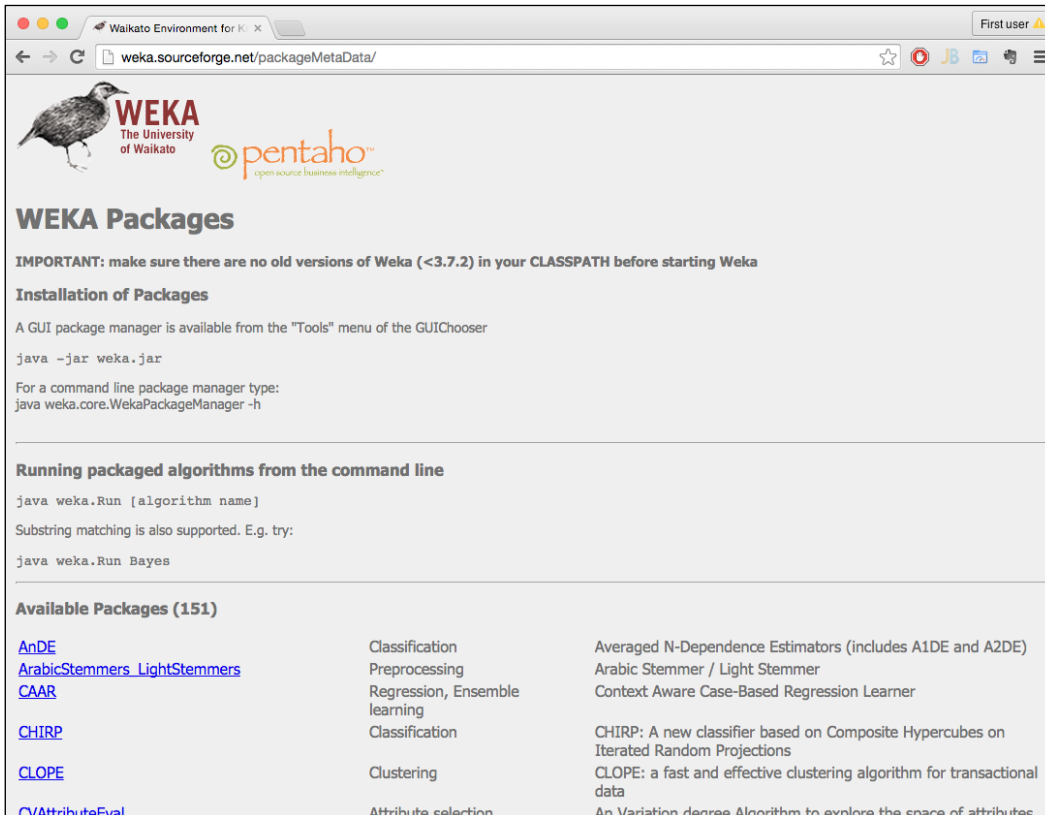
In this section, we will loosely follow the steps as described in their report. Note, this is not an exact implementation of their approach, but rather a solution overview that will include the necessary steps to dive deeper.

The general overview of steps is as follows:

1. First, we will preprocess the data by removing attributes that clearly do not bring any value, for example, all the missing or constant values; fixing missing values in order to help machine learning algorithms, which cannot deal with them; and converting categorical attributes to numerical.
2. Next, we will run attributes selection algorithm to select only a subset of attribute that can help in prediction of tasks.
3. In the third step, we will instantiate the Ensemble Selection algorithms with a wide variety of models, and, finally, evaluate the performance.

Before we start

For this task, we will need an additional Weka package, `ensembleLibrary`. Weka 3.7.2 or higher versions support external packages developed mainly by the academic community. A list of **WEKA Packages** is available at <http://weka.sourceforge.net/packageMetaData/> as shown at the following screenshot:



Find and download the latest available version of the `ensembleLibrary` package at <http://prdownloads.sourceforge.net/weka/ensembleLibrary1.0.5.zip?do=download>.

After you unzip the package, locate `ensembleLibrary.jar` and import it to your code, as follows:

```
import weka.classifiers.meta.EnsembleSelection;
```

Data pre-processing

First, we will utilize Weka's built-in `weka.filters.unsupervised.attribute.RemoveUseless` filter, which works exactly as its name suggests. It removes the attributes that do not vary much, for instance, all constant attributes are removed, and attributes that vary too much, almost at random. The maximum variance, which is applied only to nominal attributes, is specified with the `-M` parameter. The default parameter is 99%, which means that if more than 99% of all instances have unique attribute values, the attribute is removed, as follows:

```
RemoveUseless removeUseless = new RemoveUseless();
removeUseless.setOptions(new String[] { "-M", "99" }); // threshold
removeUseless.setInputFormat(data);
data = Filter.useFilter(data, removeUseless);
```

Next, we will replace all the missing values in the dataset with the modes (nominal attributes) and means (numeric attributes) from the training data by using the `weka.filters.unsupervised.attribute.ReplaceMissingValues` filter. In general, missing values replacement should be proceeded with caution while taking into consideration the meaning and context of the attributes:

```
ReplaceMissingValues fixMissing = new ReplaceMissingValues();
fixMissing.setInputFormat(data);
data = Filter.useFilter(data, fixMissing);
```

Finally, we will discretize numeric attributes, that is, we transform numeric attributes into intervals using the `weka.filters.unsupervised.attribute.Discretize` filter. With the `-B` option, we set to split numeric attributes into four intervals, and the `-R` option will specify the range of attributes (only numeric attributes will be discretized):

```
Discretize discretizeNumeric = new Discretize();
discretizeNumeric.setOptions(new String[] {
    "-B", "4", // no of bins
    "-R", "first-last"}); //range of attributes
fixMissing.setInputFormat(data);
data = Filter.useFilter(data, fixMissing);
```

Attribute selection

In the next step, we will select only informative attributes, that is, attributes that more likely help with prediction. A standard approach to this problem is to check the information gain carried by each attribute. We will use the `weka.attributeSelection.AttributeSelection` filter, which requires two additional methods: `evaluator`, that is, how attribute usefulness is calculated, and `search algorithms`, that is, how to select a subset of attributes.

In our case, we first initialize `weka.attributeSelection.InfoGainAttributeEval` that implements calculation of information gain:

```
InfoGainAttributeEval eval = new InfoGainAttributeEval();
Ranker search = new Ranker();
```

To select only top attributes above some threshold, we initialize `weka.attributeSelection.Ranker` to rank the attributes with information gain above a specific threshold. We specify this with the `-T` parameter, while keeping the value of the threshold low to keep the attributes with at least some information:

```
search.setOptions(new String[] { "-T", "0.001" });
```



The general rule for settings this threshold is to sort the attributes by information gain and pick the threshold where the information gain drops to negligible value.



Next, we can initialize the `AttributeSelection` class, set the evaluator and ranker, and apply the attribute selection to our dataset:

```
AttributeSelection attSelect = new AttributeSelection();
attSelect.setEvaluator(eval);
attSelect.setSearch(search);

// apply attribute selection
attSelect.SelectAttributes(data);
```

Finally, we remove the attributes that were not selected in the last run by calling the `reduceDimensionality(Instances)` method.

```
// remove the attributes not selected in the last run
data = attSelect.reduceDimensionality(data);
```

At the end, we are left with 214 out of 230 attributes.

Model selection

Over the years, practitioners in the field of machine learning have developed a wide variety of learning algorithms and improvements to the existing ones. There are so many unique supervised learning methods that it is challenging to keep track of all of them. As characteristics of the datasets vary, no one method is the best in all the cases, but different algorithms are able to take advantage of different characteristics and relationships of a given dataset. The property the Ensemble Selection algorithm is to try to leverage (Jung, 2005):

Intuitively, the goal of ensemble selection algorithm is to automatically detect and combine the strengths of these unique algorithms to create a sum that is greater than the parts. This is accomplished by creating a library that is intended to be as diverse as possible to capitalize on a large number of unique learning approaches. This paradigm of overproducing a huge number of models is very different from more traditional ensemble approaches. Thus far, our results have been very encouraging.

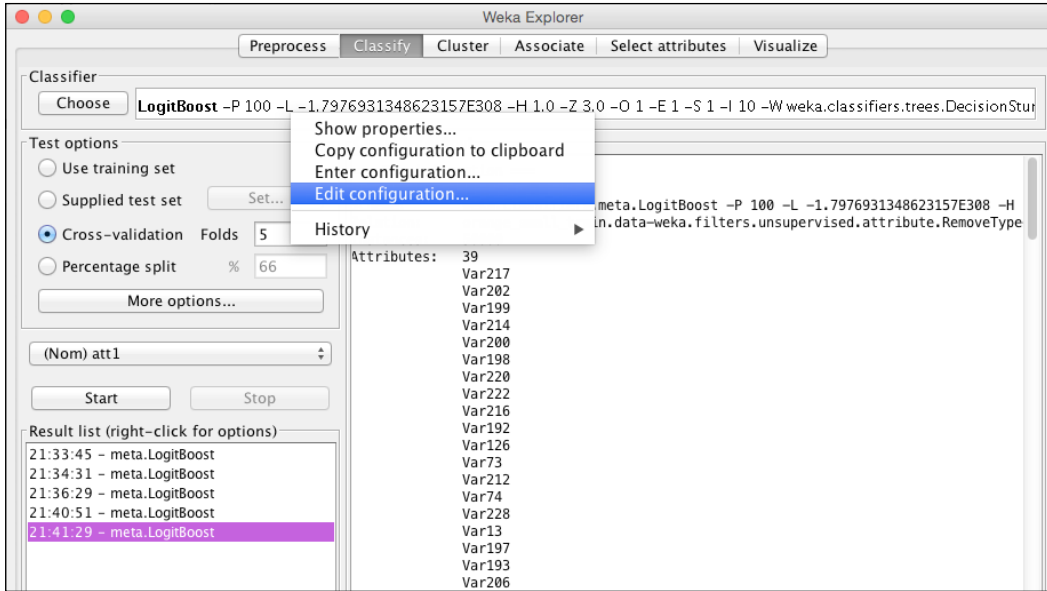
First, we need to create the model library by initializing the `weka.classifiers.EnsembleLibrary` class, which will help us define the models:

```
EnsembleLibrary ensembleLib = new EnsembleLibrary();
```

Next, we add the models and their parameters as strings to the library as string values, for example, we can add three decision tree learners with different parameters, as follows:

```
ensembleLib.addModel("weka.classifiers.trees.J48 -S -C 0.25 -B -M  
2");  
ensembleLib.addModel("weka.classifiers.trees.J48 -S -C 0.25 -B -M  
2 -A");
```


If you are familiar with the Weka graphical interface, you can also explore the algorithms and their configurations there and copy the configuration as shown in the following screenshot: right-click on the algorithm name and navigate to **Edit configuration** | **Copy configuration string**:



To complete the example, we added the following algorithms and their parameters:

- Naive Bayes that was used as default baseline:
`ensembleLib.addModel("weka.classifiers.bayes.NaiveBayes");`
- k-nearest neighbors based on lazy models?:
`ensembleLib.addModel("weka.classifiers.lazy.IBk");`
- Logistic regression as simple logistic with default parameters:
`ensembleLib.addModel("weka.classifiers.functions.SimpleLogistic");`
- Support vector machines with default parameters:
`ensembleLib.addModel("weka.classifiers.functions.SMO");`
- AdaBoost, which is an ensemble method itself:
`ensembleLib.addModel("weka.classifiers.meta.AdaBoostM1");`

- LogitBoost, an ensemble method based on logistic regression:

```
ensembleLib.addModel("weka.classifiers.meta.LogitBoost");
```
- Decision stump, an ensemble method based on one-level decision trees:

```
ensembleLib.addModel("classifiers.trees.DecisionStump");
```

As the `EnsembleLibrary` implementation is primarily focused on GUI and console users, we have to save the models into a file by calling the `saveLibrary(File, EnsembleLibrary, JComponent)` method, as follows:

```
EnsembleLibrary.saveLibrary(new
    File(path+"ensembleLib.model.xml"), ensembleLib, null);
System.out.println(ensembleLib.getModels());
```

Next, we can initialize the Ensemble Selection algorithm by instantiating the `weka.classifiers.meta.EnsembleSelection` class. Let's first review the following method options:

- `-L </path/to/modelLibrary>`: This specifies the `modelLibrary` file, continuing the list of all models.
- `-W </path/to/working/directory>`: This specifies the working directory, where all models will be stored.
- `-B <numModelBags>`: This sets the number of bags, that is, the number of iterations to run the Ensemble Selection algorithm.
- `-E <modelRatio>`: This sets the ratio of library models that will be randomly chosen to populate each bag of models.
- `-V <validationRatio>`: This sets the ratio of the training data set that will be reserved for validation.
- `-H <hillClimbIterations>`: This sets the number of hill climbing iterations to be performed on each model bag.
- `-I <sortInitialization>`: This sets the ratio of the ensemble library that the sort initialization algorithm will be able to choose from, while initializing the ensemble for each model bag.
- `-X <numFolds>`: This sets the number of cross validation folds.
- `-P <hillclimbMetric>`: This specifies the metric that will be used for model selection during the hill climbing algorithm. Valid metrics are accuracy, rmse, roc, precision, recall, fscore, and all.

- `-A <algorithm>`: This specifies the algorithm to be used for ensemble selection. Valid algorithms are forward (default) for forward selection, backward for backward elimination, both for both forward and backward elimination, best to simply print the top performer from the ensemble library, and library to only train the models in the ensemble library.
- `-R`: This flags whether or not the models can be selected more than once for an ensemble.
- `-G`: This states whether the sort initialization greedily stops adding models when the performance degrades.
- `-O`: This is a flag for verbose output. This prints the performance of all the selected models.
- `-S <num>`: This is a random number seed (default 1).
- `-D`: If set, the classifier is run in the debug mode and may output additional information to the console.

We initialize the algorithm with the following initial parameters, where we specified optimizing the ROC metric:

```
EnsembleSelection ensembleSel = new EnsembleSelection();
ensembleSel.setOptions(new String[]{
    "-L", path+"ensembleLib.model.xml", // </path/to/modelLibrary>
    "-W", path+"esTmp", // </path/to/working/directory> -
    "-B", "10", // <numModelBags>
    "-E", "1.0", // <modelRatio>.
    "-V", "0.25", // <validationRatio>
    "-H", "100", // <hillClimbIterations>
    "-I", "1.0", // <sortInitialization>
    "-X", "2", // <numFolds>
    "-P", "roc", // <hillclimbMetric>
    "-A", "forward", // <algorithm>
    "-R", "true", // - Flag to be selected more than once
    "-G", "true", // - stops adding models when performance degrades
    "-O", "true", // - verbose output.
    "-S", "1", // <num> - Random number seed.
    "-D", "true" // - run in debug mode
});
```

Performance evaluation

The evaluation is heavy both computationally and memory-wise, so make sure that you initialize the JVM with extra heap space—for instance, `java -Xmx16g`—while the computation can take a couple of hours or days, depending on the number of algorithms you include in the model library. This example took 4 hours and 22 minutes on 12-core Intel Xeon E5-2420 CPU with 32 GB of memory, utilizing 10% CPU and 6 GB of memory on average.

We call our evaluation method and output the results, as follows:

```
double resES[] = evaluate(ensembleSel);
System.out.println("Ensemble Selection\n"
+ "\tchurn:      " + resES[0] + "\n"
+ "\tappetency:  " + resES[1] + "\n"
+ "\tup-sell:    " + resES[2] + "\n"
+ "\toverall:   " + resES[3] + "\n");
```

The specific set of classifiers in the model library achieved the following result:

```
Ensamble
churn:      0.7109874158176481
appetency:  0.786325687118347
up-sell:    0.8521363243575182
overall:    0.7831498090978378
```

Overall, the approach has brought us to a significant improvement of more than 15 percentage points compared to the initial baseline that we designed at the beginning of the chapter. While it is hard to give a definite answer, the improvement was mainly due to three factors: data pre-processing and attribute selection, exploration of a large variety of learning methods, and use of an ensemble-building technique that is able to take advantage of the variety of base classifiers without overfitting. However, the improvement requires a significant increase in processing time, as well as working memory.

Summary

In this chapter, we tackled the KDD Cup 2009 challenge on customer-relationship prediction, where we implemented the data pre-processing steps, addressing the missing values and redundant attributes. We followed the winning KDD Cup solution, studying how to leverage ensemble methods using a basket of learning algorithms, which can significantly boost the classification performance.

In the next chapter, we will tackle another problem addressing the customer behavior, that is, the analysis of purchasing behavior, where you will learn how to use algorithms that detect frequently occurring patterns.

5

Affinity Analysis

Affinity analysis is the heart of **Market basket analysis (MBA)**. It can discover co-occurrence relationships among activities performed by specific users or groups. In retail, affinity analysis can help you understand the purchasing behavior of customers. These insights can drive revenue through smart cross-selling and upselling strategies and can assist you in developing loyalty programs, sales promotions, and discount plans.

In this chapter, we will look into the following topics:

- Market basket analysis
- Association rule learning
- Other applications in various domains

First, we will revise the core association rule learning concepts and algorithms, such as support, lift, **Apriori algorithm**, and **FP-growth algorithm**. Next, we will use Weka to perform our first affinity analysis on supermarket dataset and study how to interpret the resulting rules. We will conclude the chapter by analyzing how association rule learning can be applied in other domains, such as **IT Operations Analytics**, medicine, and others.

Market basket analysis

Since the introduction of electronic point of sale, retailers have been collecting an incredible amount of data. To leverage this data in order to produce business value, they first developed a way to consolidate and aggregate the data to understand the basics of the business. What are they selling? How many units are moving? What is the sales amount?

Recently, the focus shifted to the lowest level of granularity – the market basket transaction. At this level of detail, the retailers have direct visibility into the market basket of each customer who shopped at their store, understanding not only the quantity of the purchased items in that particular basket, but also how these items were bought in conjunction with each other. This can be used to drive decisions about how to differentiate store assortment and merchandise, as well as effectively combine offers of multiple products, within and across categories, to drive higher sales and profits. These decisions can be implemented across an entire retail chain, by channel, at the local store level, and even for a specific customer with so-called personalized marketing, where a unique product offering is made for each customer.



MBA covers a wide variety of analysis:

- **Item affinity:** This defines the likelihood of two (or more) items being purchased together
- **Identification of driver items:** This enables the identification of the items that drive people to the store and always need to be in stock
- **Trip classification:** This analyzes the content of the basket and classifies the shopping trip into a category: weekly grocery trip, special occasion, and so on

- **Store-to-store comparison:** Understanding the number of baskets allows any metric to be divided by the total number of baskets, effectively creating a convenient and easy way to compare the stores with different characteristics (units sold per customer, revenue per transaction, number of items per basket, and so on)
- **Revenue optimization:** This helps in determining the magic price points for this store, increasing the size and value of the market basket
- **Marketing:** This helps in identifying more profitable advertising and promotions, targeting offers more precisely in order to improve ROI, generating better loyalty card promotions with longitudinal analysis, and attracting more traffic to the store
- **Operations optimization:** This helps in matching the inventory to the requirement by customizing the store and assortment to trade area demographics, optimizing store layout

Predictive models help retailers to direct the right offer to the right customer segments/profiles, as well as gain understanding of what is valid for which customer, predict the probability score of customers responding to this offer, and understand the customer value gain from the offer acceptance.

Affinity analysis

Affinity analysis is used to determine the likelihood that a set of items will be bought together. In retail, there are natural product affinities, for example, it is very typical for people who buy hamburger patties to buy hamburger rolls, along with ketchup, mustard, tomatoes, and other items that make up the burger experience.

While there are some product affinities that might seem trivial, there are some affinities that are not very obvious. A classic example is toothpaste and tuna. It seems that people who eat tuna are more prone to brush their teeth right after finishing their meal. So, why is it important for retailers to get a good grasp of the product affinities? This information is critical to appropriately plan promotions as reducing the price for some items may cause a spike on related high-affinity items without the need to further promote these related items.

In the following section, we'll look into the algorithms for association rule learning: Apriori and FP-growth.

Association rule learning

Association rule learning has been a popular approach to discover interesting relationships between items in large databases. It is most commonly applied in retail to reveal regularities between products.

Association rule learning approaches find patterns as interesting strong rules in the database using different measures of interestingness. For example, the following rule would indicate that if a customer buys onions and potatoes together, they are likely to also buy hamburger meat: $\{onions, potatoes\} \rightarrow \{burger\}$

Another classic story probably told in every machine learning class is the beer and diaper story. An analysis of supermarket shoppers' behavior showed that customers, presumably young men, who buy diapers tend also to buy beer. It immediately became a popular example of how an unexpected association rule might be found from everyday data; however, there are varying opinions as to how much of the story is true. Daniel Powers says (DSS News, 2002):

"In 1992, Thomas Blischok, manager of a retail consulting group at Teradata, and his staff prepared an analysis of 1.2 million market baskets from about 25 Osco Drug stores. Database queries were developed to identify affinities. The analysis did discover that between 5:00 and 7:00 p.m. consumers bought beer and diapers". Osco managers did NOT exploit the beer and diapers relationship by moving the products closer together on the shelves."

In addition to the preceding example from MBA, association rules are today employed in many application areas, including web usage mining, intrusion detection, continuous production, and bioinformatics. We'll take a closer look at these areas later in this chapter.

Basic concepts

Before we dive into algorithms, let's first review the basic concepts.

Database of transactions

In association rule mining, the dataset is structured a bit differently than the approach presented in the first chapter. First, there is no class value, as this is not required for learning association rules. Next, the dataset is presented as a transactional table, where each supermarket item corresponds to a binary attribute. Hence, the feature vector could be extremely large.

Consider the following example. Suppose we have four receipts as shown in the following image. Each receipt corresponds to a purchasing transaction:

<p>WWW.MACHINE-LEARNING-JAVA.COM</p> <p>GROCERY STORE 921 JAVA AVENUE NEW YORK NY 9999</p> <p>----- PURCHASE:</p> <table> <tr><td>POTATEOS</td><td>\$4.12</td></tr> <tr><td>BURGER</td><td>\$12.04</td></tr> </table> <p>VAT +11X TAX: \$1.77</p> <p>----- TOTAL: \$17.98</p> <p>PAYMENT METHOD: CREDIT CARD TRANSACTION #145829367 -001 DATE:18/09/2016 9:29:27 AM</p> <p>THANK YOU</p>	POTATEOS	\$4.12	BURGER	\$12.04	<p>WWW.MACHINE-LEARNING-JAVA.COM</p> <p>GROCERY STORE 921 JAVA AVENUE NEW YORK NY 9999</p> <p>----- PURCHASE:</p> <table> <tr><td>POTATEOS</td><td>\$4.12</td></tr> <tr><td>BURGER</td><td>\$12.04</td></tr> <tr><td>ONIONS</td><td>\$8.14</td></tr> <tr><td>BEER</td><td>\$27.55</td></tr> </table> <p>VAT +11X TAX: \$5.15</p> <p>----- TOTAL: \$52.00</p> <p>PAYMENT METHOD: CREDIT CARD TRANSACTION #1458293428 -001 DATE:18/09/2016 9:30:28 AM</p> <p>THANK YOU</p>	POTATEOS	\$4.12	BURGER	\$12.04	ONIONS	\$8.14	BEER	\$27.55	<p>WWW.MACHINE-LEARNING-JAVA.COM</p> <p>GROCERY STORE 921 JAVA AVENUE NEW YORK NY 9999</p> <p>----- PURCHASE:</p> <table> <tr><td>DIPPERS</td><td>\$29.95</td></tr> <tr><td>BEER</td><td>\$27.55</td></tr> </table> <p>VAT +11X TAX: \$6.38</p> <p>----- TOTAL: \$63.88</p> <p>PAYMENT METHOD: CREDIT CARD TRANSACTION #1458293588 -001 DATE:18/09/2016 9:31:48 AM</p> <p>THANK YOU</p>	DIPPERS	\$29.95	BEER	\$27.55	<p>WWW.MACHINE-LEARNING-JAVA.COM</p> <p>GROCERY STORE 921 JAVA AVENUE NEW YORK NY 9999</p> <p>----- PURCHASE:</p> <table> <tr><td>BURGER</td><td>\$12.04</td></tr> <tr><td>ONIONS</td><td>\$8.14</td></tr> <tr><td>BEER</td><td>\$27.55</td></tr> </table> <p>VAT +11X TAX: \$4.70</p> <p>----- TOTAL: \$47.43</p> <p>PAYMENT METHOD: CREDIT CARD TRANSACTION #1458293459 -001 DATE:18/09/2016 9:30:59 AM</p> <p>THANK YOU</p>	BURGER	\$12.04	ONIONS	\$8.14	BEER	\$27.55
POTATEOS	\$4.12																								
BURGER	\$12.04																								
POTATEOS	\$4.12																								
BURGER	\$12.04																								
ONIONS	\$8.14																								
BEER	\$27.55																								
DIPPERS	\$29.95																								
BEER	\$27.55																								
BURGER	\$12.04																								
ONIONS	\$8.14																								
BEER	\$27.55																								

To write these receipts in the form of a transactional database, we first identify all the possible items that appear in the receipts. These items are onions, potatoes, burger, beer, and dippers. Each purchase, that is, transaction, is presented in a row, and there is 1 if an item was purchased within the transaction and 0 otherwise, as shown in the following table:

Transaction ID	Onions	Potatoes	Burger	Beer	Dippers
1	0	1	1	0	0
2	1	1	1	1	0
3	0	0	0	1	1
4	1	0	1	1	0

This example is really small. In practical applications, the dataset often contains thousands or millions of transactions, which allow learning algorithm the discovery of statistically significant patterns.

Itemset and rule

Itemset is simply a set of items, for example, $\{\text{onions, potatoes, burger}\}$. A rule consists of two itemsets, X and Y , in the following format $X \rightarrow Y$.

This indicates a pattern that when the X itemset is observed, Y is also observed. To select interesting rules, various measures of significance can be used.

Support

Support, for an itemset, is defined as the proportion of transactions that contain the itemset. The $\{potatoes, burger\}$ itemset in the previous table has the following support as it occurs in 50% of transactions (2 out of 4 transactions) $supp(\{potatoes, burger\}) = 2/4 = 0.5$.

Intuitively, it indicates the share of transactions that support the pattern.

Confidence

Confidence of a rule indicates its accuracy. It is defined as

$$Conf(X \rightarrow Y) = supp(X \cup Y) / supp(X).$$

For example, the $\{onions, burger\} \rightarrow \{beer\}$ rule has the confidence $0.5/0.5 = 1.0$ in the previous table, which means that 100% of the times when *onions* and *burger* are bought together, *beer* is bought as well.

Apriori algorithm

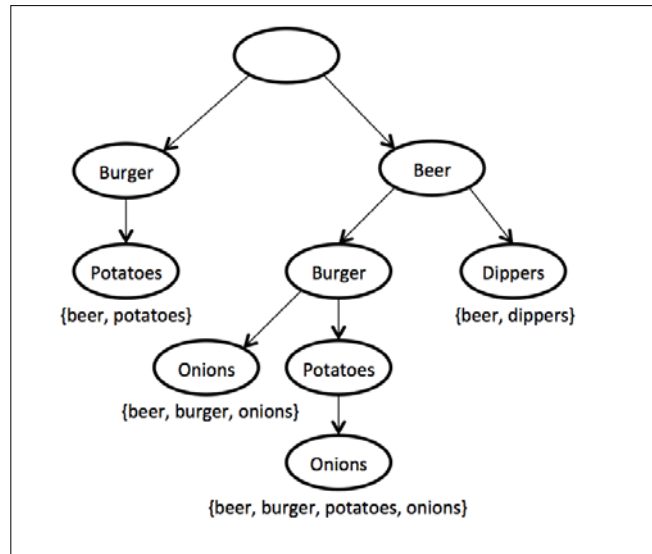
Apriori algorithm is a classic algorithm used for frequent pattern mining and association rule learning over transactional. By identifying the frequent individual items in a database and extending them to larger itemsets, Apriori can determine the association rules, which highlight general trends about a database.

Apriori algorithm constructs a set of itemsets, for example, $itemset1 = \{Item A, Item B\}$, and calculates support, which counts the number of occurrences in the database. Apriori then uses a bottom-up approach, where frequent itemsets are extended, one item at a time, and it works by eliminating the largest sets as candidates by first looking at the smaller sets and recognizing that a large set cannot be frequent unless all its subsets are. The algorithm terminates when no further successful extensions are found.

Although, Apriori algorithm is an important milestone in machine learning, it suffers from a number of inefficiencies and tradeoffs. In the following section, we'll look into a more recent FP-growth technique.

FP-growth algorithm

FP-growth, where **frequent pattern (FP)**, represents the transaction database as a prefix tree. First, the algorithm counts the occurrence of items in the dataset. In the second pass, it builds a prefix tree, an ordered tree data structure commonly used to store a string. An example of prefix tree based on the previous example is shown in the following diagram:



If many transactions share most frequent items, prefix tree provides high compression close to the tree root. Large itemsets are grown directly, instead of generating candidate items and testing them against the entire database. Growth starts at the bottom of the tree, by finding all the itemsets matching minimal support and confidence. Once the recursive process has completed, all large itemsets with minimum coverage have been found and association rule creation begins.

FP-growth algorithms have several advantages. First, it constructs an FP-tree, which encodes the original dataset in a substantially compact presentation. Second, it efficiently builds frequent itemsets, leveraging the **FP-tree structure** and **divide-and-conquer** strategy.

The supermarket dataset

The supermarket dataset, located in `datasets/chap5/supermarket.arff`, describes the shopping habits of supermarket customers. Most of the attributes stand for a particular item group, for example, dairy foods, beef, potatoes; or department, for example, department 79, department 81, and so on. The following image shows an excerpt of the database, where the value is *t* if the customer had bought an item and missing otherwise. There is one instance per customer. The dataset contains no class attribute, as this is not required to learn association rules. A sample of data is shown in the following table:

coffee	saucers-gravy-pkle	confectionary	puddings-deserts	dishcloths-scour	deod-disinfectan1	frozen foods	razor blades	fuels-garden aids	spices	jams-spreads
1	1	1	0	1	0	1	1	0	0	0
0	1	0	0	0	1	1	0	0	0	0
0	1	0	1	0	0	0	0	0	0	0
0	0	0	0	0	0	1	0	0	0	1
1	1	0	0	0	0	1	0	0	0	1
0	0	1	0	0	0	1	0	0	1	0
0	1	0	0	0	0	0	0	0	0	0
1	1	0	0	0	0	1	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	1	0	1	0	0	0	0	0	0	1
1	0	0	0	1	0	1	0	0	0	0
0	0	0	1	0	0	0	0	0	0	1
1	1	0	0	0	0	1	0	0	0	1
0	0	0	1	0	0	0	0	0	0	1
0	1	0	0	0	0	0	0	0	0	0
0	1	0	0	0	0	1	0	1	0	0
0	0	0	0	0	0	0	0	0	0	0
1	1	0	0	0	0	0	0	0	0	0
0	1	1	0	0	0	1	0	0	0	0
0	1	0	0	0	0	1	0	0	0	0
0	1	1	1	1	0	1	0	0	0	1
0	1	1	0	0	0	1	0	0	0	0
0	1	0	0	0	0	1	0	0	0	0
0	0	0	1	1	0	1	0	0	0	0
0	0	0	0	0	0	1	0	0	0	0
0	0	0	1	1	0	1	0	0	0	0
0	0	0	0	0	0	1	0	0	0	0
0	1	1	1	1	1	1	0	0	0	0

Discover patterns

To discover shopping patterns, we will use the two algorithms that we have looked into before, Apriori and FP-growth.

Apriori

We will use the Apriori algorithm as implemented in Weka. It iteratively reduces the minimum support until it finds the required number of rules with the given minimum confidence:

```
import java.io.BufferedReader;
import java.io.FileReader;
import weka.core.Instances;
import weka.associations.Apriori;
```

First, we will load the supermarket dataset:

```
Instances data = new Instances(
  new BufferedReader(
    new FileReader("datasets/chap5/supermarket.arff")));
```

Next, we will initialize an Apriori instance and call the `buildAssociations(Instances)` function to start frequent pattern mining, as follows:

```
Apriori model = new Apriori();
model.buildAssociations(data);
```

Finally, we can output the discovered itemsets and rules, as shown in the following code:

```
System.out.println(model);
```

The output is as follows:

Apriori

=====

Minimum support: 0.15 (694 instances)

Minimum metric <confidence>: 0.9

Number of cycles performed: 17

Generated sets of large itemsets:

Size of set of large itemsets L(1): 44

Size of set of large itemsets L(2): 380

Size of set of large itemsets L(3): 910

Size of set of large itemsets L(4): 633

Size of set of large itemsets L(5): 105

Size of set of large itemsets L(6): 1

Best rules found:

1. biscuits=t frozen foods=t fruit=t total=high 788 ==> bread and cake=t
723 <conf:(0.92)> lift:(1.27) lev:(0.03) [155] conv:(3.35)

2. baking needs=t biscuits=t fruit=t total=high 760 ==> bread and cake=t
696 <conf:(0.92)> lift:(1.27) lev:(0.03) [149] conv:(3.28)

3. baking needs=t frozen foods=t fruit=t total=high 770 ==> bread and
cake=t 705 <conf:(0.92)> lift:(1.27) lev:(0.03) [150] conv:(3.27)

...

The algorithm outputs ten best rules according to confidence. Let's look the first rule and interpret the output, as follows:

```
biscuits=t frozen foods=t fruit=t total=high 788 ==> bread and cake=t 723
<conf:(0.92)> lift:(1.27) lev:(0.03) [155] conv:(3.35)
```

It says that when biscuits, frozen foods, and fruits are bought together and the total purchase price is high, it is also very likely that bread and cake are purchased as well. The {biscuits, frozen foods, fruit, total high} itemset appears in 788 transactions, while the {bread, cake} itemset appears in 723 transactions. The confidence of this rule is 0.92, meaning that the rule holds true in 92% of transactions where the {biscuits, frozen foods, fruit, total high} itemset is present.

The output also reports additional measures such as lift, leverage, and conviction, which estimate the accuracy against our initial assumptions, for example, the 3.35 conviction value indicates that the rule would be incorrect 3.35 times as often if the association was purely a random chance. Lift measures the number of times X and Y occur together than expected if they were statistically independent (lift=1). The 2.16 lift in the X -> Y rule means that the probability of X is 2.16 times greater than the probability of Y.

FP-growth

Now, let's try to get the same results with more efficient FP-growth algorithm. FP-growth is also implemented in the `weka.associations` package:

```
import weka.associations.FPGrowth;
```

The FP-growth is initialized similarly as we did earlier:

```
FPGrowth fpgModel = new FPGrowth();
fpgModel.buildAssociations(data);
System.out.println(fpgModel);
```

The output reveals that FP-growth discovered 16 rules:

```
FPGrowth found 16 rules (displaying top 10)
```

```
1. [fruit=t, frozen foods=t, biscuits=t, total=high]: 788 ==> [bread and
cake=t]: 723 <conf:(0.92)> lift:(1.27) lev:(0.03) conv:(3.35)
2. [fruit=t, baking needs=t, biscuits=t, total=high]: 760 ==> [bread and
cake=t]: 696 <conf:(0.92)> lift:(1.27) lev:(0.03) conv:(3.28)
...
```

We can observe that FP-growth found the same set of rules as Apriori; however, the time required to process larger datasets can be significantly shorter.

