A Notebook on Probability, Statistics, and Data Science

To my family, friends and communities members who have been dedicating to the presentation of this notebook, and to all students, researchers and faculty members who might find this notebook helpful.

Contents

Fo	orewo	ord	xi
P	refac	e	xiii
Li	st of	Figures	$\mathbf{x}\mathbf{v}$
Li	st of	Tables	xvii
Ι	Pr	obability	1
1	Bas	ics	3
_	1.1	Randomness and Stochasticity	3
		1.1.1 Random Experiments	4
		1.1.2 Sample Space	$\overline{4}$
		1.1.3 Events	5
	1.2	Probability	5
		1.2.1 Classical and Empirical Probability	5
		1.2.2 Axioms	6
		1.2.3 Conditional Probability	7
	1.3	Permutation and Combination	7
2	Rar	ndom Variables and Probability Distributions	9
	2.1	Discrete and Continuous Random Variables	9
		2.1.1 Discrete Random Variables	9
		2.1.2 Continuous Random Variables	10
	2.2	Joint Distribution	10
		2.2.1 Joint Probability Density Function	11
		2.2.2 Conditional Distribution	11
		2.2.3 Parameter Estimation with Conditional Distribution .	12
		2.2.4 Geometric Probability	13
	2.3	Probability Density Function of Derived Variables	14
3	Mea	asures of Distribution	17
	3.1	Expectation	17
	3.2	Variance and Standard Deviation	18
	3.3	Moments	19
	3.4	Covariance and Correlation	21

vi	Contents	

		3.4.1 One Variable	21
		3.4.2 Covariance	22
		3.4.3 Correlation	23
	3.5	Important Theorem	24
			24
		3.5.2 Central Limit Theorem	24
4	Cor	nmonly Seen Distributions	25
	4.1	·	25
	4.2		27
			27
		4.2.2 Multivariate Normal Distribution	27
	4.3	Poisson Distribution	28
	4.4	Uniform Distribution	30
	4.5	Cauchy Distribution	30
	4.6		30
		4.6.1 The Γ Distribution	31
		4.6.2 The χ^2 Distribution	32
		4.6.3 The β Distribution	33
	4.7	Student's t -Distribution	34
	4.8	F-Distribution	34
	4.9	Weibull Distribution	34
II	\mathbf{S}	tatistics	35
5	San	npling	37
	5.1		37
	5.2		40
	5.3	Sample Statistics	4.1
		Sample Statistics	41
6	Est		41 43
		imation	43
6 7		imation	
	Reg	imation	43
7	Reg Hyl	imation gression pothesis Testing	43 45
7 8 9	Reg Hyl Bay	imation gression pothesis Testing vesian Methods	43 45 47 49
7 8	Reg Hyl Bay	imation gression pothesis Testing vesian Methods	43 45 47
7 8 9 II	Reg Hy _I Bay	imation gression pothesis Testing vesian Methods Fools	43 45 47 49
7 8 9 II	Reg Hyj Bay I Z	imation gression pothesis Testing vesian Methods Fools	43 45 47 49 51
7 8 9 II	Reg Hyl Bay I 7 R E 10.1	imation gression pothesis Testing vesian Methods Fools R and RStudio Installation	43 45 47 49 51
7 8 9 II	Reg Hyl Bay I 7 R E 10.1	imation gression pothesis Testing vesian Methods Fools R and RStudio Installation R Packages Management	43 45 47 49 51 53
7 8 9 II	Reg Hyl Bay I 7 R E 10.1	imation gression pothesis Testing resian Methods Fools R and RStudio Installation R Packages Management 10.2.1 Manage Packages with Built-in Functions	43 45 47 49 51 53 54 54

Contents	vii
10.3 R Programming Basics	56
10.3.1 Data Types	57
10.3.2 Conditionals and Loops	59
10.3.3 User-Defined Functions	63
10.4 Vector and Matrix	63
10.4.1 Vector	63
10.4.2 Matrix	66
10.4.3 Matrix Visualization Using matplot()	69
10.5 Data Frames	71
10.5.1 Import Data into Data Frame	73
10.5.2 Access Data in Data Frame	73
10.5.3 Filter Data from Data Frame	76
10.5.4 Create Data Frames	77
10.6 Basic Data Visualizations Using qplot()	78
10.7 Advanced Data Visualizations Using ggplot()	79
10.7.1 Grammar of Graphics	81
10.7.2 Data, Aesthetics and Geometries Layers	81
10.7.2 Data, Aesthetics and Geometries Layers	83
	85
10.7.4 Facets Layers	89
10.7.5 Coordinates Layers	
10.7.6 Themes Layers	89
11 R Advanced	93
11.1 Data Preparation	93
11.1.1 Data Type Conversion	93
11.1.2 Handling Missing Data	95
11.2 Connectivity with Data Sources	98
12 R Practice	99
13 Python Basics	101
13.1 NumPy	101
· ·	101
13.2 SciPy	104
13.3 Matplotlib and Seaborn	104
13.3.1 Matplotlib	-
13.3.2 Seaborn	105
13.4 Pandas	107
13.4.1 Data Importing	108
13.4.2 Series and Data Frame	110
14 Python Practice: Artificial Intelligence	111
14.1 Quick Review	112
14.1.1 AI Pipeline	112
14.1.2 Data Preparation and Model Evaluation	113
14.1.3 Commonly Seen ANN Use Cases	114
14.1.4 Computer Vision	114

Contents

	14.1.5	Natural Language Processing	114
14.2	Tensor	Flow	114
	14.2.1	TensorFlow Basics	115
	14.2.2	Classification and Regression	115
	14.2.3	Computer Vision	119
		General Sequential Data Processing	120
	14.2.5	Natural Language Processing	120
		TensorFlow on Different Platforms	120
14.3	PyTore	ch	120
	14.3.1	PyTorch Basics	120
		Classification and Regression	120
	14.3.3	Computer Vision	120
		General Sequential Data Processing	120
		Natural Language Processing	120
		PyTorch on Different Platforms	120
		•	
IV S	eman	tic Web	121
15 Som	antic 1	Web Basics	123
		f Data	124
10.1		Web 1.0 and 2.0	124
		Web 3.0	125
		Semantic Web Vision	126
		Semantic Web Stack	126
		Semantic Web Limitations and Challenges	130
15.2		egy	131
10.2		Philosophy Perspective	131
		Semantic Web Perspective	131
		Ontology Types and Categories	132
15.3			133
10.0	_	Syntax and Semantics	133
		Logic Framework	134
		Logical Expression	136
		Logical Equivalence	137
		Logical Reasoning	138
10 D	-		
		Description Framework	141
		m Resource Identifier	
16.2		rce Description Framework (RDF)	143
		Triple Representation	143
			144
		Lists	146
		Reification	147
		Converting RDB to RDF	147
	10.2.0	Conclusion	148

Contents	ix

16.3 RDF Schema (RDFS)	
16.3.1 RDFS Motivation	149
16.3.2 RDF Versus RDF/RDFS via a	n Example 150
16.3.3 RDFS Expanded Class and Pro	operties 151
16.3.4 Semantics inside RDF/RDFS	
16.4 SPARQL Protocol and RDF Query La	anguage (SPARQL) 153
16.4.1 SPARQL for Basic Query	
16.4.2 SPARQL for Advanced Operat	ions 155
16.4.3 Default Graph and Named Gra	ph 157
16.4.4 SPARQL Programming Return	
16.4.5 Underlying Data Structure of T	
17 Web Ontology Language	161
17.1 RDF/RDFS Limitations	
17.2 OWL Vision	
17.3 OWL Basic Syntax	
17.4 OWL Advanced Syntax	
17.5 Semantic Web with Rules	
18 Semantic Web Practice	171
18.1 Ontological Engineering	
18.2 Ontology Design	
18.2.1 Ontology Management Activity	
18.2.2 Ontology Design Basics	
18.3 Linked Data Engineering	
18.4 Semantic Search	
18.5 Triplestore	
18.6 Example: Semantic Web for Home Ass	
18.6.1 Define Classes Hierarchy	
18.6.2 Define Properties Hierarchy	
18.6.3 Add OWL	
18.6.4 Data Retrieval Examples	
18.7 Reference: Commonly Used Namespac	
Bibliography	179
Diningi apily	118

Foreword

If software or e-books can be made completely open-source, why not a note-book?

This brings me back to the summer of 2009 when I started my third year as a high school student in Harbin No. 3 High School. In around the end of August when the results of Gaokao (National College Entrance Examination of China, annually held in July) are released, people from photocopy shops would start selling notebooks photocopies that they claim to be from the top scorers of the exam. Much curious as I was about what these notebooks look like, never have I expected myself to actually learn anything from them, mainly for the following three reasons.

First of all, some (in fact many) of these notebooks were more difficult to understand than the textbooks. I guess we cannot blame the top scorers for being so smart that they sometimes make things extremely brief or overwhelmingly complicated.

Secondly, why would I want to adapt to notebooks of others when I had my own notebooks which in my opinion should be just as good as theirs.

And lastly, as a student in the top-tier high school myself, I knew that the top scorers of the coming year would probably be a schoolmate or a classmate. Why would I want to pay that much money to a complete stranger in a photocopy shop for my friend's notebook, rather than requesting a copy from him or her directly?

However, things had changed after my becoming an undergraduate student in 2010. There were so many modules and materials to learn in a university, and as an unfortunate result, students were often distracted from digging deeply into a module (For those who were still able to do so, you have my highest respect). The situation became even worse as I started pursuing my Ph.D. in 2014. As I had to focus on specific research areas entirely, I could hardly split much time on other irrelevant but still important and interesting contents.

This motivated me to start reading and taking notebooks for selected books and articles, just to force myself to spent time learning new subjects out of my comfort zone. I used to take hand-written notebooks. My very first notebook was on *Numerical Analysis*, an entrance level module for engineering background graduate students. Till today I still have on my hand dozens of these notebooks. Eventually, one day it suddenly came to me: why not digitalize them, and make them accessible online and open-source, and let everyone read and edit it?

xii Foreword

As most of the open-source software, this notebook (and it applies to the other notebooks in this series as well) does not come with any "warranty" of any kind, meaning that there is no guarantee for the statement and knowledge in this notebook to be absolutely correct as it is not peer reviewed. **Do NOT cite this notebook in your academic research paper or book!** Of course, if you find anything helpful with your research, please trace back to the origin of the citation and double confirm it yourself, then on top of that determine whether or not to use it in your research.

This notebook is suitable as:

- a quick reference guide;
- a brief introduction for beginners of the module;
- a "cheat sheet" for students to prepare for the exam (Don't bring it to the exam unless it is allowed by your lecturer!) or for lecturers to prepare the teaching materials.

This notebook is NOT suitable as:

- a direct research reference;
- a replacement to the textbook;

because as explained the notebook is NOT peer reviewed and it is meant to be simple and easy to read. It is not necessary brief, but all the tedious explanation and derivation, if any, shall be "fold into appendix" and a reader can easily skip those things without any interruption to the reading experience.

Although this notebook is open-source, the reference materials of this notebook, including textbooks, journal papers, conference proceedings, etc., may not be open-source. Very likely many of these reference materials are licensed or copyrighted. Please legitimately access these materials and properly use them.

Some of the figures in this notebook is drawn using Excalidraw, a very interesting tool for machine to emulate hand-writing. The Excalidraw project can be found in GitHub, *excalidraw/excalidraw*.

Preface

This notebook introduces probability, statistics, data science and engineering. They are the "must-have" ability in most, if not all, academic and industrial projects.

In Part I of the notebook, probability theory is introduced. Probability theory studies how likely an event is to occur. It offers rich models and tools to describe random values and stochastic events.

In Part II of the notebook, statistics is introduced. Statistics is a collection of methods to analyze and observe insights from data, verify statistics hypothesis and draw conclusions and predictions.

In Part III of the notebook, commonly used software and toolkits for statistics analysis and data science are introduced. Different from Parts I and II which focus more on theory, Part III focuses more on the tools to solve practical problems. Notice that Artificial Intelligence (AI) has become increasingly popular in recent years for data analysis. There is a separate notebook on introducing AI. In this notebook, the tools to process data using AI is only briefly introduced, again focusing only on the usage of software, not theory.

As a bonus, in Part IV, semantic web, the database framework defined under Web 3.0, is introduced. Semantic web does not necessarily contribute to solving a specific problem using a specific statistics theory. It is rather a concept or organizing and exchanging data efficiently.

Key references of this notebook are summarized as follows. Notice that these materials are so very widely used here and there in the entire notebook that it becomes improbable to address them each time they are used.

For probability and statistics parts:

- Spiegel, Murray, John Schiller, and Alu Srinivasan. Probability and statistics. 2020.
- Dekking, Frederik Michel, et al., A Modern Introduction to Probability and Statistics: Understanding why and how. Vol. 488. London: Springer, 2005.

For data science part:

- Kirill Eremenko, *R Programming A-Z: R For Data Science With Real Exercises*, Udemy Course.
- Lakshmanan, Valliappa, Martin Görner, and Ryan Gillard. Practical Machine Learning for Computer Vision. "O'Reilly Media, Inc.", 2021.

xiv Preface

 $\bullet\,$ Jose Portilla, Complete TensorFlow 2 and Keras Deep Learning Bootcamp, Udemy

Online materials such as tutorials from YouTube, Bilibili, etc., are also used in forming this notebook. ChatGPT-4 is used as a consultant in forming this notebook.

List of Figures

2.1	Sample space of two people arriving at part from 8:00 AM to 9:00 AM	14
3.1	Demonstration of PDF with different skewness	20
3.2	Demonstration of PDF with different excess kurtosis	21
4.1	Poisson distribution with different λ	28
4.2	Poisson distribution Approximation of Binomial Distributions.	29
4.3	Cauchy Distribution	31
4.4	Gamma Distribution	32
4.5	The χ^2 Distribution	33
5.1	Sample with replacement, $N=100,M=500.$	38
5.2	Sample without replacement, $N = 100$, $M = 500$	39
5.3	Sample with replacement, $N = 10000$, $M = 500$	39
5.4	Sample without replacement, $N = 10000$, $M = 500$	40
	RStudio's graphical interface for package management	56
	A demonstration of a matrix in R	67
	A demonstration of using matplot to plot trends	70
	Plot of penalty success rate of the 3 players in 10 matches Plot of average point gained per throw attempt for the 3 players	72
	in 10 matches	72
10.6	A demonstration of qplot	79
	A demonstration of qplot on mortgage price data frame A second demonstration of qplot on mortgage price data	80
	frame.	80
10.9	Multiple layers in chart design	82
	OAn example of box plot of the mortgage price data frame using	
	<pre>ggplot() and geom_boxplot()</pre>	84
10.1	1An example of using geom_smooth() for scatter point fitting.	85
10.12	2An example of histogram plot of house price per unit area in	
	different regions in a single plot	87
10.13	BUse facets to plot the histogram of price per unit are of the	
	house in different regions (subplots in rows)	88
10.14	4Use facets to plot the histogram of price per unit are of the	
	house in different regions (subplots in columns)	88

10.15Add coordinates layer using xlim() and ylim()	90 90 92
 13.1 Plot Fibonacci series as scatter plot. 13.2 Histogram plot using Seaborn. 13.3 Count plot using Seaborn. 13.4 Box plot using Seaborn. The box gives IQR. The bars below and above the box give Q₁ - 1.5 × IQR and Q₃ + 1.5 × IQR, 	105 106 106
respectively. The dots are outliers	107 109
ing pandas	109
15.1 Semantic web stack [2]	127 132
16.1 Semiotic triangle of Apple	142
born in Ulm"	144
place	145
nodes	145
16.5 An example of a container	146
16.6 An example of a collection	147
16.7 An example of RDF reification	148
16.8 Semantic web of a few books, and their authors and publishers	. 149
16.9 Semantic web of fruits and their colors, with RDF implementation in green RDF/RDFS in red	150
17.1 An RDF model that demonstrates "animal eats food"	162
18.1 An example of an RDF model in GraphDB that describes house assets. This is only a demonstration graph and the information inside is artificial and not true.	176

List of Tables

10.1	Commonly used data types	57
10.2	Numerical calculations	59
10.3	Logical comparisons	60
10.4	String operations	60
10.5	Probability related operations	6.
10.6	Statistics related functions	6.
10.7	Commonly used commands for data frame exploration	74
10.8	Commonly used commands for data frame exploration	85
10.9	Functions that fit smooth lines to scatter points	80
15.1	Numerical calculations	138
18.1	Commonly used name spaces in RDF models. URI is neglected since they can be easily found online	17'

Part I Probability

1

Basics

CONTENTS

1.1	Rando	omness and Stochasticity	3
	1.1.1	Random Experiments	4
	1.1.2	Sample Space	4
		Events	
1.2	Proba	bility	5
		Classical and Empirical Probability	
		Axioms	
	1.2.3	Conditional Probability	7
1.3		itation and Combination	

This chapter introduces the basic concepts, axioms, theorems, and fundamental calculations in probability theory.

1.1 Randomness and Stochasticity

People use *random* and *stochastic* to describe a situation or a model whose outcome is not precisely predictable.

By saying randomness, we often refer to the case where the value of a variable (as in "random variable") is not entirely predictable. The value is not known precisely until it is measured. The chance of it taking particular values may follow some pattern which can be described as the statistic property of the variable. An example of a random variable would be the result of tossing a coin. The pattern of the result is clear. It can be either head (X=1) or tail (X=0), each with a 50% chance. It remains unknown until the coin is tossed.

By saying *stochasticity*, we often refer to the case where the outcome of a process (as in "stochastic process") cannot be modeled or predicted precisely. This might be a result of incomplete modeling or random disturbance to the system. As a result, the process becomes non-deterministic, i.e., the output of the system cannot be uniquely determined by the model and the input.

A stochastic process can still be described by a parametric model, but with some unknowns (usually described by random variables) in the equations.

This notebook mainly studies random variables and its impact on data science. Stochastic process and control of stochastic systems are introduced elsewhere in control system related notebooks.

1.1.1 Random Experiments

"Experiment" is one of the most important concepts in science and engineering. Very often, experiments are used to verify a theory. In these cases, experiments are carefully designed so that its outcome is deterministic and predictable as long as the theory is correct. By observing the results of the experiments matching the prediction of the theory, we build confidence in the theory.

However, there are other experiments where we do not have full control over the result due to the lack of information or incomplete modeling. Such experiments include tossing a coin, predicting the GDP of a country next year, etc. These experiments are known as *random experiments*.

1.1.2 Sample Space

A set S that consists of all possible outcomes of a random experiment is called a $sample\ space$.

If a sample space has a finite number of elements, it is called a *finite sample space*. Otherwise, it is called a *infinite sample space*. If the elements in an infinite sample space can be mapped to natural numbers, the sample space is also called a *countably infinite sample space*. Otherwise, the sample space is called a *noncountably infinite sample space*. Examples are given below.

- Finite sample space. Randomly pick 5 balls out of a bag that contains 50 red balls and 50 green balls. The number of red balls in the picked 5 balls forms a finite sample space.
- Countably infinite sample space. The number of products a new machine can manufacture before it is broken forms a countably infinite sample space.
- Noncountably infinite sample space. The precise timestamp someone gets home after work forms a noncountably infinite sample space.

Basics 5

Infinity VS Infinity

Though both infinity, the total number of natural numbers, i.e., the cardinality of \mathbb{N} , is less than the total number of real numbers, i.e., the cardinality of \mathbb{R} .

The cardinality of \mathbb{N} is known as \aleph_0 . It is also the cardinality of all rational numbers. In computability theory, it is also the cardinality of all computable numbers (i.e., numbers that can be computed to any desired precision by a finite terminating algorithm) and computable functions (i.e., algorithms). The total number of real numbers is known as \aleph_1 . It is also the cardinality of all irrational numbers and complex numbers, the number of points on an axis, on an axis interval, or in a high dimensional space \mathbb{R}^n where n is a finite integer.