Other applications in various areas

We looked into affinity analysis to demystify shopping behavior patterns in supermarkets. Although the roots of association rule learning are in analyzing point-of-sale transactions, they can be applied outside the retail industry to find relationships among other types of baskets. The notion of a basket can easily be extended to services and products, for example, to analyze items purchased using a credit card, such as rental cars and hotel rooms, and to analyze information on value-added services purchased by telecom customers (call waiting, call forwarding, DSL, speed call, and so on), which can help the operators determine the ways to improve their bundling of service packages.

Additionally, we will look into the following examples of potential cross-industry applications:

- Medical diagnosis
- Protein sequences
- Census data
- Customer relationship management
- IT Operations Analytics

Medical diagnosis

Applying association rules in medical diagnosis can be used to assist physicians while curing patients. The general problem of the induction of reliable diagnostic rules is hard as, theoretically, no induction process can guarantee the correctness of induced hypotheses by itself. Practically, diagnosis is not an easy process as it involves unreliable diagnosis tests and the presence of noise in training examples.

Nevertheless, association rules can be used to identify likely symptoms appearing together. A transaction, in this case, corresponds to a medical case, while symptoms correspond to items. When a patient is treated, a list of symptoms is recorded as one transaction.

Protein sequences

A lot of research has gone into understanding the composition and nature of proteins; yet many things remain to be understood satisfactorily. It is now generally believed that amino-acid sequences of proteins are not random.

With association rules, it is possible to identify associations between different amino acids that are present in a protein. A protein is a sequences made up of 20 types of amino acids. Each protein has a unique three-dimensional structure, which depends on the amino-acid sequence; slight change in the sequence may change the functioning of protein. To apply association rules, a protein corresponds to a transaction, while amino acids and their structure corespond to the items.

Such association rules are desirable for enhancing our understanding of protein composition and hold the potential to give clues regarding the global interactions amongst some particular sets of amino acids occurring in the proteins. Knowledge of these association rules or constraints is highly desirable for synthesis of artificial proteins.

Census data

Censuses make a huge variety of general statistical information about the society available to both researchers and general public. The information related to population and economic census can be forecasted in planning public services (education, health, transport, and funds) as well as in business (for setting up new factories, shopping malls, or banks and even marketing particular products).

To discover frequent patterns, each statistical area (for example, municipality, city, and neighborhood) corresponds to a transaction, and the collected indicators correspond to the items.

Customer relationship management

The **customer relationship management (CRM)**, as we briefly discussed in the previous chapters, is a rich source of data through which companies hope to identify the preference of different customer groups, products, and services in order to enhance the cohesion between their products and services and their customers.

Association rules can reinforce the knowledge management process and allow the marketing personnel to know their customers well in order to provide better quality services. For example, association rules can be applied to detect a change of customer behavior at different time snapshots from customer profiles and sales data. The basic idea is to discover changes from two datasets and generate rules from each dataset to carry out rule matching.

IT Operations Analytics

Based on records of a large number of transactions, association rule learning is well-suited to be applied to the data that is routinely collected in day-to-day IT operations, enabling IT Operations Analytics tools to detect frequent patterns and identify critical changes. IT specialists need to see the big picture and understand, for example, how a problem on a database could impact an application server.

For a specific day, IT operations may take in a variety of alerts, presenting them in a transactional database. Using an association rule learning algorithm, IT Operations Analytics tools can correlate and detect the frequent patterns of alerts appearing together. This can lead to a better understanding about how a component impacts another.

With identified alert patterns, it is possible to apply predictive analytics. For example, a particular database server hosts a web application and suddenly an alert about a database is triggered. By looking into frequent patterns identified by an association rule learning algorithm, this means that the IT staff needs to take action before the web application is impacted.

Association rule learning can also discover alert events originating from the same IT event. For example, every time a new user is added, six changes in the Windows operating system are detected. Next, in the **Application Portfolio Management (APM)**, IT may face multiple alerts, showing that the transactional time in a database as high. If all these issues originate from the same source (such as getting hundreds of alerts about changes that are all due to a Windows update), this frequent pattern mining can help to quickly cut through a number of alerts, allowing the IT operators to focus on truly critical changes.

Summary

In this chapter, you learned how to leverage association rule learning on transactional datasets to gain insight about frequent patterns. We performed an affinity analysis in Weka and learned that the hard work lies in the analysis of results – careful attention is required when interpreting rules, as association (that is, correlation) is not the same as causation.

In the next chapter, we'll look at how to take the problem of item recommendation to the next level using scalable machine learning library, Apache Mahout, which is able to handle big data.

6

Recommendation Engine with Apache Mahout

Recommendation engines are probably one of the most applied data science approaches in startups today. There are two principal techniques for building a recommendation system: **content-based filtering** and **collaborative filtering**. The content-based algorithm uses the properties of the items to find items with similar properties. Collaborative filtering algorithms take user ratings or other user behavior and make recommendations based on what users with similar behavior liked or purchased.

This chapter will first explain the basic concepts required to understand recommendation engine principles and then demonstrate how to utilize Apache Mahout's implementation of various algorithms to quickly get a scalable recommendation engine. This chapter will cover the following topics:

- How to build a recommendation engine
- Getting Apache Mahout ready
- Content-based approach
- Collaborative filtering approach

By the end of the chapter, you will learn the kind of recommendation engine that is appropriate for our problem and how to quickly implement one.

Basic concepts

Recommendation engines aim to show user items of interest. What makes them different from search engines is that the relevant content usually appears on a website without requesting it and users don't have to build queries as recommendation engines observe user's actions and construct query for users without their knowledge.

Arguably, the most well-known example of recommendation engine is www.amazon.com, providing personalized recommendation in a number of ways. The following image shows an example of **Customers Who Bought This Item Also Bought**. As we will see later, this is an example of collaborative item-based recommendation, where items similar to a particular item are recommended:



An example of recommendation engine from www.amazon.com.

In this section, we will introduce key concepts related to understanding and building recommendation engines.

Key concepts

Recommendation engine requires the following four inputs to make recommendations:

- Item information described with attributes
- User profile such as age range, gender, location, friends, and so on
- User interactions in form of ratings, browsing, tagging, comparing, saving, and emailing
- Context where the items will be displayed, for example, item category and item's geographical location

These inputs are then combined together by the recommendation engine to help us answer the following questions:

- Users who bought, watched, viewed, or bookmarked this item also bought, watched, viewed, or bookmarked...
- Items similar to this item...

- Other users you may know...
- Other users who are similar to you...

Now let's have a closer look at how this combining works.

User-based and item-based analysis

Building a recommendation engine depends on whether the engine searches for related items or users when trying to recommend a particular item.

In item-based analysis, the engine focuses on identifying items that are similar to a particular item; while in user-based analysis, users similar to the particular user are first determined. For example, users with the same profile information (age, gender, and so on) or actions history (bought, watched, viewed, and so on) are determined and then the same items are recommended to other similar users.

Both approaches require us to compute a similarity matrix, depending on whether we're analyzing item attributes or user actions. Let's take a deeper look at how this is done.

Approaches to calculate similarity

There are three fundamental approaches to calculate similarity, as follows:

- Collaborative filtering algorithms take user ratings or other user behavior and make recommendations based on what users with similar behavior liked or purchased
- The content-based algorithm uses the properties of the items to find items with similar properties
- A hybrid approach combining collaborative and content-based filtering

Let's take a look at each approach in detail.

Collaborative filtering

Collaborative filtering is based solely on user ratings or other user behavior, making recommendations based on what users with similar behavior liked or purchased.

A key advantage of collaborative filtering is that it does not rely on item content, and therefore, it is capable of accurately recommending complex items such as movies, without understanding the item itself. The underlying assumption is that people who agreed in the past will agree in the future, and that they will like similar kinds of items as they liked in the past.

A major disadvantage of this approach is the so-called cold start, meaning that if we want to build an accurate collaborative filtering system, the algorithm often needs a large amount of user ratings. This usually takes collaborative filtering out of the first version of the product and it is introduced later when a decent amount of data is collected.

Content-based filtering

Content-based filtering, on the other hand, is based on a description of items and a profile of user's preferences combined as follows. First, the items are described with attributes, and to find similar items, we measure the distance between items using a distance measure such as **cosine distance** or **Pearson coefficient** (more about distance measures is in *Chapter 1, Applied Machine Learning Quick Start*). Now, the user profile enters the equation. Given the feedback about the kind of items the user likes, we can introduce weights specifying the importance of a specific item attribute. For instance, Pandora Radio streaming service applies content-based filtering to create stations using more than 400 attributes. A user initially picks a song with specific attributes, and by providing feedback, important song attributes are emphasized.

This approach initially needs very little information on user feedback, thus it effectively avoids the cold-start issue.

Hybrid approach

Now collaborative versus content-based to choose? Collaborative filtering is able to learn user preferences from user's actions regarding one content source and use them across other content types. Content-based filtering is limited to recommending content of the same type that the user is already using. This provides value to different use cases, for example, recommending news articles based on news browsing is useful, but it is much more useful if different sources such as books and movies can be recommended based on news browsing.

Collaborative filtering and content-based filtering are not mutually exclusive; they can be combined to be more effective in some cases. For example, **Netflix** uses collaborative filtering to analyze searching and watching patterns of similar users, as well as content-based filtering to offer movies that share characteristics with films that the user has rated highly.

There is a wide variety of hybridization techniques such as weighted, switching, mixed, feature combination, feature augmentation, cascade, meta-level, and so on. Recommendation systems are an active area in machine learning and data mining community with special tracks on data science conferences. A good overview of techniques is summarized in the paper *Toward the next generation of recommender systems: a survey of the state-of-the-art and possible extensions* by Adomavicius and Tuzhilin (2005), where the authors discuss different approaches and underlying algorithms and provide references to further papers. To get more technical and understand all the tiny details when a particular approach makes sense, you should look at the book edited by Ricci et al. (2010) *Recommender Systems Handbook* (1st ed.), Springer-Verlag New York.

Exploitation versus exploration

In recommendation system, there is always a tradeoff between recommending items that fall into the user's sweet spot based on what we already know about the user (**exploitation**) and recommending items that don't fall into user's sweet spot with the aim to expose user to some novelties (**exploration**). Recommendation systems with little exploration will only recommend items consistent with the previous user ratings, preventing showing items outside their current bubble. In practice, serendipity of getting new items out of user's sweet spot is often desirable, leading to pleasant surprise and potential discovery of new sweet spots.

In this section, we discussed the essential concepts required to start building recommendation engines. Now, let's take a look at how to actually build one with Apache Mahout.

Getting Apache Mahout

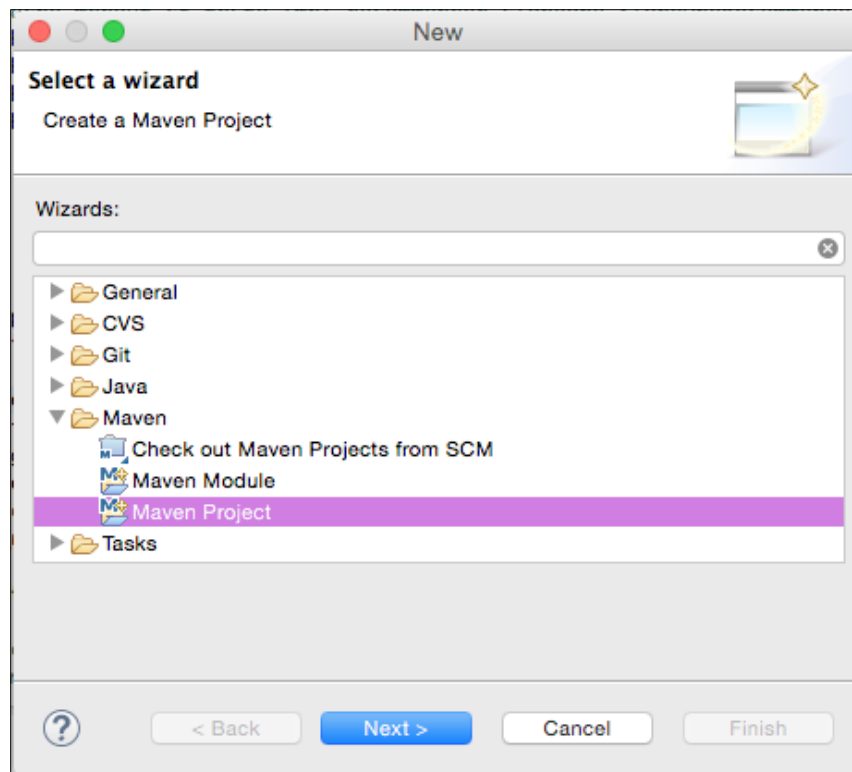
Mahout was introduced in *Chapter 2, Java Tools and Libraries for Machine Learning*, as a scalable machine learning library. It provides a rich set of components with which you can construct a customized recommendation system from a selection of algorithms. The creators of Mahout say it is designed to be enterprise-ready; it's designed for performance, scalability, and flexibility.

Mahout can be configured to run in two flavors: with or without Hadoop for a single machine and distributed processing, correspondingly. We will focus on configuring Mahout without Hadoop. For more advanced configurations and further uses of Mahout, I would recommend two recent books: *Learning Apache Mahout* (Tiwary, 2015) and *Learning Apache Mahout Classification* (Gupta, 2015).

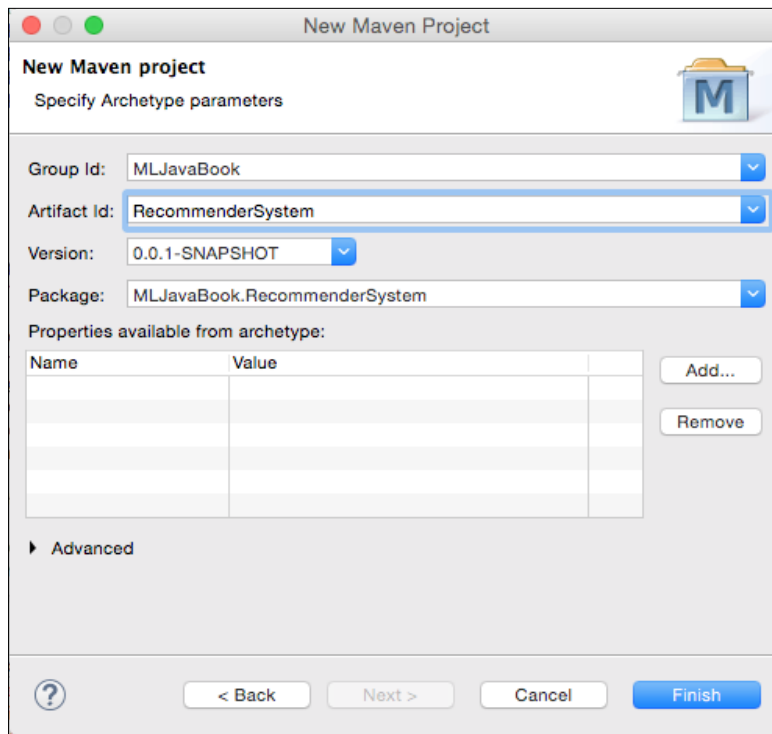
As Apache Mahout's build and release system is based on Maven, we will need to learn how to install it. We will look at the most convenient approach using Eclipse with Maven plugin.

Configuring Mahout in Eclipse with the Maven plugin

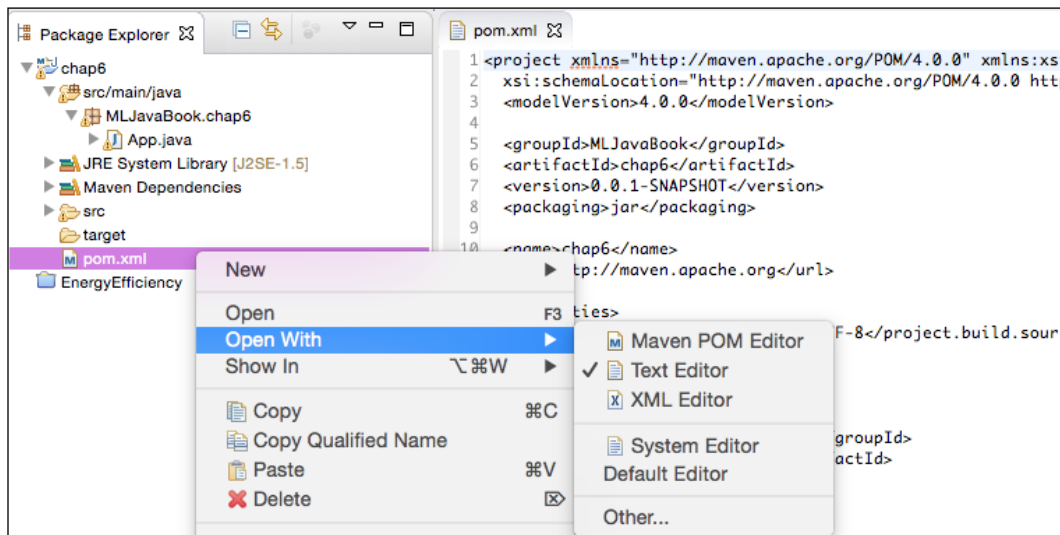
We will need a recent version of Eclipse, which can be downloaded from its home page. We use Eclipse Luna in this book. Open Eclipse and start a new **Maven Project** with default settings as shown in the following screenshot:



The **New Maven project** screen will appear as shown in the following image:



Now, we need to tell the project to add Mahout jar and its dependencies to the project. Locate the `pom.xml` file and open it with the text editor (left click on **Open With | Text Editor**), as shown in the following screenshot:



Locate the line starting with `<dependencies>` and add the following code in the next line:

```
<dependency>
  <groupId>org.apache.mahout</groupId>
  <artifactId>mahout-mr</artifactId>
  <version>0.10.0</version>
</dependency>
```

That's it, Mahout is added and we are ready to begin now.

Building a recommendation engine

To demonstrate both the content-based filtering and collaborative filtering approaches, we'll build a book-recommendation engine.

Book ratings dataset

In this chapter, we will work with book ratings dataset (Ziegler et al, 2005) collected in a four-week crawl. It contains data on 278,858 members of the **Book-Crossing** website and 1,157,112 ratings, both implicit and explicit, referring to 271,379 distinct ISBNs. User data is anonymized, but with demographic information. The dataset is available at:

<http://www2.informatik.uni-freiburg.de/~chiegler/BX/>.

The Book-Crossing dataset comprises three files described at their website as follows:

- **BX-Users:** This contains the users. Note that user IDs (User-ID) have been anonymized and mapped to integers. Demographic data is provided (Location and Age) if available. Otherwise, these fields contain NULL-values.
- **BX-Books:** Books are identified by their respective ISBN. Invalid ISBNs have already been removed from the dataset. Moreover, some content-based information is given (Book-Title, Book-Author, Year-Of-Publication, and Publisher), obtained from Amazon Web Services. Note that in case of several authors, only the first author is provided. URLs linking to cover images are also given, appearing in three different flavors (Image-URL-S, Image-URL-M, and Image-URL-L), that is, small, medium, and large. These URLs point to the Amazon website.
- **BX-Book-Ratings:** This contains the book rating information. Ratings (Book-Rating) are either explicit, expressed on a scale of 1-10 (higher values denoting higher appreciation), or implicit, expressed by 0.

Loading the data

There are two approaches for loading the data according to where the data is stored: file or database. First, we will take a detailed look at how to load the data from the file, including how to deal with custom formats. At the end, we quickly take a look at how to load the data from a database.

Loading data from file

Loading data from file can be achieved with the `FileDataModel` class, expecting a comma-delimited file, where each line contains a `userID`, `itemID`, optional preference value, and optional `timestamp` in the same order, as follows:

```
userID,itemID[,preference[,timestamp]]
```

Optional preference accommodates applications with binary preference values, that is, user either expresses a preference for an item or not, without degree of preference, for example, with like/dislike.

A line that begins with hash, #, or an empty line will be ignored. It is also acceptable for the lines to contain additional fields, which will be ignored.

The `DataModel` class assumes the following types:

- `userID`, `itemID` can be parsed as `long`
- preference value can be parsed as `double`
- `timestamp` can be parsed as `long`

If you are able to provide the dataset in the preceding format, you can simply use the following line to load the data:

```
DataModel model = new FileDataModel(new File(path));
```

This class is not intended to be used for very large amounts of data, for example, tens of millions of rows. For that, a JDBC-backed `DataModel` and a database are more appropriate.

In real world, however, we cannot always ensure that the input data supplied to us contain only integer values for `userID` and `itemID`. For example, in our case, `itemID` correspond to ISBN book numbers uniquely identifying items, but these are not integers and the `FileDataModel` default won't be suitable to process our data.

Now, let's consider how to deal with the case where our `itemID` is a string. We will define our custom data model by extending `FileDataModel` and overriding the `long readItemIDFromString(String)` method in order to read `itemID` as a string and convert it into long and return a unique long value. To convert `String` to unique long, we'll extend another Mahout `AbstractIDMigrator` helper class, which is designed exactly for this task.

Now, let's first look at how `FileDataModel` is extended:

```
class StringItemIdFileDataModel extends FileDataModel {

    //initialize migrator to covert String to unique long
    public ItemMemIDMigrator memIdMigtr;

    public StringItemIdFileDataModel(File dataFile, String regex)
        throws IOException {
        super(dataFile, regex);
    }

    @Override
    protected long readItemIDFromString(String value) {

        if (memIdMigtr == null) {
            memIdMigtr = new ItemMemIDMigrator();
        }

        // convert to long
        long retValue = memIdMigtr.toLongID(value);
        //store it to cache
        if (null == memIdMigtr.toStringID(retValue)) {
            try {
                memIdMigtr.singleInit(value);
            } catch (TasteException e) {
                e.printStackTrace();
            }
        }
        return retValue;
    }

    // convert long back to String
    String getItemIDAsString(long itemId) {
        return memIdMigtr.toStringID(itemId);
    }
}
```

Other useful methods that can be overridden are as follows:

- `readUserIDFromString(String value)` if user IDs are not numeric
- `readTimestampFromString(String value)` to change how timestamp is parsed

Now, let's take a look how `AbstractIDMigrator` is extended:

```
class ItemMemIDMigrator extends AbstractIDMigrator {

    private FastByIDMap<String> longToString;

    public ItemMemIDMigrator() {
        this.longToString = new FastByIDMap<String>(10000);
    }

    public void storeMapping(long longID, String stringID) {
        longToString.put(longID, stringID);
    }

    public void singleInit(String stringID) throws TasteException {
        storeMapping(toLongID(stringID), stringID);
    }

    public String toStringID(long longID) {
        return longToString.get(longID);
    }
}
```

Now, we have everything in place and we can load our dataset with the following code:

```
StringItemIdFileDataModel model = new StringItemIdFileDataModel(
    new File("datasets/chap6/BX-Book-Ratings.csv"), ";");
System.out.println(
    "Total items: " + model.getNumItems() +
    "\nTotal users: " + model.getNumUsers());
```

This outputs the total number of users and items:

Total items: 340556

Total users: 105283

We are ready to move on and start making recommendations.

Loading data from database

Alternately, we can also load the data from database using one of the JDBC data models. In this chapter, we will not dive into the detailed instructions about how to set up database, connection, and so on, but just give a sketch on how this can be done.

Database connectors have been moved to a separate package `mahout-integration`, hence we have to first add the package to our dependency list. Open the `pom.xml` file and add the following dependency:

```
<dependency>
  <groupId>org.apache.mahout</groupId>
  <artifactId>mahout-integration</artifactId>
  <version>0.7</version>
</dependency>
```

Consider that we want to connect a MySQL database. In this case, we will also need a package that handles database connections. Add the following to the `pom.xml` file:

```
<dependency>
  <groupId>mysql</groupId>
  <artifactId>mysql-connector-java</artifactId>
  <version>5.1.35</version>
</dependency>
```

Now, we have all the packages, so we can create a connection. First, let's initialize a `DataSource` class with connection details, as follows:

```
MysqlDataSource dbsource = new MysqlDataSource();
dbsource.setUser("user");
dbsource.setPassword("pass");
dbsource.setServerName("hostname.com");
dbsource.setDatabaseName("db");
```

Mahout integration implements `JDBCDataModel` to various databases that can be accessed via JDBC. By default, this class assumes that there is `DataSource` available under the JNDI name `jdbc/taste`, which gives access to a database with a `taste_preferences` table with the following schema:

```
CREATE TABLE taste_preferences (
  user_id BIGINT NOT NULL,
  item_id BIGINT NOT NULL,
  preference REAL NOT NULL,
  PRIMARY KEY (user_id, item_id)
)
```

```
CREATE INDEX taste_preferences_user_id_index ON taste_preferences
(user_id);
CREATE INDEX taste_preferences_item_id_index ON taste_preferences
(item_id);
```

A database-backed data model is initialized as follows. In addition to the DB connection object, we can also specify the custom table name and table column names, as follows:

```
DataModel dataModel = new MySQLJDBCDataModel(dbsource,
"taste_preferences",
"user_id", "item_id", "preference", "timestamp");
```

In-memory database

Last, but not least, the data model can be created on the fly and held in memory. A database can be created from an array of preferences holding user ratings for a set of items.

We can proceed as follows. First, we create a `FastByIdMap` hash map of preference arrays, `PreferenceArray`, which stores an array of preferences:

```
FastByIdMap <PreferenceArray> preferences = new FastByIdMap
<PreferenceArray> ();
```

Next, we can create a new preference array for a user that will hold their ratings. The array must be initialized with a size parameter that reserves that many slots in memory:

```
PreferenceArray prefsForUser1 =
new GenericUserPreferenceArray (10);
```

Next, we set user ID for current preference at position 0. This will actually set the user ID for all preferences:

```
prefsForUser1.setUserID (0, 1L);
```

Set item ID for current preference at position 0:

```
prefsForUser1.setItemID (0, 101L);
```

Set preference value for preference at 0:

```
prefsForUser1.setValue (0, 3.0f);
```


Continue for other item ratings:

```
prefsForUser1.setItemID (1, 102L);  
prefsForUser1.setValue (1, 4.5F);
```

Finally, add user preferences to the hash map:

```
preferences.put (1L, prefsForUser1); // use userID as the key
```

The preference hash map can be now used to initialize `GenericDataModel`:

```
DataModel dataModel = new GenericDataModel(preferences);
```

This code demonstrates how to add two preferences for a single user; while in practical application, you'll want to add multiple preferences for multiple users.

Collaborative filtering

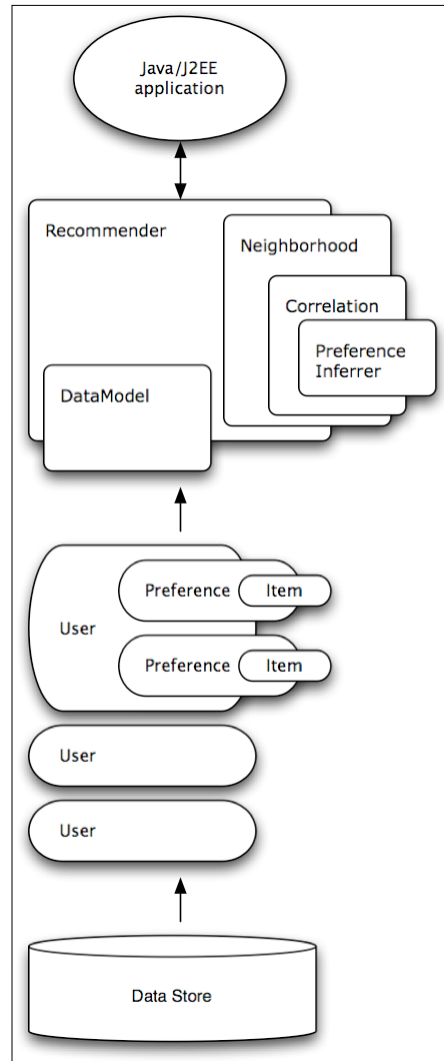
Recommendation engines in Mahout can be built with the `org.apache.mahout.cf.taste` package, which was formerly a separate project called `Taste` and has continued development in Mahout.

A Mahout-based collaborative filtering engine takes the users' preferences for items (tastes) and returns the estimated preferences for other items. For example, a site that sells books or CDs could easily use Mahout to figure out the CDs that a customer might be interested in listening to with the help of the previous purchase data.

Top-level packages define the Mahout interfaces to the following key abstractions:

- **DataModel**: This represents a repository of information about users and their preferences for items
- **UserSimilarity**: This defines a notion of similarity between two users
- **ItemSimilarity**: This defines a notion of similarity between two items
- **UserNeighborhood**: This computes neighborhood users for a given user
- **Recommender**: This recommends items for user

A general structure of the concepts is shown in the following diagram:



User-based filtering

The most basic user-based collaborative filtering can be implemented by initializing the previously described components as follows.

First, load the data model:

```
StringItemIdFileDataModel model = new StringItemIdFileDataModel(
    new File("/datasets/chap6/BX-Book-Ratings.csv", ";");
```

Next, define how to calculate how the users are correlated, for example, using Pearson correlation:

```
UserSimilarity similarity =  
    new PearsonCorrelationSimilarity(model);
```

Next, define how to tell which users are similar, that is, users that are close to each other according to their ratings:

```
UserNeighborhood neighborhood =  
    new ThresholdUserNeighborhood(0.1, similarity, model);
```

Now, we can initialize a `GenericUserBasedRecommender` default engine with data model, neighborhood, and similar objects, as follows:

```
UserBasedRecommender recommender =  
    new GenericUserBasedRecommender(model, neighborhood, similarity);
```

That's it. Our first basic recommendation engine is ready. Let's discuss how to invoke recommendations. First, let's print the items that the user already rated along with ten recommendations for this user:

```
long userID = 80683;  
int noItems = 10;  
  
List<RecommendedItem> recommendations = recommender.recommend(  
    userID, noItems);  
  
System.out.println("Rated items by user:");  
for(Preference preference : model.getPreferencesFromUser(userID)) {  
    // convert long itemID back to ISBN  
    String itemISBN = model.getItemIDAsString(  
        preference.getItemID());  
    System.out.println("Item: " + books.get(itemISBN) +  
        " | Item id: " + itemISBN +  
        " | Value: " + preference.getValue());  
}  
  
System.out.println("\nRecommended items:");  
for (RecommendedItem item : recommendations) {  
    String itemISBN = model.getItemIDAsString(item.getItemID());  
    System.out.println("Item: " + books.get(itemISBN) +  
        " | Item id: " + itemISBN +  
        " | Value: " + item.getValue());  
}
```

This outputs the following recommendations along with their scores:

Rated items:

Item: The Handmaid's Tale | Item id: 0395404258 | Value: 0.0
Item: Get Clark Smart : The Ultimate Guide for the Savvy Consumer | Item id: 1563526298 | Value: 9.0
Item: Plum Island | Item id: 0446605409 | Value: 0.0
Item: Blessings | Item id: 0440206529 | Value: 0.0
Item: Edgar Cayce on the Akashic Records: The Book of Life | Item id: 0876044011 | Value: 0.0
Item: Winter Moon | Item id: 0345386108 | Value: 6.0
Item: Sarah Bishop | Item id: 059032120X | Value: 0.0
Item: Case of Lucy Bending | Item id: 0425060772 | Value: 0.0
Item: A Desert of Pure Feeling (Vintage Contemporaries) | Item id: 0679752714 | Value: 0.0
Item: White Abacus | Item id: 0380796155 | Value: 5.0
Item: The Land of Laughs : A Novel | Item id: 0312873115 | Value: 0.0
Item: Nobody's Son | Item id: 0152022597 | Value: 0.0
Item: Mirror Image | Item id: 0446353957 | Value: 0.0
Item: All I Really Need to Know | Item id: 080410526X | Value: 0.0
Item: Dreamcatcher | Item id: 0743211383 | Value: 7.0
Item: Perplexing Lateral Thinking Puzzles: Scholastic Edition | Item id: 0806917695 | Value: 5.0
Item: Obsidian Butterfly | Item id: 0441007813 | Value: 0.0

Recommended items:

Item: Keeper of the Heart | Item id: 0380774933 | Value: 10.0
Item: Bleachers | Item id: 0385511612 | Value: 10.0
Item: Salem's Lot | Item id: 0451125452 | Value: 10.0
Item: The Girl Who Loved Tom Gordon | Item id: 0671042858 | Value: 10.0
Item: Mind Prey | Item id: 0425152898 | Value: 10.0
Item: It Came From The Far Side | Item id: 0836220730 | Value: 10.0
Item: Faith of the Fallen (Sword of Truth, Book 6) | Item id: 081257639X | Value: 10.0
Item: The Talisman | Item id: 0345444884 | Value: 9.86375
Item: Hamlet | Item id: 067172262X | Value: 9.708363
Item: Untamed | Item id: 0380769530 | Value: 9.708363

Item-based filtering

The `ItemSimilarity` is the most important point to discuss here. Item-based recommenders are useful as they can take advantage of something very fast: they base their computations on item similarity, not user similarity, and item similarity is relatively static. It can be precomputed, instead of recomputed in real time.

Thus, it's strongly recommended that you use `GenericItemSimilarity` with precomputed similarities if you're going to use this class. You can use `PearsonCorrelationSimilarity` too, which computes similarities in real time, but you will probably find this painfully slow for large amounts of data:

```
StringItemIdFileDataModel model = new StringItemIdFileDataModel(
    new File("datasets/chap6/BX-Book-Ratings.csv"), ";");

ItemSimilarity itemSimilarity = new
    PearsonCorrelationSimilarity(model);

ItemBasedRecommender recommender = new
    GenericItemBasedRecommender(model, itemSimilarity);

String itemISBN = "0395272238";
long itemID = model.readItemIDFromString(itemISBN);
int noItems = 10;
List<RecommendedItem> recommendations =
    recommender.mostSimilarItems(itemID, noItems);

System.out.println("Recommendations for item:
    "+books.get(itemISBN));

System.out.println("\nMost similar items:");
for (RecommendedItem item : recommendations) {
    itemISBN = model.getItemIDAsString(item.getItemID());
    System.out.println("Item: " + books.get(itemISBN) + " | Item id:
        " + itemISBN + " | Value: " + item.getValue());
}
```

Recommendations for item: Close to the Bone

Most similar items:

```
Item: Private Screening | Item id: 0345311396 | Value: 1.0
Item: Heartstone | Item id: 0553569783 | Value: 1.0
Item: Clockers / Movie Tie In | Item id: 0380720817 | Value: 1.0
Item: Rules of Prey | Item id: 0425121631 | Value: 1.0
Item: The Next President | Item id: 0553576666 | Value: 1.0
Item: Orchid Beach (Holly Barker Novels (Paperback)) | Item id:
0061013412 | Value: 1.0
```

```
Item: Winter Prey | Item id: 0425141233 | Value: 1.0
Item: Night Prey | Item id: 0425146413 | Value: 1.0
Item: Presumed Innocent | Item id: 0446359866 | Value: 1.0
Item: Dirty Work (Stone Barrington Novels (Paperback)) | Item id:
0451210158 | Value: 1.0
```

The resulting list returns a set of items similar to particular item that we selected.

Adding custom rules to recommendations

It often happens that some business rules require us to boost the score of the selected items. In the book dataset, for example, if a book is recent, we want to give it a higher score. That's possible using the `IDRescorer` interface implementing, as follows:

- `rescore(long, double)` that takes `itemId` and original score as an argument and returns a modified score
- `isFiltered(long)` that may return `true` to exclude a specific item from recommendation or `false` otherwise

Our example could be implemented as follows:

```
class MyRescorer implements IDRescorer {

    public double rescore(long itemId, double originalScore) {
        double newScore = originalScore;
        if (bookIsNew(itemId)) {
            originalScore *= 1.3;
        }
        return newScore;
    }

    public boolean isFiltered(long arg0) {
        return false;
    }

}
```

An instance of `IDRescorer` is provided when invoking `recommender.recommend`:

```
IDRescorer rescorer = new MyRescorer();
List<RecommendedItem> recommendations =
recommender.recommend(userID, noItems, rescorer);
```

Evaluation

You might wonder how to make sure that the returned recommendations make any sense? The only way to be really sure about how effective recommendations are is to use A/B testing in a live system with real users. For example, the A group receives a random item as a recommendation, while the B group receives an item recommended by our engine.

As this is neither always possible nor practical, we can get an estimate with offline statistical evaluation. One way to proceed is to use the k-fold cross validation introduced in *Chapter 1, Applied Machine Learning Quick Start*. We partition dataset into multiple sets, some are used to train our recommendation engine and the rest to test how well it recommends items to unknown users.

Mahout implements the `RecommenderEvaluator` class that splits a dataset in two parts. The first part, 90% by default, is used to produce recommendations, while the rest of the data is compared against estimated preference values to test the match. The class does not accept a recommender object directly, you need to build a class implementing the `RecommenderBuilder` interface instead, which builds a recommender object for a given `DataModel` object that is then used for testing. Let's take a look at how this is implemented.

First, we create a class that implements the `RecommenderBuilder` interface. We need to implement the `buildRecommender` method, which will return a recommender, as follows:

```
public class BookRecommender implements RecommenderBuilder {
    public Recommender buildRecommender(DataModel dataModel) {
        UserSimilarity similarity =
            new PearsonCorrelationSimilarity(model);
        UserNeighborhood neighborhood =
            new ThresholdUserNeighborhood(0.1, similarity, model);
        UserBasedRecommender recommender =
            new GenericUserBasedRecommender(
                model, neighborhood, similarity);
        return recommender;
    }
}
```

Now that we have class that returns a recommender object, we can initialize a `RecommenderEvaluator` instance. Default implementation of this class is the `AverageAbsoluteDifferenceRecommenderEvaluator` class, which computes the average absolute difference between the predicted and actual ratings for users. The following code shows how to put the pieces together and run a hold-out test.

First, load a data model:

```
DataModel dataModel = new FileDataModel(  
    new File("/path/to/dataset.csv"));
```

Next, initialize an evaluator instance, as follows:

```
RecommenderEvaluator evaluator =  
    new AverageAbsoluteDifferenceRecommenderEvaluator();
```

Initialize the `BookRecommender` object, implementing the `RecommenderBuilder` interface:

```
RecommenderBuilder builder = new MyRecommenderBuilder();
```

Finally, call the `evaluate()` method, which accepts the following parameters:

- `RecommenderBuilder`: This is the object implementing `RecommenderBuilder` that can build recommender to test
- `DataModelBuilder`: `DataModelBuilder` to use, or if null, a default `DataModel` implementation will be used
- `DataModel`: This is the dataset that will be used for testing
- `trainingPercentage`: This indicates the percentage of each user's preferences to use to produced recommendations; the rest are compared to estimated preference values to evaluate the recommender performance
- `evaluationPercentage`: This is the percentage of users to be used in evaluation

The method is called as follows:

```
double result = evaluator.evaluate(builder, null, model, 0.9,  
    1.0);  
System.out.println(result);
```

The method returns a `double`, where 0 presents the best possible evaluation, meaning that the recommender perfectly matches user preferences. In general, lower the value, better the match.

Online learning engine

What about the online aspect? The above will work great for existing users; but what about new users which register in the service? For sure, we want to provide some reasonable recommendations for them as well. Creating a recommendation instance is expensive (it definitely takes longer than a usual network request), so we can't just create a new recommendation each time.

Luckily, Mahout has a possibility of adding temporary users to a data model. The general set up is as follows:

- Periodically recreate the whole recommendation using current data (for example, each day or hour, depending on how long it takes)
- When doing a recommendation, check whether the user exists in the system
- If yes, complete the recommendation as always
- If no, create a temporary user, fill in the preferences, and do the recommendation

The first part (periodically recreating the recommender) may be actually quite tricky if you are limited on memory: when creating the new recommender, you need to hold two copies of the data in memory (in order to still be able to server requests from the old one). However, as this doesn't really have anything to do with recommendations, I won't go into details here.

As for the temporary users, we can wrap our data model with a `PlusAnonymousConcurrentUserDataModel` instance. This class allows us to obtain a temporary user ID; the ID must be later released so that it can be reused (there's a limited number of such IDs). After obtaining the ID, we have to fill in the preferences, and then, we can proceed with the recommendation as always:

```
class OnlineRecommendation{

    Recommender recommender;
    int concurrentUsers = 100;
    int noItems = 10;

    public OnlineRecommendation() throws IOException {

        DataModel model = new StringItemIdFileDataModel(
            new File (/chap6/BX-Book-Ratings.csv"), "");
        PlusAnonymousConcurrentUserDataModel plusModel = new
            PlusAnonymousConcurrentUserDataModel
                (model, concurrentUsers);
        recommender = ...;

    }

    public List<RecommendedItem> recommend(long userId,
        PreferenceArray preferences){

        if (userExistsInDataModel (userId)) {
```

```
        return recommender.recommend(userId, noItems);
    }

    else{

        PlusAnonymousConcurrentUserDataModel plusModel =
            (PlusAnonymousConcurrentUserDataModel)
                recommender.getDataModel();

        // Take an available anonymous user form the poll
        Long anonymousUserID = plusModel.takeAvailableUser();

        // Set temporary preferences
        PreferenceArray tempPrefs = preferences;
        tempPrefs.setUserID(0, anonymousUserID);
        tempPrefs.setItemID(0, itemID);
        plusModel.setTempPrefs(tempPrefs, anonymousUserID);

        List<RecommendedItem> results =
            recommender.recommend(anonymousUserID, noItems);

        // Release the user back to the poll
        plusModel.releaseUser(anonymousUserID);

        return results;
    }
}
}
```

Content-based filtering

Content-based filtering is out of scope in the Mahout framework, mainly because it is up to you to decide how to define similar items. If we want to do a content-based item-item similarity, we need to implement our own `ItemSimilarity`. For instance, in our book's dataset, we might want to make up the following rule for book similarity:

- If genres are the same, add 0.15 to similarity
- If author is the same, add 0.50 to similarity

We could now implement our own similarity measure as follows:

```
class MyItemSimilarity implements ItemSimilarity {
    ...
    public double itemSimilarity(long itemID1, long itemID2) {
        MyBook book1 = lookupMyBook (itemID1);
        MyBook book2 = lookupMyBook (itemID2);
        double similarity = 0.0;
        if (book1.getGenre().equals(book2.getGenre()))
            similarity += 0.15;
        }
        if (book1.getAuthor().equals(book2.getAuthor ())) {
            similarity += 0.50;
        }
        return similarity;
    }
    ...
}
```

We then use this `ItemSimilarity` instead of something like `LogLikelihoodSimilarity` or other implementations with a `GenericItemBasedRecommender`. That's about it. This is as far as we have to go to perform content-based recommendation in the Mahout framework.

What we saw here is one of the simplest forms of content-based recommendation. Another approach could be to create a content-based profile of users, based on a weighted vector of item features. The weights denote the importance of each feature to the user and can be computed from individually-rated content vectors.

Summary

In this chapter, you learned the basic concepts of recommendation engines, the difference between collaborative and content-based filtering, and how to use Apache Mahout, which is a great basis to create recommenders as it is very configurable and provides many extension points. We looked at how to pick the right configuration parameter values, set up rescoring, and evaluate the recommendation results.

With this chapter, we completed data science techniques to analyze customer behavior that started with customer-relationship prediction in *Chapter 4, Customer Relationship Prediction with Ensembles*, and continued with affinity analytics in *Chapter 5, Affinity Analysis*. In the next chapter, we will move on to other topics, such as fraud and anomaly detection.

7

Fraud and Anomaly Detection

Outlier detection is used to identify exceptions, rare events, or other anomalous situations. Such anomalies may be hard-to-find needles in a haystack, but their consequences may nonetheless be quite dramatic, for instance, credit card fraud detection, identifying network intrusion, faults in a manufacturing processes, clinical trials, voting activities, and criminal activities in e-commerce. Therefore, discovered anomalies represent high value when they are found or high costs if they are not found. Applying machine learning to outlier detection problems brings new insight and better detection of outlier events. Machine learning can take into account many disparate sources of data and find correlations that are too obscure for human analysis to identify.

Take the example of e-commerce fraud detection. With machine learning algorithm in place, the purchaser's online behavior, that is, website browsing history, becomes a part of the fraud detection algorithm rather than simply considering the history of purchases made by the cardholder. This involves analyzing a variety of data sources, but it is also a far more robust approach to e-commerce fraud detection.

In this chapter, we will cover the following topics:

- Problems and challenges
- Suspicious pattern detection
- Anomalous pattern detection
- Working with unbalanced datasets
- Anomaly detection in time series

Suspicious and anomalous behavior detection

The problem of learning patterns from sensor data arises in many applications, including e-commerce, smart environments, video surveillance, network analysis, human-robot interaction, ambient assisted living, and so on. We focus on detecting patterns that deviate from regular behaviors and might represent a security risk, health problem, or any other abnormal behavior contingency.

In other words, deviant behavior is a data pattern that either does not conform to the expected behavior (anomalous behavior) or matches a previously defined unwanted behavior (suspicious behavior). Deviant behavior patterns are also referred to as outliers, exceptions, peculiarities, surprise, misuse, and so on. Such patterns relatively occur infrequently; however, when they do occur, their consequences can be quite dramatic, and often negatively so. Typical examples include credit card fraud detection, cyber-intrusions, and industrial damage. In e-commerce, fraud is estimated to cost merchants more than \$200 billion a year; in healthcare, fraud is estimated to cost taxpayers \$60 billion a year; for banks, the cost is over \$12 billion.

Unknown-unknowns

When Donald Rumsfeld, US Secretary of Defense, had a news briefing on February 12, 2002, about the lack of evidence linking the government of Iraq to the supply of weapons of mass destruction to terrorist groups, it immediately became a subject of much commentary. Rumsfeld stated (DoD News, 2012):

"Reports that say that something hasn't happened are always interesting to me, because as we know, there are known knowns; there are things we know we know. We also know there are known unknowns; that is to say we know there are some things we do not know. But there are also unknown unknowns – the ones we don't know we don't know. And if one looks throughout the history of our country and other free countries, it is the latter category that tend to be the difficult ones."

The statement might seem confusing at first, but the idea of unknown unknowns was well studied among scholars dealing with risk, NSA, and other intelligence agencies. What the statement basically says is the following:

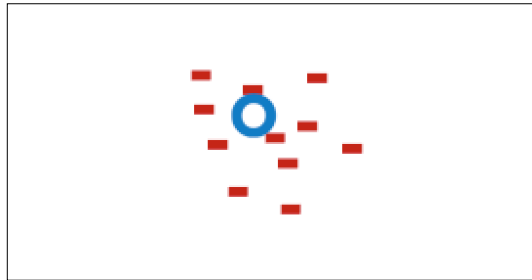
- **Known-knowns:** These are well-known problems or issues we know how to recognize them and how deal with them
- **Known-unknowns:** These are expected or foreseeable problems, which can be reasonably anticipated, but have not occurred before

- **Unknown-unknowns:** These are unexpected and unforeseeable problems, which pose significant risk as they cannot be anticipated, based on previous experience

In the following sections, we will look into two fundamental approaches dealing with the first two types of knowns and unknowns: suspicious pattern detection dealing with known-knowns and anomalous pattern detection targeting known-unknowns.

Suspicious pattern detection

The first approach assumes a behavior library that encodes negative patterns shown as red minus signs in the following image, and thus recognizing that observed behavior corresponds to identifying a match in the library. If a new pattern (blue circle) can be matched against negative patterns, then it is considered suspicious:



For example, when you visit a doctor, she inspects various health symptoms (body temperature, pain levels, affected areas, and so on) and matches the symptoms to a known disease. In machine learning terms, the doctor collects attributes and performs classifications.

An advantage of this approach is that we immediately know what is wrong; for example, assuming we know the disease, we can select appropriate treatment procedure.

A major disadvantage of this approach is that it can detect only suspicious patterns that are known in advance. If a pattern is not inserted into a negative pattern library, then we will not be able to recognize it. This approach is, therefore, appropriate for modeling known-knowns.

Anomalous pattern detection

The second approach uses the pattern library in an inverse fashion, meaning that the library encodes only positive patterns marked with green plus signs in the following image. When an observed behavior (blue circle) cannot be matched against the library, it is considered anomalous:



This approach requires us to model only what we have seen in the past, that is, normal patterns. If we return to the doctor example, the main reason we visited the doctor in the first place was because we did not feel fine. Our perceived state of feelings (for example, headache, sore skin) did not match our usual feelings, therefore, we decided to seek doctor. We don't know which disease caused this state nor do we know the treatment, but we were able to observe that it doesn't match the usual state.

A major advantage of this approach is that it does not require us to say anything about non-normal patterns; hence, it is appropriate for modeling known-unknowns and unknown-unknowns. On the other hand, it does not tell us what exactly is wrong.

Analysis types

Several approaches have been proposed to tackle the problem either way. We broadly classify anomalous and suspicious behavior detection in the following three categories: pattern analysis, transaction analysis, and plan recognition. In the following sections, we will quickly look into some real-life applications.

Pattern analysis

An active area of anomalous and suspicious behavior detection from patterns is based on visual modalities such as camera. Zhang et al (2007) proposed a system for a visual human motion analysis from a video sequence, which recognizes unusual behavior based on walking trajectories; Lin et al (2009) described a video surveillance system based on color features, distance features, and a count feature, where evolutionary techniques are used to measure observation similarity. The system tracks each person and classifies their behavior by analyzing their trajectory patterns. The system extracts a set of visual low-level features in different parts of the image, and performs a classification with SVMs to detect aggressive, cheerful, intoxicated, nervous, neutral, and tired behavior.

Transaction analysis

Transaction analysis assumes discrete states/transactions in contrast to continuous observations. A major research area is **Intrusion Detection (ID)** that aims at detecting attacks against information systems in general. There are two types of ID systems, signature-based and anomaly-based, that broadly follow the suspicious and anomalous pattern detection as described in the previous sections. A comprehensive review of ID approaches was published by Gyanchandani et al (2012).

Furthermore, applications in ambient-assisted living that are based on wearable sensors also fit to transaction analysis as sensing is typically event-based. Lymberopoulos et al (2008) proposed a system for automatic extraction of the users' **spatio-temporal** patterns encoded as sensor activations from the sensor network deployed inside their home. The proposed method, based on location, time, and duration, was able to extract frequent patterns using the Apriori algorithm and encode the most frequent patterns in the form of a Markov chain. Another area of related work includes **Hidden Markov Models (HMMs)** (Rabiner, 1989) that are widely used in traditional activity recognition for modeling a sequence of actions, but these topics are already out of scope of this book.

Plan recognition

Plan recognition focuses on a mechanism for recognizing the unobservable state of an agent, given observations of its interaction with its environment (Avrahami-Zilberbrand, 2009). Most existing investigations assume discrete observations in the form of activities. To perform anomalous and suspicious behavior detection, plan recognition algorithms may use a hybrid approach, a symbolic plan recognizer is used to filter consistent hypotheses, passing them to an evaluation engine, which focuses on ranking.

These were advanced approaches applied to various real-life scenarios targeted at discovering anomalies. In the following sections, we'll dive into more basic approaches for suspicious and anomalous pattern detection.

Fraud detection of insurance claims

First, we'll take a look at suspicious behavior detection, where the goal is to learn known patterns of frauds, which correspond to modeling known-knowns.

Dataset

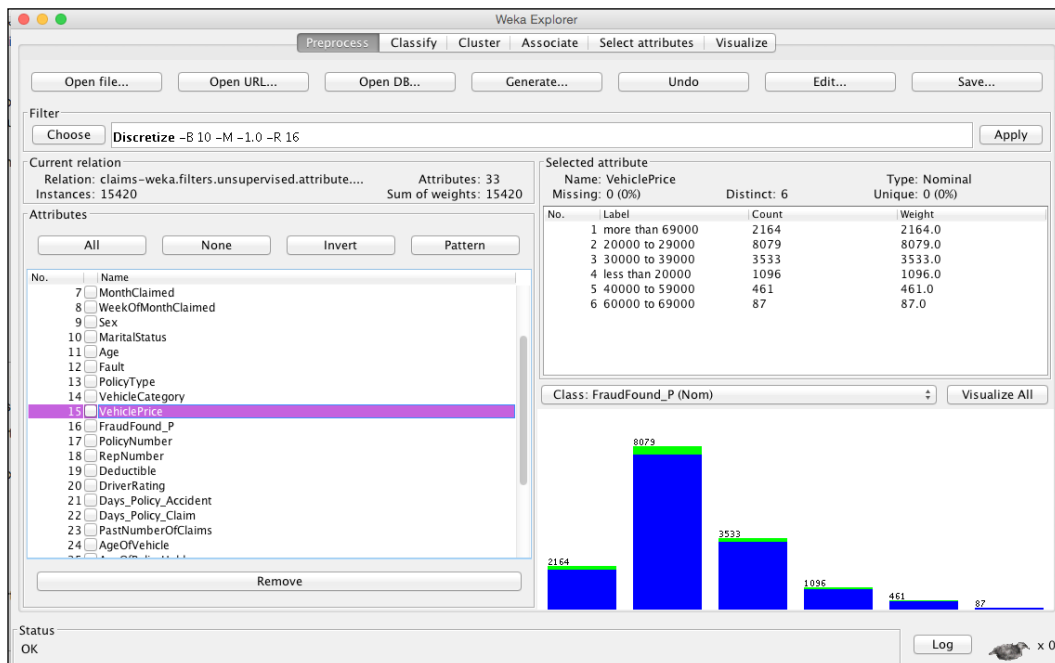
We'll work with a dataset describing insurance transactions publicly available at **Oracle Database Online Documentation** (2015), as follows:

http://docs.oracle.com/cd/B28359_01/datamine.111/b28129/anomalies.htm

The dataset describes insurance vehicle incident claims for an undisclosed insurance company. It contains 15,430 claims; each claim comprises 33 attributes describing the following components:

- Customer demographic details (**Age**, **Sex**, **MaritalStatus**, and so on)
- Purchased policy (**PolicyType**, **VehicleCategory**, number of supplements, agent type, and so on)
- Claim circumstances (day/month/week claimed, policy report filed, witness present, past days between incident-policy report, incident-claim, and so on)
- Other customer data (number of cars, previous claims, **DriverRating**, and so on)
- Fraud found (yes and no)

A sample of the database shown in the following screenshot depicts the data loaded into Weka:



Now the task is to create a model that will be able to identify suspicious claims in future. The challenging thing about this task is the fact that only 6% of claims are suspicious. If we create a dummy classifier saying no claim is suspicious, it will be accurate in 94% cases. Therefore, in this task, we will use different accuracy measures: precision and recall.

Recall the outcome table from *Chapter 1, Applied Machine Learning Quick Start*, where there are four possible outcomes denoted as true positive, false positive, false negative, and true negative:

		Classified as	
Actual		Fraud	No fraud
	Fraud	TP – true positive	FN – false negative
	No fraud	FP – false positive	TN – true negative

Precision and recall are defined as follows:

- Precision is equal to the proportion of correctly raised alarms, as follows:

$$\text{Pr} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

- Recall is equal to the proportion of deviant signatures, which are correctly identified as such:

$$\text{Re} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

- With these measures, our dummy classifier scores $\text{Pr} = 0$ and $\text{Re} = 0$ as it never marks any instance as fraud ($\text{TP} = 0$). In practice, we want to compare classifiers by both numbers, hence we use *F-measure*. This is a de-facto measure that calculates a harmonic mean between precision and recall, as follows:

$$F - \text{measure} = \frac{2 * \text{Pr} * \text{Re}}{\text{Pr} + \text{Re}}$$

Now let's move on to designing a real classifier.

Modeling suspicious patterns

To design a classifier, we can follow the standard supervised learning steps as described in *Chapter 1, Applied Machine Learning Quick Start*. In this recipe, we will include some additional steps to handle unbalanced dataset and evaluate classifiers based on precision and recall. The plan is as follows:

- Load the data in the .csv format
- Assign the class attribute
- Convert all the attributes from numeric to nominal in order to make sure there are no incorrectly loaded numerical values
- **Experiment 1:** Evaluate models with k-fold cross validation
- **Experiment 2:** Rebalance dataset to a more balanced class distribution and manually perform cross validation
- Compare classifiers by recall, precision, and f-measure

First, let's load the data using the CSVLoader class, as follows:

```
String filePath = "/Users/bostjan/Dropbox/ML Java Book/book/datasets/chap07/claims.csv";

CSVLoader loader = new CSVLoader();
loader.setFieldSeparator(",");
loader.setSource(new File(filePath));
Instances data = loader.getDataSet();
```

Next, we need to make sure all the attributes are nominal. During the data import, Weka applies some heuristics to guess the most probable attribute type, that is, numeric, nominal, string, or date. As heuristics cannot always guess the correct type, we can set types manually, as follows:

```
NumericToNominal toNominal = new NumericToNominal();
toNominal.setInputFormat(data);
data = Filter.useFilter(data, toNominal);
```

Before we continue, we need to specify the attribute that we will try to predict. We can achieve this by calling the `setClassIndex(int)` function:

```
int CLASS_INDEX = 15;
data.setClassIndex(CLASS_INDEX);
```

Next, we need to remove an attribute describing the policy number as it has no predictive value. We simply apply the `Remove` filter, as follows:

```
Remove remove = new Remove();
remove.setInputFormat(data);
remove.setOptions(new String[]{"-R", "+POLICY_INDEX});
data = Filter.useFilter(data, remove);
```

Now we are ready to start modeling.

Vanilla approach

The vanilla approach is to directly apply the lesson as demonstrated in *Chapter 3, Basic Algorithms – Classification, Regression, Clustering*, without any pre-processing and not taking into account dataset specifics. To demonstrate drawbacks of vanilla approach, we will simply build a model with default parameters and apply k-fold cross validation.

First, let's define some classifiers that we want to test:

```
ArrayList<Classifier>models = new ArrayList<Classifier>();
models.add(new J48());
models.add(new RandomForest());
models.add(new NaiveBayes());
models.add(new AdaBoostM1());
models.add(new Logistic());
```

Next, we create an `Evaluation` object and perform k-fold cross validation by calling the `crossValidate(Classifier, Instances, int, Random, String[])` method, outputting precision, recall, and fMeasure:

```
int FOLDS = 3;
Evaluation eval = new Evaluation(data);

for(Classifier model : models){
    eval.crossValidateModel(model, data, FOLDS,
        new Random(1), new String[] {});
    System.out.println(model.getClass().getName() + "\n"+
        "\tRecall: "+eval.recall(FRAUD) + "\n"+
        "\tPrecision: "+eval.precision(FRAUD) + "\n"+
        "\tF-measure: "+eval.fMeasure(FRAUD));
}
```

The evaluation outputs the following scores:

```
weka.classifiers.trees.J48
  Recall:    0.03358613217768147
  Precision: 0.9117647058823529
  F-measure: 0.06478578892371996
...
weka.classifiers.functions.Logistic
  Recall:    0.037486457204767065
  Precision: 0.2521865889212828
  F-measure: 0.06527070364082249
```

We can see the results are not very promising. Recall, that is, the share of discovered frauds among all frauds is only 1-3%, meaning that only 1-3/100 frauds are detected. On the other hand, precision, that is, the accuracy of alarms is 91%, meaning that in 9/10 cases, when a claim is marked as fraud, the model is correct.

Dataset rebalancing

As the number of negative examples, that is, frauds, is very small, compared to positive examples, the learning algorithms struggle with induction. We can help them by giving them a dataset, where the share of positive and negative examples is comparable. This can be achieved with dataset rebalancing.

Weka has a built-in filter, **Resample**, which produces a random subsample of a dataset using either sampling with replacement or without replacement. The filter can also bias distribution towards a uniform class distribution.

We will proceed by manually implementing k-fold cross validation. First, we will split the dataset into k equal folds. Fold k will be used for testing, while the other folds will be used for learning. To split dataset into folds, we'll use the `StratifiedRemoveFolds` filter, which maintains the class distribution within the folds, as follows:

```
StratifiedRemoveFolds kFold = new StratifiedRemoveFolds();
kFold.setInputFormat(data);

double measures[][] = new double[models.size()][3];

for(int k = 1; k <= FOLDS; k++){

    // Split data to test and train folds
    kFold.setOptions(new String[]{
```

```

"-N", ""+FOLDS, "-F", ""+k, "-S", "1"});
Instances test = Filter.useFilter(data, kFold);

kFold.setOptions(new String[]{
    "-N", ""+FOLDS, "-F", ""+k, "-S", "1", "-V"});
// select inverse "-V"
Instances train = Filter.useFilter(data, kFold);

```

Next, we can rebalance train dataset, where the `-z` parameter specifies the percentage of dataset to be resampled, and `-B` bias the class distribution towards uniform distribution:

```

Resample resample = new Resample();
resample.setInputFormat(data);
resample.setOptions(new String[]{"-Z", "100", "-B", "1"}); //with
replacement
Instances balancedTrain = Filter.useFilter(train, resample);

```

Next, we can build classifiers and perform evaluation:

```

for(ListIterator<Classifier>it = models.listIterator();
    it.hasNext();){
    Classifier model = it.next();
    model.buildClassifier(balancedTrain);
    eval = new Evaluation(balancedTrain);
    eval.evaluateModel(model, test);

// save results for average
    measures[it.previousIndex()][0] += eval.recall(FRAUD);
    measures[it.previousIndex()][1] += eval.precision(FRAUD);
    measures[it.previousIndex()][2] += eval.fMeasure(FRAUD);
}

```

Finally, we calculate the average and output the best model:

```

// calculate average
for(int i = 0; i < models.size(); i++){
    measures[i][0] /= 1.0 * FOLDS;
    measures[i][1] /= 1.0 * FOLDS;
    measures[i][2] /= 1.0 * FOLDS;
}

// output results and select best model
Classifier bestModel = null; double bestScore = -1;
for(ListIterator<Classifier> it = models.listIterator();
    it.hasNext();){

```

```
Classifier model = it.next();
double fMeasure = measures[it.previousIndex()][2];
System.out.println(
    model.getClass().getName() + "\n"+
    "\tRecall:    "+measures[it.previousIndex()][0] + "\n"+
    "\tPrecision: "+measures[it.previousIndex()][1] + "\n"+
    "\tF-measure: "+fMeasure);
if(fMeasure > bestScore){
    bestScore = fMeasure;
    bestModel = model;
}
}
System.out.println("Best model:"+bestModel.getClass().getName());
```

Now the performance of the models has significantly improved, as follows:

```
weka.classifiers.trees.J48
  Recall:    0.44204845100610574
  Precision: 0.14570766048577555
  F-measure: 0.21912423640160392
...
weka.classifiers.functions.Logistic
  Recall:    0.7670657247204478
  Precision: 0.13507459756495374
  F-measure: 0.22969038530557626
Best model: weka.classifiers.functions.Logistic
```

What we can see is that all the models have scored significantly better; for instance, the best model, Logistic Regression, correctly discovers 76% of frauds, while producing a reasonable amount of false alarms—only 13% of claims marked as fraud are indeed fraudulent. If an undetected fraud is significantly more expensive than investigation of false alarms, then it makes sense to deal with an increased number of false alarms.

The overall performance has most likely still some room for improvement; we could perform attribute selection and feature generation and apply more complex model learning that we discussed in *Chapter 3, Basic Algorithms – Classification, Regression, Clustering*.

Anomaly detection in website traffic

In the second example, we'll focus on modeling the opposite of the previous example. Instead of discussing what typical fraud-less cases are, we'll discuss the normal expected behavior of the system. If something cannot be matched against our expected model, it will be considered anomalous.

Dataset

We'll work with a publicly available dataset released by Yahoo Labs that is useful for discussing how to detect anomalies in time series data. For Yahoo, the main use case is in detecting unusual traffic on Yahoo servers.

Even though Yahoo announced that their data is publicly available, you have to apply to use it, and it takes about 24 hours before the approval is granted. The dataset is available here:

<http://webscope.sandbox.yahoo.com/catalog.php?datatype=s&did=70>

The data set comprises real traffic to Yahoo services, along with some synthetic data. In total, the dataset contains 367 time series, each of which contain between 741 and 1680 observations, recorded at regular intervals. Each series is written in its own file, one observation per line. A series is accompanied by a second column indicator with a one if the observation was an anomaly, and zero otherwise. The anomalies in real data were determined by human judgment, while those in the synthetic data were generated algorithmically. A snippet of the synthetic times series data is shown in the following table:

timestamp	value	anomaly	change point	trend	noise	12 hour seasonality	daily seasonality	weekly seasonality
1422237600	4333.43	0	0	4599	1.81	-190.95	-128.86	52.44
1422241200	4316.14	0	0	4602	-14.65	-220.5	-105.21	54.51
1422244800	4403.20	0	0	4605	7.04	-190.95	-74.39	56.51
1422248400	4531.20	0	0	4608	13.52	-110.25	-38.51	58.43
1422252000	4967.50	1	0	4911	-3.77	-6.91	-2.33	60.27

Snippet of the synthetic time-series data

In the following section, we'll learn how to transform time series data to attribute presentation that allows us to apply machine learning algorithms.

Anomaly detection in time series data

Detecting anomalies in raw, streaming time series data requires some data transformations. The most obvious way is to select a time window and sample time series with fixed length. In the next step, we want to compare a new time series to our previously collected set to detect if something is out of the ordinary.

The comparison can be done with various techniques, as follows:

- Forecasting the most probable following value, as well as confidence intervals (for example, Holt-Winters exponential smoothing). If a new value is out of forecasted confidence interval, it is considered anomalous.
- Cross correlation compares new sample to library of positive samples, it looks for exact match. If the match is not found, it is marked as anomalous.
- Dynamic time wrapping is similar to cross correlation, but allows signal distortion in comparison.
- Discretizing signal to bands, where each band corresponds to a letter. For example, $A=[min, mean/3]$, $B=[mean/3, mean*2/3]$, and $C=[mean*2/3, max]$ transforms the signal to a sequence of letters such as *aAABAACAABBA....*. This approach reduces the storage and allows us to apply text-mining algorithms that we will discuss in *Chapter 10, Text Mining with Mallet – Topic Modeling and Spam Detection*.
- Distribution-based approach estimates distribution of values in a specific time window. When we observe a new sample, we can compare whether distribution matches to the previously observed one.

This list is, by no means, exhaustive. Different approaches are focused on detecting different anomalies (for example, in value, frequency, and distribution). We will focus on a version of distribution-based approaches.

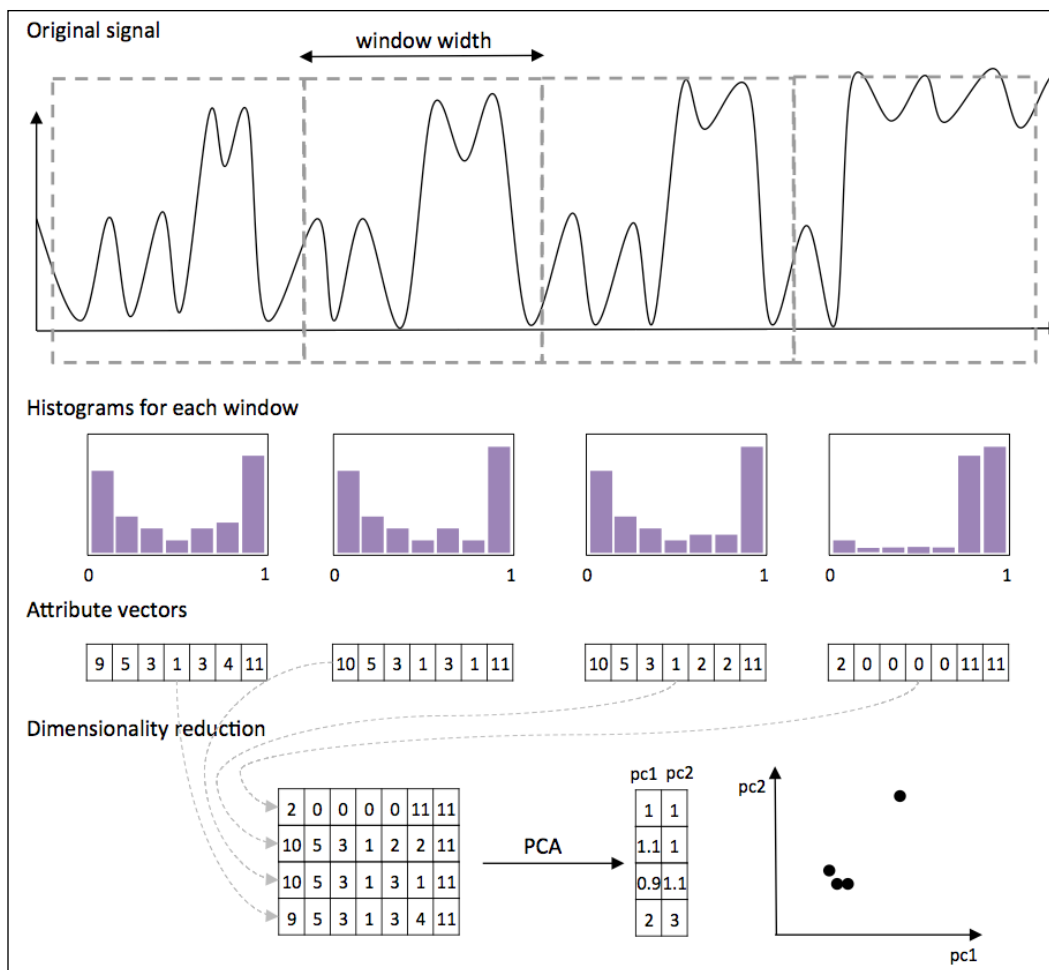
Histogram-based anomaly detection

In histogram-based anomaly detection, we split signals by some selected time window as shown in the following image.