Even larger infinity quantities, such as $\aleph_2, \aleph_3, \ldots, \aleph_{\omega}, \ldots$ are also defined, though they may not have an intuitive explanation as \aleph_0 and \aleph_1 .

Deeper discussion of this topic requires advanced set theory and is not given in this notebook.

Finite sample space and countably infinite sample space are also known as discrete sample space, whereas noncountably infinite sample space, nondiscrete sample space.

1.1.3 Events

An event is a subset A of the sample space S, i.e., it is a portion of possible outcomes. An event may or may not occur depending on the outcome of an experiment. In the special case where A has only one element, the event is also called an *elementary event*, or just a *sample*. In the special case where A = S, the event is also called a certain event.

1.2 Probability

Given an experiment and an event, there is always uncertainty as to whether the event will occur, and *probability* is a measurement of likelihood that the event is going to occur. For example, by saying a probability of 100% or 1, we mean that we are certain that the event would occur.

1.2.1 Classical and Empirical Probability

There are different ways to calculate or estimate the probability of an event. In the *classical approach* where we know the total number of outcomes n, all of which have a equal chance of happening, and there are h ways for the event to

occur, then the probability of the event is h/p. In the frequency approach, we can repeat the experiment n times where n is a large number, and record the number of instants where the event happened as h. The empirical probability of the event is obtained by calculating h/p which should converge to the true probability as n goes large.

1.2.2 Axioms

Both the classical approach and the frequency approach have some draw-backs. It is often difficult to define "equal chance" and "large number" in the two approaches respectively. Therefore, *axiomatic approach* is introduced as follows.

Let the sample space be denoted by S, and events by A_i . For simplicity, assume that S is discrete, or at least part of S associated by the events is discrete. Define $P(\cdot)$ as the probability function and P(A) the probability of an event A subject to following axioms:

- 1. For every event A_i , $P(A_i) \ge 0$.
- 2. For the certain event S, P(S) = 1.
- 3. For mutually exclusive events A_1 and A_2 , $P(A_1 \cup A_2) = P(A_1) + P(A_2)$.

Using the above axioms, a bunch of well-known theorems can be derived, such as

- If $A_1 \subseteq A_2$, $P(A_2 A_1) = P(A_2) P(A_1)$.
- For every event A_i , $0 \le P(A_i) \le 1$.
- For the impossible event \emptyset , $P(\emptyset) = 0$.
- For the complement of an event A', P(A') = 1 P(A).
- For mutually exclusive events A_i , i=1,...,n, if $A=\bigcup_{i=1}^n A_i$, $P(A)=\sum_{i=1}^n P(A_i)$.
- For two events A_1 and A_2 , $P(A_1 \cup A_2) = P(A_1) + P(A_2) P(A_1 \cap A_2)$.

With the above axioms, we can revisit classical probability as follows.

Assume that a discrete and finite sample space S consists of the following elementary events (elementary events are always mutually exclusive) A_i , i = 1, ..., n, i.e.,

$$S = \bigcup_{i=1}^{n} A_i$$

and assume equal probabilities for all the elementary events, i.e., the probability of each event is given by

$$P(A_i) = \frac{1}{n}, i = 1, ..., n$$

Basics 7

Define an event A that is made up of h such elementary events out of A_i . The probability of A can then be calculated by

$$P(A) = \frac{h}{n}$$

where h, n are nothing but the cardinality of A and S with respect to the elementary events.

1.2.3 Conditional Probability

Assume two events A and B. The conditional probability P(B|A) describes the probability of B given that A has occurred. The definition is given below. Think of this definition as S replaced by A since A has been confirmed occurred.

$$P(B|A) \equiv \frac{P(A \cap B)}{P(A)}$$

Or equivalently,

$$P(A \cap B) = P(A)P(B|A)$$

From the above,

$$P(B|A) = \frac{P(B)P(A|B)}{P(A)}$$

which is known as the Bayes' rule.

In the special case where P(B|A) = P(B), or equivalently $P(A \cap B) = P(A)P(B)$, the two events A and B are said to be *independent events*.

1.3 Permutation and Combination

In classical probability, calculating the cardinalities of event A and sample space S is the key to solving the problem. Permutations and combinations are widely used for such calculations.

Suppose that there are n distinct objects, and we would like to select r objects from them, and arrange them in a particular order. The permutation

$$nPr = n(n-1)(n-2)\dots(n-r+1)$$
 (1.1)

gives the number of possible outcomes. In the special case r = n,

$$nPn = n(n-1)(n-2)\dots 1$$

where $n(n-1)(n-2)\dots 1$ is often denoted as n!. Therefore, (1.1) becomes

$$nPr = \frac{n!}{(n-r)!}$$

In permutation, the arrangement of the selected r objects matters. When the order does not matter, combination

$$nCr = \frac{nPr}{r!}$$
$$= \frac{n!}{r!(n-r)!}$$

is used to calculate the total number of outcomes. Notice that nCr is also denoted by $\binom{n}{r}$.

Random Variables and Probability Distributions

CONTENTS

2.1	Discre	te and Continuous Random Variables	9
	2.1.1	Discrete Random Variables	9
	2.1.2	Continuous Random Variables	10
2.2	Joint I	Distribution	10
	2.2.1	Joint Probability Density Function	10
	2.2.2	Conditional Distribution	11
	2.2.3	Parameter Estimation with Conditional Distribution	12
	2.2.4	Geometric Probability	13
2.3	Probal	bility Density Function of Derived Variables	13

The definition of random variable has been introduced in earlier chapter. This chapter digs deeper into the different types and properties of random variables, and how we describe a random variable.

2.1 Discrete and Continuous Random Variables

A random variable may be discrete or continuous depending on the sample space of the variable.

2.1.1 Discrete Random Variables

If a random variable X takes only discrete values $x_1, x_2, ...$, it is called a discrete random variable. The probability of X taking a particular value x is denoted by P(X = x). Sometimes for simplicity, it is simply denoted by function f(x) = P(X = x) when there is no ambiguities. In this case, P(X = x) and f(x) are called the *probability function* (also known as *probability mass function*) of X.

Furthermore, define $P(X \leq x)$ or F(x) as the *cumulative distribution* function. It is easy to prove that F(x) is nondecreasing, and $\lim_{x\to 0} F(x) = 0$,

 $\lim_{x\to\infty} F(x) = 1$. Also, F(x) "jumps" at each $P(X = x_k) > 0$ and it is continuous from the right.

2.1.2 Continuous Random Variables

A random variable X may also take continuous values in many applications. For example, let X denote the time consumption to finish a task, which can be any value above 0 hours.

In this case, the chance for X to take a precise value, say for the student to finish his homework using precisely 25 minutes 13 seconds and 750 milliseconds, is very small (in fact, mathematically zero). The probability function P(X=x) becomes useless. The cumulative distribution function $F(x) = P(X \le x)$ still makes sense, as it essentially calculates the probability of X given by not a precise value but within a range.

Inspired by this, define probability density function (PDF) or f(x) for continuous random variable as follows. The probability density function f(x) is such that

$$P(a < X < b) = \int_{a}^{b} f(x)dx$$

Therefore,

$$F(x) = P(X \le x)$$
$$= \int_{-\infty}^{x} f(\epsilon) d\epsilon$$

Notice that f(x) itself is not probability. It is f(x)dx accumulating in range $x \in (a, b)$ that forms the probability, hence the name "probability density".

In science and engineering problems, continuous random variables are more common than discrete random variables. Therefore, the main focus of this notebook from this point on would be continuous signals. Discrete random variables can sometimes be taken as a special case of continuous random variables with impulse PDF.

2.2 Joint Distribution

Joint distribution is used to study the correlation of multiple random variables. It is widely used in system identification where insights of unknown variables are obtained by observing the known variables, assuming their correlation is known.

2.2.1 Joint Probability Density Function

Consider two random variables X and Y. In the case of discrete variables, define joint probability function and joint cumulative distribution function as follows.

$$\begin{array}{rcl} f(x,y) & = & P\left(X=x,Y=y\right) \\ F(x,y) & = & P\left(X\leq x,Y\leq y\right) \\ & = & \sum_{u\leq x} \sum_{v\leq y} f(u,v) \end{array}$$

In the case of continuous random variables, joint PDF of X and Y is such that

$$\int_{x=a}^{b} \int_{y=c}^{d} f(x,y) dx dy = P(a < X < b, c < y < d)$$
 (2.1)

Integrate (2.1) w.r.t an axis gives the PDF of the other variable. For example,

$$\int_{y=-\infty}^{\infty} f(x,y) = f_X(x)$$
 (2.2)

which becomes a function of a single variable, x. Here $f_X(x)$ is nothing but the PDF of X, completely ignoring the other variable Y.

We can define cumulative distribution function for joint distribution likewise as follows.

$$\begin{split} F(x,y) &=& P(X \leq x, Y \leq y) \\ &=& \int_{u=-\infty}^{x} \int_{v=-\infty}^{y} f(u,v) du dv \\ F_X(x) &=& P(X \leq x) \\ &=& \int_{u=-\infty}^{x} \int_{v=-\infty}^{\infty} f(u,v) du dv \end{split}$$

2.2.2 Conditional Distribution

As introduced earlier, (2.2) gives the PDF of X by completely ignoring Y. In other words, it is the best guess of X (in the form of PDF) when there is no measurement of Y.

Conditional distribution, on the other hand, try to find the best guess of X when Y is measured. For example, from (2.1), let Y be measured by a fixed value y, and we would like to calculate $f_{X|Y}(x|Y=y)$, i.e., the PDF of X given the known information Y=y. In many literatures, $f_{X|Y}(x|Y=y)$ is denoted by $f_{X|Y}(x|y)$ for simplicity.

The conditional PDF $f_{X|Y}(x|y)$ can be obtained as follows. It is essentially Bayes' rule applied on continuous variables.

$$f_{X|Y}(x|y) = \frac{f(x,y)}{f_Y(y)} \tag{2.3}$$

Equation (2.3) is given as an analytical equation of both x and y. Just substitute the measured value of y into (2.3) to get a function of x alone. Since (2.3) itself is a PDF of X, its integration w.r.t x is certainly 1. This can be easily verified as follows.

$$\int_{x=-\infty}^{\infty} f_{X|Y}(x|y)dx = \int_{x=-\infty}^{\infty} \frac{f(x,y)}{f_Y(y)}dx$$
$$= \frac{\int_{x=-\infty}^{\infty} f(x,y)dx}{f_Y(y)}$$
$$= \frac{f_Y(y)}{f_Y(y)}$$
$$= 1$$

for any y. Notice that given a particular value of y, $f_Y(y)$ is a constant value independent of x, hence can be taken out from the integration in the intermediate derivation.

If X and Y are independent variables, $f(x,y) = f_X(x)f_Y(y)$. In this case, (2.3) becomes

$$f_{X|Y}(x|y) = f_X(x)$$

which implies that the information of Y = y does not affect our understanding of X, just as if the information is completely ignored.

2.2.3 Parameter Estimation with Conditional Distribution

Equation (2.3) is widely used in parameter estimation. Let θ be the parameter(s) to be estimated, and x the measurement(s) that reflects θ via measurement model $f(\theta, x)$ which is given in the form of joint distribution of θ and X. The model is known as priori knowledge. Notice that parameters and measurements are correlated via joint PDF but not deterministic due to measurement noise and other uncertain factors.

Given measurement(s) X = x, the posteriori estimation of θ (in the form of PDF) can be obtained as follows. From (2.3),

$$f_{\theta|X}(\theta|x) = \frac{f_{X|\theta}(x|\theta)f_{\theta}(\theta)}{f_{X}(x)}$$
 (2.4)

where on the right side $f_{X|\theta}(x|\theta)$ is known as the *likelihood function* that describes the likelihood of measuring X = x if the actual parameter(s) is θ ,

 $f_{\theta}(\theta)$ the *priori* knowledge of the distribution of θ , and finally $f_X(x)$ which is known as *evidence*, in this case a constant known value given X = x.

Equation (2.4) can be memorized as follows.

$$posteriori = \frac{likelihood \times priori}{evidence}$$

An estimation of θ can then be obtained by simply calculate its expectation as follows

$$E(\theta) = \int_{-\infty}^{\infty} f_{\theta|X}(\theta|x)d\theta$$

Terms *priori* and *prior*, *posteriori* and *posterior* can be used interchangeably, respectively.

2.2.4 Geometric Probability

Geometric probability is an extension of classical probability to continuous random variables. Similar with classical probability where it is prerequisite that all elementary events share the same probability, in geometric probability it is assumed that all random variables (usually 2 or 3 variables in total) are uncorrelated and uniformly distributed. The probability of an event happening then can be obtained by measuring the area or volume of the event.

Theory wise, geometric probability is just one of the many joint PDF applications that is simple, intuitive and useful. An example is used to illustrate the use of geometric probability.

Geometric Probability Example

Consider two people trying to meet up at the park, but they forgot to tell each other the time to meet. Both of them will arrive at the park at anytime between 8:00 AM and 9:00 AM with equal chance, and wait for the other for 15 minutes or until 9:00 AM, whichever is earlier, and then will leave the park if the other does not show up.

Calculate the chance of the two people successfully meeting up.

Draw the sample space in a 2-D plot as shown in Fig. 2.1, where the x-axis and y-axis are the arriving time of the two people, respectively. The shaded area represents the samples where the two would meet up successfully. Divide the shared area by the total sample space area to get P = 7/16 = 43.75% which is the probability that the two would meet up successfully.

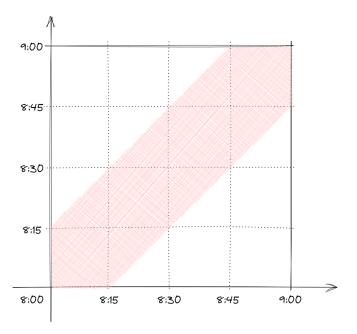


FIGURE 2.1 Sample space of two people arriving at part from 8:00 AM to 9:00 AM.

2.3 Probability Density Function of Derived Variables

Let X be a random variable with PDF $f_X(x)$. Let U be another random variable which is a function of X via $U = \phi(X)$. The PDF of U can be calculated as follows. For simplicity, assume that $U = \phi(X)$ is a injective function (one-to-one function), and $X = \phi^{-1}(U) = \psi(U)$. In that case, the PDF of U, g(u), can be obtained as follows.

$$g(u) = |\psi'(u)| f(\psi(u))$$

For example, let U = aX, $X = \frac{U}{a}$.

$$g(u) = \frac{1}{a} f\left(\frac{u}{a}\right)$$

Let X, Y be two random variables with joint distribution f(x,y). Let U = X + Y. The PDF of U can be obtained as follows.

$$g(u) = \int_{-\infty}^{\infty} f(x, u - x) dx$$
 (2.5)

In the special case where X and Y are independent, $f(x,y) = f_X(x)f_Y(y)$,

and (2.5) becomes

$$g(u) = \int_{-\infty}^{\infty} f_X(x) f_Y(u - x) dx$$
$$= f_X * f_Y$$

where * denotes the convolution operator.

Measures of Distribution

CONTENTS

3.1	Expect	ation	17
3.2		ce and Standard Deviation	
3.3	Momen	its	19
3.4	Covaria	ance and Correlation	20
	3.4.1	One Variable	21
	3.4.2	Covariance	21
	3.4.3	Correlation	23
3.5		ant Theorem	
	3.5.1	Law of Large Numbers	24
	3.5.2	Central Limit Theorem	24

Given the probability or probability density of a random variable, a lot of information can be abstracted such as the expectation which describes the average value of that random variable if the experiment is repeated many times.

Commonly seen measures of a distribution are introduced in this chapter.

3.1 Expectation

In probability and statistics sense, expectation (also known as mean) describes the average value of a random variable, if the variable is generated many times. For discrete random variable X, the expectation is given below.

$$E(X) = \sum_{i=1}^{n} x_i P(x_i)$$
(3.1)

where $E(\cdot)$ is used to denote the expectation, and n the cardinality of the sample space. For continuous random variable, it is

$$E(X) = \int_{-\infty}^{\infty} x f(x) dx \tag{3.2}$$

Expectation is sometimes denoted as μ in literatures.

Some features of expectation calculation are given below.

$$E(cX) = cE(X)$$

$$E(X+Y) = E(X) + E(Y)$$

where c is a constant and X, Y are two random variables. These features can be easily derived from (3.1) and (3.2). Furthermore, if X, Y are independent, recall $f(x,y) = f_X(x)f_Y(y)$,

$$\begin{split} \mathbf{E}(XY) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy f(x,y) dx dy \\ &= \int_{-\infty}^{\infty} x f_X(x) dx \times \int_{-\infty}^{\infty} y f_Y(y) dy \\ &= \mathbf{E}(X) \mathbf{E}(Y) \end{split}$$

3.2 Variance and Standard Deviation

Variance and standard deviation describe how spread samples are from its expectation. It is defined as follows.

$$Var(X) = E\left((X - E(X))^{2}\right)$$

$$= E\left(X^{2} - 2XE(X) + E(X)^{2}\right)$$

$$= E\left(X^{2}\right) - E(X)^{2}$$

$$\sigma_{X} = \sqrt{Var(X)}$$
(3.3)
$$(3.4)$$

Variance and standard deviation are sometimes denoted as σ^2 and σ respectively.

For continuous random variables, from (3.4) the variance is

$$Var(X) = \int_{-\infty}^{\infty} (x - E(X))^2 f(x) dx$$

Think of X as an estimation of a parameter θ . If the estimation is non biased, $E(X) = \theta$. Therefore, from the above equation, for a non-biased estimation, the variance of the estimates (left) equals to the MSE of the estimate (right) with sample numbers approaching infinity.

Some features of variance calculation are given below.

$$Var(cX) = c^2 Var(X)$$

For independent random variables X and Y,

$$Var(X \pm Y) = Var(X) + Var(Y)$$

Mean and standard deviation can be used to standardize a random variable as follows.

$$X^* = \frac{X - \mu}{\sigma}$$

where μ , σ are the mean and standard deviation of random variable X respectively. The standardized random variable, X^* , has a mean of 0 and standard deviation of 1.

3.3 Moments

In mathematics, the r-th moment of a continuous function f(x) about c is defined as follows.

$$\mu_n = \int_{-\infty}^{\infty} (x-c)^n f(x) dx$$

By simply saying *moment* without specifying c, it is by default that c=0. Let f(x) be a PDF. In this sense, the 0-th order and the 1-st order moment of a probability density function can be calculated as follows.

$$\mu_0 = \int_{-\infty}^{\infty} f(x)dx = 1$$

$$\mu_1 = \int_{-\infty}^{\infty} x f(x)dx = E(X)$$

where it can be seen that the 0-th and 1-st moment of a PDF is 1 and its mean, respectively.

Further more, let $c = \mu_1$ be the mean of the random variable to calculate the second *central moment* μ_2 as follows.

$$\mu_2 = \int_{-\infty}^{\infty} (x - \mu_1)^2 f(x) dx = \operatorname{Var}(X)$$

which is the variance of the random variable.

Using mean and variance to further define standardized moments as shown below. Define $\bar{\mu}_k$, μ_k and σ_k for $k \geq 3$ as follows.

$$\bar{\mu}_k = \frac{\mu_k}{\sigma^k}$$

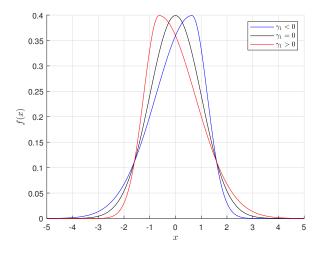


FIGURE 3.1

Demonstration of PDF with different skewness.

where

$$\mu_k = \mathbf{E}\left((X - \mu_1)^k\right) = \int_{-\infty}^{\infty} (x - \mu_1)^k f(x) dx$$
$$\sigma^k = (\mu_2)^{\frac{k}{2}} = \left(\int_{-\infty}^{\infty} (x - \mu_1)^2 f(x) dx\right)^{\frac{k}{2}}$$

with μ_1 and μ_2 used to denote mean (1-st order moment) and variance (2-nd order central moment) of the random variable, respectively.

The 3-rd and 4-th order standardized moments are known as the *skewness* and *kurtosis* of the PDF, respectively. In some literatures, skewness is denoted by $\gamma_1 = \bar{\mu}_3$, and kurtosis $\gamma_2 = \bar{\mu}_4$.

The skewness γ_1 is a measure of asymmetry of the PDF. When $\gamma_1 > 0$ or positive skew, the distribution has a long tail on the right side of the PDF. When $\gamma_1 < 0$ or negative skew, the distribution has a long tail on the left side. When $\gamma_1 = 0$, the PDF might be symmetric (but not necessarily so). Examples are given in Fig. 3.1.

The kurtosis γ_2 measures the "tailedness" of a probability distribution, i.e., whether the PDF has heavy tail or thin tail. The normal distribution, which has $\gamma_2 = 3$, is often used as a benchmark. Excess kurtosis is kurtosis subtracting 3, making the normal distribution having the excess kurtosis of 0. A positive excess kurtosis would mean a "heavier" tail than the normal distribution. Examples are given in Fig. 3.2.

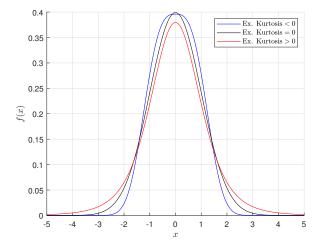


FIGURE 3.2 Demonstration of PDF with different excess kurtosis.

3.4 Covariance and Correlation

Covariance and correlation are defined for a joint distribution with multiple random variables. For simplicity, consider only two random variables X, Y whose joint distribution is given by $f_{XY}(x,y)$. The idea derived from here can be generated to more variables.

3.4.1 One Variable

The PDF of one variable can be derived from the joint distribution using (2.2). It is straight forward to get the expectation and variance for that variable as follows.

$$E(X) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x f(x, y) dx dy$$

$$Var(X) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - E(X))^{2} f(x, y) dx dy$$

3.4.2 Covariance

The covariance of two variables is defined and calculated as follows.

$$Cov(X,Y) = E((x - E(X))(y - E(Y)))$$
(3.5)

$$= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - E(X))(y - E(Y))f(x, y)dxdy \qquad (3.6)$$

where Cov(X, Y) is sometimes denoted by σ_{XY} . Notice that unlike variance that is always positive, covariance can be zero or negative. If X, Y are independent variables, $f(x, y) = f_X(x)f_Y(y)$. From (3.6)

$$Cov(X,Y) = \int_{-\infty}^{\infty} (x - E(X)) f_X(x) dx \times \int_{-\infty}^{\infty} (y - E(Y)) f_Y(y) dy$$
$$= 0$$

If the covariance of the two variables is zero, the two variables are called *uncorrelated*. Independent variables are always uncorrelated. However, uncorrelated variables are not necessarily independent.

Furthermore,

$$Cov(X, Y)^2 \le Var(X)Var(Y)$$

The proof is as follows. Notice that lemma (3.8) is used in the proof. Using (3.8) on (3.5), (3.3) gives

$$Cov(X,Y)^{2} = E((x - E(X))(y - E(Y)))^{2}$$

$$\leq E((x - E(X)))^{2} E((y - E(Y)))^{2}$$

$$= Var(X)Var(Y)$$

or equivalently

$$\sigma_{XY}^2 \quad \leq \quad \sigma_X^2 \sigma_Y^2$$

Noticing that while σ_X , σ_Y are always nonnegative while σ_{XY} is not,

$$-\sigma_X \sigma_Y \le \sigma_{XY} \le \sigma_X \sigma_Y \tag{3.7}$$

Lemma

For two random variables X and Y,

$$E(XY)^2 \le E(X^2)E(Y^2) \tag{3.8}$$

Proof

$$0 \leq \mathbb{E}\left(\left(X - Y \frac{\mathbb{E}(XY)}{\mathbb{E}(Y^2)}\right)^2\right)$$

$$= \mathbb{E}\left(X^2 - 2XY \frac{\mathbb{E}(XY)}{\mathbb{E}(Y^2)} + \mathbb{E}(Y^2) \frac{\mathbb{E}(XY)^2}{\mathbb{E}(Y^2)^2}\right)$$

$$= \mathbb{E}(X^2) - 2\mathbb{E}(XY) \frac{\mathbb{E}(XY)}{\mathbb{E}(Y^2)} + \mathbb{E}(Y^2) \frac{\mathbb{E}(XY)^2}{\mathbb{E}(Y^2)^2}$$

$$= \mathbb{E}(X^2) - 2 \frac{\mathbb{E}(XY)^2}{\mathbb{E}(Y^2)} + \frac{\mathbb{E}(XY)^2}{\mathbb{E}(Y^2)}$$

$$= \mathbb{E}(X^2) - \frac{\mathbb{E}(XY)^2}{\mathbb{E}(Y^2)}$$

Therefore

$$\frac{\mathrm{E}(XY)^2}{\mathrm{E}(Y^2)} \leq \mathrm{E}(X^2)$$
$$\mathrm{E}(XY)^2 \leq \mathrm{E}(X^2)\mathrm{E}(Y^2)$$

3.4.3 Correlation

The *correlation* of two variables is defined and calculated as follows.

$$\rho = \frac{\operatorname{Cov}(X, Y)}{\sqrt{\operatorname{Var}(X)}\sqrt{\operatorname{Var}(Y)}}$$

or simply

$$\rho = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

From (3.7), apparently $-1 \le \rho \le 1$. When two variables are uncorrelated or independent, $\rho=0$. If $\rho=1$, it is called *perfect positive correlation*. This happens usually because the two variables are positively linearly depended, for example, X=2Y or X=Y+1. If $\rho=-1$, it is called *perfect negative correlation*, and the similar idea applies.

3.5 Important Theorem

There are a few important theorems frequently used in the study of probability and statistics. They are introduced here.

3.5.1 Law of Large Numbers

The Law of Large Numbers (LLN) is a theorem that basically says if performing the same experiment a large number of times, the average of the outcome of the experiments should eventually converges to an expected value. The larger the number of times the experiment is repeated, the closer the average to the expected value.

In mathematics form, let X be a random variable which represents the outcome of an experiment. Let X_i be a sample of the outcome. According to LLN,

$$\lim_{n \to \infty} \sum_{i=1}^{n} \frac{X_i}{n} = \bar{X}$$

3.5.2 Central Limit Theorem

Central Limit Theorem (CLT) states the following observation. For independent and identically distributed (i.i.d) random variables not necessarily normal distribution, the sampling distribution of the standardized sample mean tends towards the standard normal distribution.

Let X be a random variable not necessarily normal distribution, and it has mean and variance of μ and $\sigma^2 < \infty$. Let X_i be samples of the random variable. The standardized sample mean of the n samples are calculated as follows.

standardized sample mean
$$=\frac{\bar{X}_n - \mu}{\sigma_{\bar{X}_n}}$$

where

$$\bar{X}_n = \sum_{i=1}^n X_i$$

$$\sigma_{\bar{X}_n} = \frac{\sigma}{\sqrt{n}}$$

For a large n, the distribution of the standardized sample mean should look like normal distribution.

Commonly Seen Distributions

CONTENTS

4.1	Bernoulli and Binomial Distribution	25
4.2	Normal Distribution	27
	4.2.1 Single Normal Distribution	27
	4.2.2 Multivariate Normal Distribution	27
4.3	Poisson Distribution	28
4.4	Uniform Distribution	29
4.5	Cauchy Distribution	30
4.6	Gamma, Chi-Squared and Beta Distribution	30
	4.6.1 The Γ Distribution	30
	4.6.2 The χ^2 Distribution	32
	4.6.3 The β Distribution	33
4.7	Student's t-Distribution	33
4.8	F-Distribution	34
4.9	Weibull Distribution	34

Commonly seen distributions, both discrete and continuous, are introduced here. Some of them are very useful in statistics analysis and are introduced in more details in later part of the notebook.

4.1 Bernoulli and Binomial Distribution

Bernoulli distribution is a discrete probability distribution of random variable X that can take only 2 values, 0 or 1. The probability of X taking 1 is denoted by p, while 0 is q = 1 - p as shown below.

$$f(x) = P(X = x) = \begin{cases} p & x = 1\\ q & x = 0 \end{cases}$$

subject to $0 \le p \le 1, \ 0 \le q \le 1$ and p+q=1. Each test is also called a Bernoulli trail.

The expectation, variance, skewness and excess kurtosis of the distribution are $p,\ pq,\ \frac{q-p}{\sqrt{pq}}$ and $\frac{1-6pq}{pq}$, respectively.

It is assumed that Bernoulli trails are tested repeatedly. Each bernoulli trail has a probability of p to take value 1, and q=1-p to take value 0. The test is carried out n times. The number of the tests returning 1 is a random variable $0 \le X \le n$. The discrete distribution of X taking values $0, \ldots, n$ follows binomial distribution, whose probability mass function is given by

$$f(x) = P(X = x)$$

$$= {n \choose x} p^{x} (1-p)^{n-x}$$

$$= \frac{n!}{x! (n-x)!} p^{x} (1-p)^{n-x}$$
(4.1)

Bernoulli distribution can be taken as a special case of binomial distribution with n=1 at x=1. The expectation, variance, skewness and excess kurtosis of the distribution are np, npq, $\frac{q-p}{\sqrt{npq}}$ and $\frac{1-6pq}{npq}$, respectively.

According to central limit theorem, when n is large, binomial distribution can be approximated by normal distribution.

Binomial distribution can be extended to multinomial distribution, where instead of a single event A happening with probability p or not happening with probability q, s.t. p + q = 1, consider multiple events A_1 , A_2 , ..., each with probability p_1 , p_2 , ..., respectively, s.t. $p_1 + p_2 + \ldots + p_m = 1$. Consider a total of n tests. The number of event A_1 happening is a random variable X_1 , event A_2 , X_2 , and so on. Multinomial distribution studies the probability of

$$P(X_1 = x_1, X_2 = x_2, \dots, X_m = p_m) = \frac{n}{x_1! x_2! \dots x_m!} p_1^{x_1} p_2^{x_2} \dots p_m^{x_m}$$
s.t.
$$\sum x_i = n$$

Do distinguish binomial distribution with hypergeometric distribution. Binomial distribution can be used to model "sampling with replacement" process. Consider picking up a marbles from a bag containing mixture of red and blue marbles whose numbers are given by r and b respectively. Repeat the test n times. After each test, put the marble back to the bag. The number of blue marbles collected is a random variable X that follows

$$f(x) = P(X = x) = \binom{n}{x} \left(\frac{b}{b+r}\right)^x \left(\frac{r}{b+r}\right)^{n-x} = \binom{n}{x} \frac{b^x r^{n-x}}{(b+r)^n}$$

which is a binomial distribution. On the other hand, if after each test the marble is not returned to the bag, it becomes a hypergeometric distribution

and the probability follows

$$f(x) = P(X = x) = \frac{\binom{b}{x} \binom{r}{n-x}}{\binom{b+r}{n}}$$

with $\max(0, n - r) \le x \le \min(n, b)$.

When b, r are far larger than n, the hypergeometric distribution can be approximated by binomial distribution. When b and r approaches infinity with constant ratio b/(b+r) and r/(b+r), hypergeometric distribution converges to binomial distribution.

4.2 Normal Distribution

Normal distribution, also known as Gaussian distribution, is a continuous distribution, and it is one of the most widely used assumption of random noise. According to central limit theorem, many aggregated values can be characterized by normal distribution. We will discuss single normal distribution first, followed by joint normal distribution.

4.2.1 Single Normal Distribution

The PDF of the normal distribution is given by

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}}e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

where μ , σ^2 are the mean and variance of the distribution, respectively. The skewness and excess kurtosis of normal distribution are zero.

Let X be a random variable following normal distribution with mean μ and variance σ^2 . This can be denoted by $X \sim \mathcal{N}(\mu, \sigma^2)$. Let $Z = \frac{X - \mu}{\sigma}$, and Z would be a standard normal distribution with mean 0 and variance 1.

For a random variable X following normal distribution, the probabilities of its value falling between $\pm \sigma$, $\pm 2\sigma$ and $\pm 3\sigma$ are 0.6827, 0.9545 and 0.9973 respectively. In statistics, sometimes we will take samples with residuals larger than $|3\sigma|$ as outliers.

4.2.2 Multivariate Normal Distribution

Multivariate normal distribution describes a vector of random variables $X = [X_1, \ldots, X_m]$ that follows joint normal distribution. The joint PDF is given

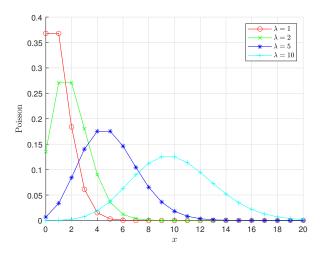


FIGURE 4.1

Poisson distribution with different λ .

by

$$f(x) = \frac{1}{\sqrt{(2\pi)^m |\Sigma|}} e^{-\frac{(x-\mu)^T \Sigma^{-1} (x-\mu)}{2}}$$

where $\mu \in \mathbb{R}^m$ is the mean of x and $\Sigma \in \mathbb{R}^{m \times m}$ the covariance matrix. The off diagonal elements in Σ describes the correlation of elements in X. The $|\Sigma|$ is the determinant of Σ .

4.3 Poisson Distribution

Poisson distribution is a discrete probability distribution that takes non-negative integer values 0, 1, 2,

The probability function of poisson distribution is given by

$$f(x) = P(X = x) = \frac{\lambda^x e^{-\lambda}}{x!}$$

where $\lambda > 0$ is a constant that describes the shape of the probability function, as shown in Fig. 4.1.

The sum of two variables following Poisson distributions with $\lambda = \lambda_1$ and $\lambda = \lambda_2$ follows Poisson distribution with $\lambda = \lambda_1 + \lambda_2$.

Poisson distribution is often used to describe the probability of an event

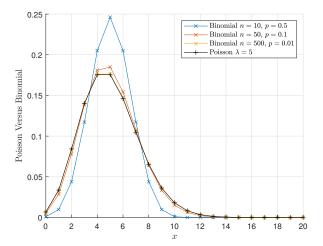


FIGURE 4.2
Poisson distribution Approximation of Binomial Distributions.

happening in the given window, for example, the number of times a machine fails in a given length of duration (say one year), assuming each failure is independent from another. Binomial distribution probability function converges to Poisson distribution when $n \to \infty$. In practice, when n is large and p is small so that np is decent, binomial distribution can be approximated by Poisson distribution with $\lambda = np$. A demonstration is given in Fig. 4.2.

From the demonstration given by Fig. 4.2, we can interpret Poisson distribution as follows. Assume that a test is carried out many times (n times) individually and independently. During each test, an event may happen with a small chance (p probability). The total number of that event happening has an expectation of $\lambda = np$, and it follows Poisson distribution when $n \to \infty$.

Poisson distribution can be used model things like the failure rates of a machine in its lifespan. The machines runs for many hours which correspond with n. In each hour of its running, it has a failure rate which correspond with p. It can also model visitors to a shopping mall in a day. Everyone living in the city may visit the shopping mall. The population of the city is correspond with n, and the chance for a person to visit the mall, p.

The expectation, variance, skewness and excess kurtosis of the distribution are λ , λ , $\frac{1}{\sqrt{\lambda}}$ and $\frac{1}{\lambda}$, respectively.

4.4 Uniform Distribution

Uniform distribution can be divided into continuous distribution and discrete uniform distribution. Under the scope of this discussion, continuous uniform distribution is studied.

The PDF of continuous uniform distribution is given by

$$f(x) = \begin{cases} \frac{1}{b-a} & a \le x \le b \\ 0 & \text{otherwise} \end{cases}$$

The expectation, variance, skewness and excess kurtosis of the distribution are $\frac{a+b}{2}$, $\frac{(b-a)^2}{12}$, 0 and $-\frac{6}{5}$, respectively.

4.5 Cauchy Distribution

Cauchy distribution, named after Augustin Cauchy, is a continuous distribution. The PDF is given by

$$f(x) = \frac{1}{\pi \gamma \left(1 + \left(\frac{x - x_0}{\gamma}\right)^2\right)}$$
$$= \frac{1}{\pi} \frac{\gamma}{(x - x_0)^2 + \gamma^2}$$

and it is plotted in Fig. 4.3.

From mathematics perspective, Cauchy distribution describes the ratio of two independent normal distributed random variables with mean zero. It is a useful distribution in physics.

4.6 Gamma, Chi-Squared and Beta Distribution

The Γ , χ^2 and β distributions are introduced as follows.

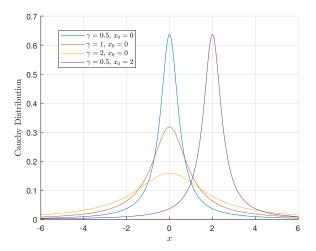


FIGURE 4.3

Cauchy Distribution.

4.6.1 The Γ Distribution

Gamma distribution is a continuous distribution with the following PDF

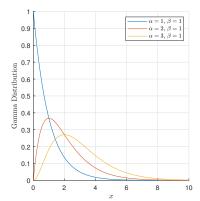
$$f(x) = \begin{cases} \frac{\beta^{\alpha}}{\Gamma(\alpha)} x^{k-1} e^{-\beta x} & x > 0 \\ 0 & \text{otherwise} \end{cases}$$
 (4.2)

where $\alpha > 0$ and $\beta > 0$ are the shape and rate parameters respectively (some literatures uses scale parameter $\theta = 1/\beta$ in the equation), and $\Gamma(\cdot)$ denotes the Gamma function

$$\Gamma(\alpha) = \int_0^\infty t^\alpha - 1e^{-t}dt$$

for $\alpha > 0$. Gamma distribution is often used to model the interval of two events in the Poisson process.

The plot of the PDF of Gamma distribution is given in Fig. 4.4.



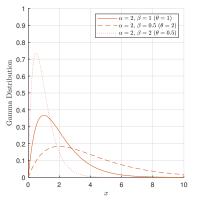


FIGURE 4.4
Gamma Distribution.

A Little Bit about Gamma Function

Gamma function is given by

$$\Gamma(\alpha) = \int_0^\infty t^\alpha - 1e^{-t}dt$$

It is clear that $\Gamma(1) = 1$. If $\alpha > 1$, an integration by parts shows that

$$\Gamma(\alpha) = (\alpha - 1)\Gamma(\alpha - 1)$$

Therefore, $\Gamma(n) = \Gamma(n-1)\Gamma(n-2)... \times 3 \times 2 \times 1 \times \Gamma(1) = (\alpha-1)!$ for $\alpha \in \mathbb{N}^+$, with 0! = 1.

The sum of Gamma distributions with the same scale parameter also follows Gamma distribution. The expectation, variance, skewness and excess kurtosis of the distribution are $\frac{\alpha}{\beta}$, $\frac{\alpha}{\beta^2}$, $\frac{2}{\sqrt{\alpha}}$ and $\frac{6}{\alpha}$, respectively.