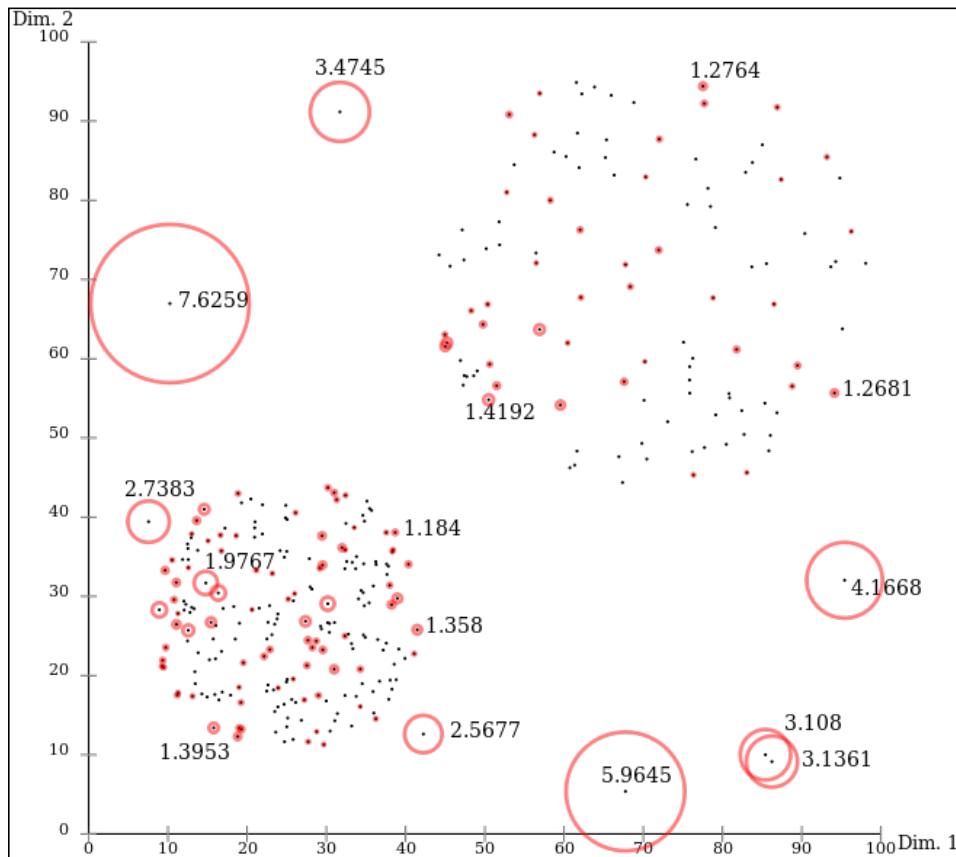
For each window, we calculate the histogram, that is, for a selected number of buckets, we count how many values fall into each bucket. The histogram captures basic distribution of values in a selected time window as shown at the center of the diagram.

Histograms can be then directly presented as instances, where each bin corresponds to an attribute. Further, we can reduce the number of attributes by applying a dimensionality-reduction technique such as **Principal Component Analysis (PCA)**, which allows us to visualize the reduced-dimension histograms in a plot as shown at the bottom-right of the diagram, where each dot corresponds to a histogram.

In our example, the idea is to observe website traffic for a couple of days and then to create histograms, for example, four-hour time windows to build a library of positive behavior. If a new time window histogram cannot be matched against positive library, we can mark it as an anomaly:



For comparing a new histogram to a set of existing histograms, we will use a density-based k-nearest neighbor algorithm, **Local Outlier Factor (LOF)** (Breunig et al, 2000). The algorithm is able to handle clusters with different densities as shown in the following image. For example, the upper-right cluster is large and widespread as compared to the bottom-left cluster, which is smaller and denser:



Let's get started!

Loading the data

In the first step, we'll need to load the data from text files to a Java object. The files are stored in a folder, each file contains one-time series with values per line. We'll load them into a list of Double:

```
String filePath = "chap07/ydata/AlBenchmark/real";  
List<List<Double>> rawData = new ArrayList<List<Double>>();
```

We will need the min and max value for histogram normalization, so let's collect them in this data pass:

```
double max = Double.MIN_VALUE;
double min = Double.MAX_VALUE;

for(int i = 1; i<= 67; i++){
    List<Double> sample = new ArrayList<Double>();
    BufferedReader reader = new BufferedReader(new
        FileReader(filePath+i+".csv"));

    boolean isAnomaly = false;
    reader.readLine();
    while(reader.ready()){
        String line[] = reader.readLine().split(",");
        double value = Double.parseDouble(line[1]);
        sample.add(value);

        max = Math.max(max, value);
        min = Double.min(min, value);

        if(line[2] == "1")
            isAnomaly = true;
    }
    System.out.println(isAnomaly);
    reader.close();

    rawData.add(sample);
}
```

The data is loaded, now let's move on to histograms.

Creating histograms

We will create a histogram for a selected time window with the `WIN_SIZE` width. The histogram will hold the `HIST_BINS` value buckets. The histograms consisting of list of doubles will be stored into an array list:

```
int WIN_SIZE = 500;
int HIST_BINS = 20;
int current = 0;

List<double[]> dataHist = new ArrayList<double[]>();
for(List<Double> sample : rawData){
```

```
double[] histogram = new double[HIST_BINS];
for(double value : sample){
    int bin = toBin(normalize(value, min, max), HIST_BINS);
    histogram[bin]++;
    current++;
    if(current == WIN_SIZE){
        current = 0;
        dataHist.add(histogram);
        histogram = new double[HIST_BINS];
    }
}
dataHist.add(histogram);
}
```

Histograms are now completed. The last step is to transform them into Weka's Instance objects. Each histogram value will correspond to one Weka attribute, as follows:

```
ArrayList<Attribute> attributes = new ArrayList<Attribute>();
for(int i = 0; i<HIST_BINS; i++){
    attributes.add(new Attribute("Hist-"+i));
}
Instances dataset = new Instances("My dataset", attributes,
    dataHist.size());
for(double[] histogram: dataHist){
    dataset.add(new Instance(1.0, histogram));
}
```

The dataset is now loaded and ready to be plugged into an anomaly-detection algorithm.

Density based k-nearest neighbors

To demonstrate how LOF calculates scores, we'll first split the dataset into training and testing set using the `testCV(int, int)` function. The first parameter specifies the number of folds, while the second parameter specifies which fold to return.

```
// split data to train and test
Instances trainData = dataset.testCV(2, 0);
Instances testData = dataset.testCV(2, 1);
```

The LOF algorithm is not a part of the default Weka distribution, but it can be downloaded through Weka's package manager:

<http://weka.sourceforge.net/packageMetaData/localOutlierFactor/index.html>

LOF algorithm has two implemented interfaces: as an unsupervised filter that calculates LOF values (known-unknowns) and as a supervised k-nn classifier (known-knowns). In our case, we want to calculate the outlier-ness factor, therefore, we'll use the unsupervised filter interface:

```
import weka.filters.unsupervised.attribute.LOF;
```

The filter is initialized the same way as a usual filter. We can specify the k number of neighbors, for example, $k=3$, with `-min` and `-max` parameters. LOF allows us to specify two different k parameters, which are used internally as the upper and lower bound to find the minimal/maximal number of values:

```
LOF lof = new LOF();
lof.setInputFormat(trainData);
lof.setOptions(new String[]{"-min", "3", "-max", "3"});
```

Next, we load training instances into the filter that will serve as a positive example library. After we complete the loading, we call the `batchFinished()` method to initialize internal calculations:

```
for(Instance inst : trainData){
    lof.input(inst);
}
lof.batchFinished();
```

Finally, we can apply the filter to test data. Filter will process the instances and append an additional attribute at the end containing the LOF score. We can simply output the score on the console:

```
Instances testDataLofScore = Filter.useFilter(testData, lof);

for(Instance inst : testDataLofScore){
    System.out.println(inst.value(inst.numAttributes()-1));
}
```

The LOF score of the first couple of test instances is as follows:

```
1.306740014927325
1.318239332210458
1.0294812291949587
1.1715039094530768
```

To understand the LOF values, we need some background on the LOF algorithm. It compares the density of an instance to the density of its nearest neighbors. The two scores are divided, producing the LOF score. The LOF score around 1 indicates that the density is approximately equal, while higher LOF values indicate that the density of the instance is substantially lower than the density of its neighbors. In such cases, the instance can be marked as anomalous.

Summary

In this chapter, we looked into detecting anomalous and suspicious patterns. We discussed the two fundamental approaches focusing on library encoding either positive or negative patterns. Next, we got our hands on two real-life datasets, where we discussed how to deal with unbalanced class distribution and perform anomaly detection in time series data.

In the next chapter, we'll dive deeper into patterns and more advanced approaches to build pattern-based classifier, discussing how to automatically assign labels to images with deep learning.

8

Image Recognition with Deeplearning4j

Images have become ubiquitous in web services, social networks, and web stores. In contrast to humans, computers have great difficulty in understanding what is in the image and what does it represent. In this chapter, we'll first look at the challenges required to teach computers how to understand images, and then focus on an approach based on deep learning. We'll look at a high-level theory required to configure a deep learning model and discuss how to implement a model that is able to classify images using a Java library, Deeplearning4j.

This chapter will cover the following topics:

- Introducing image recognition
- Discussing deep learning fundamentals
- Building an image recognition model

Introducing image recognition

A typical goal of image recognition is to detect and identify an object in a digital image. Image recognition is applied in factory automation to monitor product quality; surveillance systems to identify potentially risky activities, such as moving persons or vehicles; security applications to provide biometric identification through fingerprints, iris, or facial features; autonomous vehicles to reconstruct conditions on the road and environment and so on.

Digital images are not presented in a structured way with attribute-based descriptions; instead, they are encoded as the amount of color in different channels, for instance, black-white and red-green-blue channels. The learning goal is to then identify patterns associated with a particular object. The traditional approach for image recognition consists of transforming an image into different forms, for instance, identify object corners, edges, same-color blobs, and basic shapes. Such patterns are then used to train a learner to distinguish between objects. Some notable examples of traditional algorithms are:

- Edge detection finds boundaries of objects within an image
- Corner detection identifies intersections of two edges or other interesting points, such as line endings, curvature maxima/minima, and so on
- Blob detection identifies regions that differ in a property, such as brightness or color, compared to its surrounding regions
- Ridge detection identifies additional interesting points at the image using smooth functions
- **Scale Invariant Feature Transform (SIFT)** is a robust algorithm that can match objects even if their scale or orientation differs from the representative samples in database
- Hough transform identifies particular patterns in the image

A more recent approach is based on deep learning. Deep learning is a form of neural network, which mimics how the brain processes information. The main advantage of deep learning is that it is possible to design neural networks that can automatically extract relevant patterns, which in turn, can be used to train a learner. With recent advances in neural networks, image recognition accuracy was significantly boosted. For instance, the **ImageNet** challenge (ImageNet, 2016), where competitors are provide more than 1.2 million images from 1,000 different object categories, reports that the error rate of the best algorithm was reduced from 28% in 2010, using SVM, to only 7% in 2014, using deep neural network.

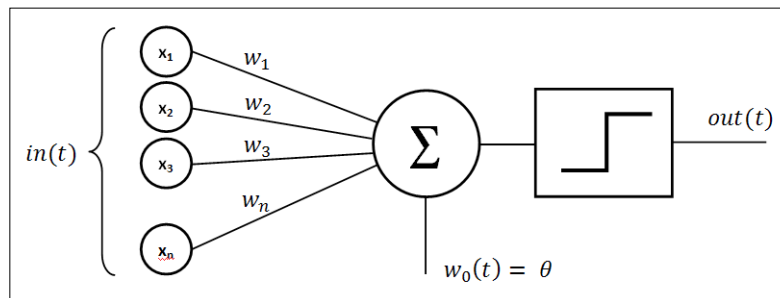
In this chapter, we'll take a quick look at neural networks, starting from the basic building block – **perceptron** – and gradually introducing more complex structures.

Neural networks

The first neural networks, introduced in the sixties, are inspired by biological neural networks. Recent advances in neural networks proved that deep neural networks fit very well in pattern recognition tasks, as they are able to automatically extract interesting features and learn the underlying presentation. In this section, we'll refresh the fundamental structures and components from a single perceptron to deep networks.

Perceptron

Perceptron is a basic neural network building block and one of the earliest supervised algorithms. It is defined as a sum of features, multiplied by corresponding weights and a bias. The function that sums all of this together is called **sum transfer function** and it is fed into an **activation function**. If the binary step activation function reaches a threshold, the output is 1, otherwise 0, which gives us a binary classifier. A schematic illustration is shown in the following diagram:

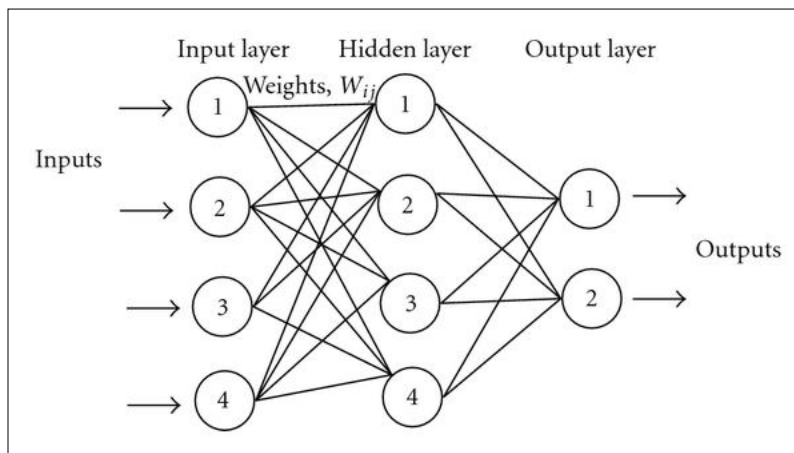


Training perceptrons involves a fairly simple learning algorithm that calculates the errors between the calculated output values and correct training output values, and uses this to create an adjustment to the weights; thus implementing a form of gradient descent. This algorithm is usually called the **delta rule**.

Single-layer perceptron is not very advanced, and nonlinearly separable functions, such as XOR, cannot be modeled using it. To address this issue, a structure with multiple perceptrons was introduced, called **multilayer perceptron**, also known as **feedforward neural network**.

Feedforward neural networks

A feedforward neural network is an artificial neural network that consists of several perceptrons, which are organized by layers, as shown in the following diagram: input layer, output layer, and one or more hidden layers. Each layer perceptron, also known as neuron, has direct connections to the perceptrons in the next layer; whereas, connections between two neurons carry a weight similar to the perceptron weights. The diagram shows a network with a four-unit **Input layer**, corresponding to the size of feature vector of length 4, a four-unit **Hidden layer**, and a two-unit **Output layer**, where each unit corresponds to one class value:



The most popular approach to train multilayer networks is backpropagation. In backpropagation, the calculated output values are compared with the correct values in the same way as in delta rule. The error is then fed back through the network by various techniques, adjusting the weights of each connection in order to reduce the value of the error. The process is repeated for sufficiently large number of training cycles, until the error is under a certain threshold.

Feedforward neural network can have more than one hidden layer; whereas, each additional hidden layer builds a new abstraction atop the previous layers. This often leads to more accurate models; however, increasing the number of hidden layers leads to the following two known issues:

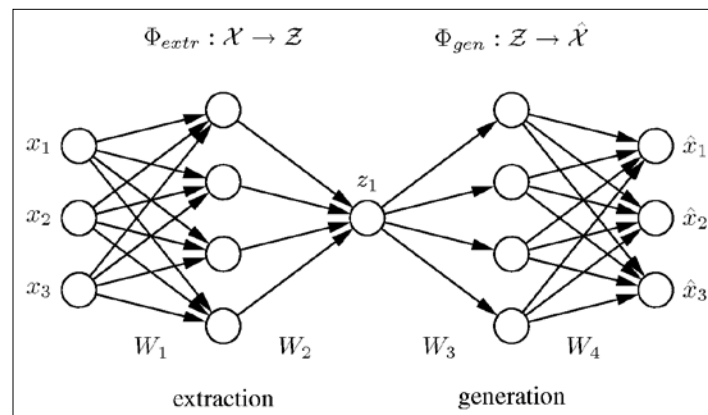
- **Vanishing gradients problem:** With more hidden layers, the training with backpropagation becomes less and less useful for passing information to the front layers, causing these layers to train very slowly
- **Overfitting:** The model fits the training data too well and performs poorly on real examples

Let's look at some other networks structures that address these issues.

Autoencoder

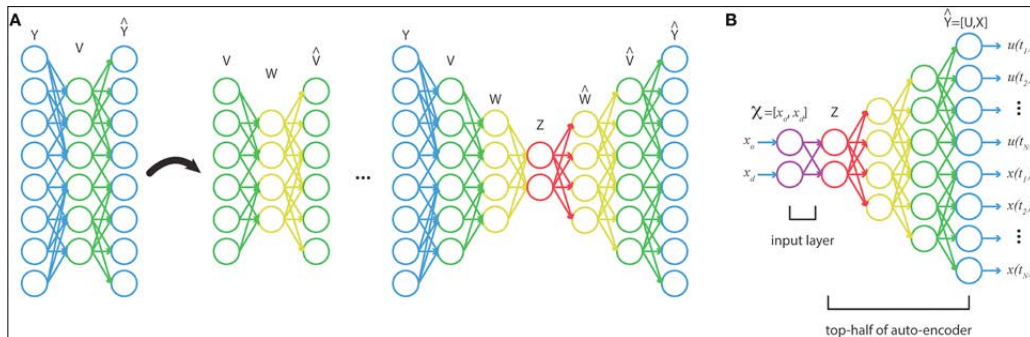
Autoencoder is a feedforward neural network that aims to learn how to compress the original dataset. Therefore, instead of mapping features to input layer and labels to output layer, we will map the features to both input and output layers. The number of units in hidden layers is usually different from the number of units in input layers, which forces the network to either expand or reduce the number of original features. This way the network will learn the important features, while effectively applying dimensionality reduction.

An example network is shown below. The three-unit input layer is first expanded into a four-unit layer and then compressed into a single-unit layer. The other side of the network restores the single-layer unit back to the four-unit layer, and then to the original three-input layer:



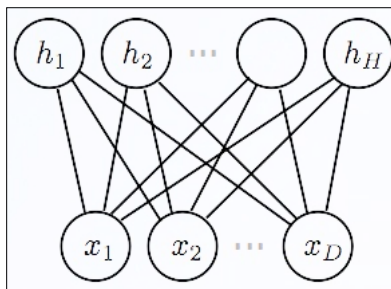
Once the network is trained, we can take the left-hand side to extract image features as we would with traditional image processing.

The autoencoders can be also combined into **stacked autoencoders**, as shown in the following image. First, we will discuss the hidden layer in a basic autoencoder, as described previously. Then we will take the learned hidden layer (green circles) and repeat the procedure, which in effect, learns a more abstract presentation. We can repeat the procedure multiple times, transforming the original features into increasingly reduced dimensions. At the end, we will take all the hidden layers and stack them into a regular feedforward network, as shown at the top-right part of the diagram:



Restricted Boltzmann machine

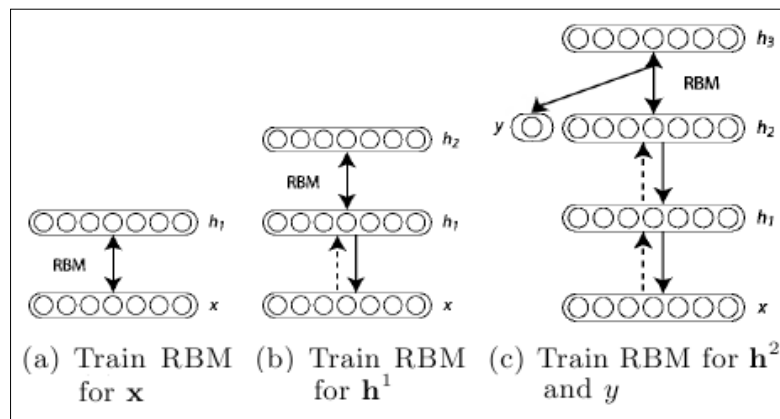
Restricted Boltzmann machine is an undirected neural network, also denoted as **Generative Stochastic Networks (GSNs)**, and can learn probability distribution over its set of inputs. As the name suggests, they originate from Boltzmann machine, a recurrent neural network introduced in the eighties. Restricted means that the neurons must form two fully connected layers – input layer and hidden layer – as show in the following diagram:



Unlike feedforward networks, the connections between the visible and hidden layers are undirected, hence the values can be propagated in both visible-to-hidden and hidden-to-visible directions.

Training Restricted Boltzmann machines is based on the **Contrastive Divergence** algorithm, which uses a gradient descent procedure, similar to backpropagation, to update weights, and **Gibbs sampling** is applied on the **Markov chain** to estimate the gradient – the direction on how to change the weights.

Restricted Boltzmann machines can also be stacked to create a class known as **Deep Belief Networks (DBNs)**. In this case, the hidden layer of RBM acts as a visible layer for the RBM layer, as shown in the following diagram:



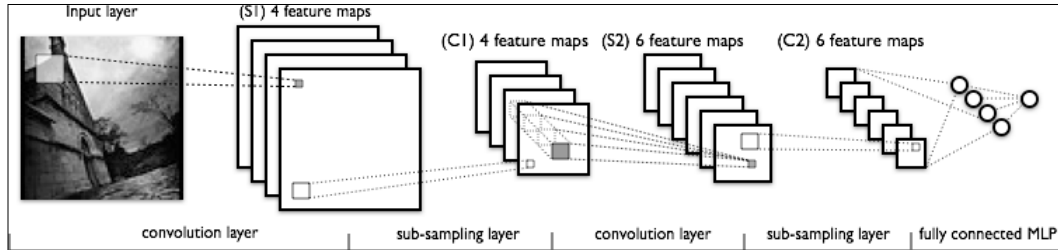
The training, in this case, is incremental; training layer by layer.

Deep convolutional networks

A network structure that recently achieves very good results at image recognition benchmarks is **Convolutional Neural Network (CNN)** or ConvNet. CNNs are a type of feedforward neural network that are structured in such a way that emulates behavior of the visual cortex, exploiting 2D structures of an input image, that is, patterns that exhibit spatially local correlation.

A CNN consists of a number of convolutional and subsampling layers, optionally followed by fully connected layers. An example is shown in the following image. The input layer reads all the pixels at an image and then we apply multiple filters. In the following image, four different filters are applied. Each filter is applied to the original image, for example, one pixel of a 6×6 filter is calculated as the weighted sum of a 6×6 square of input pixels and corresponding 6×6 weights. This effectively introduces filters similar to the standard image processing, such as smoothing, correlation, edge detection, and so on. The resulting image is called **feature map**. In the example in the image, we have four feature maps, one for each filter.


The next layer is the subsampling layer, which reduces the size of the input. Each feature map is subsampled typically with mean or max pooling over a contiguous region of 2×2 (up to 5×5 for large images). For example, if the feature map size is 16×16 and the subsampling region is 2×2 , the reduced feature map size is 8×8 , where 4 pixels (2×2 square) are combined into a single pixel by calculating max, min, mean, or some other functions:



The network may contain several consecutive convolution and subsampling layers, as shown in the preceding diagram. A particular feature map is connected to the next reduced/convoluted feature map, while feature maps at the same layer are not connected to each other.

After the last subsampling or convolutional layer, there is usually a fully connected layer, identical to the layers in a standard multilayer neural network, which represents the target data.

CNN is trained using a modified backpropagation algorithm that takes the subsampling layers into account and updates the convolutional filter weights based on all the values where this filter is applied.

[ Some good CNN designs can be found at the ImageNet competition results page:
<http://www.image-net.org/>
An example is AlexNet, described in the *ImageNet Classification with Deep Convolutional Neural Networks* paper by A. Krizhevsky et al.]

This concludes our review of main neural network structures. In the following section, we'll move to the actual implementation.

Image classification

In this section, we will discuss how to implement some of the neural network structures with the `deeplearning4j` library. Let's start.

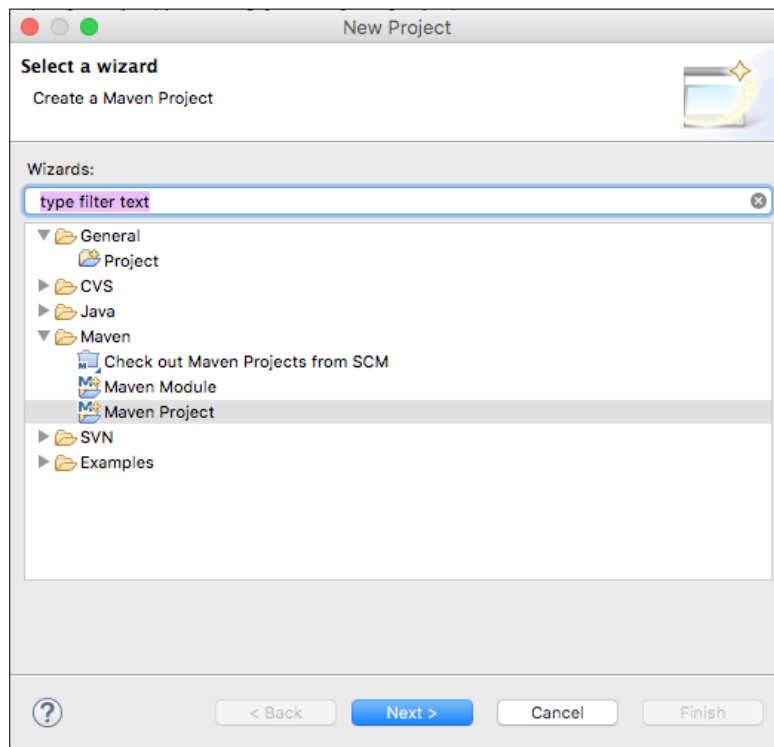
Deeplearning4j

As we discussed in *Chapter 2, Java Libraries and Platforms for Machine Learning*, `deeplearning4j` is an open source, distributed deep learning project in Java and Scala. `Deeplearning4j` relies on Spark and Hadoop for MapReduce, trains models in parallel, and iteratively averages the parameters they produce in a central model. A detailed library summary is presented in *Chapter 2, Java Libraries and Platforms for Machine Learning*.

Getting DL4J

The most convenient way to get `deeplearning4j` is through the Maven repository:

1. Start a new Eclipse project and pick **Maven Project**, as shown in the following screenshot:



2. Open the `pom.xml` file and add the following dependencies under the `<dependencies>` section:

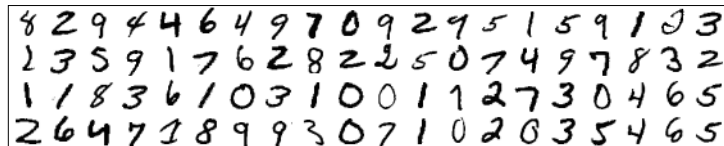
```
<dependency>
  <groupId>org.deeplearning4j</groupId>
  <artifactId>deeplearning4j-nlp</artifactId>
  <version>${dl4j.version}</version>
</dependency>

<dependency>
  <groupId>org.deeplearning4j</groupId>
  <artifactId>deeplearning4j-core</artifactId>
  <version>${dl4j.version}</version>
</dependency>
```

3. Finally, right-click on **Project**, select **Maven**, and pick **Update project**.

MNIST dataset

One of the most famous datasets is MNIST dataset, which consists of handwritten digits, as shown in the following image. The dataset comprises 60,000 training and 10,000 testing images:



The dataset is commonly used in image recognition problems to benchmark algorithms. The worst recorded error rate is 12%, with no preprocessing and using a SVM in one-layer neural network. Currently, as of 2016, the lowest error rate is only 0.21%, using the **DropConnect** neural network, followed by **deep convolutional network** at 0.23%, and deep feedforward network at 0.35%.

Now, let's see how to load the dataset.

Loading the data

Deeplearning4j provides the MNIST dataset loader out of the box. The loader is initialized as `DataSetIterator`. Let's first import the `DataSetIterator` class and all the supported datasets that are part of the `impl` package, for example, Iris, MNIST, and others:

```
import org.deeplearning4j.datasets.iterator.DataSetIterator;
import org.deeplearning4j.datasets.iterator.impl.*;
```

Next, we'll define some constants, for instance, the images consist of 28 x 28 pixels and there are 10 target classes and 60,000 samples. Initialize a new `MnistDataSetIterator` class that will download the dataset and its labels. The parameters are iteration batch size, total number of examples, and whether the datasets should be binarized or not:

```
int numRows = 28;
int numColumns = 28;
int outputNum = 10;
int numSamples = 60000;
int batchSize = 100;
DataSetIterator iter = new MnistDataSetIterator(batchSize,
numSamples,true);
```

Having an already-implemented data importer is really convenient, but it won't work on your data. Let's take a quick look at how is it implemented and what needs to be modified to support your dataset. If you're eager to start implementing neural networks, you can safely skip the rest of this section and return to it when you need to import your own data.



To load the custom data, you'll need to implement two classes: `DataSetIterator` that holds all the information about the dataset and `BaseDataFetcher` that actually pulls the data either from file, database, or web. Sample implementations are available on GitHub at <https://github.com/deeplearning4j/deeplearning4j/tree/master/deeplearning4j-core/src/main/java/org/deeplearning4j/datasets/iterator/impl>.

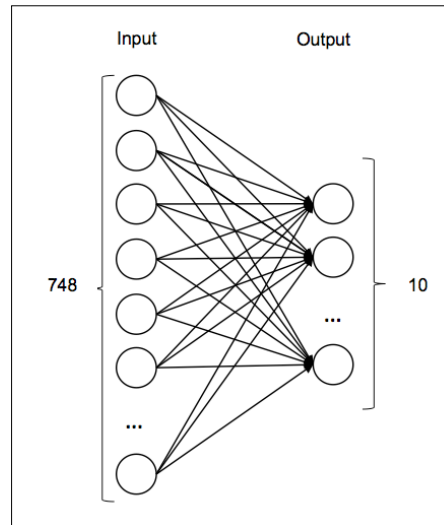
Another option is to use the **Canova** library, which is developed by the same authors, at <http://deeplearning4j.org/canovadoc/>.

Building models

In this section, we'll discuss how to build an actual neural network model. We'll start with a basic single-layer neural network to establish a benchmark and discuss the basic operations. Later, we'll improve this initial result with DBN and **Multilayer Convolutional Network**.

Building a single-layer regression model

Let's start by building a single-layer regression model based on the softmax activation function, as shown in the following diagram. As we have a single layer, **Input** to the neural network will be all the figure pixels, that is, $28 \times 28 = 748$ neurons. The number of **Output** neurons is **10**, one for each digit. The network layers are fully connected, as shown in the following diagram:



A neural network is defined through a `NeuralNetConfiguration Builder` object as follows:

```
MultiLayerConfiguration conf = new  
NeuralNetConfiguration.Builder()
```

We will define the parameters for gradient search in order to perform iterations with the **conjugate gradient optimization** algorithm. The momentum parameter determines how fast the optimization algorithm converges to an local optimum – the higher the momentum, the faster the training; but higher speed can lower model's accuracy, as follows:

```
.seed(seed)  
.gradientNormalization(GradientNormalization.ClipElementWiseAbsoluteValue)  
.gradientNormalizationThreshold(1.0)  
.iterations(iterations)  
.momentum(0.5)  
.momentumAfter(Collections.singletonMap(3, 0.9))  
.optimizationAlgo(OptimizationAlgorithm.CONJUGATE_GRADIENT)
```

Next, we will specify that the network will have one layer and define the error function (`NEGATIVELOGLIKELIHOOD`), internal perceptron activation function (`softmax`), and the number of input and output layers that correspond to total image pixels and the number of target variables:

```
.list(1)
.layer(0, new
OutputLayer.Builder(LossFunction.NEGATIVELOGLIKELIHOOD)
.activation("softmax")
.nIn(numRows*numColumns).nOut(outputNum).build())
```

Finally, we will set the network to `pretrain`, disable backpropagation, and actually build the untrained network structure:

```
.pretrain(true).backprop(false)
.build();
```

Once the network structure is defined, we can use it to initialize a new `MultiLayerNetwork`, as follows:

```
MultiLayerNetwork model = new MultiLayerNetwork(conf);
model.init();
```

Next, we will point the model to the training data by calling the `setListeners` method, as follows:

```
model.setListeners(Collections.singletonList((IterationListener)
new ScoreIterationListener(listenerFreq)));
```

We will also call the `fit(int)` method to trigger an end-to-end network training:

```
model.fit(iter);
```

To evaluate the model, we will initialize a new `Evaluation` object that will store batch results:

```
Evaluation eval = new Evaluation(outputNum);
```

We can then iterate over the dataset in batches in order to keep the memory consumption at a reasonable rate and store the results in an `eval` object:

```
DataSetIterator testIter = new MnistDataSetIterator(100,10000);
while(testIter.hasNext()) {
    DataSet testMnist = testIter.next();
    INDArray predict2 =
    model.output(testMnist.getFeatureMatrix());
    eval.eval(testMnist.getLabels(), predict2);
}
```

Finally, we can get the results by calling the `stats()` function:

```
log.info(eval.stats());
```

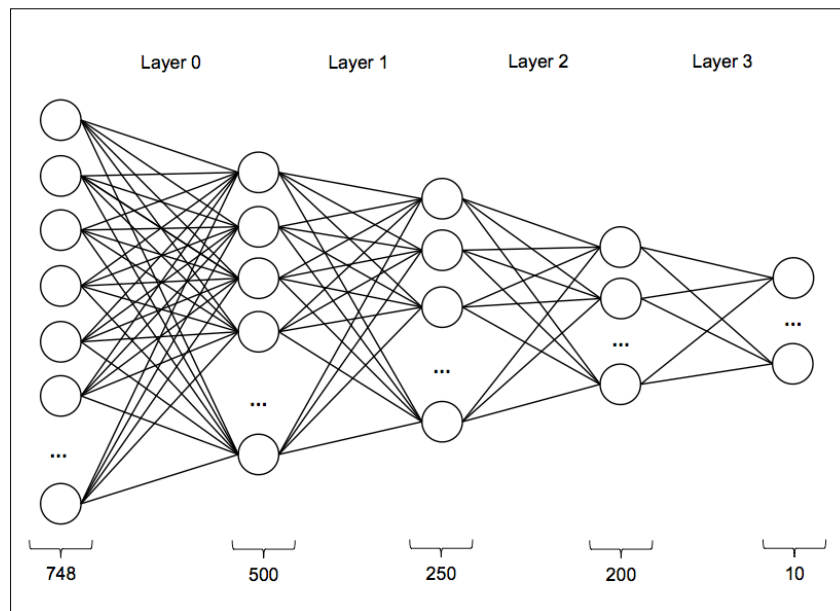
A basic one-layer model achieves the following accuracy:

```
Accuracy: 0.8945  
Precision: 0.8985  
Recall: 0.8922  
F1 Score: 0.8953
```

Getting 89.22% accuracy, that is, 10.88% error rate, on MNIST dataset is quite bad. We'll improve that by going from a simple one-layer network to the moderately sophisticated deep belief network using Restricted Boltzmann machines and Multilayer Convolutional Network.

Building a deep belief network

In this section, we'll build a deep belief network based on Restricted Boltzmann machine, as shown in the following diagram. The network consists of four layers: the first layer reduces the **748** inputs to **500** neurons, then to **250**, followed by **200**, and finally to the last **10** target values:



As the code is the same as in the previous example, let's take a look at how to configure such a network:

```
MultiLayerConfiguration conf = new
    NeuralNetConfiguration.Builder()
```

We defined the gradient optimization algorithm, as shown in the following code:

```
    .seed(seed)
    .gradientNormalization(
        GradientNormalization.ClipElementWiseAbsoluteValue)
    .gradientNormalizationThreshold(1.0)
    .iterations(iterations)
    .momentum(0.5)
    .momentumAfter(Collections.singletonMap(3, 0.9))
    .optimizationAlgo(OptimizationAlgorithm.CONJUGATE_GRADIENT)
```

We will also specify that our network will have four layers:

```
    .list(4)
```

The input to the first layer will be 748 neurons and the output will be 500 neurons. We'll use the root mean squared-error cross entropy, Xavier algorithm, to initialize weights by automatically determining the scale of initialization based on the number of input and output neurons, as follows:

```
    .layer(0, new RBM.Builder()
        .nIn(numRows*numColumns)
        .nOut(500)
        .weightInit(WeightInit.XAVIER)
        .lossFunction(LossFunction.RMSE_XENT)
        .visibleUnit(RBM.VisibleUnit.BINARY)
        .hiddenUnit(RBM.HiddenUnit.BINARY)
        .build())
```

The next two layers will have the same parameters, except the number of input and output neurons:

```
    .layer(1, new RBM.Builder()
        .nIn(500)
        .nOut(250)
        .weightInit(WeightInit.XAVIER)
        .lossFunction(LossFunction.RMSE_XENT)
        .visibleUnit(RBM.VisibleUnit.BINARY)
        .hiddenUnit(RBM.HiddenUnit.BINARY)
        .build())
```

```
.layer(2, new RBM.Builder()
.nIn(250)
.nOut(200)
.weightInit(WeightInit.XAVIER)
.lossFunction(LossFunction.RMSE_XENT)
.visibleUnit(RBM.VisibleUnit.BINARY)
.hiddenUnit(RBM.HiddenUnit.BINARY)
.build())
```

Now the last layer will map the neurons to outputs, where we'll use the `softmax` activation function, as follows:

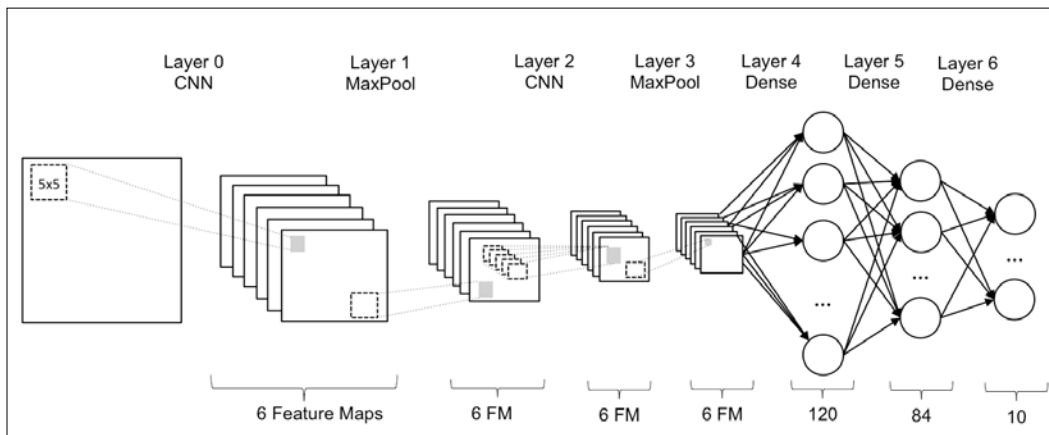
```
.layer(3, new OutputLayer.Builder()
.nIn(200)
.nOut(outputNum)
.lossFunction(LossFunction.NEGATIVELOGLIKELIHOOD)
.activation("softmax")
.build())
.pretrain(true).backprop(false)
.build();
```

The rest of the training and evaluation is the same as in the single-layer network example. Note that training deep network might take significantly more time compared to a single-layer network. The accuracy should be around 93%.

Now let's take a look at another deep network.

Build a Multilayer Convolutional Network

In the final example, we'll discuss how to build a convolutional network, as shown in the following diagram. The network will consist of seven layers: first, we'll repeat two pairs of convolutional and subsampling layers with max pooling. The last subsampling layer is then connected to a densely connected feedforward neuronal network, comprising 120 neurons, 84 neurons, and 10 neurons in the last three layers, respectively. Such a network effectively forms the complete image recognition pipeline, where the first four layers correspond to feature extraction and the last three layers correspond to the learning model:



Network configuration is initialized as we did earlier:

```
MultiLayerConfiguration.Builder conf = new
    NeuralNetConfiguration.Builder()
```

We will specify the gradient descent algorithm and its parameters, as follows:

```
.seed(seed)
.iterations(iterations)
.activation("sigmoid")
.weightInit(WeightInit.DISTRIBUTION)
.dist(new NormalDistribution(0.0, 0.01))
.learningRate(1e-3)
.learningRateScoreBasedDecayRate(1e-1)
.optimizationAlgo(
    OptimizationAlgorithm.STOCHASTIC_GRADIENT_DESCENT)
```

We will also specify the seven network layers, as follows:

```
.list(7)
```

The input to the first convolutional layer is the complete image, while the output is six feature maps. The convolutional layer will apply a 5 x 5 filter, and the result will be stored in a 1 x 1 cell:

```
.layer(0, new ConvolutionLayer.Builder(
    new int[]{5, 5}, new int[]{1, 1})
    .name("cnn1")
    .nIn(numRows*numColumns)
    .nOut(6)
    .build())
```


The second layer is a subsampling layer that will take a 2 x 2 region and store the max result into a 2 x 2 element:

```
.layer(1, new SubsamplingLayer.Builder(  
    SubsamplingLayer.PoolingType.MAX,  
    new int[] {2, 2}, new int[] {2, 2})  
    .name("maxpool1")  
    .build())
```

The next two layers will repeat the the previous two layers:

```
.layer(2, new ConvolutionLayer.Builder(new int[] {5, 5}, new  
    int[] {1, 1})  
    .name("cnn2")  
    .nOut(16)  
    .biasInit(1)  
    .build())  
.layer(3, new SubsamplingLayer.Builder  
    (SubsamplingLayer.PoolingType.MAX, new  
    int[] {2, 2}, new int[] {2, 2})  
    .name("maxpool2")  
    .build())
```


Now we will wire the output of the subsampling layer into a dense feedforward network, consisting of 120 neurons, and then through another layer, into 84 neurons, as follows:

```
.layer(4, new DenseLayer.Builder()  
    .name("ffn1")  
    .nOut(120)  
    .build())  
.layer(5, new DenseLayer.Builder()  
    .name("ffn2")  
    .nOut(84)  
    .build())
```

The final layer connects 84 neurons with 10 output neurons:

```
.layer(6, new OutputLayer.Builder  
    (LossFunctions.LossFunction.NEGATIVELOGLIKELIHOOD)  
    .name("output")  
    .nOut(outputNum)  
    .activation("softmax") // radial basis function required  
    .build())  
.backprop(true)  
.pretrain(false)  
.cnnInputSize(numRows, numColumns, 1);
```

To train this structure, we can reuse the code that we developed in the previous two examples. Again, the training might take some time. The network accuracy should be around 98%.

 As model training significantly relies on linear algebra, training can be significantly sped up by using **Graphics Processing Unit (GPU)** for an order of magnitude. As GPU backend is at the time of writing undergoing a rewrite, please check the latest documentation at <http://deeplearning4j.org/documentation>

As we saw in different examples, increasingly more complex neural networks allow us to extract relevant features automatically, thus completely avoiding traditional image processing. However, the price we pay for this is an increased processing time and a lot of learning examples to make this approach efficient.

Summary

In this chapter, we discussed how to recognize patterns in images in order to distinguish between different classes by covering fundamental principles of deep learning and discussing how to implement them with the deeplearning4j library. We started by refreshing the basic neural network structure and discussed how to implement them to solve handwritten digit recognition problem.

In the next chapter, we'll look further into patterns; however, instead of patterns in images, we'll tackle patterns with temporal dependencies that can be found in sensor data.

9

Activity Recognition with Mobile Phone Sensors

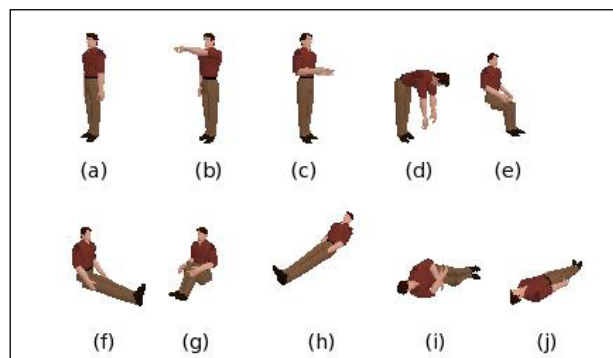
While the previous chapter focused on pattern recognition in images, this chapter is all about recognizing patterns in sensor data, which, in contrast to images, has temporal dependencies. We will discuss how to recognize granular daily activities such as walking, sitting, and running using mobile phone inertial sensors. The chapter also provides references to related research and emphasizes best practices in the activity recognition community.

The topics covered in this chapter will include the following:

- Introducing activity recognition, covering mobile phone sensors and activity recognition pipeline
- Collecting sensor data from mobile devices
- Discussing activity classification and model evaluation
- Deploying activity recognition model

Introducing activity recognition

Activity recognition is an underpinning step in behavior analysis, addressing healthy lifestyle, fitness tracking, remote assistance, security applications, elderly care, and so on. Activity recognition transforms low-level sensor data from sensors, such as accelerometer, gyroscope, pressure sensor, and GPS location, to a higher-level description of behavior primitives. In most cases, these are basic activities, for example, walking, sitting, lying, jumping, and so on, as shown in the following image, or they could be more complex behaviors, such as going to work, preparing breakfast, shopping, and so on:



In this chapter, we will discuss how to add the activity recognition functionality into a mobile application. We will first look at what does an activity recognition problem looks like, what kind of data do we need to collect, what are the main challenges are, and how to address them?

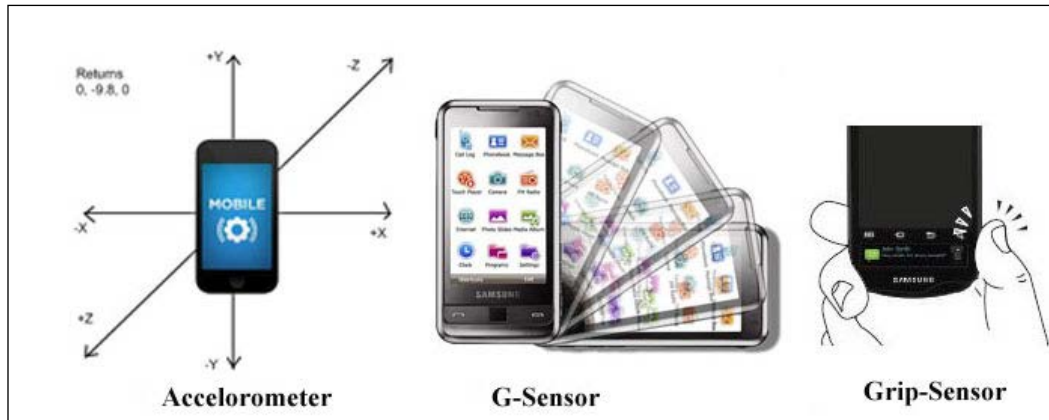
Later, we will follow an example to see how to actually implement activity recognition in an Android application, including data collection, data transformation, and building a classifier.

Let's start!

Mobile phone sensors

Let's first review what kinds of mobile phone sensors there are and what they report. Most smart devices are now equipped with a several built-in sensors that measure the motion, position, orientation, and conditions of the ambient environment. As sensors provide measurements with high precision, frequency, and accuracy, it is possible to reconstruct complex user motions, gestures, and movements. Sensors are often incorporated in various applications; for example, gyroscope readings are used to steer an object in a game, GPS data is used to locate the user, and accelerometer data is used to infer the activity that the user is performing, for example, cycling, running, or walking.

The next image shows a couple of examples what kind of interactions the sensors are able to detect:



Mobile phone sensors can be classified into the following three broad categories:

- **Motion sensors** measure acceleration and rotational forces along the three perpendicular axes. Examples of sensors in this category include accelerometers, gravity sensors, and gyroscopes.
- **Environmental sensors** measure a variety of environmental parameters, such as illumination, air temperature, pressure, and humidity. This category includes barometers, photometers, and thermometers.
- **Position sensors** measure the physical position of a device. This category includes orientation sensors and magnetometers.

More detailed descriptions for different mobile platforms are available at the following links:



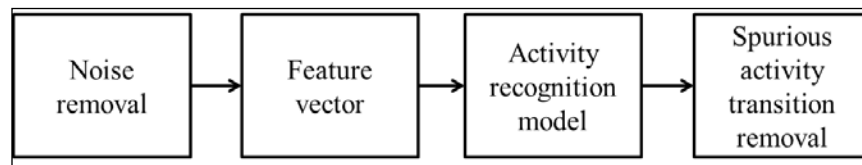
- Android sensors framework: http://developer.android.com/guide/topics/sensors/sensors_overview.html
- iOS Core Motion framework: https://developer.apple.com/library/ios/documentation/CoreMotion/Reference/CoreMotion_Reference/
- Windows Phone: [https://msdn.microsoft.com/en-us/library/windows/apps/hh202968\(v=vs.105\).aspx](https://msdn.microsoft.com/en-us/library/windows/apps/hh202968(v=vs.105).aspx)

In this chapter, we will work only with Android's sensors framework.

Activity recognition pipeline

Classifying multidimensional time-series sensor data is inherently more complex compared to classifying traditional, nominal data as we saw in the previous chapters. First, each observation is temporally connected to the previous and following observations, making it very difficult to apply a straightforward classification of a single set of observations only. Second, the data obtained by sensors at different time points stochastic, that is unpredictable due to influence of sensor noise, environmental disturbances, and many other reasons. Moreover, an activity can comprise various sub-activities executed in different manner and each person performs the activity a bit differently, which results in high intraclass differences. Finally, all these reasons make an activity recognition model imprecise, resulting in new data being often misclassified. One of the highly desirable properties of an activity recognition classifier is to ensure continuity and consistency in the recognized activity sequence.

To deal with these challenges, activity recognition is applied to a pipeline as shown in the following:



In the first step, we attenuate as much noise as we can, for example, by reducing sensor sampling rate, removing outliers, applying high/low-pass filters, and so on. In the next phase, we construct a feature vector, for instance, we convert sensor data from time domain to frequency domain by applying **Discrete Fourier Transform (DFT)**. DFT is a method that takes a list of samples as an input and returns a list of sinusoid coefficients ordered by their frequencies. They represent a combination of frequencies that are present in the original list of samples.



An gentle introduction of Fourier transform is written by Pete Bevelacqua at <http://www.thefouriertransform.com/>.
If you want to get more technical and theoretical background on the Fourier transform, take a look at the eighth and ninth lectures in the class by Robert Gallanger and Lizhong Zheng at MIT open course:
<http://theopenacademy.com/content/principles-digital-communication>

Next, based on the feature vector and set of training data, we can build an activity recognition model that assigns an atomic action to each observation. Therefore, for each new sensor reading, the model will output the most probable activity label. However, models make mistakes. Hence, the last phase smooths the transitions between activities by removing transitions that cannot occur in reality, for example, it is not physically feasible that the transition between activities lying-standing-lying occurs in less than half a second, hence such transition between activities is smoothed as lying-lying-lying.

The activity recognition model is constructed with a supervised learning approach, which consists of training and classification steps. In the training step, a set of labeled data is provided to train the model. The second step is used to assign a label to the new unseen data by the trained model. The data in both phases must be pre-processed with the same set of tools, such as filtering and feature-vector computation.

The post-processing phase, that is, spurious activity removal, can also be a model itself and, hence, also requires a learning step. In this case, the pre-processing step also includes activity recognition, which makes such arrangement of classifiers a meta-learning problem. To avoid overfitting, it is important that the dataset used for training the post-processing phase is not the same as that used for training the activity recognition model.

We will roughly follow a lecture on smartphone programming by professor Andrew T. Campbell from Dartmouth University and leverage data collection mobile app that they developed in the class (Campbell, 2011).

The plan

The plan consists of training phase and deployment phase. Training phase shown in the following image boils down to the following steps:

1. Install Android Studio and import `MyRunsDataCollector.zip`.
2. Load the application in your Android phone.
3. Collect your data, for example, standing, walking, and running, and transform the data to a feature vector comprising of FFT transforms. Don't panic, low-level signal processing functions such as FFT will not be written from scratch as we will use existing code to do that. The data will be saved on your phone in a file called `features.arff`.
4. Create and evaluate an activity recognition classifier using exported data and implement filter for spurious activity transitions removal.
5. Plug the classifier back into the mobile application.

If you don't have an Android phone, or if you want to skip all the steps related to mobile application, just grab an already-collected dataset located in `data/features.arff` and jump directly to the *Building a classifier* section.

Collecting data from a mobile phone

This section describes the first three steps from the plan. If you want to directly work with the data, you can just skip this section and continue to the following *Building a classifier* section. There are many open source mobile apps for sensor data collection, including an app by Prof. Campbell that we will use in this chapter. The application implements the essentials to collect sensor data for different activity classes, for example, standing, walking, running, and others.

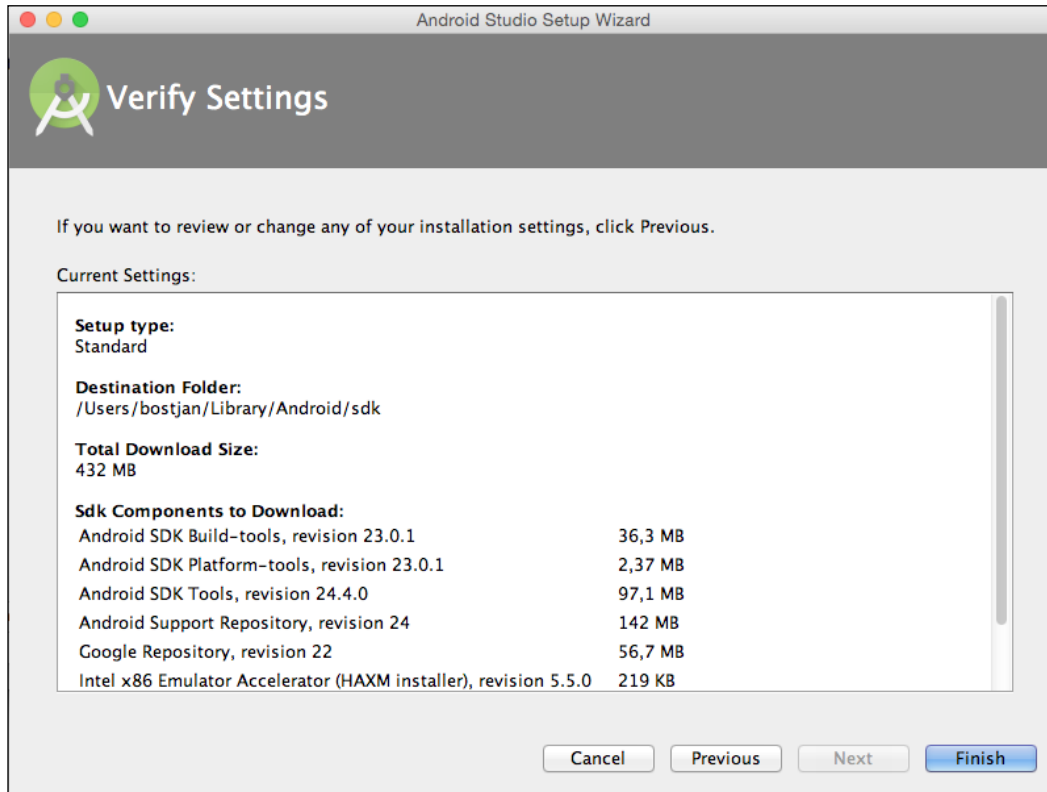
Let's start by preparing the Android development environment. If you have already installed it, jump to the *Loading the data collector* section.

Installing Android Studio

Android Studio is a development environment for Android platform. We will quickly review installation steps and basic configurations required to start the app on a mobile phone. For more detailed introduction to Android development, I would recommend an introductory book, *Android 5 Programming by Example* by Kyle Mew.

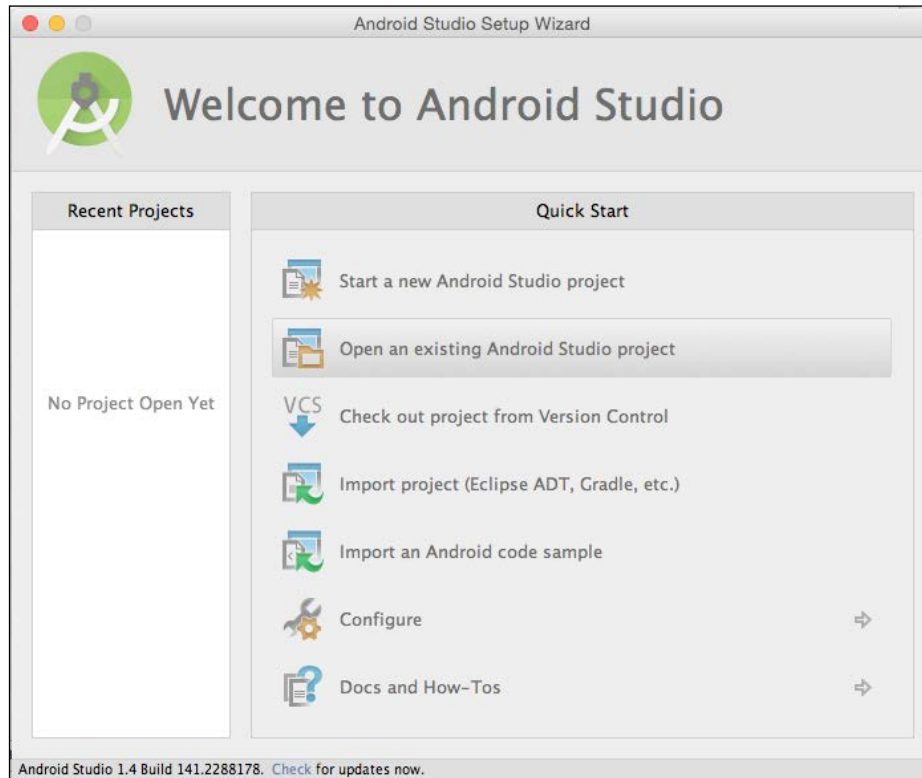
Grab the latest Android Studio for developers at <http://developer.android.com/sdk/installing/index.html?pkg=studio> and follow the installation instructions. The installation will take over 10 minutes, occupying approximately 0.5 GB of space:

Then:

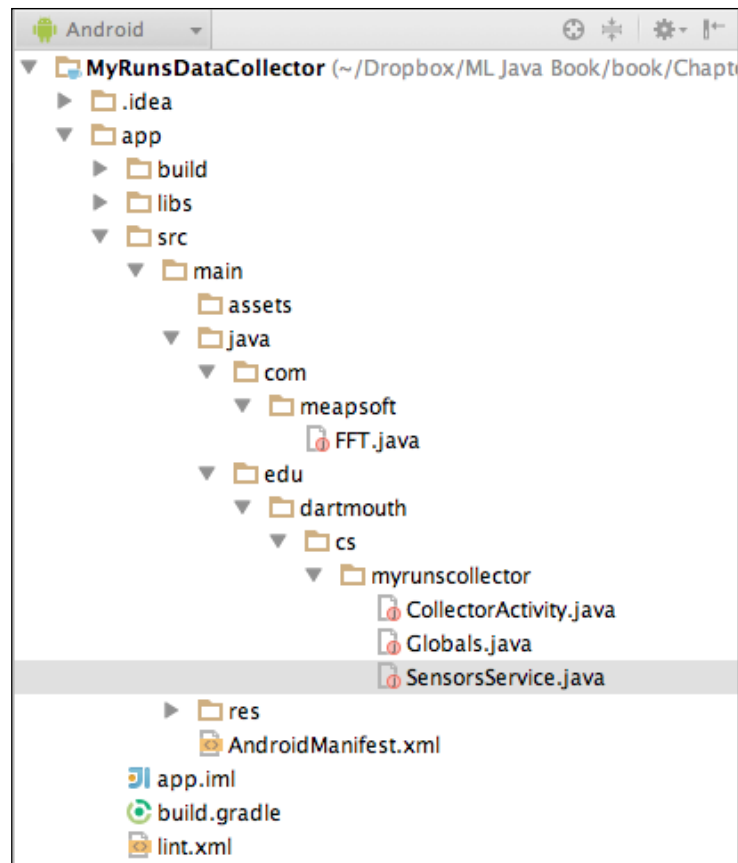


Loading the data collector

First, grab source code of MyRunsDataCollector from <http://www.cs.dartmouth.edu/~campbell/cs65/code/myrunsdatacollector.zip>. Once the Android Studio is installed, choose to open an existing Android Studio project as shown in the following image and select the MyRunsDataCollector folder. This will import the project to Android Studio:



After the project import is completed, you should be able to see the project files structure, as shown in the following image. As shown in the following, the collector consists of `CollectorActivity.java`, `Globals.java`, and `SensorsService.java`. The project also shows `FFT.java` implementing low-level signal processing:



The main `myrunscollector` package contains the following classes:

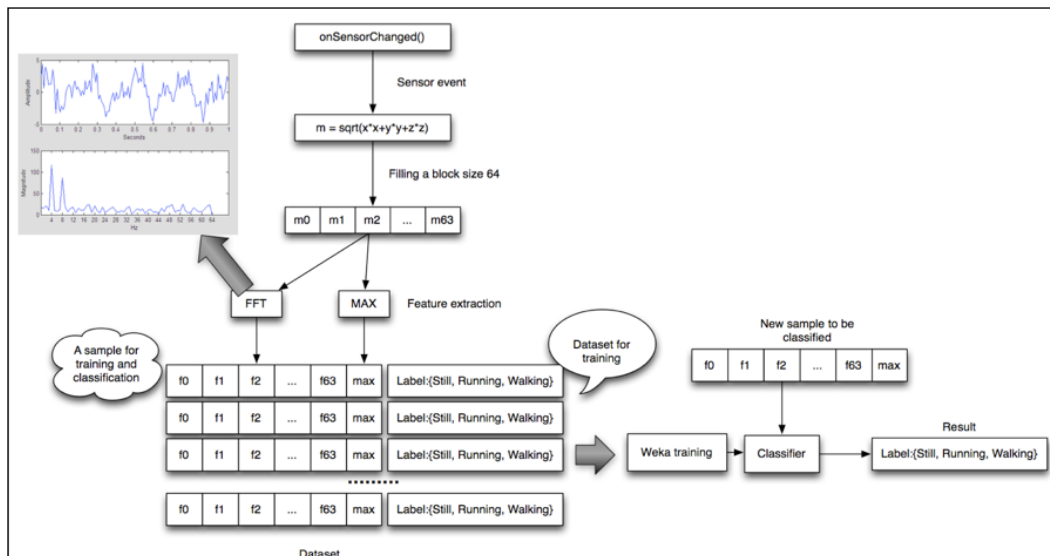
- `Globals.java`: This defines global constants such as activity labels and IDs, data filenames, and so on
- `CollectorActivity.java`: This implements user interface actions, that is, what happens when specific button is pressed
- `SensorsService.java`: This implements a service that collects data, calculates the feature vector as we will discuss in the following sections, and stores the data into a file on the phone

The next question that we will address is how to design features.

Feature extraction

Finding an appropriate representation of the person's activities is probably the most challenging part of activity recognition. The behavior needs to be represented with simple and general features so that the model using these features will also be general and work well on behaviors different from those in the learning set.

In fact, it is not difficult to design features specific to the captured observations in a training set; such features would work well on them. However, as the training set captures only a part of the whole range of human behavior, overly specific features would likely fail on general behavior:



Let's see how it is implemented in `MyRunsDataCollector`. When the application is started, a method called `onSensorChanged()` gets a triple of accelerometer sensor readings (x , y , and z) with a specific time stamp and calculates the magnitude from the sensor readings. The methods buffers up to 64 consecutive magnitudes marked before computing the FFT coefficients (Campbell, 2015):

"As shown in the upper left of the diagram, FFT transforms a time series of amplitude over time to magnitude (some representation of amplitude) across frequency; the example shows some oscillating system where the dominant frequency is between 4-8 cycles/second called Hertz (H) – imagine a ball attached to an elastic band that this stretched and oscillates for a short period of time, or your gait while walking, running -- one could look at these systems in the time and frequency domains. The x,y,z accelerometer readings and the magnitude are time domain variables. We transform these time domain data into the frequency domain because they can represent the distribution in a nice compact manner that the classifier will use to build a decision tree model. For example, the rate of the amplitude transposed to the frequency domain may look something like the figure bottom plot -- the top plot is time domain and the bottom plot a transformation of the time to the frequency domain.

The training phase also stores the maximum (MAX) magnitude of the ($m0..m63$) and the user supplied label (e.g., walking) using the collector. The individual features are computed as magnitudes ($f0..f63$), the MAX magnitude and the class label."

Now let's move on to the actual data collection.

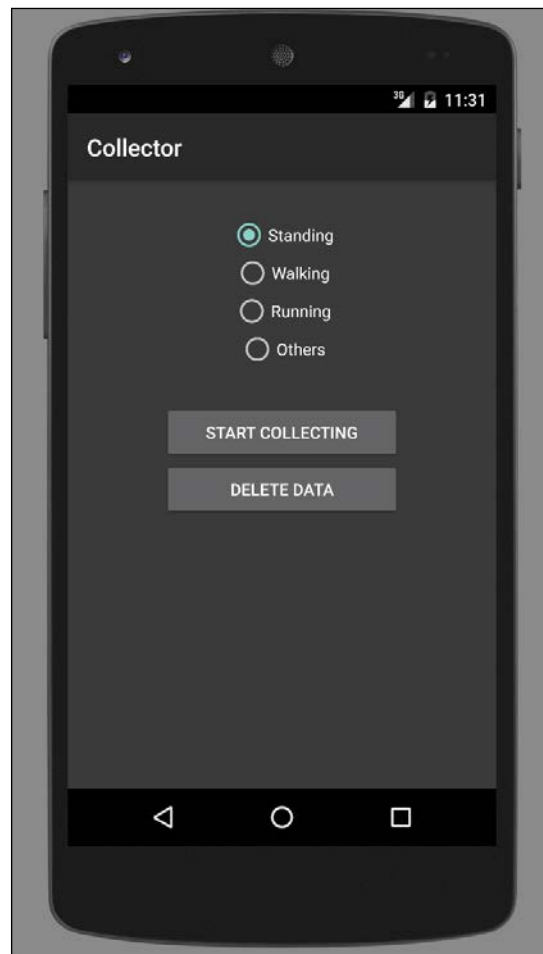
Collecting training data

We can now use the collector to collect training data for activity recognition. The collector supports three activities by default: standing, walking, and running, as shown in the application screenshot in the following figure.

You can select an activity, that is, target class value, and start recording the data by clicking the **START COLLECTING** button. Make sure that each activity is recorded for at least three minutes, for example, if the **Walking** activity is selected, press **START COLLECTING** and walk around for at least three minutes. At the end of the activity, press stop collecting. Repeat this for each of the activities.

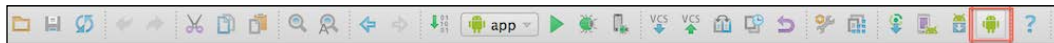
You could also collect different scenarios involving these activities, for example, walking in the kitchen, walking outside, walking in a line, and so on. By doing so, you will have more data for each activity class and a better classifier. Makes sense, right? The more data, the less confused the classifier will be. If you only have a little data, overfitting will occur and the classifier will confuse classes – standing with walking, walking with running. However, the more data, the less they get confused. You might collect less than three minutes per class when you are debugging, but for your final polished product, the more data, the better it is. Multiple recording instances will simply be accumulated in the same file.

Note, the delete button removes the data that is stored in a file on the phone. If you want to start over again, hit delete before starting otherwise, the new collected data will be appended at the end of the file:

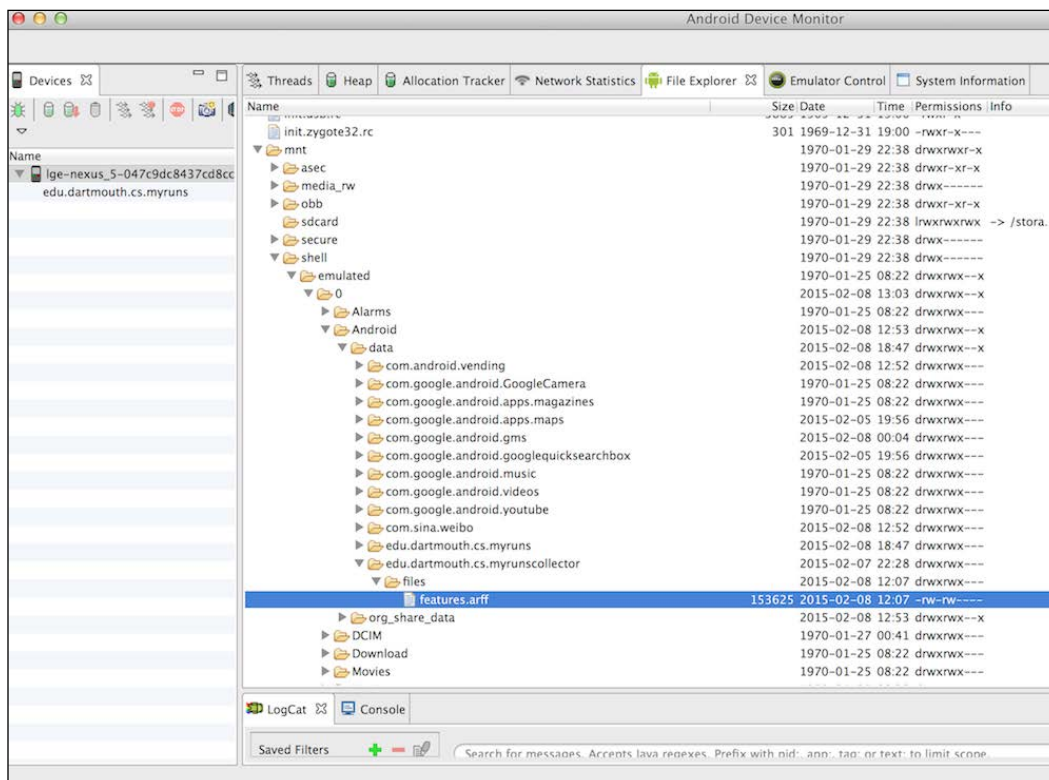


The collector implements the diagram as discussed in the previous sections: it collects accelerometer samples, computes the magnitudes, uses the `FFT.java` class to compute the coefficients, and produces the feature vectors. The data is then stored in a Weka formatted `features.arff` file. The number of feature vectors will vary as you will collect a small or large amount of data. The longer you collect the data, the more feature vectors are accumulated.