4.6.2 The χ^2 Distribution

The χ^2 estimation is a special case of Gamma distribution. In (4.2), let $\alpha = \frac{r}{2}$ and $\beta = \frac{1}{2}$ to get the PDF of χ^2 distribution as follows.

$$f(x) = \begin{cases} \frac{1}{\Gamma(r/2)2^{r/2}} x^{r/2-1} e^{-x/2} & x > 0\\ 0 & \text{otherwise} \end{cases}$$
 (4.3)

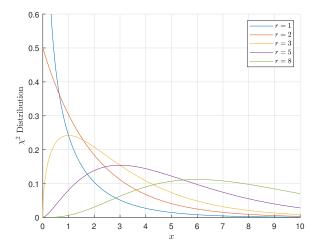


FIGURE 4.5 The χ^2 Distribution.

Equation (4.3) is also denoted by $\chi^2(r)$, where r is known as the degrees of freedom of the chi-squared distribution. The plot of PDF are shown in Fig. 4.5.

The χ^2 distribution can be used to model the sum of the squares of r independent standard normal distributions, making it an important distribution is statistics, such as in outlier test.

The expectation, variance, skewness and excess kurtosis of the distribution are r, 2r, $\sqrt{\frac{8}{r}}$ and $\frac{12}{r}$, respectively.

4.6.3 The β Distribution

Let X_1 and X_2 be two independent random variables following Gamma distributions with shape parameters α_1 and α_2 respectively, and with rate parameters 1. Denote $X = \frac{X_1}{X_1 + X_2}$. In this case, X would be a random variable following β distribution whose PDF can be derived from (4.2) and it is given by

$$f(x) = \begin{cases} \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} x^{\alpha_1 - 1} (1 - x)^{\alpha_2 - 1} & 0 < x < 1 \\ 0 & \text{otherwise} \end{cases}$$
(4.4)

In some literatures, notation (a,b) or (α,β) are used to denote (α_1,α_2) in (4.4) respectively.

4.7 Student's t-Distribution

"nobreak

4.8 *F*-Distribution

"nobreak

4.9 Weibull Distribution

Part II Statistics

Sampling

CONTENTS

5.1	Sampling Methods	37
5.2	Model of Population	38
	Sample Statistics	

There is an important underlying assumption in statistics: the features observed from samples can somehow reflect the features of the entire population, where the samples are (randomly) selected from the population. Intuitively, this assumption shall hold true in many occasions.

Statistics studies how to select samples, how to make the best use of the samples information, and how accurate and reliable the information derived from the samples is when applying them on the population. In this chapter, we start from study sampling.

5.1 Sampling Methods

When selecting elements from the population, make sure that all elements have a equal probability of being selected, hence, random sampling. Depending on how many times a member can be sampled, we have

- Sampling with replacement: a member can be chosen more than once.
- Sampling without replacement: a member can be chosen no more than once.

Sometimes it is interesting to compare the differences of the two methods, especially when the population is finite. An obvious difference is that by using sampling with replacement all the samples can be considered as "independent event", while by using sampling without replacement, previous samples may change the distributions in the remaining population, thus making the samples relevant. In this case, using sampling with replacement can theoretically be considered as sampling from an infinite population (by thinking that the population is duplicated as many times as necessary).

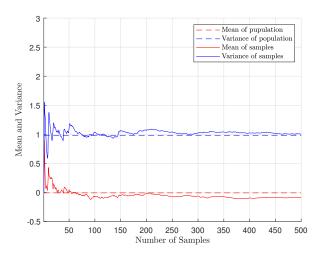


FIGURE 5.1 Sample with replacement, N = 100, M = 500.

In practice, the population is usually so large, that sampling from a finite population can be considered as sampling from an infinite population, and the two methods would make no differences as far as it is concerned.

Consider the following examples. A set of N random variables are generated from a Gaussian distribution as the population. Sample the population M times using sampling with replacement and sampling without replacement, respectively. Calculate the sampled mean and variance after each sampling instance, and see how it converges to the mean and variance of the population.

In the first example, let N=100 and M=500. Figures 5.1 and 5.2 gives the cumulative mean and variance of sampling with and without replacement, respectively. The mean and variance are given by red and blue curves, respectively. The statistics obtained from the cumulative samples and from the population are given by the solid and dashed curves, respectively. Notice that in Fig. 5.2, after number of samples exceeding 100, the entire population has been sampled, and thus the sampling stops. This explains why its mean and variance stop fluctuating and converge to the population mean and variance, respectively.

In practice, however, the population size is often orders of magnitudes larger than the number of samples. In the second example, let N=10000 and M=500. The corresponding figures are given in Figs. 5.3 and 5.4. There is no obvious differences of the two figures from statistics perspective.

Sampling 39

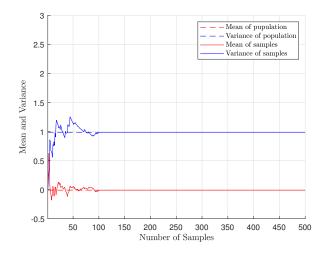


FIGURE 5.2 Sample without replacement, $N=100,\,M=500.$

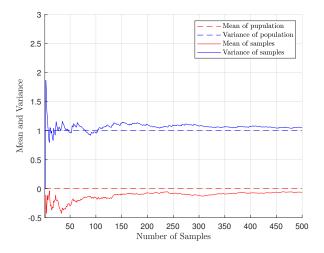


FIGURE 5.3 Sample with replacement, $N=10000,\,M=500.$

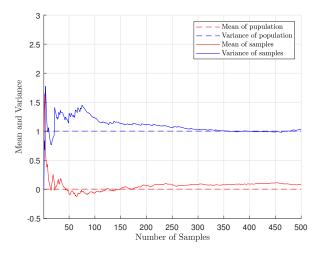


FIGURE 5.4 Sample without replacement, N = 10000, M = 500.

5.2 Model of Population

The features of the population is often not known, or at least not known entirely. It is possible to make some preliminary assumptions to the distribution of population, with parameters to be further confirmed using the samples.

For example, let X be a variable of the population. It could be, for example, the heights of all teenagers in a city. We can make an assumption that X follows some distribution f(x). A widely used assumption, in this scenario, is that f(x) is a Gaussian distribution with mean μ and standard deviation σ , and each element in the population, X_i , can be taken as a random variable generated from f(x). In the case of Gaussian distribution, since it is uniquely characterized by μ and σ , other quantities such as the median, moments, skewness, etc., can be derived once μ and σ is calibrated.

The questions rise sequentially are:

- What are the parameters in the assumed distribution?
- Does it indeed follow the assumed distribution?

The answers to the above questions need to be found out via the samples, or more precisely, from the sample statistics.

Sampling 41

5.3 Sample Statistics

Estimation

Regression

Hypothesis Testing

$Bayesian\ Methods$

Part III

Tools

R Basics

CONTENTS

10.1	R and	RStudio Installation	53
10.2	R Pacl	kages Management	54
	10.2.1	Manage Packages with Built-in Functions	54
	10.2.2	Manage Packages with Third-Party Packages	55
	10.2.3	Manage Packages with RStudio IDE	55
10.3	R Prog	gramming Basics	55
	10.3.1	Data Types	57
	10.3.2	Conditionals and Loops	59
	10.3.3	User-Defined Functions	63
10.4	Vector	and Matrix	63
	10.4.1	Vector	63
	10.4.2	Matrix	66
	10.4.3	Matrix Visualization Using matplot()	69
10.5	Data I	Frames	71
	10.5.1	Import Data into Data Frame	71
	10.5.2	Access Data in Data Frame	73
	10.5.3	Filter Data from Data Frame	75
	10.5.4	Create Data Frames	77
10.6	Basic 1	Data Visualizations Using qplot()	78
10.7	Advan	ced Data Visualizations Using ggplot()	79
	10.7.1	Grammar of Graphics	79
	10.7.2	Data, Aesthetics and Geometries Layers	81
	10.7.3	Statistics Layers	83
	10.7.4	Facets Layers	85
	10.7.5	Coordinates Layers	87
	10.7.6	Themes Layers	89

This chapter and the sequential chapters introduce R language, a widely used programming language for statistical computing and data visualization.

This chapter focuses on the fundamental setup and basic operations such as the installation of R and RStudio, data types, basic data frames manipulation and visualization methods. Chapter 11 introduces advanced skills and tools widely used in data preparation and processing with R language. Finally, Chapter 12 puts the use of R language in data science into practice.

10.1 R and RStudio Installation

R is a programming language for statistical computing and visualization. It is widely used among statisticians and data miners for developing statistical software and carrying out data analysis. R is free and can be downloaded from [4], from which more details about R can also be found.

RStudio, also known as Posit, is an IDE widely used for R programming and testing. RStudio IDE is open-source and free of charge for personal use. It can be downloaded from [5].

Download R and RStudio from the aforementioned web sites, and install them following the instructions.

10.2 R Packages Management

R uses packages to boost its capability. R packages, both built-in and third-party, provide powerful features for data analysis and visualization. Some packages may also come with demonstrative sample data frames. The packages can be published and shared online. CRAN is by far the most popular platform to store and share R packages.

10.2.1 Manage Packages with Built-in Functions

To install or remove a package, use

```
install.packages("<package>")
remove.packages("<package>")
```

respectively. For example,

install.packages("pacman")

Built-in packages come along with R installation. Only third-party libraries need additional installation.

To load a package, use

library(<package>)

For example

library(pacman)

After loading a package, the data frames and functions defined in that package can be used normally. Alternatively, use prefix <package>::<function>, <package>::<dataframe> to call the functions and data frames defined in a package without loading it. The later approach can become inconvenient when the function or data frame is used frequently.

To unload a package, use

```
detach("package:<package>", unload = TRUE)
For example,
detach("package:pacman", unload = TRUE)
```

10.2.2 Manage Packages with Third-Party Packages

There are third-party packages that provide package management functions, for example pacman. With pacman installed and loaded, use the following commands to install, load and unload packages respectively.

```
p_install(<package>, ...) # install
p_load(<package>, ...) # install and load
p_unload(<package>, ...) # unload
p_unload(all) # unload all
```

An example of using pacman to load packages are given as follows.

```
pacman::p_load(
    pacman, # package management
    dplyr, # data manipulation
    GGally, # data visualization
    ggplot2, # data visualization
    ggthemes, # data visualization
    ggvis, # data visualization
    httr, # url and http
    lubridate, # date and time manipulation
    plotly, # data visualization
    rio, # io
    rmarkdown, # documentation
    shiny, # web apps development
    stringr, # string operation
    tidyr # data tidying
)
```

The above command can be executed without loading pacman in advance since pacman::p_load() prefix is used. The listed packages are commonly used in R projects, and a brief explanation to them is given as comments in the code listing following #. Notice that if a package is not installed in the computer, the above command will try to download and install the package.

10.2.3 Manage Packages with RStudio IDE

RStudio provides a graphical interface to manage packages as shown in Fig. 10.1. A package can be loaded or unloaded simply by checking and unchecking the package.

File	s Plots Packages	Help Vie	er Presentation			-0
0	Install ① Update				Q,	
	Name		Description		Version	
User	r Library					4
	askpass		Safe Password Entry for R, Git, and SSH		1.1	⊕ ⊗
	assertthat		Easy Pre and Post Assertions		0.2.1	0 0
	base64enc		Tools for base64 encoding		0.1-3	⊕ ⊗
	bit		Classes and Methods for Fast Memory-Efficient	Boolean Selections	4.0.5	⊕ ⊗
	bit64		A S3 Class for Vectors of 64bit Integers		4.0.5	0 0
	bslib		Custom 'Bootstrap' 'Sass' Themes for 'shiny' an	d 'rmarkdown'	0.4.2	⊕ ⊗
	cachem		Cache R Objects with Automatic Pruning		1.0.6	0 0
	cellranger		Translate Spreadsheet Cell Ranges to Rows and	Columns	1.1.0	⊕ ⊗
	cli Helpers for Developing Command Line		Helpers for Developing Command Line Interfac	es	3.5.0	⊕ ⊗
	clipr		Read and Write from the System Clipboard		0.8.0	⊕ ⊗
	colorspace		A Toolbox for Manipulating and Assessing Colors and Palettes		2.0-3	⊕ ⊗
	commonmark High Per		High Performance CommonMark and Github N	larkdown Rendering in R	1.8.1	⊕ ⊗
	cpp11		A C++11 Interface for R's C Interface		0.4.3	⊕ ⊗
	crayon		Colored Terminal Output		1.5.2	⊕ ⊗
	crosstalk		Inter-Widget Interactivity for HTML Widgets		1.2.0	⊕ ⊗
	curl		A Modern and Flexible Web Client for R		4.3.3	● ⊗
	data.table Extension of `data.frame` 1.14.6			⊕ ⊗		
	digest Create Compact Hash Digests of R Objects 0.6		0.6.31	⊕ ⊗		
	dplyr		A Grammar of Data Manipulation		1.0.10	0 0
	ellipsis		Tools for Working with		0.3.2	⊕ ⊗
	evaluate		Parsing and Evaluation Tools that Provide More	Details than the Default	0.19	0 0

FIGURE 10.1

RStudio's graphical interface for package management.

10.3 R Programming Basics

This section introduces the basics of R programming, including its data types and basic syntax.

It is worth mentioning in the beginning that R is case sensitive. Use # to lead a comment. Use print(<variable>) or simply type the name of the variable to print the value of a variable to the console. Use ?<function>, ?<dataframe> to access the manual for a function or data frame.

As a quick demonstrative example, the following is a small piece of code written in R. It generates generalized t-distribution noise and save the noise samples into a local CSV file.

TABLE 10.1

Commonly used data types.

Data Type	Syntax (Example)	Description	
integer	n <- 2L	An integer. Define an integer by a value	
		followed by L.	
double	x <- 2	An double float value.	
complex	z <- 3+2i	A complex value.	
character	a <- "a"	A character or a string.	
logical	q <- T	A boolean value. Use T, TRUE and F,	
		FALSE to represent true and false respec-	
		tively.	

```
# Generate pseudo random numbers
library('optimx')
library('numDeriv')
library('sgt')
set.seed(sgt_seed)
x = rsgt(n = sgt_n, mu = sgt_mu, sigma = sgt_sigma, lambda = sgt_lambda
, p = sgt_p, q = sgt_q, mean.cent = TRUE, var.adj = FALSE)
# Write data
write.table(x, 'generate_sgt_data.csv', sep=",", row.names = FALSE, col
.names=FALSE)
```

10.3.1 Data Types

R supports many data types. Commonly used data types are summarized in Table 10.1, where notice that <- is used to assign a value to a variable. Use typeof() to check the type of a variable. Alternatively, use is.numeric(), is.integer(), is.double(), is.character(), etc., to check whether a variable belongs to a particular data type.

Examples of assigning variables and checking their types are given as follows. Notice that > is the prompt indicating that the commands are executed in a console.

```
> n <- 2L
> typeof(n)
[1] "integer"

> x <- 2
> typeof(x)
[1] "double"

> z <- 3+2i
> typeof(z)
```

```
[1] "complex"
> a <- "h"
> typeof(a)
[1] "character"
> q <- T
> typeof(q)
[1] "logical"
```

To transform data from one type to another, use as.<datatype>(). Examples of transforming data types are given as follows.

```
> n1 <- as.integer(2)
> typeof(n1)
[1] "integer"
> n2 <- as.integer("2")
> typeof(n2)
[1] "integer"
> x1 <- as.double(2L)
> typeof(x1)
[1] "double"
> x2 <- as.double("2")
> typeof(x2)
[1] "double"
> z1 <- as.complex("3+2i")</pre>
> typeof(z1)
[1] "complex"
> a1 <- as.character(2L)</pre>
> typeof(a1)
[1] "character"
> a2 <- as.character(2)</pre>
> typeof(a2)
[1] "character"
```

R supports arithmetic calculations of variables, including +, -, *, /, % (integer division), % (modulus) and $\hat{}$ (exponential). Examples of arithmetic calculations are given as follows.

```
> a <- 16
> b <- 3
> add <- a + b
> sub <- a - b
> multi <- a * b
> division <- a / b</pre>
```

TABLE 10.2

Numerical calculations Description Syntax (Example) Absolute value. abs(x)sqrt(x) Square root. Smallest larger/equal integer. ceiling(x) Largest smaller/equal integer. floor(x) trunc(x) Integer part of a variable. Round to n digit after decimal. round(x, n=0)sin(x)Trigonometric sin function. cos(x)Trigonometric cos function. tan(x)Trigonometric tan function. log(x)Natural logarithm. log10(x) Common logarithm. exp(x)Exponent.

```
> int_division <- a %/% b</pre>
 modulus <- a %% b
 exponent <- a ^ b
 add
[1] 19
> sub
[1] 13
> multi
[1] 48
> division
[1] 5.333333
> int_division
[1] 5
> modulus
[1] 1
> exponent
[1] 4096
```

R uses built-in and third-party functions which extend its capability. There is a rich set of functions for numerical calculations, string operations, probability calculations and statistics analysis. Some of them are summarized in Tables 10.2, 10.3, 10.4, 10.5 and 10.6.

10.3.2 Conditionals and Loops

The if statement syntax is given as follows.

TABLE 10.3

Logical comparisons.

Syntax (Example)	Description
х == у	Equal.
x != y	Not equal.
x > y, x < y	Greater than; less than.
x >= y, x <= y	Greater than or equal to; less than or equal to.
! x	Not.
x & y	And.
x y	Or.
isTRUE(x)	Is true.

TABLE 10.4

String operations.

Syntax (Example)	Description
substr(s, n1, n2)	Segment of a string, from the n_1 -th charac-
	ter to n_2 -th character, both characters in-
	cluded.
<pre>grep(p, s)</pre>	Searching of a pattern in a string.
sub(s1, s2, s)	Find and replace patterns in a string.
paste(s1, s2,, p="")	Concatenate strings.
strsplit(s, p)	Split a string.
tolower(s)	Convert to lower case.
toupper(s)	Convert to upper case.

61

TABLE 10.5

Prob	ability	v related	operations.
1 101	<i>j</i> abiii (v rerateu	operanons.

is.
Description
Calculate the PDF of Gaussian distribution.
Calculate the CDF of Gaussian distribution.
Inverse function of pnorm().
Generate Gaussian distribution samples.
Calculate the probability of a binominal dis-
tribution.
Calculate the comulative probability of a bi-
nominal distribution.
Inverse function of pbinom().
Generate binominal distribution samples.
Calculate the probability of a Poisson distri-
bution.
Calculate the comulative probability of a Pois-
son distribution.
Inverse function of ppois().
Generate Poisson distribution samples.
Calculate the PDF of uniform distribution.
Calculate the CDF of uniform distribution.
Inverse function of punif().
Generate uniform distribution samples.

TABLE 10.6
Statistics relate

Statistics	related	functions.

Syntax (Example)	Description
mean(1)	Mean.
sd(1)	Standard deviation.
median(1)	Median.
range(1)	Minimum and maximum.
min(l)	Minimum.
max(1)	Maximum.
sum(1)	Sum

An example of using if statement is givne below.

where + indicates the splitting of commands in multiple lines in the console. The for loop syntax is given as follows.

where **<vector>** is a list of numbers or characters. Examples are given below.

where c() creates a list from items.

The while loop syntax is given as follows.

```
while(<condition>) {
      <command>
}
```

An example of using while loop is given below.

```
> counter <- 0
> while(counter < 5){
+    print(counter)
+    counter <- counter + 1
+ }
[1] 0
[1] 1
[1] 2
[1] 3
[1] 4</pre>
```

An example using the above to estimate the ratio of data of a standard normal distribution between ± 1 is given below.

10.3.3 User-Defined Functions

Define a simple function as follows. Run the codes where the function is described before calling the function.

10.4 Vector and Matrix

Vector and matrix are the fundamental forms that R stores and processes data.

10.4.1 Vector

There are different types of vectors in R depending on the data type of the elements it stores. The commonly used vector types include numeric vector and character vector. All elements in a vector must have the same data type.

When different data type values are stored in a vector, they will be transferred to the most general data type.