Once you stop collecting the training data using the collector tool, we need to grab the data to carry on the workflow. We can use the file explorer in **Android Device Monitor** to upload the `features.arff` file from the phone and to store it on the computer. You can access your Android Device Monitor by clicking on the Android robot icon as shown in the following image:



By selecting your device on the left, your phone storage content will be shown on the right-hand side. Navigate through `mnt/shell/emulated/Android/data/edu.dartmouth.cs.myrunscollector/files/features.arff`:



To upload this file to your computer, you need to select the file (it is highlighted) and click **Upload**.

Now we are ready to build a classifier.

Building a classifier

Once sensor samples are represented as feature vectors having the class assigned, it is possible to apply standard techniques for supervised classification, including feature selection, feature discretization, model learning, k-fold cross validation, and so on. The chapter will not delve into the details of the machine learning algorithms. Any algorithm that supports numerical features can be applied, including SVMs, random forest, AdaBoost, decision trees, neural networks, multi-layer perceptrons, and others.

Therefore, let's start with a basic one, decision trees: load the dataset, build set class attribute, build a decision tree model, and output the model:

```
String databasePath = "/Users/bostjan/Dropbox/ML Java Book/book/
datasets/chap9/features.arff";

// Load the data in arff format
Instances data = new Instances(new BufferedReader(new
    FileReader(databasePath)));

// Set class the last attribute as class
data.setClassIndex(data.numAttributes() - 1);

// Build a basic decision tree model
String[] options = new String[]{};
J48 model = new J48();
model.setOptions(options);
model.buildClassifier(data);

// Output decision tree
System.out.println("Decision tree model:\n"+model);
```

The algorithm first outputs the model, as follows:

Decision tree model:

J48 pruned tree

max <= 10.353474

```

|   fft_coef_0000 <= 38.193106: standing (46.0)
|   fft_coef_0000 > 38.193106
|   |   fft_coef_0012 <= 1.817792: walking (77.0/1.0)
|   |   fft_coef_0012 > 1.817792
|   |   |   max <= 4.573082: running (4.0/1.0)
|   |   |   max > 4.573082: walking (24.0/2.0)
max > 10.353474: running (93.0)

```

Number of Leaves : 5

Size of the tree : 9

The tree is quite simplistic and seemingly accurate as majority class distributions in the terminal nodes are quite high. Let's run a basic classifier evaluation to validate the results:

```

// Check accuracy of model using 10-fold cross-validation
Evaluation eval = new Evaluation(data);
eval.crossValidateModel(model, data, 10, new Random(1), new
String[] {});
System.out.println("Model performance:\n"+
eval.toSummaryString());

```

This outputs the following model performance:

Correctly Classified Instances	226	92.623 %
Incorrectly Classified Instances	18	7.377 %
Kappa statistic	0.8839	
Mean absolute error	0.0421	
Root mean squared error	0.1897	
Relative absolute error	13.1828 %	
Root relative squared error	47.519 %	
Coverage of cases (0.95 level)	93.0328 %	
Mean rel. region size (0.95 level)	27.8689 %	
Total Number of Instances	244	

The classification accuracy scores very high, 92.62%, which is an amazing result. One important reason why the result is so good lies in our evaluation design. What I mean here is the following: sequential instances are very similar to each other, if we split them randomly during a 10-fold cross validation, there is a high chance that we use almost identical instances for both training and testing; hence, straightforward k-fold cross validation produces an optimistic estimate of model performance.

A better approach is to use folds that correspond to different sets of measurements or even different people. For example, we can use the application to collect learning data of five people. Then, it makes sense to run k-person cross validation, where the model is trained on four people and tested on the fifth person. The procedure is repeated for each person and the results are averaged. This will give us a much more realistic estimate of the model performance.

Leaving evaluation comment aside, let's look at how to deal with classifier errors.

Reducing spurious transitions

At the end of the activity recognition pipeline, we want to make sure that the classifications are not too volatile, that is, we don't want activities to change every millisecond. A basic approach is to design a filter that ignores quick changes in the activity sequence.

We build a filter that remembers the last window activities and returns the most frequent one. If there are multiple activities with the same score, it returns the most recent one.

First, we create a new `SpuriousActivityRemoval` class that will hold a list of activities and the window parameter:

```
class SpuriousActivityRemoval{

    List<Object> last;
    int window;

    public SpuriousActivityRemoval(int window){
        this.last = new ArrayList<Object>();
        this.window = window;
    }
}
```

Next, we create the `Object filter(Object)` method that will take an activity and return a filtered activity. The method first checks whether we have enough observations. If not, it simply stores the observation and returns the same value, as shown in the following code:

```
public Object filter(Object obj){
    if(last.size() < window){
        last.add(obj);
        return obj;
    }
}
```

If we already collected window observations, we simply return the most frequent observation, remove the oldest observation, and insert the new observation:

```

    Object o = getMostFrequentElement(last);
    last.add(obj);
    last.remove(0);
    return o;
}

```

What is missing here is a function that returns the most frequent element from a list of objects. We implement this with a hash map, as follows:

```

private Object getMostFrequentElement(List<Object> list){

    HashMap<String, Integer> objectCounts = new HashMap<String,
        Integer>();
    Integer frequentCount = 0;
    Object frequentObject = null;

```

Now, we iterate over all the elements in the list, insert each unique element into a hash map, or update its counter if it is already in the hash map. At the end of the loop, we store the most frequent element that we found so far, as follows:

```

    for(Object obj : list){
        String key = obj.toString();
        Integer count = objectCounts.get(key);
        if(count == null){
            count = 0;
        }
        objectCounts.put(key, ++count);

        if(count >= frequentCount){
            frequentCount = count;
            frequentObject = obj;
        }
    }

    return frequentObject;
}
}

```

Let's run a simple example:

```

String[] activities = new String[]{"Walk", "Walk", "Walk", "Run",
    "Walk", "Run", "Run", "Sit", "Sit", "Sit"};
SpuriousActivityRemoval dlpFilter = new
    SpuriousActivityRemoval(3);

```

```
for(String str : activities){
    System.out.println(str + " -> " + dlpFilter.filter(str));
}
```

The example outputs the following activities:

```
Walk -> Walk
Walk -> Walk
Walk -> Walk
Run -> Walk
Walk -> Walk
Run -> Walk
Run -> Run
Sit -> Run
Sit -> Run
Sit -> Sit
```

The result is a continuous sequence of activities, that is, we do not have quick changes. This adds some delay, but unless this is absolutely critical for the application, it is acceptable.

Activity recognition may be enhanced by appending n previous activities as recognized by the classifier to the feature vector. The danger of appending previous activities is that the machine learning algorithm may learn that the current activity is always the same as the previous one, as this will often be the case. The problem may be solved by having two classifiers, A and B: the classifier B's attribute vector contains n previous activities as recognized by the classifier A. The classifier A's attribute vector does not contain any previous activities. This way, even if B gives a lot of weight to the previous activities, the previous activities as recognized by A will change as A is not burdened with B's inertia.

All that remains to do is to embed the classifier and filter into our mobile application.

Plugging the classifier into a mobile app

There are two ways to incorporate a classifier into a mobile application. The first one involves exporting a model in the Weka format, using the Weka library as a dependency in our mobile application, loading the model, and so on. The procedure is identical to the example we saw in *Chapter 3, Basic Algorithms – Classification, Regression, and Clustering*. The second approach is more lightweight; we export the model as a source code, for example, we create a class implementing the decision tree classifier. Then we can simply copy and paste the source code into our mobile app, without even importing any Weka dependencies.

Fortunately, some Weka models can be easily exported to source code by the `toSource(String)` function:

```
// Output source code implementing the decision tree
System.out.println("Source code:\n" +
    model.toSource("ActivityRecognitionEngine"));
```

This outputs an `ActivityRecognitionEngine` class that corresponds to our model. Now, let's take a closer look at the outputted code:

```
class ActivityRecognitionEngine {

    public static double classify(Object[] i)
        throws Exception {

        double p = Double.NaN;
        p = ActivityRecognitionEngine.N17a7cec20(i);
        return p;
    }
    static double N17a7cec20(Object []i) {
        double p = Double.NaN;
        if (i[64] == null) {
            p = 1;
        } else if (((Double) i[64]).doubleValue() <= 10.353474) {
            p = ActivityRecognitionEngine.N65b3120a1(i);
        } else if (((Double) i[64]).doubleValue() > 10.353474) {
            p = 2;
        }
        return p;
    }
    ...
}
```

The outputted `ActivityRecognitionEngine` class implements the decision tree that we discussed earlier. The machine-generated function names, such as `N17a7cec20(Object [])`, correspond to decision tree nodes. The classifier can be called by the `classify(Object [])` method, where we should pass a feature vector obtained by the same procedure as we discussed in the previous sections. As usual, it returns a double, indicating a class label index.

Summary

In this chapter, we discussed how to implement an activity recognition model for mobile applications. We looked into the completed process, including data collection, feature extraction, model building, evaluation, and model deployment.

In the next chapter, we will move on to another Java library targeted at text analysis – Mallet.

10

Text Mining with Mallet – Topic Modeling and Spam Detection

In this chapter, we will first discuss what text mining is, what kind of analysis is it able to offer, and why you might want to use it in your application. We will then discuss how to work with Mallet, a Java library for natural language processing, covering data import and text pre-processing. Afterwards, we will look into two text mining applications: topic modeling, where we will discuss how text mining can be used to identify topics found in the text documents without reading them individually; and spam detection, where we will discuss how to automatically classify text documents into categories.

This chapter will cover the following topics:

- Introducing text mining
- Installing and working with Mallet
- Topic modeling
- Spam detection

Introducing text mining

Text mining, or text analytics, refers to the process of automatically extracting high-quality information from text documents, most often written in natural language, where high-quality information is considered to be relevant, novel, and interesting.

While a typical text-analytics application is to scan a set of documents to generate a search index, text mining can be used in many other applications, including text categorization into specific domains; text clustering to automatically organize a set of documents; sentiment analysis to identify and extract subjective information in documents; concept/entity extraction that is capable of identifying people, places, organizations, and other entities from documents; document summarization to automatically provide the most important points in the original document; and learning relations between named entities.

The process based on statistical pattern mining usually involves the following steps:

1. Information retrieval and extraction.
2. Transforming unstructured text data into structured; for example, parsing, removing noisy words, lexical analysis, calculating word frequencies, deriving linguistic features, and so on.
3. Discovery of patterns from structured data and tagging/annotation.
4. Evaluation and interpretation of the results.

Later in this chapter, we will look at two application areas: topic modeling and text categorization. Let's examine what they bring to the table.

Topic modeling

Topic modeling is an unsupervised technique and might be useful if you need to analyze a large archive of text documents and wish to understand what the archive contains, without necessarily reading every single document by yourself. A text document can be a blog post, e-mail, tweet, document, book chapter, diary entry, and so on. Topic modeling looks for patterns in a corpus of text; more precisely, it identifies topics as lists of words that appear in a statistically meaningful way. The most well-known algorithm is **Latent Dirichlet Allocation** (Blei et al, 2003), which assumes that author composed a piece of text by selecting words from possible baskets of words, where each basket corresponds to a topic. Using this assumption, it becomes possible to mathematically decompose text into the most likely baskets from where the words first came. The algorithm then iterates over this process until it converges to the most likely distribution of words into baskets, which we call topics.

For example, if we use topic modeling on a series of news articles, the algorithm would return a list of topics and keywords that most likely comprise of these topics. Using the example of news articles, the list might look similar to the following:

- Winner, goal, football, score, first place
- Company, stocks, bank, credit, business
- Election, opponent, president, debate, upcoming

By looking at the keywords, we can recognize that the news articles were concerned with sports, business, upcoming election, and so on. Later in this chapter, we will learn how to implement topic modeling using the news article example.

Text classification

In text classification, or text categorization, the goal is to assign a text document according to its content to one or more classes or categories, which tend to be a more general subject area such as vehicles or pets. Such general classes are referred to as topics, and the classification task is then called text classification, text categorization, topic classification, or topic spotting. While documents can be categorized according to other attributes such as document type, author, printing year, and so on, the focus in this chapter will be on the document content only. Examples of text classification include the following components:

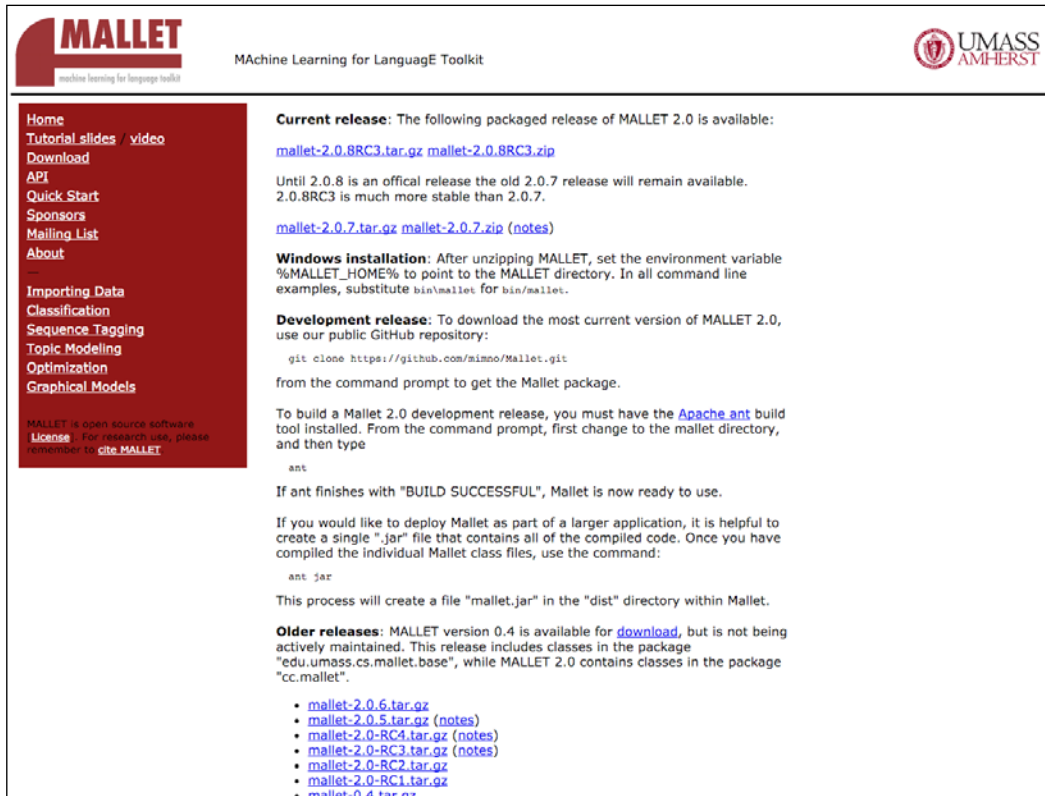
- Spam detection in e-mail messages, user comments, webpages, and so on
- Detection of sexually-explicit content
- Sentiment detection, which automatically classifies a product/service review as positive or negative
- E-mail sorting according to e-mail content
- Topic-specific search, where search engines restrict searches to a particular topic or genre, thus providing more accurate results

These examples show how important text classification is in information retrieval systems, hence most modern information retrieval systems use some kind of text classifier. The classification task that we will use as an example in this book is text classification for detecting e-mail spam.

We continue this chapter with an introduction to Mallet, a Java-based package for statistical natural language processing, document classification, clustering, topic modeling, information extraction, and other machine learning applications to text. We will then cover two text-analytics applications, namely, topics modeling and spam detection as text classification.

Installing Mallet

Mallet is available for download at UMass Amherst University website at <http://mallet.cs.umass.edu/download.php>. Navigate to the **Download** section as shown in the following image and select the latest stable release (2.0.8, at the time of writing this book):



The screenshot shows the Mallet website's download page. At the top left is the Mallet logo with the tagline "Machine Learning for Language Toolkit". At the top right is the UMass Amherst logo. A dark red sidebar on the left contains navigation links: Home, Tutorial slides / video, Download, API, Quick Start, Sponsors, Mailing List, About, Importing Data, Classification, Sequence Tagging, Topic Modeling, Optimization, and Graphical Models. Below the sidebar is a small note: "Mallet is open source software. License: For research use, please remember to cite Mallet". The main content area has a white background and contains the following text:

Current release: The following packaged release of Mallet 2.0 is available:
[mallet-2.0.8RC3.tar.gz](#) [mallet-2.0.8RC3.zip](#)

Until 2.0.8 is an official release the old 2.0.7 release will remain available. 2.0.8RC3 is much more stable than 2.0.7.
[mallet-2.0.7.tar.gz](#) [mallet-2.0.7.zip](#) (notes)

Windows installation: After unzipping Mallet, set the environment variable %MALLET_HOME% to point to the Mallet directory. In all command line examples, substitute bin\mallet for bin/mallet.

Development release: To download the most current version of Mallet 2.0, use our public GitHub repository:
`git clone https://github.com/mimno/Mallet.git`
from the command prompt to get the Mallet package.

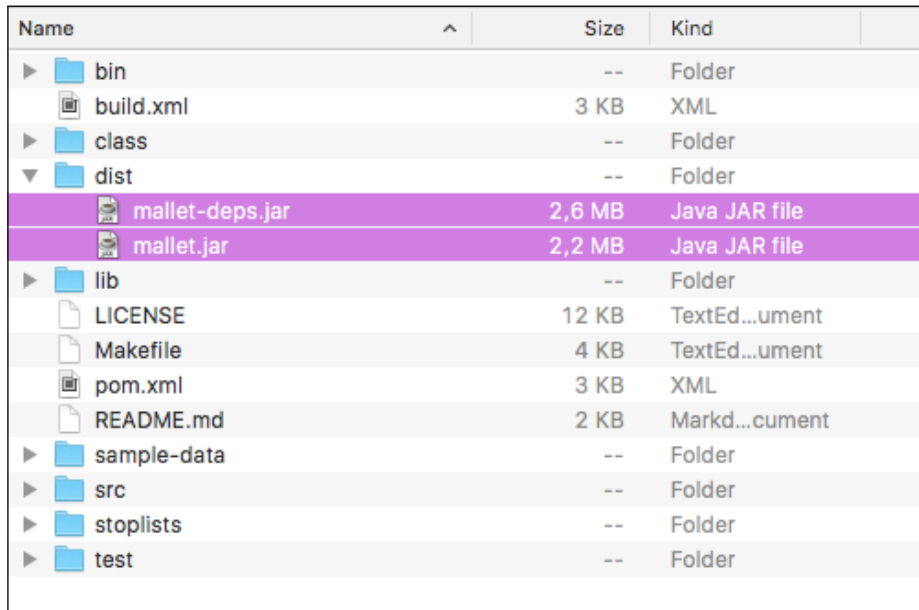
To build a Mallet 2.0 development release, you must have the [Apache ant](#) build tool installed. From the command prompt, first change to the mallet directory, and then type
`ant`
If ant finishes with "BUILD SUCCESSFUL", Mallet is now ready to use.

If you would like to deploy Mallet as part of a larger application, it is helpful to create a single ".jar" file that contains all of the compiled code. Once you have compiled the individual Mallet class files, use the command:
`ant jar`
This process will create a file "mallet.jar" in the "dist" directory within Mallet.

Older releases: Mallet version 0.4 is available for [download](#), but is not being actively maintained. This release includes classes in the package "edu.umass.cs.mallet.base", while Mallet 2.0 contains classes in the package "cc.mallet".

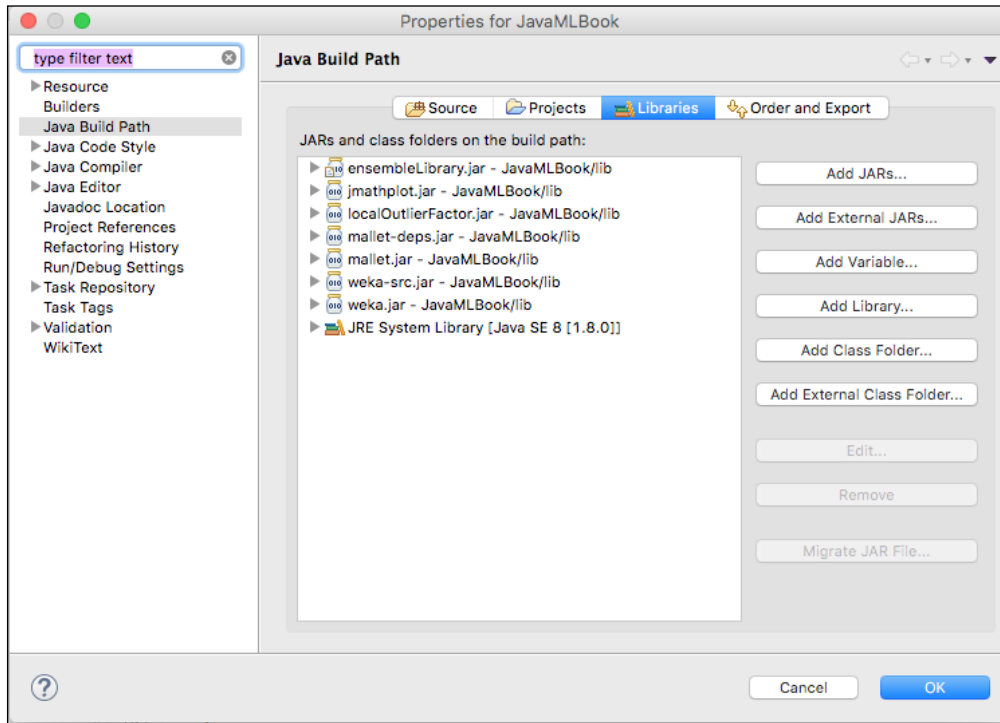
- [mallet-2.0.6.tar.gz](#)
- [mallet-2.0.5.tar.gz](#) (notes)
- [mallet-2.0-RC4.tar.gz](#) (notes)
- [mallet-2.0-RC3.tar.gz](#) (notes)
- [mallet-2.0-RC2.tar.gz](#)
- [mallet-2.0-RC1.tar.gz](#)
- [mallet-0.4.tar.gz](#)

Download the ZIP file and extract the content. In the extracted directory, you should find a folder named `dist` with two JAR files: `mallet.jar` and `mallet-deps.jar`. The first one contains all the packaged Mallet classes, while the second one packs all the dependencies. Include both JARs in your project as referenced libraries, as shown in the following image:



Name	Size	Kind
▶ bin	--	Folder
build.xml	3 KB	XML
▶ class	--	Folder
▼ dist	--	Folder
mallet-deps.jar	2,6 MB	Java JAR file
mallet.jar	2,2 MB	Java JAR file
▶ lib	--	Folder
LICENSE	12 KB	TextEd...ument
Makefile	4 KB	TextEd...ument
pom.xml	3 KB	XML
README.md	2 KB	Markd...cument
▶ sample-data	--	Folder
▶ src	--	Folder
▶ stoplists	--	Folder
▶ test	--	Folder

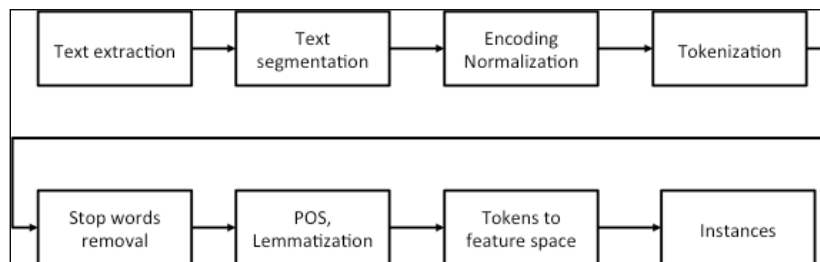
If you are using Eclipse, right click on **Project**, select **Properties**, and pick **Java Build Path**. Select the **Libraries** tab and click **Add External JARs**. Now, select the two JARs and confirm, as shown in the following screenshot:



Now we are ready to start using Mallet.

Working with text data

One of the main challenges in text mining is transforming unstructured written natural language into structured attribute-based instances. The process involves many steps as shown in the following image:



First, we extract some text from the Internet, existing documents, or databases. At the end of the first step, the text could still be presented in the XML format or some other proprietary format. The next step is to, therefore, extract the actual text only and segment it into parts of the document, for example, title, headline, abstract, body, and so on. The third step is involved with normalizing text encoding to ensure the characters are presented the same way, for example, documents encoded in formats such as ASCII, ISO 8859-1, and Windows-1250 are transformed into Unicode encoding. Next, tokenization splits the document into particular words, while the following step removes frequent words that usually have low predictive power, for example, the, a, I, we, and so on.

The **part-of-speech (POS)** tagging and lemmatization step could be included to transform each token (that is, word) to its basic form, which is known as lemma, by removing word endings and modifiers. For example, running becomes run, better becomes good, and so on. A simplified approach is stemming, which operates on a single word without any context of how the particular word is used, and therefore, cannot distinguish between words having different meaning, depending on the part of speech, for example, axes as plural of axe as well as axis.

The last step transforms tokens into a feature space. Most often feature space is a **bag-of-words (BoW)** presentation. In this presentation, a set of all words appearing in the dataset is created, that is, a bag of words. Each document is then presented as a vector that counts how many times a particular word appears in the document.

Consider the following example with two sentences:

- Jacob likes table tennis. Emma likes table tennis too.
- Jacob also likes basketball.

The bag of words in this case consists of {Jacob, likes, table, tennis, Emma, too, also, basketball}, which has eight distinct words. The two sentences could be now presented as vectors using the indexes of the list, indicating how many times a word at a particular index appears in the document, as follows:

- [1, 2, 2, 2, 1, 0, 0, 0]
- [1, 1, 0, 0, 0, 0, 1, 1]

Such vectors finally become instances for further learning.



Another very powerful presentation based on the BoW model is **word2vec**. Word2vec was introduced in 2013 by a team of researchers led by Tomas Mikolov at Google. Word2vec is a neural network that learns distributed representations for words. An interesting property of this presentation is that words appear in clusters, such that some word relationships, such as analogies, can be reproduced using vector math. A famous example shows that king - man + woman returns queen.

Further details and implementation are available at the following link:

<https://code.google.com/archive/p/word2vec/>

Importing data

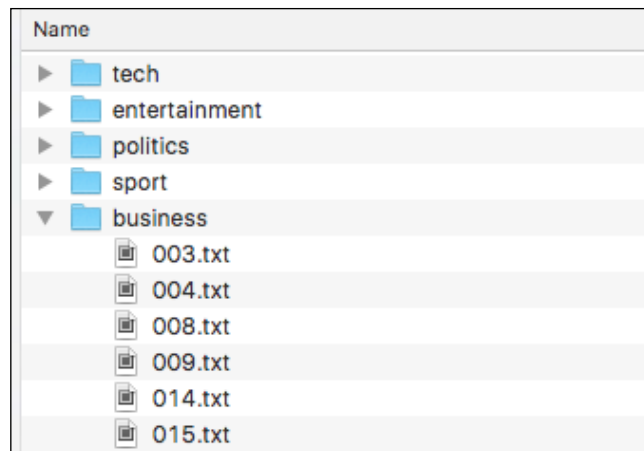
In this chapter, we will not look into how to scrap a set of documents from a website or extract them from database. Instead, we will assume that we already collected them as set of documents and store them in the `.txt` file format. Now let's look at two options how to load them. The first option addresses the situation where each document is stored in its own `.txt` file. The second option addresses the situation where all the documents are stored in a single file, one per line.

Importing from directory

Mallet supports reading from directory with the `cc.mallet.pipe.iterator.FileIterator` class. File iterator is constructed with the following three parameters:

- A list of `File[]` directories with text files
- File filter that specifies which files to select within a directory
- A pattern that is applied to a filename to produce a class label

Consider the data structured into folders as shown in the following image. We have documents organized in five topics by folders (`tech`, `entertainment`, `politics`, and `sport`, `business`). Each folder contains documents on particular topics, as shown in the following image:



In this case, we initialize `iterator` as follows:

```
FileIterator iterator =  
    new FileIterator(new File[] {new File("path-to-my-dataset") },  
        new TxtFilter(),  
        FileIterator.LAST_DIRECTORY);
```

The first parameter specifies the path to our root folder, the second parameter limits the iterator to the `.txt` files only, while the last parameter asks the method to use the last directory name in the path as class label.

Importing from file

Another option to load the documents is through `cc.mallet.pipe.iterator.CsvIterator.CsvIterator(Reader, Pattern, int, int, int)`, which assumes all the documents are in a single file and returns one instance per line extracted by a regular expression. The class is initialized by the following components:

- `Reader`: This is the object that specifies how to read from a file
- `Pattern`: This is a regular expression, extracting three groups: data, target label, and document name
- `int, int, int`: These are the indexes of data, target, and name groups as they appear in a regular expression

Consider a text document in the following format, specifying document name, category and content:

```
AP881218 local-news A 16-year-old student at a private
  Baptist...
AP880224 business The Bechtel Group Inc. offered in 1985 to...
AP881017 local-news A gunman took a 74-year-old woman hostage...
AP900117 entertainment Cupid has a new message for lovers
  this...
AP880405 politics The Reagan administration is weighing w...
```

To parse a line into three groups, we can use the following regular expression:

```
^(\\S*) [\\s,]* (\\S*) [\\s,]* (.*)$
```

There are three groups that appear in parenthesis, (), where the third group contains the data, the second group contains the target class, and the first group contains the document ID. The iterator is initialized as follows:

```
CsvIterator iterator = new CsvIterator (
  fileReader,
  Pattern.compile("^ (\\S*) [\\s,]* (\\S*) [\\s,]* (.*)$"),
  3, 2, 1);
```

Here the regular expression extracts the three groups separated by an empty space and their order is 3, 2, 1.


Now let's move to data pre-processing pipeline.

Pre-processing text data

Once we initialized an iterator that will go through the data, we need to pass the data through a sequence of transformations as described at the beginning of this section. Mallet supports this process through a pipeline and a wide variety of steps that could be included in a pipeline, which are collected in the `cc.mallet.pipe` package. Some examples are as follows:

- `Input2CharSequence`: This is a pipe that can read from various kinds of text sources (either URI, File, or Reader) into `CharSequence`
- `CharSequenceRemoveHTML`: This pipe removes HTML from `CharSequence`
- `MakeAmpersandXMLFriendly`: This converts `&` to `&` in tokens of a token sequence
- `TokenSequenceLowercase`: This converts the text in each token in the token sequence in the data field to lower case

- `TokenSequence2FeatureSequence`: This converts the token sequence in the data field of each instance to a feature sequence
- `TokenSequenceNGrams`: This converts the token sequence in the data field to a token sequence of ngrams, that is, combination of two or more words


 The full list of processing steps is available in the following Mallet documentation:
<http://mallet.cs.umass.edu/api/index.html?cc/mallet/pipe/iterator/package-tree.html>

Now we are ready to build a class that will import our data.

First, let's build a pipeline, where each processing step is denoted as a pipeline in Mallet. Pipelines can be wired together in a serial fashion with a list of `ArrayList<Pipe>` objects:

```
ArrayList<Pipe> pipeList = new ArrayList<Pipe>();
```

Begin by reading data from a file object and converting all the characters into lower case:

```
pipeList.add(new Input2CharSequence("UTF-8"));
pipeList.add( new CharSequenceLowercase() );
```

Next, tokenize raw strings with a regular expression. The following pattern includes Unicode letters and numbers and the underscore character:

```
Pattern tokenPattern =
    Pattern.compile("[\\p{L}\\p{N}_]+");

pipeList.add(new CharSequence2TokenSequence(tokenPattern));
```

Remove stop words, that is, frequent words with no predictive power, using a standard English stop list. Two additional parameters indicate whether stop word removal should be case-sensitive and mark deletions instead of just deleting the words. We'll set both of them to `false`:

```
pipeList.add(new TokenSequenceRemoveStopwords(false, false));
```

Instead of storing the actual words, we can convert them into integers, indicating a word index in the bag of words:

```
pipeList.add(new TokenSequence2FeatureSequence());
```

We'll do the same for the class label; instead of label string, we'll use an integer, indicating a position of the label in our bag of words:

```
pipeList.add(new Target2Label());
```

We could also print the features and the labels by invoking the `PrintInputAndTarget` pipe:

```
pipeList.add(new PrintInputAndTarget());
```

Finally, we store the list of pipelines in a `SerialPipes` class that will covert an instance through a sequence of pipes:

```
SerialPipes pipeline = new SerialPipes(pipeList);
```

Now let's take a look at how apply this in a text mining application!

Topic modeling for BBC news

As discussed earlier, the goal of topic modeling is to identify patterns in a text corpus that correspond to document topics. In this example, we will use a dataset originating from BBC news. This dataset is one of the standard benchmarks in machine learning research, and is available for non-commercial and research purposes.

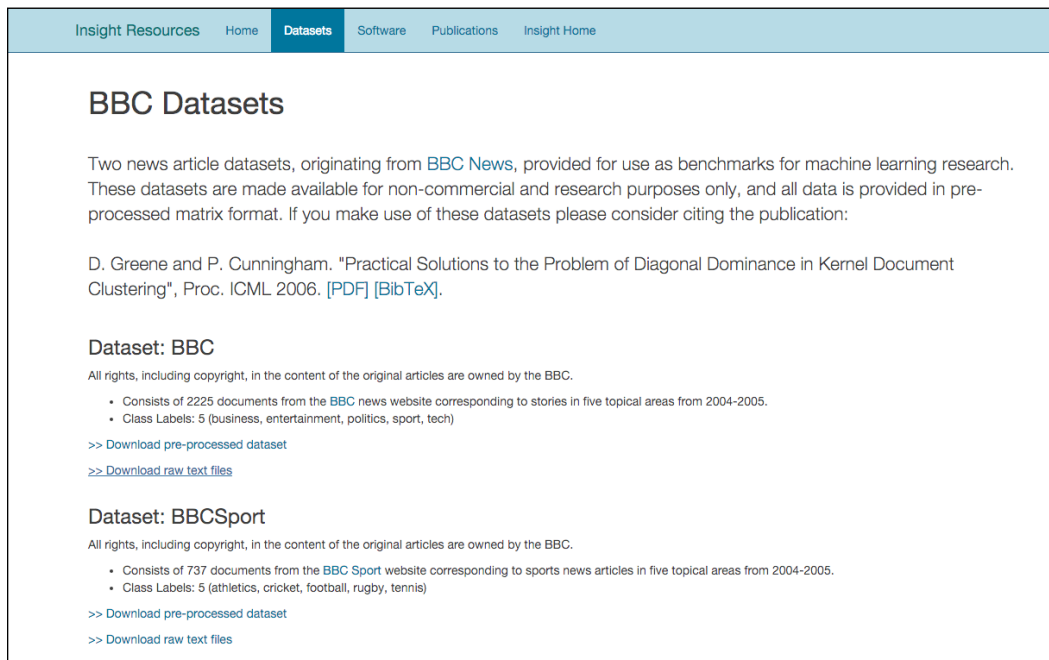
The goal is to build a classifier that is able to assign a topic to an uncategorized document.

BBC dataset

Greene and Cunningham (2006) collected the BBC dataset to study a particular document-clustering challenge using support vector machines. The dataset consists of 2,225 documents from the BBC News website from 2004 to 2005, corresponding to the stories collected from five topical areas: business, entertainment, politics, sport, and tech. The dataset can be grabbed from the following website:

<http://mlg.ucd.ie/datasets/bbc.html>

Download the raw text files under the **Dataset: BBC** section. You will also notice that the website contains already processed dataset, but for this example, we want to process the dataset by ourselves. The ZIP contains five folders, one per topic. The actual documents are placed in the corresponding topic folder, as shown in the following screenshot:



Insight Resources Home **Datasets** Software Publications Insight Home

BBC Datasets

Two news article datasets, originating from [BBC News](#), provided for use as benchmarks for machine learning research. These datasets are made available for non-commercial and research purposes only, and all data is provided in pre-processed matrix format. If you make use of these datasets please consider citing the publication:

D. Greene and P. Cunningham. "Practical Solutions to the Problem of Diagonal Dominance in Kernel Document Clustering", Proc. ICML 2006. [\[PDF\]](#) [\[BibTeX\]](#).

Dataset: BBC

All rights, including copyright, in the content of the original articles are owned by the BBC.

- Consists of 2225 documents from the [BBC news website](#) corresponding to stories in five topical areas from 2004-2005.
- Class Labels: 5 (business, entertainment, politics, sport, tech)

>> [Download pre-processed dataset](#)
>> [Download raw text files](#)

Dataset: BBCSport

All rights, including copyright, in the content of the original articles are owned by the BBC.

- Consists of 737 documents from the [BBC Sport website](#) corresponding to sports news articles in five topical areas from 2004-2005.
- Class Labels: 5 (athletics, cricket, football, rugby, tennis)

>> [Download pre-processed dataset](#)
>> [Download raw text files](#)

Now, let's build a topic classifier.

Modeling

Start by importing the dataset and processing the text:

```
import cc.mallet.types.*;
import cc.mallet.pipe.*;
import cc.mallet.pipe.iterator.*;
import cc.mallet.topics.*;

import java.util.*;
import java.util.regex.*;
import java.io.*;
```

```
public class TopicModeling {  
  
    public static void main(String[] args) throws Exception {  
  
        String dataFolderPath = args[0];  
        String stopListFilePath = args[1];
```

Create a default pipeline as previously described:

```
ArrayList<Pipe> pipeList = new ArrayList<Pipe>();  
pipeList.add(new Input2CharSequence("UTF-8"));  
Pattern tokenPattern = Pattern.compile("[\\p{L}\\p{N}_]+");  
pipeList.add(new CharSequence2TokenSequence(tokenPattern));  
pipeList.add(new TokenSequenceLowercase());  
pipeList.add(new TokenSequenceRemoveStopwords(new  
    File(stopListFilePath), "utf-8", false, false, false));  
pipeList.add(new TokenSequence2FeatureSequence());  
pipeList.add(new Target2Label());  
SerialPipes pipeline = new SerialPipes(pipeList);
```

Next, initialize folderIterator:

```
FileIterator folderIterator = new FileIterator(  
    new File[] {new File(dataFolderPath)},  
    new TxtFilter(),  
    FileIterator.LAST_DIRECTORY);
```

Construct a new instance list with the pipeline that we want to use to process the text:

```
InstanceList instances = new InstanceList(pipeline);
```

Finally, process each instance provided by the iterator:

```
instances.addThruPipe(folderIterator);
```

Now let's create a model with five topics using the `cc.mallet.topics.ParallelTopicModel.ParallelTopicModel` class that implements a simple threaded **Latent Dirichlet Allocation (LDA)** model. LDA is a common method for topic modeling that uses Dirichlet distribution to estimate the probability that a selected topic generates a particular document. We will not dive deep into the details in this chapter; the reader is referred to the original paper by D. Blei et al. (2003). Note that there is another classification algorithm in machine learning with the same acronym that refers to **Linear Discriminant Analysis (LDA)**. Beside the common acronym, it has nothing in common with the LDA model.

The class is instantiated with parameters `alpha` and `beta`, which can be broadly interpreted, as follows:

- High `alpha` value means that each document is likely to contain a mixture of most of the topics, and not any single topic specifically. A low `alpha` value puts less of such constraints on documents, and this means that it is more likely that a document may contain mixture of just a few, or even only one, of the topics.
- A high `beta` value means that each topic is likely to contain a mixture of most of the words, and not any word specifically; while a low value means that a topic may contain a mixture of just a few of the words.

In our case, we initially keep both parameters low (`alpha_t = 0.01`, `beta_w = 0.01`) as we assume topics in our dataset are not mixed much and there are many words for each of the topics:

```
int numTopics = 5;
ParallelTopicModel model =
    new ParallelTopicModel(numTopics, 0.01, 0.01);
```

Next, add instances to the model, and as we are using parallel implementation, specify the number of threads that will run in parallel, as follows:

```
model.addInstances(instances);
model.setNumThreads(4);
```

Run the model for a selected number of iterations. Each iteration is used for better estimation of internal LDA parameters. For testing, we can use a small number of iterations, for example, 50; while in real applications, use 1000 or 2000 iterations. Finally, call the `void estimate()` method that will actually build an LDA model:

```
model.setNumIterations(1000);
model.estimate();
```

The model outputs the following result:

```
0 0,06654  game england year time win world 6
1 0,0863  year 1 company market growth economy firm
2 0,05981  people technology mobile mr games users music
3 0,05744  film year music show awards award won
4 0,11395  mr government people labour election party blair
```

```
[beta: 0,11328]
<1000> LL/token: -8,63377
```

```
Total time: 45 seconds
```

LL/token indicates the model's log-likelihood, divided by the total number of tokens, indicating how likely the data is given the model. Increasing values mean the model is improving.

The output also shows the top words describing each topic. The words correspond to initial topics really well:

- **Topic 0:** game, England, year, time, win, world, 6 → sport
- **Topic 1:** year, 1, company, market, growth, economy, firm → finance
- **Topic 2:** people, technology, mobile, mr, games, users, music → tech
- **Topic 3:** film, year, music, show, awards, award, won → entertainment
- **Topic 4:** mr, government, people, labor, election, party, blair → politics

There are still some words that don't make much sense, for instance, mr, 1, and 6. We could include them in the stop word list. Also, some words appear twice, for example, award and awards. This happened because we didn't apply any stemmer or lemmatization pipe.

In the next section, we'll take a look to check whether the model is of any good.

Evaluating a model

As statistical topic modeling has unsupervised nature, it makes model selection difficult. For some applications, there may be some extrinsic tasks at hand, such as information retrieval or document classification, for which performance can be evaluated. However, in general, we want to estimate the model's ability to generalize topics regardless of the task.

Wallach et al. (2009) introduced an approach that measures the quality of a model by computing the log probability of held-out documents under the model. Likelihood of unseen documents can be used to compare models – higher likelihood implies a better model.

First, let's split the documents into training and testing set (that is, held-out documents), where we use 90% for training and 10% for testing:

```
// Split dataset
InstanceList[] instanceSplit= instances.split(new Randoms(), new
double[] {0.9, 0.1, 0.0});
```

Now, let's rebuild our model using only 90% of our documents:

```
// Use the first 90% for training
model.addInstances(instanceSplit[0]);
```

```

model.setNumThreads(4);
model.setNumIterations(50);
model.estimate();

```

Next, initialize an estimator that implements Wallach's log probability of held-out documents, `MarginalProbEstimator`:

```

// Get estimator
MarginalProbEstimator estimator = model.getProbEstimator();

```



An intuitive description of LDA is summarized by Annalyn Ng in her blog:

<https://annalyzin.wordpress.com/2015/06/21/laymans-explanation-of-topic-modeling-with-lda-2/>

To get deeper insight into the LDA algorithm, its components, and it working, take a look at the original paper LDA by David Blei et al. (2003) at <http://jmlr.csail.mit.edu/papers/v3/blei03a.html> or take a look at the summarized presentation by D. Santhanam of Brown University at http://www.cs.brown.edu/courses/csci2950-p/spring2010/lectures/2010-03-03_santhanam.pdf.

The class implements many estimators that require quite deep theoretical knowledge of how the LDA method works. We'll pick the left-to-right evaluator, which is appropriate for a wide range of applications, including text mining, speech recognition, and others. The left-to-right evaluator is implemented as the `double evaluateLeftToRight` method, accepting the following components:

- `Instances heldOutDocuments`: This test the instances
- `int numParticles`: This algorithm parameter indicates the number of left-to-right tokens, where default value is 10
- `boolean useResampling`: This states whether to resample topics in left-to-right evaluation; resampling is more accurate, but leads to quadratic scaling in the length of documents
- `PrintStream docProbabilityStream`: This is the file or `stdout` in which we write the inferred log probabilities per document

Let's run the estimator, as follows:

```

double loglike = estimator.evaluateLeftToRight(
    instancesSplit[1], 10, false, null);
System.out.println("Total log likelihood: "+loglike);

```


In our particular case, the estimator outputs the following log likelihood, which makes sense when it is compared to other models that are either constructed with different parameters, pipelines, or data – the higher the log likelihood, the better the model is:

```
Total time: 3 seconds
Topic Evaluator: 5 topics, 3 topic bits, 111 topic mask
Total log likelihood: -360849.4240795393
Total log likelihood
```

Now let's take a look at how to make use of this model.

Reusing a model

As we are usually not building models on the fly, it often makes sense to train a model once and use it repeatedly to classify new data.

Note that if you'd like to classify new documents, they need go through the same pipeline as other documents – the pipe needs to be the same for both training and classification. During training, the pipe's data alphabet is updated with each training instance. If you create a new pipe with the same steps, you don't produce the same pipeline as its data alphabet is empty. Therefore, to use the model on new data, save/load the pipe along with the model and use this pipe to add new instances.

Saving a model

Mallet supports a standard method for saving and restoring objects based on serialization. We simply create a new instance of `ObjectOutputStream` class and write the object into a file as follows:

```
String modelPath = "myTopicModel";

//Save model
ObjectOutputStream oos = new ObjectOutputStream(
    new FileOutputStream (new File(modelPath+".model")));
oos.writeObject(model);
oos.close();

//Save pipeline
oos = new ObjectOutputStream(
    new FileOutputStream (new File(modelPath+".pipeline")));
oos.writeObject(pipeline);
oos.close();
```

Restoring a model

Restoring a model saved through serialization is simply an inverse operation using the `ObjectInputStream` class:

```
String modelPath = "myTopicModel";

//Load model
ObjectInputStream ois = new ObjectInputStream(
    new FileInputStream (new File(modelPath+".model")));
ParallelTopicModel model = (ParallelTopicModel) ois.readObject();
ois.close();

// Load pipeline
ois = new ObjectInputStream(
    new FileInputStream (new File(modelPath+".pipeline")));
SerialPipes pipeline = (SerialPipes) ois.readObject();
ois.close();
```

We discussed how to build an LDA model to automatically classify documents into topics. In the next example, we'll look into another text mining problem – text classification.

E-mail spam detection

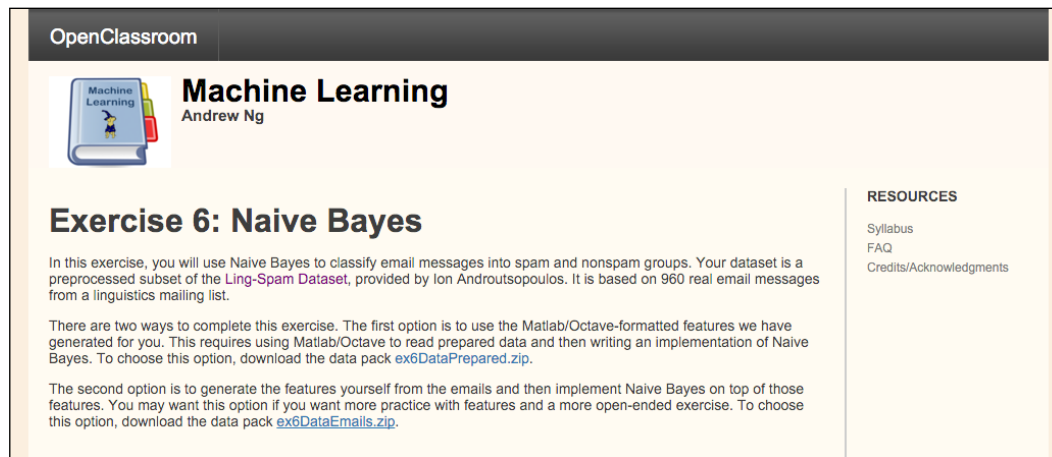
Spam or electronic spam refers to unsolicited messages, typically carrying advertising content, infected attachments, links to phishing or malware sites, and so on. While the most widely recognized form of spam is e-mail spam, spam abuses appear in other media as well: website comments, instant messaging, Internet forums, blogs, online ads, and so on.

In this chapter, we will discuss how to build naive Bayesian spam filtering, using bag-of-words representation to identify spam e-mails. The naive Bayes spam filtering is one of the basic techniques that was implemented in the first commercial spam filters; for instance, **Mozilla Thunderbird** mail client uses native implementation of such filtering. While the example in this chapter will use e-mail spam, the underlying methodology can be applied to other type of text-based spam as well.

E-mail spam dataset

Androutsopoulos et al. (2000) collected one of the first e-mail spam datasets to benchmark spam-filtering algorithms. They studied how the naive Bayes classifier can be used to detect spam, if additional pipes such as stop list, stemmer, and lemmatization contribute to better performance. The dataset was reorganized by Andrew Ng in OpenClassroom's machine learning class, available for download at <http://openclassroom.stanford.edu/MainFolder/DocumentPage.php?course=MachineLearning&doc=exercises/ex6/ex6.html>.

Select and download the second option, `ex6DataEmails.zip`, as shown in the following image:



OpenClassroom

Machine Learning
Andrew Ng

Exercise 6: Naive Bayes

In this exercise, you will use Naive Bayes to classify email messages into spam and nonspam groups. Your dataset is a preprocessed subset of the [Ling-Spam Dataset](#), provided by Ion Androutsopoulos. It is based on 960 real email messages from a linguistics mailing list.

There are two ways to complete this exercise. The first option is to use the Matlab/Octave-formatted features we have generated for you. This requires using Matlab/Octave to read prepared data and then writing an implementation of Naive Bayes. To choose this option, download the data pack [ex6DataPrepared.zip](#).

The second option is to generate the features yourself from the emails and then implement Naive Bayes on top of those features. You may want this option if you want more practice with features and a more open-ended exercise. To choose this option, download the data pack [ex6DataEmails.zip](#).

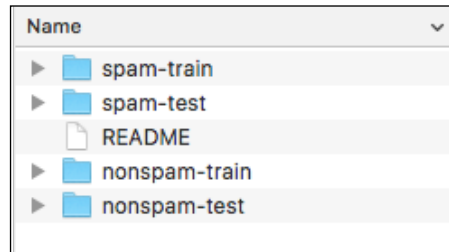
RESOURCES

- Syllabus
- FAQ
- Credits/Acknowledgments

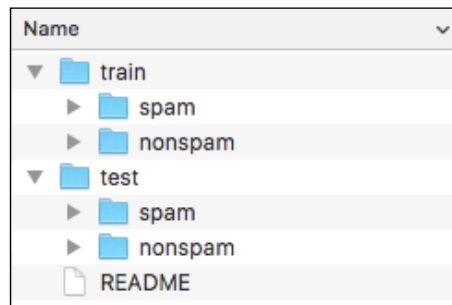
The ZIP contains four folders (Ng, 2015):

- The `nonspam-train` and `spam-train` folders contain the pre-processed e-mails that you will use for training. They have 350 e-mails each.
- The `nonspam-test` and `spam-test` folders constitute the test set, containing 130 spam and 130 nonspam e-mails. These are the documents that you will make predictions on. Notice that even though separate folders tell you the correct labeling, you should make your predictions on all the test documents without this knowledge. After you make your predictions, you can use the correct labeling to check whether your classifications were correct.

To leverage Mallet's folder iterator, let's reorganize the folder structure as follows. Create two folders, `train` and `test`, and put the `spam/nospam` folders under the corresponding folders. The initial folder structure is as shown in the following image:



The final folder structure will be as shown in the following image:



The next step is to transform e-mail messages to feature vectors.

Feature generation

Create a default pipeline as described previously:

```
ArrayList<Pipe> pipeList = new ArrayList<Pipe>();
pipeList.add(new Input2CharSequence("UTF-8"));
Pattern tokenPattern = Pattern.compile("[\\p{L}\\p{N}_]+");
pipeList.add(new CharSequence2TokenSequence(tokenPattern));
pipeList.add(new TokenSequenceLowercase());
pipeList.add(new TokenSequenceRemoveStopwords(new
    File(stopListFilePath), "utf-8", false, false, false));
pipeList.add(new TokenSequence2FeatureSequence());
pipeList.add(new FeatureSequence2FeatureVector());
pipeList.add(new Target2Label());
SerialPipes pipeline = new SerialPipes(pipeList);
```

Note that we added an additional `FeatureSequence2FeatureVector` pipe that transforms a feature sequence into a feature vector. When we have data in a feature vector, we can use any classification algorithm as we saw in the previous chapters. We'll continue our example in Mallet to demonstrate how to build a classification model.

Next, initialize a folder iterator to load our examples in the `train` folder comprising e-mail examples in the `spam` and `nonspam` subfolders, which will be used as example labels:

```
FileIterator folderIterator = new FileIterator(  
    new File[] {new File(dataFolderPath)},  
    new TxtFilter(),  
    FileIterator.LAST_DIRECTORY);
```

Construct a new instance list with the pipeline that we want to use to process the text:

```
InstanceList instances = new InstanceList(pipeline);
```

Finally, process each instance provided by the iterator:

```
instances.addThruPipe(folderIterator);
```

We have now loaded the data and transformed it into feature vectors. Let's train our model on the training set and predict the `spam/nonspam` classification on the test set.

Training and testing

Mallet implements a set of classifiers in the `cc.mallet.classify` package, including decision trees, naive Bayes, AdaBoost, bagging, boosting, and many others. We'll start with a basic classifier, that is, a naive Bayes classifier. A classifier is initialized by the `ClassifierTrainer` class, which returns a classifier when we invoke its `train(Instances)` method:

```
ClassifierTrainer classifierTrainer = new NaiveBayesTrainer();  
Classifier classifier = classifierTrainer.train(instances);
```

Now let's see how this classier works and evaluate its performance on a separate dataset.

Model performance

To evaluate the classifier on a separate dataset, let's start by importing the e-mails located in our test folder:

```
InstanceList testInstances = new
    InstanceList(classifier.getInstancePipe());
folderIterator = new FileIterator(
    new File[] {new File(testFolderPath)},
    new TxtFilter(),
    FileIterator.LAST_DIRECTORY);
```

We will pass the data through the same pipeline that we initialized during training:

```
testInstances.addThruPipe(folderIterator);
```

To evaluate classifier performance, we'll use the `cc.mallet.classify.Trial` class, which is initialized with a classifier and set of test instances:

```
Trial trial = new Trial(classifier, testInstances);
```

The evaluation is performed immediately at initialization. We can then simply take out the measures that we care about. In our example, we'd like to check the precision and recall on classifying spam e-mail messages, or F-measure, which returns a harmonic mean of both values, as follows:

```
System.out.println(
    "F1 for class 'spam': " + trial.getF1("spam"));
System.out.println(
    "Precision:" + trial.getPrecision(1));
System.out.println(
    "Recall:" + trial.getRecall(1));
```

The evaluation object outputs the following results:

```
F1 for class 'spam': 0.9731800766283524
Precision: 0.9694656488549618
Recall: 0.9769230769230769
```

The results show that the model correctly discovers 97.69% of spam messages (recall), and when it marks an e-mail as spam, it is correct in 96.94% cases. In other words, it misses approximately 2 per 100 spam messages and marks 3 per 100 valid messages as spam. Not really perfect, but it is more than a good start!

Summary

In this chapter, we discussed how text mining is different from traditional attribute-based learning, requiring a lot of pre-processing steps in order to transform written natural language into feature vectors. Further, we discussed how to leverage Mallet, a Java-based library for natural language processing by applying it to two real life problems. First, we modeled topics in news corpus using the LDA model to build a model that is able to assign a topic to new document. We also discussed how to build a naive Bayesian spam-filtering classifier using the bag-of-words representation.

This chapter concludes the technical demonstrations of how to apply various libraries to solve machine learning tasks. As we were not able to cover more interesting applications and give further details at many points, the next chapter gives some further pointers on how to continue learning and dive deeper into particular topics.

11

What is Next?

This chapter brings us to the end of our journey of reviewing machine learning Java libraries and discussing how to leverage them to solve real-life problems. However, this should not be the end of your journey by all means. This chapter will give you some practical advice on how to start deploying your models in the real world, what are the catches, and where to go to deepen your knowledge. It also gives you further pointers about where to find additional resources, materials, venues, and technologies to dive deeper into machine learning.

This chapter will cover the following topics:

- Important aspects of machine learning in real life
- Standards and markup languages
- Machine learning in the cloud
- Web resources and competitions

Machine learning in real life

Papers, conference presentations, and talks often don't discuss how the models were actually deployed and maintained in production environment. In this section, we'll look into some aspects that should be taken into consideration.

Noisy data

In practice, data typically contains errors and imperfections due to various reasons such as measurement errors, human mistakes, errors of expert judgment in classifying training examples, and so on. We refer to all of these as noise. Noise can also come from the treatment of missing values when an example with unknown attribute value is replaced by a set of weighted examples corresponding to the probability distribution of the missing value. The typical consequences of noise in learning data are low prediction accuracy of learned model in new data and complex models that are hard to interpret and to understand by the user.

Class unbalance

Class unbalance is a problem we come across in *Chapter 7, Fraud and Anomaly Detection*, where the goal was to detect fraudulent insurance claims. The challenge is that a very large part of the dataset, usually more than 90%, describes normal activities and only a small fraction of the dataset contains fraudulent examples. In such a case, if the model always predicts normal, then it is correct 90% of the time. This problem is extremely common in practice and can be observed in various applications, including fraud detection, anomaly detection, medical diagnosis, oil spillage detection, facial recognition, and so on.

Now knowing what the class unbalance problem is and why is it a problem, let's take a look at how to deal with this problem. The first approach is to focus on measures other than classification accuracy, such as recall, precision, and f-measure. Such measures focus on how accurate a model is at predicting minority class (recall) and what is the share of false alarms (precision). The other approach is based on resampling, where the main idea is to reduce the number of overrepresented examples in such way that the new set contains a balanced ratio of both the classes.

Feature selection is hard

Feature selection is arguably the most challenging part of modeling that requires domain knowledge and good insights into the problem at hand. Nevertheless, properties of well-behaved features are as follows:

- **Reusability:** Features should be available for reuse in different models, applications, and teams
- **Transformability:** You should be able to transform a feature with an operation, for example, `log()`, `max()`, or combine multiple features together with a custom calculation

- **Reliability:** Features should be easy to monitor and appropriate unit tests should exist to minimize bugs/issues
- **Interpretability:** In order to perform any of the previous actions, you need to be able to understand the meaning of features and interpret their values

The better you are able to capture the features, the more accurate your results will be.

Model chaining

Some models might produce an output, which is used as the feature in another model. Moreover, we can use multiple models—ensembles—turning any model into a feature. This is a great way to get better results, but this can lead to problems too. Care must be taken that the output of your model is ready to accept dependencies. Also, try to avoid feedback loops, as they can create dependencies and bottlenecks in pipeline.

Importance of evaluation

Another important aspect is model evaluation. Unless you apply your models to actual new data and measure a business objective, you're not doing predictive analytics. Evaluation techniques, such as cross-validation and separated train/test set, simply split your test data, which can give only you an estimate of how your model will perform. Life often doesn't hand you a train dataset with all the cases defined, so there is a lot of creativity involved in defining these two sets in a real-world dataset.

At the end of the day, we want to improve a business objective, such as improve ad conversion rate, get more clicks on recommended items, and so on. To measure the improvement, execute A/B tests, measure differences in metrics across statistically identical populations that each experience a different algorithm. Decisions on the product are always data-driven.



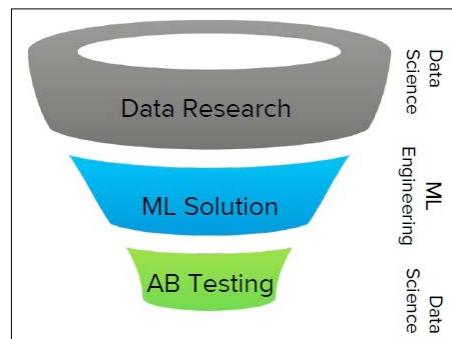
A/B testing is a method for a randomized experiment with two variants: A, which corresponds to the original version, controlling the experiment; and B, which corresponds to a variation. The method can be used to determine whether the variation outperforms the original version. It can be used to test everything from website changes to sales e-mails to search ads.

Udacity offers a free course, covering design and analysis of A/B tests at <https://www.udacity.com/course/ab-testing--ud257>.

Getting models into production

The path from building an accurate model in a lab to deploying it in a product involves collaboration of data science and engineering, as shown in the following three steps and diagram:

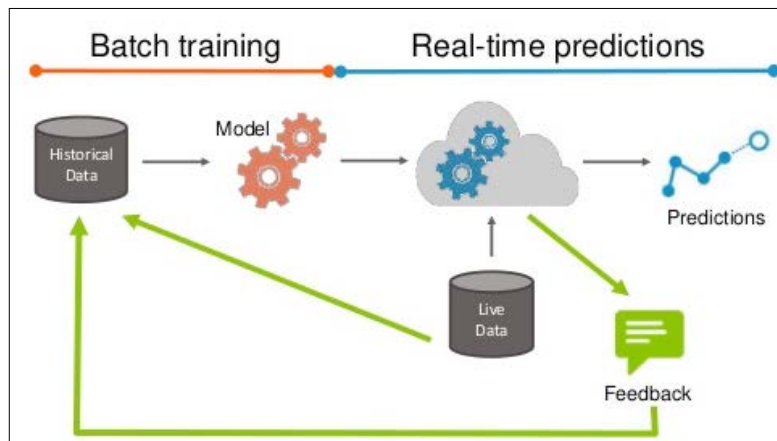
1. **Data research and hypothesis building** involves modeling the problem and executing initial evaluation.
2. **Solution building and implementation** is where your model finds its way into the product flow by rewriting it into more efficient, stable, and scalable code.
3. **Online evaluation** is the last stage where the model is evaluated with live data using A/B testing on business objectives.



Model maintenance

Another aspect that we need to address is how the model will be maintained. Is this a model that will not change over time? Is it modeling a dynamic phenomenon requiring the model to adjust its prediction over time?

The model is usually built in an of offline batch training and then used on live data to serve predictions as shown in the following figure. If we are able to receive feedback on model predictions; for instance, whether the stock went up as model predicted, whether the candidate responded to campaign, and so on, the feedback should be used to improve the initial model.



The feedback could be really useful to improve the initial model, but make sure to pay attention to the data you are sampling. For instance, if you have a model that predicts who will respond to a campaign, you will initially use a set of randomly contacted clients with specific responded/not responded distribution and feature properties. The model will focus only on a subset of clients that will most likely respond and your feedback will return you a subset of clients that responded. By including this data, the model is more accurate in a specific subgroup, but might completely miss some other group. We call this problem exploration versus exploitation. Some approaches to address this problem can be found in Osugi et al (2005) and Bondu et al (2010).

Standards and markup languages

As predictive models become more pervasive, the need for sharing the models and completing the modeling process leads to formalization of development process and interchangeable formats. In this section, we'll review two de facto standards, one covering data science processes and the other specifying an interchangeable format for sharing models between applications.

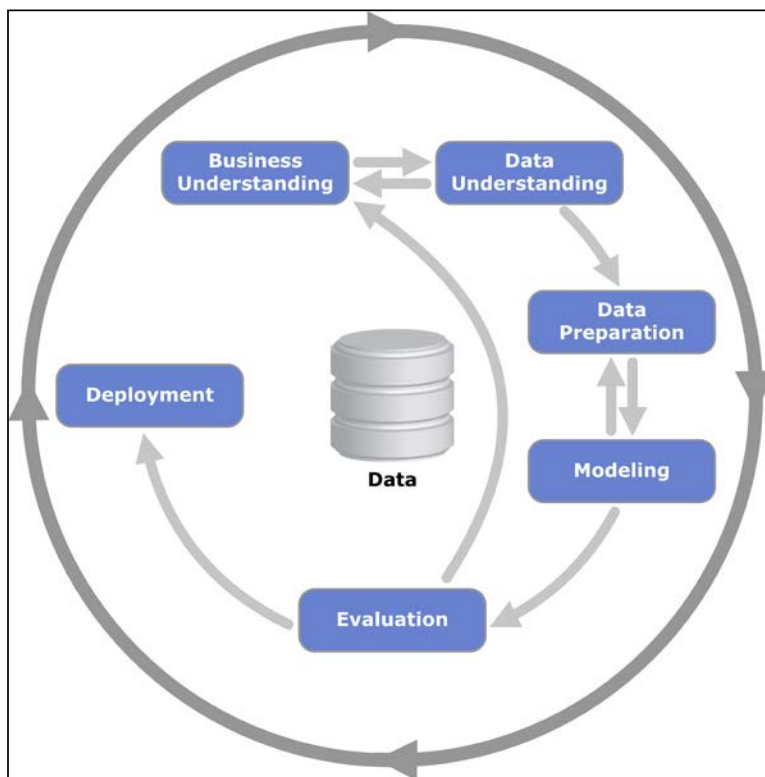
CRISP-DM

Cross Industry Standard Process for Data Mining (CRISP-DM) describing a data mining process commonly used by data scientists in industry. CRISP-DM breaks the data mining science process into the following six major phases:

- **Business understanding**
- **Data understanding**

- **Data preparation**
- **Modeling**
- **Evaluation**
- **Deployment**

In the following diagram, the arrows indicate the process flow, which can move back and forth through the phases. Also, the process doesn't stop with model deployment. The outer arrow indicates the cyclic nature of data science. Lessons learned during the process can trigger new questions and repeat the process while improving previous results:



SEMMA methodology

Another methodology is **Sample, Explore, Modify, Model, and Assess (SEMMA)**. SEMMA describes the main modeling tasks in data science, while leaving aside business aspects such as data understanding and deployment. SEMMA was developed by SAS institute, which is one of the largest vendors of statistical software, aiming to help the users of their software to carry out core tasks of data mining.

Predictive Model Markup Language

Predictive Model Markup Language (PMML) is an XML-based interchange format that allows machine learning models to be easily shared between applications and systems. Supported models include logistic regression, neural networks, decision trees, naïve Bayes, regression models, and many others. A typical PMML file consists of the following sections:

- Header containing general information
- Data dictionary, describing data types
- Data transformations, specifying steps for normalization, discretization, aggregations, or custom functions
- Model definition, including parameters
- Mining schema listing attributes used by the model
- Targets allowing post-processing of the predicted results
- Output listing fields to be outputted and other post-processing steps

The generated PMML files can be imported to any PMML-consuming application, such as Zementis **Adaptive Decision and Predictive Analytics (ADAPA)** and **Universal PMML Plug-in (UPPI)** scoring engines; Weka, which has built-in support for regression, general regression, neural network, **TreeModel**, **RuleSetModel**, and **Support Vector Machine (SVM)** model; Spark, which can export k-means clustering, linear regression, ridge regression, lasso model, binary logistic model, and SVM; and cascading, which can transform PMML files into an application on Apache Hadoop.

The next generation of PMML is an emerging format called **Portable Format for Analytics (PFA)**, providing a common interface to deploy the complete workflows across environments.

Machine learning in the cloud

Setting up a complete machine learning stack that is able to scale with the increasing amount of data could be challenging. Recent wave of **Software as a Service (SaaS)** and **Infrastructure as a Service (IaaS)** paradigm was spilled over to machine learning domain as well. The trend today is to move the actual data preprocessing, modeling, and prediction to cloud environments and focus on modeling task only.

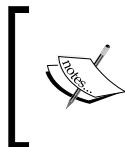
In this section, we'll review some of the promising services offering algorithms, predictive models already train in specific domain, and environments empowering collaborative workflows in data science teams.

Machine learning as a service

The first category is algorithms as a service, where you are provided with an API or even graphical user interface to connect pre-programmed components of data science pipeline together:

- **Google Prediction API** was one of the first companies that introduced prediction services through its web API. The service is integrated with **Google Cloud Storage** serving as data storage. The user can build a model and call an API to get predictions.
- **BigML** implements a user-friendly graphical interface, supports many storage providers (for instance, Amazon S3) and offers a wide variety of data processing tools, algorithms, and powerful visualizations.
- **Microsoft Azure Machine Learning** provides a large library of machine learning algorithms and data processing functions, as well as graphical user interface, to connect these components to an application. Additionally, it offers a fully-managed service that you can use to deploy your predictive models as ready-to-consume web services.
- **Amazon Machine Learning** entered the market quite late. It's main strength is seamless integration with other Amazon services, while the number of algorithms and user interface needs further improvements.
- **IBM Watson Analytics** focuses on providing models that are already hand-crafted to a particular domain such as speech recognition, machine translations, and anomaly detection. It targets a wide range of industries by solving specific use cases.
- **Prediction.IO** is a self-hosted open source platform, providing the full stack from data storage to modeling to serving the predictions. Prediction.IO can talk to Apache Spark to leverage its learning algorithms. In addition, it is shipped with a wide variety of models targeting specific domains, for instance, recommender system, churn prediction, and others.

Predictive API is an emerging new field, so these are just some of the well-known examples; **KDnuggets** compiled a list of 50 machine learning APIs at <http://www.kdnuggets.com/2015/12/machine-learning-data-science-apis.html>.



To learn more about it, you can visit PAPI, the International Conference on Predictive APIs and Apps at <http://www.papi.io> or take a look at a book by *Louis Dorard, Bootstrapping Machine Learning* (L. Dorard, 2014).

Web resources and competitions

In this section, we'll review where to find additional resources for learning, discussing, presenting, or sharpening our data science skills.

Datasets

One of the most well-known repositories of machine learning datasets is hosted by the University of California, Irvine. The UCI repository contains over 300 datasets covering a wide variety of challenges, including poker, movies, wine quality, activity recognition, stocks, taxi service trajectories, advertisements, and many others. Each dataset is usually equipped with a research paper where the dataset was used, which can give you a hint on how to start and what is the prediction baseline.