Notice that the index of a vector in R starts from 1 instead of 0. This is different with many other computer languages such as C/C++, JavaScript and Python.

Use the following syntax to create a vector.

```
<vector> <- c(<value>, ...)
```

where <value> can be single element or a vector. For example,

```
> x <- c(1,2,3,4,5)
> x
[1] 1 2 3 4 5
> y <- c(c(1,2,3),4,5,6,7)
> y
[1] 1 2 3 4 5 6 7
> z = c(x, 1,2,3,4,5)
> z
[1] 1 2 3 4 5 1 2 3 4 5
> typeof(z)
[1] "double"
```

It is possible to assign names to each element as shown below.

```
> v <- c(1, 2, 3, 4, 5)
> names(v) <- c("e1", "e2", "e3", "e4", "e5")
> print(v)
e1 e2 e3 e4 e5
1 2 3 4 5
```

Alternative ways to create a vector are given as follows. Use sequence to create a vector as follows.

```
<vector> <- seq(<from>, <to>, <by=1>)
<vector> <- <from>:<to> # equivalent to seq() with by=1
```

Use replica to create a vector as follows.

```
<vector> <- rep(<value>, <repeate>)
```

where **<value>** can be a numeric number, a character, or a vector. For example,

```
> 1 <- rep(c("a", "b", "cde"), 2)
> print(1)
[1] "a" "b" "cde" "a" "b" "cde"
```

Replica can also be used to create empty vector vy rep(NA, n).

A character vector can also be created by splitting strings using strsplit(). For example,

```
> a <- "Hello World!"
> b <- strsplit(a, "")
```

```
> print(b)
[[1]]
[1] "H" "e" "l" "l" "o" " " "W" "o" "r" "l" "d" "!"
```

To access the element in a vector, use <vector>[<index>] or <vector>["<name>"]. Notice that the first element in a vector has the index of 1 instead of 0. The index can be an integer, or an integer vector. Examples are given below.

It is also possible to address an element by its name

```
> v <- c(1, 2, 3, 4, 5)
> names(v) <- c("e1", "e2", "e3", "e4", "e5")
> print(v)
e1 e2 e3 e4 e5
1  2  3  4  5
> print(v[3])
e3
3
> print(v["e3"])
e3
3
```

Notice that it is inefficient to use deep-layer loops to access and process data in vector or matrix one-by-one because most operations in R are done by the vector basis. Vectorization operation significantly speeds up the calculation, hence quite commonly seen in high-layer languages such as R and Python.

Most numerical calculations, including +, -, *, /, %, %, $^{\circ}$, support vectorization operation. Examples are given below.

```
> a <- c(1,2,3,4,5)
> b <- c(5,4,3,2,1)
> a + b
[1] 6 6 6 6 6
> a - b
[1] -4 -2 0 2 4
> a * b
[1] 5 8 9 8 5
> a / b
```

```
[1] 0.2 0.5 1.0 2.0 5.0

> a %/% b

[1] 0 0 1 2 5

> a %% b

[1] 1 2 0 0 0

> a ^ b

[1] 1 16 27 16 5
```

It is also possible to apply logic operations using vectors. Examples are given below.

```
> a <- c(1,2,3,4,5)
> b <- c(5,4,3,2,1)
> a < b
[1] TRUE TRUE FALSE FALSE
> a > b
[1] FALSE FALSE FALSE TRUE
> a == b
[1] FALSE FALSE TRUE FALSE
```

When the sizes of the vectors are not consistent, the shorter vector will repeat and populate to align with the longer vector. Examples are given below.

```
> a <- c(1,10)
> b \leftarrow c(1,2,3,4)
> a + b
[1] 2 12 4 14
> a - b
[1] 0 8 -2 6
> a * b
[1] 1 20 3 40
> a / b
[1] 1.0000000 5.0000000 0.3333333 2.5000000
> a %/% b
[1] 1 5 0 2
> a %% b
[1] 0 0 1 2
> a ^ b
[1]
                   1 10000
```

A vector can also be used as an input argument or output return of a function.

10.4.2 Matrix

Similar with vector, all elements in a matrix must have consistent type. A demonstration of matrix indexing in R is shown in Fig. 10.2. Let the matrix be named A. The elements in the matrix can be accessed by the name of the table followed by the index coordinates. For example, in the figure, A[1,1] refers to the first element and A[4,6] the last element. It is also possible to

	L,13	[,2]	[,3]	L,4]	L,5]	L,6]
[1,]	E1,13					
[2,]						
[3,]			[3,3]			
[4,]						[4,6]

FIGURE 10.2

A demonstration of a matrix in R.

index an entire row or column. For example, use A[1,] to represent the first row of the matrix.

A matrix can be created from scratch by stacking rows as follows. First, create rows in the matrix. Then, use rbind() to bind rows. Notice that the columns and rows of a matrix can have names in addition to their associated numeric indices, and the names can also be used to for index. This is different from MATLAB or Python numpy.

```
# build rows
<row1> <- c(<value11>, ..., <value1n>)
...
<rowm> <- c(<valuem1>, ..., <valuemn>)
# build matrix
<matrix> <- rbind(<row1>, ..., <rowm>)
# (optional) clean rows
rm(<row1>, ..., <rowm>)
# give names
colnames(<matrix>) <- c("<column-name1>", ..., "<column-namen>")
rownames(<matrix>) <- c("<row-name1>", ..., "<row-namem>")
```

There are alternative ways, other than rbind(), to create a matrix. For example, matrix() convert a vector into a matrix. Similar with rbind(), cbind() binds the columns to form a matrix. Examples to create matrices using different methods are given below.

```
> A <- matrix(1:9, 3, 3)
> print(A)
[,1] [,2] [,3]
[1,] 1 4 7
[2,] 2 5 8
```

```
[3,]
            6 9
> B <- rbind(c(1, 4, 7), c(2, 5, 8), c(3, 6, 9))
> print(B)
[,1] [,2] [,3]
[1,]
            4
                7
       1
[2,]
       2
                8
            5
[3,]
       3
            6
> C \leftarrow cbind(c(1, 2, 3), c(4, 5, 6), c(7, 8, 9))
> print(C)
[,1] [,2]
          [,3]
[1,]
       1
            4
                 7
[2,]
       2
                 8
            5
[3,]
       3
            6
```

To check or assign names for a vector or matrix, use names(<vector>), rownames(<matrix>) and colnames(<matrix>), depending on dealing with either a vector or a matrix. These commands can also be used to assign names. Examples are given below.

```
> v <- c(1, 2, 3, 4, 5)
> names(v) <- c("e1", "e2", "e3", "e4", "e5")
> v
e1 e2 e3 e4 e5
1 2 3 4 5
> names(v)
[1] "e1" "e2" "e3" "e4" "e5"
> A <- matrix(1:9, 3, 3)
> colnames(A) <- c("col1", "col2", "col3")</pre>
> rownames(A) <- c("row1", "row2", "row3")</pre>
> A
col1 col2 col3
row1
       1
            4
                 7
row2
       2
            5
                 8
row3
       3
            6
> colnames(A)
[1] "col1" "col2" "col3"
> rownames(A)
[1] "row1" "row2" "row3"
> A[1,1]
[1] 1
> A["row1", "col1"]
[1] 1
> A[1,2]
[1] 4
> A["row1", "col2"]
[1] 4
```

To remove the names, simply assign NULL to the name.

Vectorization operators are defined for matrix level as well. For example,

for two matrices with the same shape, numerical operations such as +, -, *, /, %, % and $\hat{}$ can be used.

10.4.3 Matrix Visualization Using matplot()

R provides flexible and powerful data visualization tools, many of which more advanced than what is to be introduced in this section. This section introduces a simple matrix visualization function called matplot(), which plots the columns of a matrix against each other.

To demonstrate matplot(), consider the following example.

```
professor <- c(1130, 1026, 893, 922, 776)
student \leftarrow c(2, 14, 24, 49, 46)
citation <- rbind(professor, student)</pre>
colnames(citation) <- c("2018", "2019", "2020", "2021", "2022")
rownames(citation) <- c("Professor", "Student")
print(citation)
citation.ratio <- citation
citation.ratio["Professor",] <- round(citation["Professor",] / mean(</pre>
    citation["Professor",]) * 100, 1)
citation.ratio["Student",] <- round(citation["Student",] / mean(</pre>
    citation["Student",]) * 100, 1)
print(citation.ratio)
matplot(
       2018:2022, # x axis
       t(citation.ratio), # y axis
       type="b", # line and point selection
       pch = 15:16, # point shape
       col = 1:2, # color
       xlab = "Year",
       ylab = "Citations Moving Ratio (%)"
legend("bottomright", inset = 0.01, legend = rownames(citation.ratio),
    pch = 15:16, col = 1:2, horiz = F)
```

where t() used inside matplot() calculates the transpose of a matrix. Save the above in a script and execute the code, to get the following Fig. 10.3.

Notice that matplot() is not widely used in practice.

As introduced earlier, a matrix or a vector can be split and segmented to form a smaller matrix or vector. It is worth mentioning that when a single column or row is selected, R will automatically treated the return as a vector instead of a matrix. An example is given below. When a matrix downgrades to a vector, the row name (if it has only one row), or the column name (if it has only one column) will be removed.

```
> A <- matrix(1:9, 3, 3)
> is.matrix(A)
[1] TRUE
> is.vector(A)
```

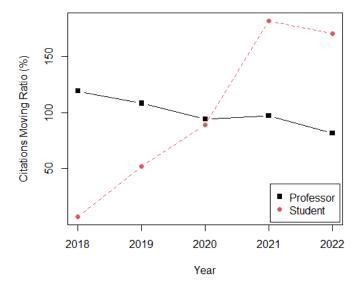


FIGURE 10.3

A demonstration of using matplot to plot trends.

```
[1] FALSE
> is.matrix(A[1,])
[1] FALSE
> is.vector(A[1,])
[1] TRUE
```

To get consistent results, when segmenting matrix to get a single row or column vector, deliberately ask R to not drop the matrix dimensions. This can be done as follows. By doing this, the names assigned to columns and rows preserve.

```
> A <- matrix(1:9, 3, 3)
> is.matrix(A[1,,drop=F]) # select a row/column
[1] TRUE
> is.matrix(A[2,3,drop=F]) # select an element
[1] TRUE
```

An example of using the above to analyze the performance of players through a series of basketball games are given below.

```
# generate table
player_name <- c("player1", "player2", "player3")
match_name <- c("match1", "match2", "match3", "match4", "match5", "
    match6", "match7", "match8", "match9", "match10")
penalty_attempt <- abs(matrix(round(rnorm(3*10, 5, 2)), 3, 10))
penalty_point <- abs(penalty_attempt - matrix(abs(round(rnorm(3*10, 1, 1))), 3, 10))</pre>
```

```
throw_attempt <- abs(matrix(round(rnorm(3*10, 15, 3)), 3, 10))
total_point <- abs(3*throw_attempt - abs(matrix(round(rnorm(3*10, 5, 1)
    ), 3, 10))) + penalty_point
rownames(penalty_attempt) <- player_name
colnames(penalty_attempt) <- match_name</pre>
rownames(penalty_point) <- player_name
colnames(penalty_point) <- match_name</pre>
rownames(throw_attempt) <- player_name
colnames(throw_attempt) <- match_name</pre>
rownames(total_point) <- player_name
colnames(total_point) <- match_name</pre>
# claim function
myplot <- function(table, xlab, ylab){</pre>
   row_name = rownames(table)
   column_name = colnames(table)
   matplot(
       1:length(column_name), # x axis
       t(table), # y axis
       type="b", # line and point selection
       pch = 1:length(row_name), # point shape
       col = 1:length(row_name), # color
       xlab = xlab,
       ylab = ylab
   legend("bottomleft", inset = 0.01, legend = row_name, pch = 1:
        length(row_name), col = 1:length(row_name), horiz = F)
# plot
myplot(penalty_point / penalty_attempt, "match", "penalty success rate
    ") # penalty successful rate
myplot((total_point - penalty_point) / throw_attempt, "match", "average
     gained point per throw") # average point gained per throw
```

The results of the above codes are given in Figs. 10.4, 10.5.

10.5 Data Frames

Data frame, just like vector and matrix, is another data structure defined in R.

Both matrix and data frame use a table structure to store data, but data frame does not require all data to have the same data type, which makes it maybe the most important and commonly used data structure in R.

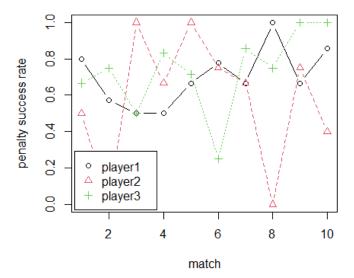


FIGURE 10.4 Plot of penalty success rate of the 3 players in 10 matches.

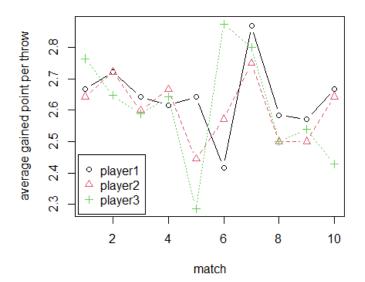


FIGURE 10.5 Plot of average point gained per throw attempt for the 3 players in 10 matches.

10.5.1 Import Data into Data Frame

One of the most common sources of data is CSV files. R provides convenient functions to read data from CSV files into data frames. Use the following commands to import data from a CSV file into a data frame.

The following command pops up a separate window that allows the user to choose a CSV file manually.

```
<data-frame> <- read.csv(file.choose()) # manual selection</pre>
```

The following commands import a specified CSV file.

```
setwd("<directory>") # navigate to the directory of the csv file
<data-frame> <- read.csv("<csv-file>.csv")
```

where notice that getwd() and setwd() are used to get and set current working directory, respectively.

10.5.2 Access Data in Data Frame

There are a few ways to access an element in a data frame as shown below.

```
<df>[<row_index>, <column_index>]
<df>[<row_index>, "<column_name>"]
<df>[, <column-name>] [<row_index>]
<df>$<column_name>[<row_index>]
```

where

```
<df>[, <column-name>]
<df>$<column-name>
```

are used to access the entire column in a data frame. Notice that different from a matrix, rows in data frame do not have names.

Table 10.7 summarizes the commonly used functions on date frames, such as checking its shape and data types. Examples of applying these functions to iris data frame from the built-in datasets package are given below.

```
> library(datasets)
> nrow(iris)
[1] 150
> ncol(iris)
[1] 5
> head(iris)
 Sepal.Length Sepal.Width Petal.Length Petal.Width Species
                     3.5
                                 1.4
                                            0.2 setosa
          5.1
          4.9
                     3.0
                                 1.4
                                             0.2 setosa
          4.7
                     3.2
                                 1.3
                                             0.2 setosa
          4.6
                     3.1
                                 1.5
                                             0.2 setosa
          5.0
                     3.6
                                 1.4
                                             0.2 setosa
          5.4
                     3.9
                                 1.7
                                             0.4 setosa
 tail(iris)
```

TABLE 10.7
Commonly used commands for data frame exploration.

Syntax (Example)	Description
nrow(df)	Number of rows.
<pre>ncol(df)</pre>	Number of columns.
head(df, n=6L)	Display the first few rows.
tail(df, n=6L)	Display the last few columns.
str(df)	A summary of the data frame, including the struc-
	ture of each column.
<pre>summary(df)</pre>	A summary of the data frame, including some of
	its statistics features.
<pre>levels(df\$<column>)</column></pre>	The level of the column.

```
Sepal.Length Sepal.Width Petal.Length Petal.Width Species
145
           6.7
                      3.3
                                            2.5 virginica
                                 5.7
146
           6.7
                      3.0
                                             2.3 virginica
                                  5.2
147
           6.3
                      2.5
                                  5.0
                                             1.9 virginica
148
           6.5
                      3.0
                                  5.2
                                             2.0 virginica
149
           6.2
                      3.4
                                  5.4
                                             2.3 virginica
150
           5.9
                      3.0
                                  5.1
                                            1.8 virginica
> str(iris)
'data.frame': 150 obs. of 5 variables:
$ Sepal.Length: num 5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 ...
$ Sepal.Width : num 3.5 3 3.2 3.1 3.6 3.9 3.4 3.4 2.9 3.1 ...
$ Petal.Length: num 1.4 1.4 1.3 1.5 1.4 1.7 1.4 1.5 1.4 1.5 ...
$ Petal.Width : num 0.2 0.2 0.2 0.2 0.2 0.4 0.3 0.2 0.2 0.1 ...
$ Species
             : Factor w/ 3 levels "setosa", "versicolor", ...: 1 1 1 1 1
    1 1 1 1 1 ...
 summary(iris)
 Sepal.Length Sepal.Width
                              Petal.Length
                                            Petal.Width
                                                               Species
Min. :4.300 Min. :2.000
                             Min. :1.000 Min. :0.100
                                                                 :50
1st Qu.:5.100 1st Qu.:2.800 1st Qu.:1.600 1st Qu.:0.300 versicolor
     :50
Median: 5.800 Median: 3.000 Median: 4.350 Median: 1.300 virginica
     :50
Mean :5.843 Mean :3.057 Mean :3.758 Mean
                                                 :1.199
3rd Qu.:6.400 3rd Qu.:3.300 3rd Qu.:5.100 3rd Qu.:1.800
Max. :7.900 Max. :4.400 Max. :6.900 Max.
> levels(iris$Species) # only works discrete-value columns
[1] "setosa"
               "versicolor" "virginica"
```

Data frame segmentation works similarly with matrix. A segmentation of multiple rows and columns of a data frame is also a data frame. Notice that when a single column is segmented, the result will be a vector by default, and drop=F can be used to preserve data frame structure. An example is given below.

```
> library(datasets)
 print(iris[1:5,])
 Sepal.Length Sepal.Width Petal.Length Petal.Width Species
         5.1
                    3.5
                                1.4
                                           0.2 setosa
         4.9
                    3.0
                                 1.4
                                           0.2 setosa
                    3.2
                                1.3
         4.7
                                           0.2 setosa
                                 1.5
         4.6
                    3.1
                                            0.2 setosa
         5.0
                    3.6
                                 1.4
                                            0.2 setosa
 print(iris[1,])
 Sepal.Length Sepal.Width Petal.Length Petal.Width Species
         5.1
                    3.5
                                 1.4
                                            0.2 setosa
> is.data.frame(iris[1,])
[1] TRUE
 print(iris[,1]) # equivalent to print(iris@Sepal.Length)
 [1] 5.1 4.9 4.7 4.6 5.0 5.4 4.6 5.0 4.4 4.9 5.4 4.8 4.8 4.3 5.8 5.7
     5.4 5.1 5.7 5.1 5.4 5.1 4.6 5.1
[25] 4.8 5.0 5.0 5.2 5.2 4.7 4.8 5.4 5.2 5.5 4.9 5.0 5.5 4.9 4.4 5.1
    5.0 4.5 4.4 5.0 5.1 4.8 5.1 4.6
[49] 5.3 5.0 7.0 6.4 6.9 5.5 6.5 5.7 6.3 4.9 6.6 5.2 5.0 5.9 6.0 6.1
     5.6 6.7 5.6 5.8 6.2 5.6 5.9 6.1
[73] 6.3 6.1 6.4 6.6 6.8 6.7 6.0 5.7 5.5 5.5 5.8 6.0 5.4 6.0 6.7 6.3
     5.6 5.5 5.5 6.1 5.8 5.0 5.6 5.7
[97] 5.7 6.2 5.1 5.7 6.3 5.8 7.1 6.3 6.5 7.6 4.9 7.3 6.7 7.2 6.5 6.4
     6.8 5.7 5.8 6.4 6.5 7.7 7.7 6.0
[121] 6.9 5.6 7.7 6.3 6.7 7.2 6.2 6.1 6.4 7.2 7.4 7.9 6.4 6.3 6.1 7.7
    6.3 6.4 6.0 6.9 6.7 6.9 5.8 6.8
[145] 6.7 6.7 6.3 6.5 6.2 5.9
> is.data.frame(iris[,1])
[1] FALSE
> print(iris[,1,drop=F]) # preserve data frame
   Sepal.Length
           5.1
2
           4.9
3
           4.7
    # WRAPPED #
148
           6.5
149
           6.2
150
           5.9
> is.data.frame(iris[,1,drop=F])
[1] TRUE
```

To add a new column to an existing data frame, just assign values to a new column name as follows.