The UCI machine learning repository can be accessed at <https://archive.ics.uci.edu>, as follows:

The screenshot shows the homepage of the UCI Machine Learning Repository. The header includes the UCI logo and navigation links for 'About', 'Citation Policy', 'Donate a Data Set', and 'Contact'. A 'View ALL Data Sets' link is also present. The main content area is divided into three columns:

- Latest News:** A list of recent updates, including the addition of new data sets and the setting up of a repository mirror.
- Newest Data Sets:** A list of the most recently added datasets, each with a UCI logo and a brief description.
- Most Popular Data Sets (hits since 2007):** A list of the most frequently accessed datasets, including Iris, Adult, Wine, Car Evaluation, Breast Cancer Wisconsin (Diagnostic), Abalone, and Wine Quality.

At the bottom left, there is a 'Featured Data Set: Movie' section with a small image of the Hollywood sign and details about the dataset's type and number of instances.

Another well-maintained collection by Xiaming Chen is hosted on GitHub:

<https://github.com/caesar0301/awesome-public-datasets>

The Awesome Public Datasets repository maintains links to more than 400 data sources from a variety of domains, ranging from agriculture, biology, economics, psychology, museums, and transportation. Datasets, specifically targeting machine learning, are collected under the image processing, machine learning, and data challenges sections.

Online courses

Learning how to become a data scientist has become much more accessible due to the availability of online courses. The following is a list of free resources to learn different skills online:

- Online courses for learning Java:
 - **Udemy: Learn Java Programming From Scratch** at <https://www.udemy.com/learn-java-programming-from-scratch>
 - **Udemy: Java Tutorial for Complete Beginners** at <https://www.udemy.com/java-tutorial>
 - LearnJavaOnline.org: **Interactive Java tutorial** at <http://www.learnjavaonline.org/>

- Online courses to learn more about machine learning:
 - **Coursera: Machine Learning (Stanford) by Andrew Ng**: This teaches you the math behind many the machine learning algorithms, explains how they work, and explores why they make sense at <https://www.coursera.org/learn/machine-learning>.
 - **Statistics 110 (Harvard) by Joe Biltzstein**: This course lets you discover the probability of related terms that you will hear many times in your data science journey. Lectures are available on YouTube at <http://projects.iq.harvard.edu/stat110/youtube>.
 - **Data Science CS109 (Harvard) by John A. Paulson**: This is a hands-on course where you'll learn about Python libraries for data science, as well as how to handle machine-learning algorithms at <http://cs109.github.io/2015/>.

Competitions

The best way to sharpen your knowledge is to work on real problems; and if you want to build a proven portfolio of your projects, machine learning competitions are a viable place to start:

- **Kaggle:** This is the number one competition platform, hosting a wide variety of challenges with large prizes, strong data science community, and lots of helpful resources. You can check it out at <https://www.kaggle.com/>.
- **CrowdANALYTIX:** This is a crowdsourced data analytics service that is focused on the life sciences and financial services industries at <https://www.crowdanalytix.com>.
- **DrivenData:** This hosts data science competitions for social good at <http://www.drivendata.org/>.

Websites and blogs

In addition to online courses and competitions, there are numerous websites and blogs publishing the latest developments in the data science community, their experience in attacking different problems, or just best practices. Some good starting points are as follows:

- **KDnuggets:** This is the de facto portal for data mining, analytics, big data, and data science, covering the latest news, stories, events, and other relevant issues at <http://www.kdnuggets.com/>.
- **Machine learning mastery:** This is an introductory-level blog with practical advice and pointers where to start. Check it out at <http://machinelearningmastery.com/>.
- **Data Science Central:** This consists of practical community articles on a variety of topics, algorithms, caches, and business cases at <http://www.datasciencecentral.com/>.
- **Data Mining Research** by Sandro Saitta at <http://www.dataminingblog.com/>.
- **Data Mining: Text Mining, Visualization and Social Media** by Matthew Hurst, covering interesting text and web mining topics, frequently with applications to Bing and Microsoft at http://datamining.typepad.com/data_mining/.
- **Geeking with Greg** by Greg Linden, inventor of Amazon recommendation engine and Internet entrepreneur. You can check it out at <http://glinden.blogspot.si/>.
- **DSGuide:** This is a collection of over 150 data science blogs at <http://dsguide.biz/reader/sources>.

Venues and conferences

The following are a few top-tier academic conferences with the latest algorithms:

- **Knowledge Discovery in Databases (KDD)**
- **Computer Vision and Pattern Recognition (CVPR)**
- **Annual Conference on Neural Information Processing Systems (NIPS)**
- **International Conference on Machine Learning (ICML)**
- **IEEE International Conference on Data Mining (ICDM)**
- **International Joint Conference on Pervasive and Ubiquitous Computing (UbiComp)**
- **International Joint Conference on Artificial Intelligence (IJCAI)**

Some business conferences are as follows:

- O'Reilly Strata Conference
- The Strata + Hadoop World Conferences
- Predictive Analytics World
- MLconf

You can also check local meetup groups.

Summary

In this chapter, we concluded the book by discussing some aspects of model deployment, we also looked into standards for data science process and interchangeable predictive model format PMML. We also reviewed online courses, competitions, web resources, and conferences that could help you in your journey towards mastering the art of machine learning.

I hope this book inspired you to dive deeper into data science and has motivated you to get your hands dirty, experiment with various libraries and get a grasp of how different problems could be attacked. Remember, all the source code and additional resources are available at the supplementary website <http://www.machine-learning-in-java.com>.

References

The following are the references for all the citations throughout the book:

- Adomavicius, G. and Tuzhilin, A.. Toward the next generation of recommender systems: a survey of the state-of-the-art and possible extensions. *IEEE Transactions on Knowledge and Data Engineering*, 17(6), 734-749. 2005.
- Bengio, Y.. Learning Deep Architectures for AI. *Foundations and Trends in Machine Learning* 2(1), 1-127. 2009. Retrieved from <http://www.iro.umontreal.ca/~bengioy/papers/ftml.pdf>.
- Blei, D. M., Ng, A. Y., and Jordan, M. I.. Latent dirichlet allocation. *Journal of Machine Learning Research*. 3, 993-1022. 2003. Retrieved from: <http://www.jmlr.org/papers/volume3/blei03a/blei03a.pdf>.
- Bondu, A., Lemaire, V., Boulle, M.. Exploration vs. exploitation in active learning: A Bayesian approach. *The 2010 International Joint Conference on Neural Networks (IJCNN)*, Barcelona, Spain. 2010.
- Breunig, M. M., Kriegel, H.-P., Ng, R. T., Sander, J.. LOF: Identifying Density-based Local Outliers (PDF). *Proceedings from the 2000 ACM SIGMOD International Conference on Management of Data*, 29(2), 93-104. 2000
- Campbell, A. T. (n.d.). Lecture 21 - Activity Recognition. Retrieved from <http://www.cs.dartmouth.edu/~campbell/cs65/lecture22/lecture22.html>.
- Chandra, N.S. Unraveling the Customer Mind. 2012. Retrieved from <http://www.cognizant.com/InsightsWhitepapers/Unraveling-the-Customer-Mind.pdf>.

- Dror, G., Boule \square , M., Guyon, I., Lemaire, V., and Vogel, D.. The 2009 Knowledge Discovery in Data Competition (KDD Cup 2009) Volume 3, Challenges in Machine Learning, Massachusetts, US. Microtome Publishing. 2009.
- Gelman, A. and Nolan, D.. Teaching Statistics a bag of tricks. Cambridge, MA. Oxford University Press. 2002.
- Goshtasby, A. A. *Image Registration Principles, Tools and Methods*. London, Springer. 2012.
- Greene, D. and Cunningham, P.. Practical Solutions to the Problem of Diagonal Dominance in Kernel Document Clustering. Proceedings from the 23rd International Conference on Machine Learning, Pittsburgh, PA. 2006. Retrieved from http://www.autonlab.org/icml_documents/camera-ready/048_Practical_Solutions.pdf.
- Gupta, A.. *Learning Apache Mahout Classification*, Birmingham, UK. Packt Publishing. 2015.
- Gutierrez, N.. Demystifying Market Basket Analysis. 2006. Retrieved from <http://www.information-management.com/specialreports/20061031/1067598-1.html>.
- Hand, D., Manilla, H., and Smith, P.. *Principles of Data Mining*. USA. MIT Press. 2001. Retrieved from ftp://gamma.sbin.org/pub/doc/books/Principles_of_Data_Mining.pdf.
- Intel. What Happens in an Internet Minute?. 2013. Retrieved from <http://www.intel.com/content/www/us/en/communications/internet-minute-infographic.html>.
- Kaluža, B.. *Instant Weka How-To*. Birmingham. Packt Publishing. 2013.
- Karkera, K. R.. *Building Probabilistic Graphical Models with Python*. Birmingham, UK. Packt Publishing. 2014.
- KDD (n.d.). KDD Cup 2009: Customer relationship prediction. Retrieved from <http://www.kdd.org/kdd-cup/view/kdd-cup-2009>.
- Koller, D. and Friedman, N.. *Probabilistic Graphical Models Principles and Techniques*. Cambridge, Mass. MIT Press. 2012.
- Kurucz, M., Siklósi, D., Bíró, I., Csizsek, P., Fekete, Z., Iwatt, R., Kiss, T., and Szabó, A.. KDD Cup 2009 @ Budapest: feature partitioning and boosting 61. JMLR W&CP 7, 65–75. 2009.
- Laptev, N., Amizadeh, S., and Billawala, Y. (n.d.). A Benchmark Dataset for Time Series Anomaly Detection. Retrieved from <http://yahoolabs.tumblr.com/post/114590420346/a-benchmark-dataset-for-time-series-anomaly>.

-
- LKurgan, L.A. and Musilek, P.. A survey of Knowledge Discovery and Data Mining process models. *The Knowledge Engineering Review*, 21(1), 1–24. 2006.
 - Lo, H.-Y., Chang, K.-W., Chen, S.-T., Chiang, T.-H., Ferng, C.-S., Hsieh, C.-J., Ko, Y.-K., Kuo, T.-T., Lai, H.-C., Lin, K.-Y., Wang, C.-H., Yu, H.-F., Lin, C.-J., Lin, H.-T., and Lin, S.-de. An Ensemble of Three Classifiers for KDD Cup 2009: Expanded Linear Model, Heterogeneous Boosting, and Selective Naive Bayes, *JMLR W&CP* 7, 57–64. 2009.
 - Magalhães, P. Incorrect information provided by your website. 2010. Retrieved from <http://www.best.eu.org/aboutBEST/helpdeskRequest.jsp?req=f5wpxc8&auth=Paulo>.
 - Mariscal, G., Marban, O., and Fernandez, C.. A survey of data mining and knowledge discovery process models and methodologies. *The Knowledge Engineering Review*, 25(2), 137–166. 2010.
 - Mew, K. (2015). *Android 5 Programming by Example*. Birmingham, UK. Packt Publishing.
 - Miller, H., Clarke, S., Lane, S., Lonie, A., Lazaridis, D., Petrovski, S., and Jones, O.. Predicting customer behavior: The University of Melbourne's KDD Cup report, *JMLR W&CP* 7, 45–55. 2009.
 - Niculescu-Mizil, A., Perlich, C., Swirszcz, G., Sindhwani, V., Liu, Y., Melville, P., Wang, D., Xiao, J., Hu, J., Singh, M., Shang, W. X., and Zhu, Y. F.. Winning the KDD Cup Orange Challenge with Ensemble Selection. *JMLR W&CP*, 7, 23–34. 2009. Retrieved from <http://jmlr.org/proceedings/papers/v7/niculescu09/niculescu09.pdf>.
 - Oracle (n.d.). Anomaly Detection. Retrieved from http://docs.oracle.com/cd/B28359_01/datamine.111/b28129/anomalies.htm.
 - Osugi, T., Deng, K., and Scott, S.. Balancing exploration and exploitation: a new algorithm for active machine learning. Fifth IEEE International Conference on Data Mining, Houston, Texas. 2005.
 - Power, D. J. (ed.). DSS News. *DSSResources.com*, 3(23). 2002. Retrieved from <http://www.dssresources.com/newsletters/66.php>.
 - Quinlan, J. R. *C4.5: Programs for Machine Learning*. San Francisco, CA. Morgan Kaufmann Publishers. 1993.
 - Rajak, A.. Association Rule Mining-Applications in Various Areas. 2008. Retrieved from https://www.researchgate.net/publication/238525379_Association_rule_mining_Applications_in_various_areas.
 - Ricci, F., Rokach, L., Shapira, B., and Kantor, P. B. (eds.). *Recommender Systems Handbook*. New York, Springer. 2010.

- Rumsfeld, D. H. and Myers, G.. DoD News Briefing – Secretary Rumsfeld and Gen. Myers. 2002. Retrieved from <http://archive.defense.gov/transcripts/transcript.aspx?transcriptid=2636>.
- Stevens, S. S.. On the Theory of Scales of Measurement. *Science*, 103 (2684), 677–680. 1946.
- Sutton, R. S. and Barto, A. G.. *Reinforcement Learning An Introduction*. Cambridge, MA: MIT Press. 1998.
- Tiwary, C.. *Learning Apache Mahout*. Birmingham, UK. Packt Publishing. 2015.
- Tsai, J., Kaminka, G., Epstein, S., Zilka, A., Rika, I., Wang, X., Ogden, A., Brown, M., Fridman, N., Taylor, M., Bowring, E., Marsella, S., Tambe, M., and Sheel, A.. ESCAPES - Evacuation Simulation with Children, Authorities, Parents, Emotions, and Social comparison. Proceedings from 10th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2011) 2 (6), 457–464. 2011. Retrieved from http://www.aamas-conference.org/Proceedings/aamas2011/papers/D3_G57.pdf.
- Tsanas, A. and Xifara. Accurate quantitative estimation of energy performance of residential buildings using statistical machine learning tools. *Energy and Buildings*, 49, 560-567. 2012.
- Utts, J.. What Educated Citizens Should Know About Statistics and Probability. *The American Statistician*, 57 (2), 74-79. 2003.
- Wallach, H. M., Murray, I., Salakhutdinov, R., and Mimno, D.. Evaluation Methods for Topic Models. Proceedings from the 26th International conference on Machine Learning, Montreal, Canada. 2009. Retrieved from <http://mimno.infosci.cornell.edu/papers/wallach09evaluation.pdf>.
- Witten, I. H. and Frank, E.. *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*. USA. Morgan Kaufmann Publishers. 2000.
- Xie, J., Rojkova, V., Pal, S., and Coggeshall, S.. A Combination of Boosting and Bagging for KDD Cup 2009. *JMLR W&CP*, 7, 35–43. 2009.
- Zhang, H.. The Optimality of Naive Bayes. Proceedings from FLAIRS 2004 conference. 2004. Retrieved from <http://www.cs.unb.ca/~hzhang/publications/FLAIRS04ZhangH.pdf>.
- Ziegler, C-N., McNee, S. M., Konstan, J. A., and Lausen, G.. Improving Recommendation Lists Through Topic Diversification. Proceedings from the 14th International World Wide Web Conference (WWW '05), Chiba, Japan. 2005. Retrieved from <http://www2.informatik.uni-freiburg.de/~chiegler/papers/WWW-05-CR.pdf>.

Index

A

A/B tests

URL 211

activation function 147

activity recognition

about 166

mobile phone sensors 166, 167

pipeline 168, 169

plan 169

AdaBoost M1 method 56

advanced modeling

attribute selection 80

data, pre-processing 79

ensembleLibrary package, using 78

model selection 81-84

performance, evaluation 85

with ensembles 77

affinity analysis

about 87, 89

cross-industry applications 97

agglomerative clustering 16

Amazon Machine Learning 216

analysis types

about 128

pattern analysis 128

transaction analysis 129

Android Device Monitor 177

Android Studio

installing 170, 171

URL 170

anomalous behavior detection

about 126

unknown-unknowns 126

anomalous pattern detection

about 128

analysis types 128

plan recognition 129

anomaly detection, in time series data

about 138

data, loading 140, 141

density based k-nearest neighbors 142-144

histogram-based anomaly

detection 138-140

histograms, creating 141, 142

anomaly detection, in website traffic

about 137

dataset, using 137

Apache Mahout

about 35, 36

configuring 105

configuring, in Eclipse with

Maven plugin 106, 107

Apache Spark

about 36-38

URL 36

Application Portfolio Management

(APM) 99

Applied Machine Learning

about 1

workflow 3, 4

Apriori algorithm

about 87, 92

used, for discovering shopping

patterns 94-96

artificial neural networks 18

association rule learning

about 90

Apriori algorithm 92

- confidence 92
- database, of transactions 90, 91
- FP-growth algorithm 93
- itemset 91
- rule 91
- support 92

autoencoder 149, 150

B

bag-of-words (BoW) 191

basic modeling
about 75
models, evaluating 75, 76
naive Bayes baseline, implementing 76, 77

basic naive Bayes classifier baseline
about 71
data, loading 73, 74
data, obtaining 72, 73

BBC dataset
URL 196

big data
dealing with 43
variety 43
velocity 43
volume 43

big data application
architecture 43

BigML 216

Book-Crossing dataset
BX-Book-Ratings file 108
BX-Books file 108
BX-Users file 108
URL 108

book-recommendation engine
book ratings dataset, using 108
building 108
collaborative filtering, implementing 114
content-based filtering,
implementing 123, 124
custom rules, adding 119
data, loading 109
data, loading from database 112, 113
data, loading from file 109-111
evaluation 120, 121
in-memory database, creating 113, 114
online learning engine 121, 122

C

Canova library
URL 155

Cassandra
about 44
URL 44

cc.mallet.pipe package
CharSequenceRemoveHTML pipeline 194
Input2CharSequence pipeline 194
MakeAmpersandXMLFriendly
pipeline 194
TokenSequence2FeatureSequence
pipeline 195
TokenSequenceLowercase pipeline 194
TokenSequenceNGrams pipeline 195

Chebyshev distance 34

classification
about 17, 46
artificial neural networks 18
classification algorithm, selecting 55, 56
confusion matrix, examining 54, 55
data, classifying 53
data, loading 48
data, using 47
decision trees learning 18
ensemble learning 19
evaluating 19, 20
evaluation 54
feature selection 49
kernel methods 18
learning algorithms, selecting 50-52
precision 20
prediction error metrics 54
probabilistic classifiers 18
recall 20
Roc curves 20, 21

classification algorithms, examples
weka.classifiers.bayes.NaiveBayes 56
weka.classifiers.functions.Multilayer
Perceptron 56
weka.classifiers.lazy.IBk 56
weka.classifiers.meta.AdaBoostM1 56
weka.classifiers.meta.Bagging 56
weka.classifiers.rules.ZeroR 55
weka.classifiers.trees.RandomForest 56
weka.classifiers.trees.RandomTree 55

- classifier**
 - building 178, 179
 - plugging, into mobile app 182, 183
 - spurious transitions, reducing 180-182
- class unbalance 210**
- clustering**
 - about 16, 17, 63, 64
 - clustering algorithms 64, 65
 - evaluation 66
- collaborative filtering**
 - about 101, 103
 - implementing, with book-recommendation engine 114
 - item-based 118, 119
 - user-based 115-117
- Comma Separated Value (CSV) 57**
- competitions 219**
- conjugate gradient optimization algorithm**
 - building 156
- content-based filtering**
 - about 101, 104
 - implementing, with book-recommendation engine 123, 124
- Contrastive Divergence algorithm 151**
- Convolutional Neural Network (CNN) 151**
- Core Motion framework, iOS**
 - URL 167
- correlation coefficient 23**
- cosine distance 104**
- cost function 17**
- cross-industry applications, of affinity analysis**
 - about 97
 - census data 98
 - customer relationship management (CRM) 98
 - IT Operations Analytics 99
 - medical diagnosis 97
 - protein sequences 97
- Cross Industry Standard Process for Data Mining (CRISP-DM) 213**
- cross-validation 26**
- CrowdANALYTIX**
 - URL 219
- CSVLoader class**
 - URL 73

- curse of dimensionality 15**
- customer relationship database**
 - about 68
 - challenge 68, 69
 - dataset 69, 70
 - evaluation 71
- Customer Relationship Management (CRM) 67**

D

- data and problem definition**
 - about 4
 - measurement scales 5, 6
- data collection**
 - about 6
 - Android Studio, installing 170, 171
 - data collector, loading 172, 173
 - data, generating 8
 - data, observing 7, 8
 - data, searching 7, 8
 - from mobile phone 170
 - training data, collecting 175-178
 - traps, sampling 9
- data collector**
 - feature extraction 174, 175
 - loading 172, 173
 - URL 172
- Data Mining Research**
 - URL 219
- data pre-processing**
 - about 9
 - data cleaning 9
 - data reduction 12, 13
 - data transformation 11, 12
 - missing values, filling 10
 - outliers, removing 11
- data science 1, 2**
- Data Science Central**
 - URL 219
- data scientist 2**
- dataset rebalancing 134-136**
- datasets 217, 218**
- Decision and Predictive Analytics (ADAPA) 215**
- decision trees 25**
- decision trees learning 18**

deep belief network
 about 18, 151
 building 158-160

deep convolutional networks 151-154

deeplearning4j (DL4J)
 about 38, 39, 153
 obtaining 153, 154
 org.deeplearning4j.base 38
 org.deeplearning4j.berkeley 38
 org.deeplearning4j.clustering 38
 org.deeplearning4j.datasets 38
 org.deeplearning4j.distributions 38
 org.deeplearning4j.eval 38
 org.deeplearning4j.exceptions 38
 org.deeplearning4j.models 38
 org.deeplearning4j.nn 39
 org.deeplearning4j.optimize 39
 org.deeplearning4j.plot 39
 org.deeplearning4j.rng 39
 org.deeplearning4j.util 39
 URL 38

delta rule 147

directory
 text data, importing 192, 193

Discrete Fourier Transform (DFT) 168

distance measures
 Euclidean distances 13
 non-Euclidean distances 14, 15

double evaluateLeftToRight method
 boolean useResampling component 201
 Instances heldOutDocuments
 component 201
 int numParticles component 201
 PrintStream docProbabilityStream
 component 201

DrivenData 219

DropConnect neural network 154

DSGuide
 URL 219

dynamic time wrapping (DTW) 34

E

Eclipse
 Apache Mahout, configuring with
 Maven plugin 106, 107

Eclipse IDE
 using 46

Edit distance 15

elbow method 64

e-mail spam detection
 about 203
 default pipeline, creating 205, 206
 e-mail spam dataset, collecting 204, 205
 model performance, evaluating 207
 testing 206
 training 206

energy efficiency dataset
 URL 56

ensembleSel.setOptions () method
 -A <algorithm> option 84
 -B <numModelBags> option 83
 -D option 84
 -E <modelRatio> option 83
 -G option 84
 -H <hillClimbIterations> option 83
 -I <sortInitialization> option 83
 -L </path/to/modelLibrary> option 83
 -O option 84
 -P <hillclimbMetric> option 83
 -R option 84
 -S <num> option 84
 -V <validationRatio> option 83
 -W </path/to/working/directory>
 option 83
 -X <numFolds> option 83

ensemble learning 19

ensembleLibrary package
 URL 78
 using 78

ensembles
 used, for advanced modeling 77

Ensemble Selection algorithm 77

environmental sensors 167

Euclidean distances 13

evaluate() method, parameters
 DataModel 121
 DataModelBuilder 121
 evaluationPercentage 121
 RecommenderBuilder 121
 trainingPercentage 121

evaluation 24

Expectation Maximization (EM)

- clustering 63

- exploitation 105

- exploration 105

F

- Feature extraction 42

- feature map 151

- feature selection 12

- feedforward neural networks 148

- file

 - text data, importing 193, 194

- Fourier transform

 - reference link 168

- FP-growth algorithm

 - about 87, 93

 - used, for discovering shopping patterns 96

- FP-tree structure 93

- fraud detection, of insurance claims

 - about 129

 - dataset, using 130, 131

 - suspicious patterns, modeling 132

- frequent pattern (FP) 93

G

- Geeking with Greg

 - URL 219

- generalization

 - about 24

 - cross-validation 26

 - leave-one-out validation 26

 - overfitting 24, 25

 - stratification 27

 - test set 26

 - train set 26

 - underfitting 24, 25

- Generalized Sequential Patterns (GSP) 32

- Generative Stochastic Networks (GSNs) 150

- Gibbs sampling 151

- GNU General Public License

 - (GNU GPL) 31

- Google Prediction API 216

- Graphics Processing Unit (GPU)

 - about 163

 - reference link 163

- GraphX 36

H

- Hadoop

 - about 44

 - URL 44

- Hadoop Distributed File System (HDFS) 37

- Hamming distance 15

- HBase

 - about 44

 - URL 44

- Hidden layer 148

- Hidden Markov Models (HMMs) 36, 129

- hierarchical clustering 16

- histogram-based anomaly detection 138-140

- Hotspot 32

- hybrid approach 104

I

- IBM Watson Analytics 216

- image classification

 - about 153

 - data, loading 154, 155

 - deeplearning4j 153

 - MNIST dataset 154

 - models, building 155

- ImageNet

 - about 146

 - URL 152

- image recognition

 - about 145, 146

 - neural networks 147

- Infrastructure as a Service (IaaS) 215

- Input layer 148

- insurance claims

 - fraud detection 129

- interval data 5

- Intrusion Detection (ID) 129

- item-based analysis 103

- item-based collaborative filtering 118, 119

- IT Operations Analytics 87

J

- Jaccard distance 14

- Java

 - need for 30

Java API packages, Weka

- weka.associations 32
- weka.classifiers 32
- weka.clusterers 33
- weka.core 33
- weka.datagenerators 33
- weka.estimators 33
- weka.experiment 33
- weka.filters 33
- weka.gui 33

Java machine learning (Java-ML)

- about 34
- URL 34

Java-ML packages

- net.sf.javaml.classification 34
- net.sf.javaml.clustering 34
- net.sf.javaml.core 34
- net.sf.javaml.distance 34
- net.sf.javaml.featureselection 35
- net.sf.javaml.filter 35
- net.sf.javaml.matrix 35
- net.sf.javaml.sampling 35
- net.sf.javaml.tools 35
- net.sf.javaml.utils 35

java -Xmx16g 85

K

Kaggle 219

KDD Cup

- about 68
- URL 68

KDnuggets

- about 216
- URL 219

kernel methods 18

k-means clustering 16

k-nearest neighbors 25

Knowledge Discovery and Data Science (KDD) 67

known-knowns 126

known-unknowns 126

L

Latent Dirichlet Allocation (LDA) 186, 198

leave-one-out validation 26

Linear Discriminant Analysis (LDA)

- about 198
- reference link 201

linear regression 59, 60

Local Outlier Factor (LOF)

- about 140
- URL 142

M

machine learning

- about 1, 2
- advantages 2, 3
- as service 216
- class unbalance 210
- evaluation 211
- feature selection 210
- in cloud 215
- in real life 209
- model chaining 211
- models, in production 212
- models, maintaining 212, 213
- noisy data 210
- online courses 218
- reinforcement learning 3
- supervised learning 2
- unsupervised learning 3
- web resources 217-219

machine learning application

- big data, dealing with 43
- building 42
- traditional machine learning 42

Machine Learning for Language Toolkit (MALLET)

- about 39, 40
- installing 188-190
- reference link 195
- URL 39

machine learning libraries

- about 30
- Apache Mahout 35, 36
- Apache Spark 36-38
- comparing 41
- deeplearning4j 38, 39
- Java machine learning (Java-ML) 34
- Machine Learning for Language Toolkit (MALLET) 39, 40

- Waikato Environment for Knowledge
 - Analysis (Weka) 30-33
- Machine learning mastery**
 - URL 219
- Mahalanobis distance** 15, 34
- Mahout interfaces, abstractions**
 - DataModel 114
 - ItemSimilarity 114
 - Recommender 114
 - UserNeighborhood 114
 - UserSimilarity 114
- Mahout libraries**
 - org.apache.mahout.cf.taste 36
 - org.apache.mahout.classifier 36
 - org.apache.mahout.clustering 36
 - org.apache.mahout.common 36
 - org.apache.mahout.ep 36
 - org.apache.mahout.math 36
 - org.apache.mahout.vectorizer 36
- MALLET, packages**
 - cc.mallet.classify 40
 - cc.mallet.cluster 40
 - cc.mallet.extract 40
 - cc.mallet.fst 40
 - cc.mallet.grmm 40
 - cc.mallet.optimize 40
 - cc.mallet.pipe 40
 - cc.mallet.topics 40
 - cc.mallet.types 40
 - cc.mallet.util 40
- Manhattan distance** 34
- market basket analysis (MBA)**
 - about 87-89
 - affinity analysis 89
 - identification, of driver items 88
 - item affinity 88
 - marketing 89
 - operations optimization 89
 - revenue optimization 89
 - store-to-store comparison 89
 - trip classification 88
- Markov chain** 151
- Maven plugin**
 - Apache Mahout, configuring with 106-108
- mean absolute error** 23
- mean squared error** 23
- measurement scales**
 - about 5, 6
 - interval data 5
 - nominal data 5
 - ordinal data 5
 - ratio data 6
- Microsoft Azure Machine Learning** 216
- Minkowski distance** 34
- missing values**
 - filling 10
- MLlib API library**
 - org.apache.spark.mllib.classification 37
 - org.apache.spark.mllib.clustering 37
 - org.apache.spark.mllib.linalg 37
 - org.apache.spark.mllib.optimization 37
 - org.apache.spark.mllib.recommendation 37
 - org.apache.spark.mllib.regression 38
 - org.apache.spark.mllib.stat 38
 - org.apache.spark.mllib.tree 38
 - org.apache.spark.mllib.util 38
- MNIST dataset** 154
- mobile app**
 - classifier, plugging into 182, 183
- mobile phone**
 - data, collecting 170
- mobile phone sensors**
 - about 166
 - environmental sensors 167
 - motion sensors 167
 - position sensors 167
 - URL, for Android 167
 - URL, for Windows Phone 167
- models**
 - building 155
 - chaining 211
 - deep belief network, building 158-160
 - in production 212
 - maintenance 212, 213
 - Multilayer Convolutional Network,
 - building 160-163
 - single layer regression model,
 - building 156-158
- MongoDB**
 - about 44
 - URL 44
- motion sensors** 167

Mozilla Thunderbird 203
Multilayer Convolutional Network
 about 155
 building 160-163
myrunscollector package
 CollectorActivity.java class 173
 Globals.java class 173
 SensorsService.java class 173

N

Naive Bayes 25
naive Bayes baseline
 implementing 76, 77
neural networks
 about 25, 147
 autoencoder 149, 150
 deep convolutional networks 151, 152
 feedforward neural networks 148
 perceptron 147
 Restricted Boltzman machine 150, 151
nominal data 5
non-Euclidean distance 14, 15

O

online learning engine 121, 122
Oracle Database Online Documentation
 URL 130
ordinal data 5
outliers
 removing 11
Output layer 148
overfitting 4, 24, 25

P

PAPI
 URL 216
part-of-speech (POS) 191
pattern analysis 128
Pearson coefficient 104
Pearson correlation coefficient 34
perceptron 18, 146, 147
plan recognition 129
p-norm distance 13
Portable Format for Analytics (PFA) 215
position sensors 167

precision 20
Prediction.IO 216
predictive apriori 32
Predictive Model Markup Language (PMML) 215
Pre-processing phase 42
Principal Component Analysis (PCA) 12, 18, 139
probabilistic classifiers 18

R

ratio data 6
recall 20
Receiver Operating Characteristics (ROC) 21
recommendation engine
 basic concepts 101, 102
 book-recommendation engine,
 building 108
 exploitation 105
 exploration 105
 item-based analysis 103
 key concepts 102, 103
 similarity, calculating 103
 user-based analysis 103
regression
 about 21, 25, 56
 attributes, analyzing 58
 correlation coefficient 23
 data, loading 56, 58
 evaluating 22
 linear regression 22
 mean absolute error 23
 mean squared error 23
 regression model, building 59
 regression model, evaluating 59
 tips 63
regression model
 building 59
 evaluating 59
 linear regression 59, 60
 regression trees 60-62
reinforcement learning 3
Resilient Distributed Dataset (RDD) 37
restricted Boltzmann machines (RBM) 18, 38, 150, 151

Roc curves 20, 21
RuleSetModel 215

S

Sample, Explore, Modify, Model, and Assess (SEMMA) 214
Scale Invariant Feature Transform (SIFT) 146
score function 17
similar items
 searching 13
similarity calculation
 about 103
 collaborative filtering 103
 content-based filtering 104
 hybrid approach 104
SimRank 15
single layer regression model
 building 156-158
Singular value decomposition (SVD) 12
Spark Streaming 36
spatio-temporal patterns 129
Spearman's footrule distance 34
stacked autoencoders 150
standards and markup languages 213
stratification 27
sum transfer function 147
supermarket dataset
 about 94
 shopping patterns, discovering 94
 shopping patterns, discovering with Apriori algorithm 94, 96
 shopping patterns, discovering with FP-growth algorithm 96
supervised learning
 about 2, 17
 classification 17
 regression 21
Support Vector Machine (SVM) 18, 25, 215
survivorship bias 9
suspicious behavior detection 126
suspicious pattern detection 127
suspicious patterns, modeling
 about 132
 dataset rebalancing 134-136
 vanilla approach 133, 134

T

target variables
 appetency probability 69
 churn probability 68
 upselling probability 69
Tertius 32
test set 26
text classification
 about 187
 examples 187
text data
 extracting 190-192
 importing 192
 importing, from directory 192
 importing, from file 193, 194
 pre-processing 194-196
text mining
 about 185, 186
 text classification 187
 topic modeling 186, 187
time series data
 anomaly detection 138
topic modeling 186, 187
topic modeling, for BBC news
 about 196
 BBC dataset, collecting 196, 197
 model, evaluating 200, 201
 modeling 197-200
 model, restoring 203
 model, reusing 202
 model, saving 202
traditional machine learning
 architecture 42
training data
 about 42
 collecting 175-178
train set 26
transaction analysis 129
TreeModel 215

U

UCI machine learning repository
 URL 217
underfitting 24, 25

Universal PMML Plug-in (UPPI) 215
unknown-unknowns 126, 127
unsupervised learning
 about 3, 13
 clustering 16, 17
 similar items, searching 13
user-based analysis 103
user-based collaborative filtering 115-117

V

vanilla approach 133, 134

W

Waikato Environment for Knowledge

Analysis (Weka)

 about 30-33

 URL 30

website traffic

 anomaly detection 137

Weka 3.6

 downloading 46

 URL 46

weka.classifiers package

 weka.classifiers.bayes 32

 weka.classifiers.evaluation 32

 weka.classifiers.functions 32

 weka.classifiers.lazy 32

 weka.classifiers.meta 33

 weka.classifiers.mi 33

 weka.classifiers.rules 33

 weka.classifiers.trees 33

WEKA Packages

 URL 78

word2vec

 about 192

 URL 192

workflow, applied machine learning

 data analysis and modeling 4

 data and problem definition 3

 data collection 4

 data preprocessing 4

 evaluation 4

X

Xiaming Chen

 URL 218

Y

Yahoo traffic dataset

 URL 137