<df>\$<new-column> <- <vector>

if there is a mismatch in size, <vector> will be cycled. To remove a column, assign NULL to the column as follows.

```
<df>$<column> <- NULL
```

10.5.3 Filter Data from Data Frame

Filtering allows selecting rows from a data frame that meet specific criteria. A true-false vector can be used as a filter as follows.

<filter-name> <- <true-false-vector> # use true-false vector as filter <df>[filter,] # implement filter on data frame

An example is given below.

```
> library(datasets)
> filter <- iris$Sepal.Length >= 7
> print(iris[filter,])
   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
51
           7.0
                      3.2
                                  4.7
                                              1.4 versicolor
103
           7.1
                       3.0
                                  5.9
                                              2.1 virginica
                                              2.1 virginica
106
           7.6
                      3.0
                                  6.6
108
           7.3
                       2.9
                                  6.3
                                              1.8 virginica
           7.2
                                  6.1
                                              2.5 virginica
110
                      3.6
118
           7.7
                      3.8
                                  6.7
                                              2.2 virginica
119
           7.7
                      2.6
                                              2.3 virginica
                                  6.9
123
           7.7
                       2.8
                                  6.7
                                              2.0 virginica
126
                      3.2
                                              1.8 virginica
           7.2
                                  6.0
130
                      3.0
           7.2
                                  5.8
                                              1.6 virginica
131
           7.4
                       2.8
                                  6.1
                                              1.9 virginica
132
           7.9
                      3.8
                                  6.4
                                              2.0 virginica
136
           7.7
                       3.0
                                   6.1
                                              2.3 virginica
> filter <- iris$Sepal.Length >= 7 & iris$Sepal.Width >= 3.5
> print(iris[filter,])
   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
110
           7.2
                       3.6
                                   6.1
                                              2.5 virginica
118
           7.7
                       3.8
                                   6.7
                                              2.2 virginica
132
           7.9
                       3.8
                                   6.4
                                              2.0 virginica
```

As shown above, it is possible to use &, | to form a more complex filter. The commands can be merged together as follows.

> p	rint(iris[iri	s\$Sepal.Leng	gth >= 7,])		
	Sepal.Length	Sepal.Width	Petal.Length	Petal.Wi	dth Species
51	7.0	3.2	4.7	1.4	versicolor
103	7.1	3.0	5.9	2.1	virginica
106	7.6	3.0	6.6	2.1	virginica
108	7.3	2.9	6.3	1.8	virginica
110	7.2	3.6	6.1	2.5	virginica
118	7.7	3.8	6.7	2.2	virginica
119	7.7	2.6	6.9	2.3	virginica
123	7.7	2.8	6.7	2.0	virginica
126	7.2	3.2	6.0	1.8	virginica
130	7.2	3.0	5.8	1.6	virginica
131	7.4	2.8	6.1	1.9	virginica
132	7.9	3.8	6.4	2.0	virginica
136	7.7	3.0	6.1	2.3	virginica

```
> nrow(iris[iris$Sepal.Length >= 7,]) # count the result number
[1] 13
```

10.5.4 Create Data Frames

So far we have been using existing iris data frame in the case study. To create a new data frame from scratch, use function data.frame() as follows.

```
<df> <- data.frame(<vector>, ...) # add a column colnames(<df>) <- c("<column-name>", ...)
```

or

```
<df> <- data.frame(<column-name> = <vector>, ...)
```

An example is given below, where a data frame of mortgage price at 3 types of areas, namely "CBD", "city" and "suburbs", are is created. Arbitrary data is used.

which gives the following result

```
> head(mortgage_price)
Region Size Price
1   CBD 81.84889 1154873.0
2   CBD 77.78946 831468.7
3   CBD 84.60477 735265.2
4   CBD 62.42625 829977.5
5   CBD 65.42723 933851.3
6   CBD 82.43867 1208589.0
```

A data frame can also be created from two existing data frames by joining them together. It works like the "JOIN" function in SQL, and it supports "INNER JOIN", "LEFT JOIN", "RIGHT JOIN" and "OUTER JOIN". The syntax is given below.

```
<df2>,
  by.x = <column-in-df1>,
  by.y = <column-in-df2>,
  all=FALSE,
  all.x = all, # left join, by default FALSE
  all,y = all, # right join, by default FALSE
  sort = TRUE # sort by the join column
```

In case the two data frames have duplicated columns other than the joining columns pair, use <df>\$<column> <- NULL to remove those columns.

10.6 Basic Data Visualizations Using qplot()

The package ggplot2 provides useful tools for visualization of a data frame. For example, both qplot() and ggplot() in ggplot2 provide plot function. Notice that in the late versions of ggplot2, qplot() is deprecated to encourage using of the more powerful ggplot(). Both functions are powerful enough to produce many different types of plots.

An example of using qplot() is given below, just to show some of its capabilities. Run the following codes to get Fig. 10.6. It can be seen that qplot() is smart enough to automatically choose plot type, background color, etc.

```
library(datasets)
library(ggplot2)
qplot(
    data=iris,
    x=Sepal.Length*Sepal.Width,
    y=Petal.Length*Petal.Width,
    color=Species,
    size=I(3),
    xlab = "Sepal Area",
    ylab = "Petal Area"
)
```

As a recap, the mortgage_price data frame created previously can be visualized as follows. Figures 10.7 and 10.8 can be obtained.

```
library(ggplot2)
rm(list=ls())
# create data frame
vec_region <- rep(c("CBD","City","Suburbs"),each = 100)
vec_size_cbd <- rnorm(100, 75, 10)
vec_size_city <- rnorm(100, 100, 15)
vec_size_suburbs <- rnorm(100, 150, 25)</pre>
```

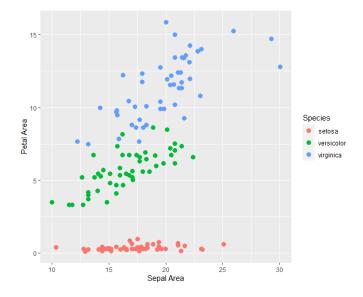
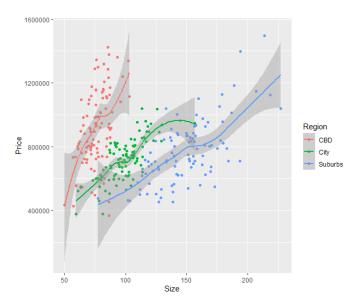


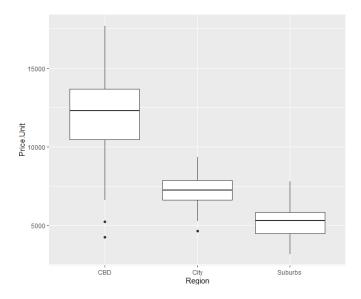
FIGURE 10.6 A demonstration of qplot.

10.7 Advanced Data Visualizations Using ggplot()

Function ggplot() is the main data visualization tool in ggolot2 package. It provides very flexible tools for data plotting.



 $\begin{tabular}{ll} FIGURE~10.7\\ A~{\rm demonstration~of~qplot~on~mortgage~price~data~frame.} \end{tabular}$



 $\begin{tabular}{ll} {\bf FIGURE~10.8} \\ {\bf A~second~demonstration~of~qplot~on~mortgage~price~data~frame.} \\ \end{tabular}$

10.7.1 Grammar of Graphics

As proposed by Leland Wilkinson's Grammar of Graphics, a chart shall contain multiple independent and reusable layers including

- Data. Original data.
- Aesthetics. How data is mapped to the graph, for example, by dots, curves, color blocks or lines/bars of different length.
- Geometries. The color or shape of each element in the graph.
- Statistics. Information derived from the data.
- Facets. Subplots of different data sets, and how they are aligned and compared.
- Coordinates. The meaning and range of axis.
- Theme. Titles, labels, legends, etc.

A demonstrative Fig. 10.9 is given to illustrate the different layers in a chart.

10.7.2 Data, Aesthetics and Geometries Layers

Function ggplot() is a very good practice of implementing the above chart design and plotting philosophy. A simple example for ggplot(), just for quick demonstration purpose, is given below.

where aes() is used to build mappings in the aesthetics.

An interesting fact when using ggplot() is that, when adding a layer to the chat, the layer is literally added to ggplot(). In the program, this step by step build up an object, where ggplot() provides the most basic layers. Therefore, the above simple example is equivalent to

and the added layers are able to inherit the aesthetics settings, if it is not overwritten. And speaking of overwriting, even the x and y axis can be overwritten. The displaying name of the labels can be overwritten by stack xlab("") and ylab("") into the chart.

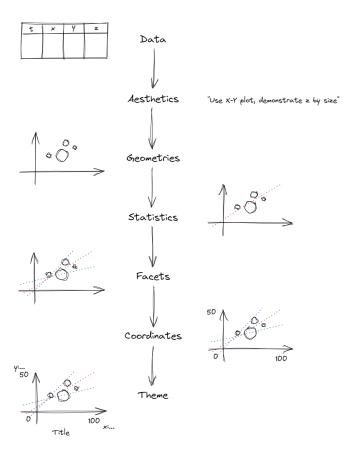


FIGURE 10.9 Multiple layers in chart design.

TABLE 10.8
Commonly used commands for data frame exploration.

Geom	Description
<pre>geom_point()</pre>	Scatter plots and dot plots.
<pre>geom_line()</pre>	Line plots.
<pre>geom_bar()</pre>	Bar plots.
<pre>geom_histogram()</pre>	Histograms.
<pre>geom_boxplot()</pre>	Box plots.
<pre>geom_violin()</pre>	Violin plots.
<pre>geom_density()</pre>	Density plots.
<pre>geom_density2d()</pre>	2-dimensional Density plots.
<pre>geom_text()</pre>	Text Annotation.
<pre>geom_label()</pre>	Label on the observations.

Function ggplot() provides many choices for geometries. The most commonly used ones are summarized in Table 10.8.

10.7.3 Statistics Layers

Similar to the case of geometries layers, statistics layers can also be stacked to ggplot(). As introduced earlier, statistics layers are often "add-on" layers that derives statistical features from the data to provide additional insights. Many functions in Table 10.8 are statistics related. More details are given below.

Consider using geom_boxplot() to visualize the mortgage_price data frame that was used in the previous section. Examples are given below.

```
library(ggplot2)
# create data frame
vec_region <- rep(c("CBD","City","Suburbs"), each = 100)</pre>
vec_size_cbd <- rnorm(100, 75, 10)</pre>
vec_size_city <- rnorm(100, 100, 15)</pre>
vec_size_suburbs <- rnorm(100, 150, 25)</pre>
vec_size = c(vec_size_cbd, vec_size_city, vec_size_suburbs)
vec_price_cbd <- vec_size_cbd*rnorm(100, 12500, 2500)</pre>
vec_price_city <- vec_size_city*rnorm(100, 7500, 1000)</pre>
vec_price_suburbs <- vec_size_suburbs*rnorm(100, 5000, 1000)</pre>
vec_price <- c(vec_price_cbd, vec_price_city, vec_price_suburbs)</pre>
mortgage_price <- data.frame(Region = vec_region, Size = vec_size,
    Price = vec_price)
rm(vec_region, vec_size_cbd, vec_size_city, vec_size_suburbs, vec_size,
     vec_price_cbd, vec_price_city, vec_price_suburbs, vec_price)
# processing
mortgage_price$Price.Unit <- mortgage_price$Price / mortgage_price$Size
p <- ggplot(data=mortgage_price, aes(x=Region, y=Price.Unit, color=
```

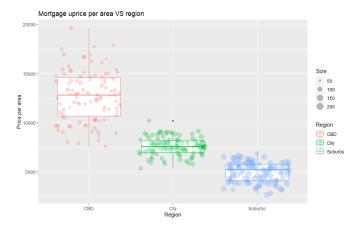


FIGURE 10.10

An example of box plot of the mortgage price data frame using ggplot() and geom_boxplot().

```
Region)) + ggtitle("Mortgage uprice per area VS region") + xlab("
Region") + ylab("Price per area")
p + geom_boxplot() + geom_jitter(aes(size=Size, color=Region), alpha
=0.25)
```

and the result is shown in Fig. 10.10. Notice that ggtitle(), xlab(), ylab(), alpha are used in the plot. They are self-explanatory. A new geometry geom_jitter() is used, which works similarly with geom_point() except the additional vibration in the horizontal axis which makes the points clearer to see.

Function geom_smooth() is widely used for curve fitting. An example is given below.

```
library(ggplot2)
# generate data
t <- 1:500
var1 <- 1.5*t + rnorm(500, 0, 100)
var2 <- 0.5*t + rnorm(500, 200, 10) + t^1.3*rnorm(500, 0, 0.1)
df <- data.frame(t=t, x=var1, y=var2)
# plot data
p <- ggplot(data=df) +
ggtitle("Plot of x and y VS t.") +
xlab("t") +
ylab("x and y") +
geom_point(aes(x=t, y=x), color="blue", shape=1, size=1.5) +
geom_smooth(aes(x=t, y=x), color="red", shape=2, size=1.5) +
geom_smooth(aes(x=t, y=y), color="red", shape=2, size=1.5) +
geom_smooth(aes(x=t, y=y), color="red")
p</pre>
```

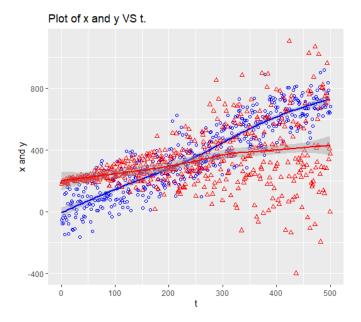


FIGURE 10.11

An example of using geom_smooth() for scatter point fitting.

Do note that aesthetics needs to be given to geom_smooth() in the above example. This is because aesthetics is not given in the base ggolot(). Notice that geom_smooth() can inherit aesthetics from the previous ggplot(), but not from the previous geom_point(). The plot is given by Fig. 10.11.

More functions similar to geom_smooth() are summarized in Table 10.9.

10.7.4 Facets Layers

The facets layer allows subplot of data. Consider the following example, where the distribution of mortgage price is studied using histogram. The following code can be used to plot the result in a single plot without the facets layer. The plot is given in Fig. 10.12.

TABLE 10.9Functions that fit smooth lines to scatter points

Function	Description
loess()	Non-parametric method for fitting a smooth line to a
<pre>smooth.spline()</pre>	scatter plot using locally weighted regression algorithm. Fits a smoothing spline to the data, which is a type of regression spline where the degree of smoothing is
	chosen automatically by cross-validation.
lm()	Linear Model, fits a linear relationship between inde-
	pendent and dependent variables by minimizing the
	residuals between the data points and the line.
glm()	Generalized Linear Model, similar to linear model, but
	it allows different distribution of error other than normal.
gam()	Generalized Additive Model, it is similar to GLM, but
	it allows non-parametric smooth functions to be added
	to the linear predictor.
<pre>geom_smooth()</pre>	A function in ggplot2 that is used to add a smooth line
	to a scatter plot, it uses method = "loess" by default
	but also allow to use other smoothing method like lm,
	gam etc.

To use facets layer, revise the code as follows. Notice that facet_grid() is added to the plot, and its input <column>~. or .~<column> (it is okay to use <column1>~<column2> as well) decide the design of the subplots (how to arrange the rows and columns of the subplots).

```
library(ggplot2)
# create data frame
Region = rep(c("CBD","City","Suburbs"), each = 500)
vec_size = list(vec_size_cbd = rnorm(500, 75, 10),
```

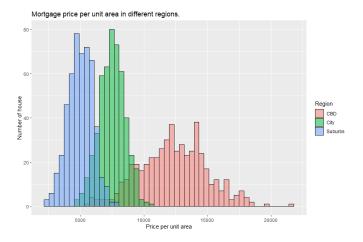


FIGURE 10.12

An example of histogram plot of house price per unit area in different regions in a single plot.

```
vec_size_city = rnorm(500, 100, 15),
               vec_size_suburbs = rnorm(500, 150, 25))
vec_price = list(vec_price_cbd = vec_size$vec_size_cbd*rnorm(500,
    12500, 2500),
               vec_price_city = vec_size$vec_size_city*rnorm(500,
                   7500, 1000),
               vec_price_suburbs = vec_size$vec_size_suburbs*rnorm
                    (500, 5000, 1000))
mortgage_price <- data.frame(Region = Region,
                          Size = unlist(vec_size),
                          Price = unlist(vec_price))
mortgage_price$Region <- as.factor(mortgage_price$Region)</pre>
mortgage_price$Price.Unit <- mortgage_price$Price / mortgage_price$Size</pre>
# plot data
p <- ggplot(data=mortgage_price, aes(x=Price.Unit))</pre>
p <- p + geom_histogram(aes(fill=Region), bins=50, color="black", alpha
    =0.5, position="identity") +
 ggtitle("Mortgage price per unit area in different regions.") +
 xlab("Price per unit area") +
 ylab("Number of house")
p + facet_grid(Region~.) # put subplots for different regions in rows
p + facet_grid(.~Region) # put subplots for different regions in
```

The results are given in Figs. 10.13 and 10.14, depending on the subplot designs.

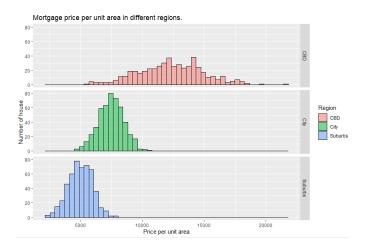


FIGURE 10.13

Use facets to plot the histogram of price per unit are of the house in different regions (subplots in rows).

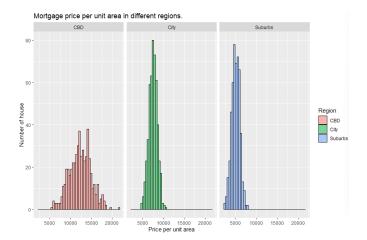


FIGURE 10.14

Use facets to plot the histogram of price per unit are of the house in different regions (subplots in columns).

10.7.5 Coordinates Layers

Coordinate control is important. The coordinate layer allows setting limits to the axis and zooming in to the chart. An example of adding coordinates layers to a plot is given as follows. The same mortgage price data frame is used for illustration.

```
library(ggplot2)
# create data frame
Region = rep(c("CBD", "City", "Suburbs"), each = 500)
vec_size = list(vec_size_cbd = rnorm(500, 75, 10),
               vec_size_city = rnorm(500, 100, 15),
               vec_size_suburbs = rnorm(500, 150, 25))
vec_price = list(vec_price_cbd = vec_size$vec_size_cbd*rnorm(500,
    12500, 2500),
               vec_price_city = vec_size$vec_size_city*rnorm(500,
                    7500, 1000),
               vec_price_suburbs = vec_size$vec_size_suburbs*rnorm
                    (500, 5000, 1000))
mortgage_price <- data.frame(Region = Region,
                          Size = unlist(vec_size),
                          Price = unlist(vec_price))
mortgage_price$Region <- as.factor(mortgage_price$Region)</pre>
mortgage_price$Price.Unit <- mortgage_price$Price / mortgage_price$Size
# plot data
p <- ggplot(data=mortgage_price, aes(x=Size, y=Price))
p <- p + geom_point(aes(color=Region))</pre>
p + xlim(50, 200) + ylim(400000, 1200000) # first chart
 + coord_cartesian(xlim=c(50, 200), ylim = c(400000, 1200000)) #
    second chart
```

where notice that two charts are generated. The first chart using xlim(), ylim removes all samples outside the boundary from the chart. While in the second chart using coord_cartesian(), all samples preserves and the chart zooms in towards the boundary. The results are given in Figs. 10.15 and 10.16, respectively. The difference can be observed near the boundary.

10.7.6 Themes Layers

Theme layers mainly refer to titles, labels, and other comments on the chart that help with understanding the content of the chart. As already demonstrated in previous examples, use xlab(), ylab() to add labels, ggtitle() to add title.

Use theme() to change the themes of the labels. An example is given below.

```
library(ggplot2)
# create data frame
Region = rep(c("CBD","City","Suburbs"), each = 500)
```

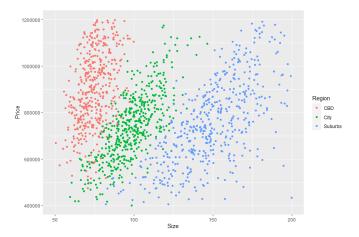


FIGURE 10.15 Add coordinates layer using xlim() and ylim().

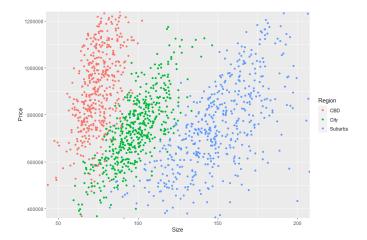


FIGURE 10.16 Add coordinates layer using coord_cartesian().

R Basics 91

```
vec_size = list(vec_size_cbd = rnorm(500, 75, 10), vec_size_city =
    rnorm(500, 100, 15), vec_size_suburbs = rnorm(500, 150, 25))
vec_price = list(vec_price_cbd = vec_size$vec_size_cbd*rnorm(500,
    12500, 2500),
               vec_price_city = vec_size$vec_size_city*rnorm(500,
                   7500, 1000),
               vec_price_suburbs = vec_size$vec_size_suburbs*rnorm
                    (500, 5000, 1000))
mortgage_price <- data.frame(Region = Region,
                          Size = unlist(vec_size),
                          Price = unlist(vec_price))
mortgage_price$Region <- as.factor(mortgage_price$Region)</pre>
mortgage_price$Price.Unit <- mortgage_price$Price / mortgage_price$Size</pre>
# plot data
p <- ggplot(data=mortgage_price, aes(x=Price.Unit))</pre>
p <- p + geom_histogram(aes(fill=Region), bins=50, color="black", alpha
    =0.5, position="identity")
p + ggtitle("Mortgage price per unit area in different regions.") +
 xlab("Price per unit area") +
 ylab("Number of house") +
 theme(axis.title.x = element_text(color = "DarkGreen", size=15),
       axis.title.y = element_text(color = "DarkRed", size=15),
       axis.text.x = element_text(size=10),
       axis.text.y = element_text(size=10),
       legend.title = element_text(size=10),
       legend.text = element_text(size=8),
       legend.position = c(1,1), # right top corner of chart
       legend.justification = c(1,1), # legend align point
       plot.title = element_text(color = "DarkBlue", size = 15)
```

The resulted chart is given in Fig. 10.17. Compare it with Fig. 10.12 to see the differences by applying theme() in the themes layer.

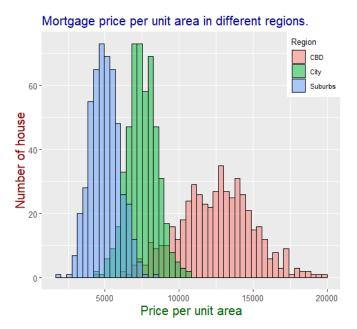


FIGURE 10.17 Mortgage price chart with theme.

11

R Advanced

CONTENTS

11.1	Data Preparation				
	11.1.1	Data Type Conversion	93		
	11.1.2	Handling Missing Data	95		
11.2	Connec	ctivity with Data Sources	98		

This chapter introduces advanced skills of using R language, including data preparation techniques, the use of list, etc.

11.1 Data Preparation

The data downloaded from sensors usually needs to go through pre-processing procedures such as filtering, normalization, etc., before it can be used by a controller, an AI engine, or for further statistics analysis. Data preparation including data tidy is one of the most tedious and time consuming parts when using R for data analysis. The section introduces useful techniques helpful with data preparation.

11.1.1 Data Type Conversion

It is important that the data types of all the columns meet expectation, especially for numeric and factor (categorical) data types. Use str(<df>) to check the column data types of a data frame, and if necessary convert data types using the following commands.

```
<df>$<column> <- factor(<df>$<column>) # character/numeric to factor <df>$<column> <- as.numeric(<df>$<column>) # character to numeric <df>$<column> <- as.numeric(as.character(<df>$<column>)) # factor to numeric
```

Notice that when converting factor type to other types, R may deal with the factor using the underlying "factorization integers" instead of the factor item names. An example is given below. It can be seen that the original 5.1, after being converted to factor then back to numeric, becomes 9. This is because the factorization integer for 5.1 is 9, as shown by printing my_factor to the console.

```
> library(datasets)
> iris$Sepal.Length
 [1] 5.1 4.9 4.7 4.6 5.0 5.4 4.6 5.0 4.4 4.9 5.4 4.8 4.8 4.3 5.8 5.7
 [17] 5.4 5.1 5.7 5.1 5.4 5.1 4.6 5.1 4.8 5.0 5.0 5.2 5.2 4.7 4.8 5.4
[33] 5.2 5.5 4.9 5.0 5.5 4.9 4.4 5.1 5.0 4.5 4.4 5.0 5.1 4.8 5.1 4.6
[49] 5.3 5.0 7.0 6.4 6.9 5.5 6.5 5.7 6.3 4.9 6.6 5.2 5.0 5.9 6.0 6.1
[65] 5.6 6.7 5.6 5.8 6.2 5.6 5.9 6.1 6.3 6.1 6.4 6.6 6.8 6.7 6.0 5.7
[81] 5.5 5.5 5.8 6.0 5.4 6.0 6.7 6.3 5.6 5.5 5.5 6.1 5.8 5.0 5.6 5.7
[97] 5.7 6.2 5.1 5.7 6.3 5.8 7.1 6.3 6.5 7.6 4.9 7.3 6.7 7.2 6.5 6.4
 [113] \ \ 6.8 \ \ 5.7 \ \ 5.8 \ \ 6.4 \ \ 6.5 \ \ 7.7 \ \ 7.7 \ \ 6.0 \ \ 6.9 \ \ 5.6 \ \ 7.7 \ \ 6.3 \ \ 6.7 \ \ 7.2 \ \ 6.2 \ \ 6.1 
[129] 6.4 7.2 7.4 7.9 6.4 6.3 6.1 7.7 6.3 6.4 6.0 6.9 6.7 6.9 5.8 6.8
[145] 6.7 6.7 6.3 6.5 6.2 5.9
> my_factor <- factor(iris$Sepal.Length)</pre>
 my_numeric <- as.numeric(my_factor)</pre>
 my_numeric
 [1] 9 7 5 4 8 12 4 8 2 7 12 6 6 1 16 15 12 9 15 9 12 9
 [23] 4 9 6 8 8 10 10 5 6 12 10 13 7 8 13 7 2 9 8 3 2 8
[45] 9 6 9 4 11 8 28 22 27 13 23 15 21 7 24 10 8 17 18 19 14 25
[67] 14 16 20 14 17 19 21 19 22 24 26 25 18 15 13 13 16 18 12 18 25 21
[89] 14 13 13 19 16 8 14 15 15 20 9 15 21 16 29 21 23 33 7 31 25 30
[111] 23 22 26 15 16 22 23 34 34 18 27 14 34 21 25 30 20 19 22 30 32 35
[133] 22 21 19 34 21 22 18 27 25 27 16 26 25 25 21 23 20 17
> typeof(my_factor)
[1] "integer"
> my_factor
 [1] 5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 5.4 4.8 4.8 4.3 5.8 5.7
[17] 5.4 5.1 5.7 5.1 5.4 5.1 4.6 5.1 4.8 5 5 5.2 5.2 4.7 4.8 5.4
[33] 5.2 5.5 4.9 5 5.5 4.9 4.4 5.1 5 4.5 4.4 5 5.1 4.8 5.1 4.6
[49] 5.3 5 7 6.4 6.9 5.5 6.5 5.7 6.3 4.9 6.6 5.2 5 5.9 6 6.1
[65] 5.6 6.7 5.6 5.8 6.2 5.6 5.9 6.1 6.3 6.1 6.4 6.6 6.8 6.7 6 5.7
[81] 5.5 5.5 5.8 6 5.4 6 6.7 6.3 5.6 5.5 5.5 6.1 5.8 5 5.6 5.7
[97] 5.7 6.2 5.1 5.7 6.3 5.8 7.1 6.3 6.5 7.6 4.9 7.3 6.7 7.2 6.5 6.4
[113] 6.8 5.7 5.8 6.4 6.5 7.7 7.7 6 6.9 5.6 7.7 6.3 6.7 7.2 6.2 6.1
[129] 6.4 7.2 7.4 7.9 6.4 6.3 6.1 7.7 6.3 6.4 6 6.9 6.7 6.9 5.8 6.8
[145] 6.7 6.7 6.3 6.5 6.2 5.9
35 Levels: 4.3 4.4 4.5 4.6 4.7 4.8 4.9 5 5.1 5.2 5.3 5.4 5.5 ... 7.9
```

When converting factor to other types, special caution is required. To convert a factor to other types such as numeric, consider converting it to character first as given in the following example.

```
> my_numeric <- as.numeric(as.character(my_factor))
> my_numeric
[1] 5.1 4.9 4.7 4.6 5.0 5.4 4.6 5.0 4.4 4.9 5.4 4.8 4.8 4.3 5.8 5.7
[17] 5.4 5.1 5.7 5.1 5.4 5.1 4.6 5.1 4.8 5.0 5.0 5.2 5.2 4.7 4.8 5.4
[33] 5.2 5.5 4.9 5.0 5.5 4.9 4.4 5.1 5.0 4.5 4.4 5.0 5.1 4.8 5.1 4.6
```

R Advanced 95

```
[49] 5.3 5.0 7.0 6.4 6.9 5.5 6.5 5.7 6.3 4.9 6.6 5.2 5.0 5.9 6.0 6.1 [65] 5.6 6.7 5.6 5.8 6.2 5.6 5.9 6.1 6.3 6.1 6.4 6.6 6.8 6.7 6.0 5.7 [81] 5.5 5.5 5.8 6.0 5.4 6.0 6.7 6.3 5.6 5.5 5.5 6.1 5.8 5.0 5.6 5.7 [97] 5.7 6.2 5.1 5.7 6.3 5.8 7.1 6.3 6.5 7.6 4.9 7.3 6.7 7.2 6.5 6.4 [113] 6.8 5.7 5.8 6.4 6.5 7.7 7.7 6.0 6.9 5.6 7.7 6.3 6.7 7.2 6.2 6.1 [129] 6.4 7.2 7.4 7.9 6.4 6.3 6.1 7.7 6.3 6.4 6.0 6.9 6.7 6.9 5.8 6.8 [145] 6.7 6.7 6.3 6.5 6.2 5.9
```

where my_factor is generated previously.

It is possible for some columns in the data frame to look like a factor and a character string, but indeed should be handled as numeric values. For example, \$\$6,125.50 in many occasions should be treated just as 6125.5. These factor or character values cannot be converted to numeric values directly.

In such cases, consider using sub() or gsub() to replace patterns in a character, then convert it into numeric values. Notice that sub() replaces only the first encounter of the pattern, while gsub() replaces all the encounters. An example of using gsub() is given below.

```
> money_character <- c("S$6,273.15", "S$215.3", "S$8,987,756.00")
> typeof(money)
[1] "character"
> a <- gsub(",", "", money) # replace "," with ""
> a <- gsub("S\\$", "", a) # replace "S$" with ""
> money_numeric <- as.numeric(a)
> money_numeric
[1] 6273.15 215.30 8987756.00
> typeof(money_numeric)
[1] "double"
```

where notice that \$ is a special character defined in R, and to escape from that \\\$ is used. Notice that applying sub() and gsub() on a factor automatically converts it to character as a hidden step.

11.1.2 Handling Missing Data

There can be missing data in the data frame. There are a few ways to deal with missing data as follows.

- If the missing data can be derived from other columns, derive the missing data and fill in the blanks.
- If the missing data does not affect the rest analysis, leave it blank.
- Delete the row.
- Use interpolations to fill in the blank.
- Use correlations and similarities to fill in the blank.
- Argument a new column add a "data-missing" flag to that row.

Flag Missing Data using NA

In R, NA is a special variable used to indicate a missing value. A general idea is to "flag" the missing data in the original data source, what format it may look like, using NA during or after the data importing. After that, use a special program in R to filter NA and deal with them separately. Sometimes a blank string "" that we would expected to be treated as NA is not treated as so. To fix that, while importing the data frame (say, from a CSV file), use the following

```
df <- read.csv("<csv-name>", na.string=c("<pattern>", ...))
```

where "<pattern>" are the patterns in the original file to be replaced by NA, for instance, "", "ERROR", etc.

Locate and Filter NA

Notice that NA is treated as a logical data type in addition to TURE and FALSE in a logical expression. These operations involving NA often return NA. Examples are given below.

```
> typeof(NA)
[1] "logical"
> TRUE == 1 # TRUE is equivalent with 1
[1] TRUE
> TRUE == 2
[1] FALSE
> FALSE == 0 # FALSE is equivalent with 0
[1] TRUE
> FALSE == -1
[1] FALSE
> TRUE == FALSE
[1] FALSE
> NA == NA
[1] NA
> NA == TRUE
[1] NA
> NA == FALSE
[1] NA
```

This applies to filtering. In filtering, when a variable of value NA is asserted with a criterion, the return is most likely NA. The filter often does not know how to deal with NA, hence it would simply return the rows with NA anyway. This can become inconvenient sometimes. An example is given below.

R Advanced 97

The return is as follows.

I	Region	Size	Price F	Price.Unit	
vec_size_cbd1	CBD	73.47779	947130.9	12890.031	
vec_size_cbd2	CBD	NA	678842.3	NA	
vec_size_cbd3	CBD	85.06748	1261029.7	14823.875	
vec_size_cbd4	CBD	92.35454	NA	NA	
vec_size_cbd5	CBD	71.06649	1276168.6	17957.388	
vec_size_city1	City	68.59874	491912.0	7170.861	
vec_size_city2	City	118.39441	804794.5	6797.572	
vec_size_city3	City	NA	534591.5	NA	
vec_size_city4	City	95.57428	583044.7	6100.436	
vec_size_city5	City	74.45356	468788.0	6296.382	
vec_size_suburbs1	Suburb	s 136.884	32 939436.	6 6862.997	
vec_size_suburbs2	Suburb	s 136.577	99 810070.	0 5931.190	
vec_size_suburbs3	Suburb	s 189.665	86 561089.	2 2958.303	
vec_size_suburbs4	Suburb	s 195.974	76 913572.	2 4661.683	
vec_size_suburbs5	Suburb	s 190.630	082 NA	NA	

```
Region
                                    Price Price.Unit
                            Size
vec_size_cbd1
                   CBD 73.47779 947130.9 12890.031
                   CBD 85.06748 1261029.7 14823.875
vec_size_cbd3
                   <NA>
                              NA
                                      NA
vec_size_cbd5
                   CBD 71.06649 1276168.6 17957.388
                   City 118.39441 804794.5 6797.572
vec_size_city2
vec_size_suburbs1 Suburbs 136.88432 939436.6 6862.997
vec_size_suburbs2 Suburbs 136.57799 810070.0 5931.190
vec_size_suburbs4 Suburbs 195.97476 913572.2 4661.683
                              NA
                                                 NA
NA.1
                   <NA>
                                       NA
```

In the above example, the intension of the program is to find all the houses with price larger than 750000. The program is able to filter out those houses

cheaper than the threshold. However, there are two rows of NA returned, as explained earlier.

Use the following to filter for all rows with/without at least one NA.

```
<df>[complete.cases(<df>),] # all complete rows
<df>[!complete.cases(<df>),] # all incomplete rows
```

where complete.cases(<df>) returns a list made up of TRUE and FALSE indicating whether the associated row is complete or now.

11.2 Connectivity with Data Sources

This section introduces the connectivity of R to the data sources, such as a file, or a database.

12

R Practice

CONTENTS

13

Python Basics

CONTENTS

13.1	NumPy	101
13.2	SciPy	103
13.3	Matplotlib and Seaborn	104
	13.3.1 Matplotlib	104
	13.3.2 Seaborn	105
13.4	Pandas	107
	13.4.1 Data Importing	108
	13.4.2 Series and Data Frame	110

Python has been increasingly popular for data science in the past few years. Many libraries and tools have been developed for Python to enhance its data analysis and visualization capabilities, just to name a few, numpy, scipy, scikit-learn, pandas, matplotlib tensorflow and pytorch.

This chapter together with a few consequent chapters introduces commonly used tools that data science adopts using Python. This part of the notebook is more application driven, and only the basic implementations are introduced. We are not digging into the theory supporting machine learning and artificial intelligence.

It has been increasingly popular today to use Python together with Conda and jupyter-lab/jupyter notebook. Conda is an open-source language-agnostic package and environment management system. Jupyter notebook is an interactive computing platform for Python and other computer programming languages. The detailed introduction to the installation and usage of Conda and Jupyter notebook is not covered here. They are used when demonstrating the examples in this chapter and the consequent chapters.

13.1 NumPy

When comes to any sort of numerical computation, one of the most popular Python packages is definitely NumPy. NumPy allows quickly deployment of

numerical vectors, matrices and tensors, as well as associated efficient numerical calculations. It is the "MATLAB" package in Python.

Details of NumPy can be found at *numpy.org*.

The following commands can be used to create NumPy arrays and matrices

```
import numpy as np
# create numpy array from python list
x = np.array([1, 2, 3, 4, 5]) # 1d
x = np.array([[1, 2], [3, 4]]) #2d
# create numpy array/matrix using built-in functions
x = np.arange(0, 5) # array([0, 1, 2, 3, 4])
x = np.linspace(0, 5, 3) # array([0, 2.5, 5])
x = np.zeros(5) # 1d zero vector
x = np.zeros([5, 5]) # 2d zero matrix
x = np.ones(5)
x = np.ones([5, 5])
x = np.eye(5)
# create random vector/matrix
np.random.seed(1) # set seed; optional
x = np.random.rand(5, 5) # uniform distribution [0, 1)
x = np.random.randn(5, 5) # standard normal distribution N(0, 1)
# reshape
x = np.random.randn(16)
y = x.reshape(4, 4)
y = x.reshape(1, 16) # return is always 2d matrix format, not 1d vector
```

Aggregation functions are defined. These functions are used to calculate maximum, minimum, sum, etc., of a vector or a matrix. Examples are given below.

```
import numpy as np

x = np.random.randn(10)

xmax = np.max(x)

xargmax = np.argmax(x)

xmin = np.min(x)

xargmin = np.argmin(x)

xsum = np.sum(x)

xprod = np.prod(x)

xmean = np.mean(x)

xstd = np.std(x)

xvar = np.var(x)

xmedian = np.median(x)
```

Python Basics 103

When some of these functions such as sum() are applied to matrix, it is important to specify the axis along which the calculation would be carried out.

```
import numpy as np
x = np.array([[1,2,3,4,5], [6,7,8,9,10]])
np.sum(x, axis=0) # array([7, 9, 11, 13, 15])
np.sum(x, axis=1) # array([15, 40])
```

Accessing (reading and modifying) values in numpy arrays or matrices using the index. Examples are given below. Notice that all examples are about vector accessing.

Matrix accessing is a bit more tricky. A matrix in NumPy is stored like a nested array. Examples are given below to access a single item, a row, and a column or a matrix.

```
import numpy as np
x = np.array([[1, 2, 3], [4, 5, 6], [7, 8, 9]])
x[1, 2] # 6 (2nd row, 3rd column)
x[0] # 1st row as a numpy array
x[:,1] # 2nd column as a numpy array
```

NumPy provides efficient and convenient vector and matrix level calculations, such as broad casting, matrix multiplication, etc. Broad casting essentially allows element-by-element basis adding, subtracting or assigning scalar value to a vector or a matrix. An example is given below.

```
import numpy as np
x = np.array([1, 2, 3, 4, 5])
x[:] = 10 # all elements become 10
```

NumPy array and matrix support boolean selection. An example is given below.

NumPy supports both element-by-element or vector-level operations, including +, -, *, /, **, numpy.sqrt(), numpy.log(), etc. Most of these operators are executed element-by-element by default.

13.2 SciPy

SciPy is a library collections of "comprehensive" algorithms widely used in scientific and technical calculations. Notice that NumPy also has built in basic and commonly algorithms in the library, such as FFT. For those algorithms not included in NumPy, there is a chance that it is in SciPy, such as K-means clustering.

A detailed documentation of SciPy functions put into different categories are given here docs.scipy.org/doc/scipy/reference/.

13.3 Matplotlib and Seaborn

Data visualization is important through out the entire data analysis process. It is about not only demonstrating the results to the audiences, but also helping the developers to understand the data and improving the inefficient designs in the pipeline. Matplotlib and Seaborn are two important data visualization libraries. They are briefly introduced in this section.

13.3.1 Matplotlib

Matplotlib is the "basic" visualization library. Many data visualization features provided by other packages such as pandas are essentially realized using Matplotlib internally.

Simple line and scatter plots using Matplotlib can be drawn easily. Examples are given below. Pandas series is used as the axis to the plots, but in reality they can be Python arrays or NumPy arrays. The scatter plot used in the example is given in Fig. 13.1.

```
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

index_s = pd.Series([1, 2, 3, 4, 5, 6, 7, 8, 9, 10])
values_s = pd.Series([1, 1, 2, 3, 5, 8, 13, 21, 34, 55])
df_dict = {
        "F_Index": index_s,
        "F_Values": values_s
}
fibonacci = pd.DataFrame(df_dict)
plt.plot(index_s, values_s) # line
plt.scatter(fibonacci["F_Index"], fibonacci["F_Values"], color="red") #
        scatter
```

Python Basics 105

```
plt.title("Fibonacci∟Series")
plt.xlabel("n")
plt.ylabel("f(n)")
```

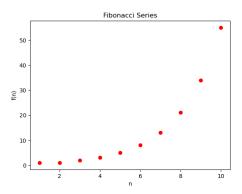


FIGURE 13.1

Plot Fibonacci series as scatter plot.

The presentation of the plot can be customized. For example, use plt.xlim(x1, x2), plt.ylim(y1, y2) to change the axis limitations, etc. It is possible to change curve color, size, style, and marker size, style, etc.

13.3.2 Seaborn

Seaborn is another visualization library built on top of the Matplotlib library. It tries to standardize the plots with a simple function call. It scarifies some flexibility that Matplotlib offers, but makes plotting easier.

An example of using Seaborn to plot a histogram is given below. The result is given in Fig. 13.2.

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

x = np.random.randn(1000)
plt.figure()
sns.histplot(x, color="red")
plt.title("Histogram_Example")
plt.xlabel("x")
plt.ylabel("Count")
```

An example of using Seaborn to plot a count plot for discrete values (usually category values) is given below. Notice that the attribute whose values are put into categories is assigned to x of ${\tt seaborn.countplot()}$. The result is given in Fig. 13.3.

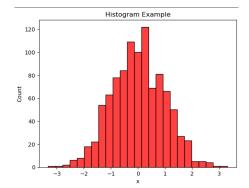


FIGURE 13.2

Histogram plot using Seaborn.

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

response = ["yes", "no", "no", "yes", "yes", "yes", "yes", "yes", "yes", "yes", "yes", "yes", "no", "yes", "NA", "NA"]
plt.figure()
sns.countplot(x=response, color="blue")
plt.title("Count_Plot_Example")
plt.xlabel("response")
plt.ylabel("Count")
```

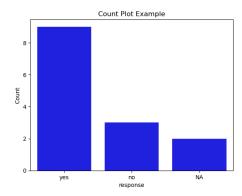


FIGURE 13.3

Count plot using Seaborn.

Bot plot gives the maximum, minimum, and interquartile range (IQR) of the data. In the box plot, the data is split into 4 parts, namely $x < Q_1$

Python Basics 107

(0%-25%), $Q_1 < x < Q_2$ (25%-50%), $Q_2 < x < Q_3$ (50%-75%), and $Q_3 < x$ (75%-100%). The IQR gives the range between Q_1 and Q_3 , i.e., the half in the middle 25%-75%. An example is given below. The result is given in Fig. 13.4.

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

x = np.random.randn(1000)
plt.figure()
sns.boxplot(x, color="blue")
plt.title("Box_Plot_Example")
plt.xlabel("Input")
plt.ylabel("Value")
```

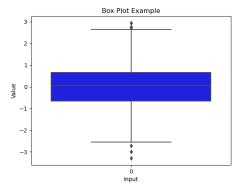


FIGURE 13.4

Box plot using Seaborn. The box gives IQR. The bars below and above the box give $Q_1-1.5\times IQR$ and $Q_3+1.5\times IQR$, respectively. The dots are outliers.

There are many more plot functions in Seaborn. For example, scatter plot seaborn.scatterplot() works similar with Matplotlib as shown by Fig. 13.1. The marker size and color can be adjusted to reflect different parameters, just like in R language. Pairplot seaborn.pairplot() generates a matrix of plots to demonstrate associations of columns in a data frame.

13.4 Pandas

Many Python packages provide functions to handle structured data such as tables, series, and data frames. Among all these packages, pandas is the all-

time star that is very widely used by developers and data scientists. With pandas, Python gains the ability to easily, flexibly and efficiently deal with data frames. The pandas package is introduced in this section.

A large portion of this chapter, including codes and examples, are from online resources such as *Data Analysis by Pandas and Python* on Udemy. My special thanks go to them.

The data frames used in the examples of this chapter may come from different public data resources. It is worth mentioning *kaggle.com*, a place where tens of thousands of data sets, code examples, and notebooks are collected and shared. Some data frames used in the examples come from Kaggle.

Two Python packages, numpy and pandas, are almost certainly used in all the examples to be presented in this chapter. Import these packages as follows.

```
import numpy as np
import pandas as pd
```

Unless otherwise mentioned, these packages are always assumed imported in this chapter. To display the packages version, use something like

```
print("Numpy version: {}.".format(np.__version__))
print("Pandas version: {}".format(pd.__version__))
```

13.4.1 Data Importing

Pandas provide variety of ways to import data from different resources, including plain texts, CSV files, databases, etc. One of the most commonly seen data sources is CSV files. Importing data from CSV files is introduced here.

Pandas provide .read_csv() function to import data from CSV files. Its basic usage is introduced below.

which reads all the information in the CSV table into a data frame as shown by Fig. 13.5. It is possible to import only selected columns as follows.

When no index column is specified, pandas will add an additional auto-incremental column and use it as the index column, as shown in Fig. 13.5 by the most-left column. When a column index is specified, pandas will use that column as index column. An example is given below. The result is shown in Fig. 13.6.

Python Basics 109

F	Rank	Title	Sales	Series	Platform(s)	Initial release date	Developer(s)	Publisher(s)
0	1	Minecraft	238000000	Minecraft	Multi-platform	November 18, 2011	Mojang Studios	Xbox Game Studios
1	2	Grand Theft Auto V	175000000	Grand Theft Auto	Multi-platform	September 17, 2013	Rockstar North	Rockstar Games
2	3	Tetris (EA)	100000000	Tetris	Multi-platform	September 12, 2006	EA Mobile	Electronic Arts
3	4	Wii Sports	82900000	Wii	Wii	November 19, 2006	Nintendo EAD	Nintendo
4	5	PUBG: Battlegrounds	75000000	PUBG Universe	Multi-platform	December 20, 2017	PUBG Corporation	PUBG Corporation
5	6	Mario Kart 8 / Deluxe	60460000	Mario Kart	Wii U / Switch	May 29, 2014	Nintendo EAD	Nintendo
6	7	Super Mario Bros.	58000000	Super Mario	Multi-platform	September 13, 1985	Nintendo R&D4	Nintendo
7	8	Red Dead Redemption 2	50000000	Red Dead	Multi-platform	October 26, 2018	Rockstar Studios	Rockstar Games
8	9	Pokémon Red / Green / Blue / Yellow	47520000	Pokémon	Multi-platform	February 27, 1996	Game Freak	Nintendo
9	10	Terraria	44500000	None	Multi-platform	May 16, 2011	Re-Logic	Re-Logic / 505 Games
10	11	Wii Fit / Plus	43800000	Wii	Wii	December 1, 2007	Nintendo EAD	Nintendo

FIGURE 13.5

The simplest data frame importing using pandas.

Publisher(s)	Sales	
		Title
Xbox Game Studios	238000000	Minecraft
Rockstar Games	175000000	Grand Theft Auto V
Electronic Arts	100000000	Tetris (EA)
Nintendo	82900000	Wii Sports
PUBG Corporation	75000000	PUBG: Battlegrounds
Nintendo	60460000	Mario Kart 8 / Deluxe
Nintendo	58000000	Super Mario Bros.
Rockstar Games	50000000	Red Dead Redemption 2
Nintendo	47520000	Pokémon Red / Green / Blue / Yellow
Re-Logic / 505 Games	44500000	Terraria
Nintendo	43800000	Wii Fit / Plus

FIGURE 13.6

Specifying index column and reading only selected columns using pandas.

It is possible to read the full-size data frame into pandas first, then subsequently select only a few (or event only one) columns to form a new subdata frame. It is possible to take only one column from the data frame, and convert it into a series. Notice that a single-column data frame is different from a series from data type perspective. Examples are given below.

best_selling_games_titles = best_selling_games_df["Title"]
best_selling_games_titles

It is worth mentioning that when generating series from data frames, the index column of the data frame is inherited by the series. This introduces an important feature of pandas series: unlike Python array where it is just single-stream sequence of data, pandas series has a separate measure of index for each element in the series, essentially making it multi-stream of data. More are introduced in later sections.

13.4.2 Series and Data Frame

14

Python Practice: Artificial Intelligence

CONTENTS

14.1	Quick	Review	112
	14.1.1	AI Pipeline	112
	14.1.2	Data Preparation and Model Evaluation	112
	14.1.3	Commonly Seen ANN Use Cases	113
	14.1.4	Computer Vision	114
	14.1.5	Natural Language Processing	114
14.2	Tensor	Flow	114
	14.2.1	TensorFlow Basics	115
	14.2.2	Classification and Regression	115
	14.2.3	Computer Vision	119
	14.2.4	General Sequential Data Processing	119
	14.2.5	Natural Language Processing	120
	14.2.6	TensorFlow on Different Platforms	120
14.3	PyTore	ch	120
	14.3.1	PyTorch Basics	120
	14.3.2	Classification and Regression	120
	14.3.3	Computer Vision	120
	14.3.4	General Sequential Data Processing	120
	14.3.5	Natural Language Processing	120
	14.3.6	PyTorch on Different Platforms	120

This chapter focuses on the introduction of commonly used Python-supported ANN engines used in data science. ANN relationships with AI, machine learning and deep learning as well as theories and mechanisms behind ANN can be found on other notebooks, hence is not given here. The introduction only contains the basic usage of these ANN engines from the implementation perspective, and it may not reflect the state-of-the-art technologies such as transformer, LLM, instruction-tuned LLM, etc. These state-of-the-art technologies are introduced on other notebooks.

There are many ANN engines for Python. Among all, TensorFlow and PyTorch are very popular and powerful generic-purpose ANN engines. They both cover a large range of supervised, reinforcement and unsupervised learning applications including classification, regression, pattern recognition, computer vision, natural language processing, clustering, abnormality detection,

and many more. They both offer variety of tools to quickly and flexibly design and deploy different types of AI models such as conventional dense networks, CNN models, RNN models, and many more. Both of them can be used to train, evaluate and run networks. Both of them provide server solutions, cloud solutions and edge computing solutions. TensorFlow and PyTorch are introduced in this chapter.

Installation of TensorFlow and PyTorch can be found in there websites. Although it is possible to run all the calculations on CPU, these ANN engines are more powerful when GPU/TPU are enabled. Depends on the OS and the GPU/TPU brands of the local system, different methods may apply to enable GPU/TPU. For example, if NVIDIA GPU is used, a software called CUDA can be used to configure and enable GPU for ANN training. The installation of TensorFlow, PyTorch, and the enabling of GPU/TPU modules are not covered here.

Alternative to running the code on a local system, consider using online platforms such as Google Colaboratory, which already have all necessary packages pre-installed and the CPU/GPU/TPU pre-configured.

14.1 Quick Review

This section briefly reviews the basic concepts used in this chapter. Details of the concepts can be found elsewhere in other notebooks.

14.1.1 AI Pipeline

AI pipeline is a set of (automated) steps used to build, train, evaluate and deploy AI models. An AI pipeline usually includes at least the following steps (for supervised learning):

- 1. Data collection
- 2. Data preparation
- 3. Model design
- 4. Model training
- 5. Model evaluation and analysis
- 6. Model deployment and testing

where notice that model design and training might need to be carried out iteratively. After the training, the performance of the model is validated using the validation set, according to which the model and its hyper parameters can be modified.

14.1.2 Data Preparation and Model Evaluation

Model training is where the magic happens. Nowadays, with the help of AI engines, it is done automatically via back propagation and other techniques. Given the same training data and model design (including training methods), the almost-the-same trained model can be reproduced by the machine. It is done in a standardized, systematic and consistent manner, hence does not distinguish the performance of the model.

It is rather the data preparation (pre-processing), model design, and model evaluation that require human guidance. Depending on the experience and skill level of the data scientist and engineer, this is where the model performance may differ. Model design and fine-tuning are closely related to model evaluation, as model evaluation results and analysis decides how the model should be tuned. Therefore, it can be concluded that good data preparation, model evaluation and analysis skills are the critical factors that affect model behavior.

To be more precise by breaking down the aforementioned critical factors:

• Data preprocessing:

- Data cleaning. This is a critical step that greatly influences the model's performance. It includes cleaning the data, handling missing values, and dealing with outliers. Often, domain knowledge plays a significant role in these steps. In practice, some of the above procedures are done using AI engine built-in functions, while others are done using other packages such as pandas.
- Feature engineering. This involves creating new features from the existing data that might help improve the model's performance. Feature selection is also crucial to reduce overfitting and improve the model's interpretability. Again, domain knowledge can be very beneficial here.

• Model design, evaluation and tuning:

- Model selection. Choosing the right model or architecture for the problem at hand requires a solid understanding of the strengths and weaknesses of different models.
- Hyperparameter tuning. While there are systematic approaches like grid search or random search, often, practical experience and intuition play a significant role in choosing the right hyperparameters.
- Model evaluation and analysis. Evaluating a model goes beyond looking at a single metric. It involves understanding the model's errors, checking its performance on different subsets of data, and considering aspects like fairness and interpretability.

Given that many, if not all, of the above steps involve a lot of human interaction, data visualization tools often play a very important role to assist humans on their tasks.

14.1.3 Commonly Seen ANN Use Cases

"nobreak

14.1.4 Computer Vision

CV, as an important part of AI, has evolved in the past decades. In the early 2010s, ANN was not used in CV. Instead, conventional deterministic approaches were widely used. With the development in deep learning, the primary approach for CV has changed to CNN. There are a few milestones along the way that together make the change happen:

- Development of GPU
- CNN with deep neural network
- Introduction of rectified linear unit (ReLU) activation function
- Regularization techniques

Recently, with the development in transformer model and large language model, CV is able to be combined with LLM for image comprehension, reasoning, and even artwork generation.

The commonly seen objectives of CV include:

- Image classification
- Object detection
- Image generation
- Image search
- Image comprehension and generation

14.1.5 Natural Language Processing

"nobreak

14.2 TensorFlow

TensorFlow is an open-source software library for machine learning developed by Google in 2015. TensorFlow 2.X is released in 2019 and it is know the official latest major updated version. TensorFlow is backed up by a large community and it supports Python as well as a few other programming languages.

Don't confuse TensorFlow with Keras, later of which is a Python library built on top of deep learning libraries such as TensorFlow, and provides a simple and useful API. In TensorFlow 2.X, Keras is officially adopted as its API. Therefore, when TensorFlow 2.X is used, there is no need to install or import Keras separately.

14.2.1 TensorFlow Basics

Unless otherwise mentioned, the following packages are imported and the command executed in the beginning of all the relevant scripts.

```
import numpy as np
import pandas as pd
import tensorflow as tf

tf.test.is_gpu_available()
```

where .is_gpu_available() tests whether GPU is enabled on the machine. The well-known numpy package defines "NumPy arrays" to store vectors, matrices and tensors. Package tensorflow also defines counterparts "Tensor-Flow tensors" for the same purpose. NumPy arrays and TensorFlow tensors can be converted from one to the other. A key difference of the two is that TensorFlow tensors related calculations are executed on GPU wherever possible, making it more efficient when comes to large-scale calculation that can be paralleled. An example is given below.

```
x = np.array([1, 2, 3, 4, 5])
y = tf.convert_to_tensor(x, dtype=tf.float64)
x = x*0.3 // cpu calculation
y = y*0.3 // gpu parallel calculation
```

14.2.2 Classification and Regression

Classification

Regression

The following data frame *kaggle.com/datasets/shree1992/housedata* is used in this example to as a demonstration to predict house pricing using regression model.

The following class SimpleRegressionModel serves as an example that can be used for the above task. Notice that this model design is only a demonstration, and it is not optimized.

```
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
```

```
from sklearn.preprocessing import MinMaxScaler
import tensorflow as tf
import seaborn as sns
class SimpleRegressionModel:
       SimpleRegressionModel trains and tests a regression model using
           the given data frame.
       def __init__(self):
              self.num_input = None
              self.num_output = None
              self.scalar = None
              self.num_training = None
              self.X_train = None
              self.Y_train = None
              self.num_validation = None
              self.X_val = None
              self.Y_val = None
              self.history = None
              self.num_test = None
              self.X_test = None
              self.Y_test = None
              self.model = None
       def import_dataset(self, df_total: pd.DataFrame, iplist: list,
           oplist: list, validation_size: float, test_size: float):
       total_X = df_total[df_total.columns.intersection(iplist)]
       object_columns = total_X.select_dtypes(include=['object'])
       if not object_columns.empty:
           dummy_columns = pd.get_dummies(object_columns)
           total_X = pd.concat([total_X.drop(object_columns, axis=1),
               dummy_columns], axis=1)
       self.num_input = len(total_X.columns)
       total_Y = df_total[df_total.columns.intersection(oplist)]
       self.num_output = len(total_Y.columns)
       train_val_X, self.X_test, train_val_Y, self.Y_test =
           train_test_split(total_X, total_Y, test_size=test_size,
           random_state=None)
       self.num_test = len(self.X_test.index)
              if validation_size == 0:
                     self.X_train = train_val_X
                     self.Y_train = train_val_Y
                     self.X_val = []
                     self.Y_val = []
                     self.num_training = len(self.X_train.index)
                     self.num_validation = 0
              else:
                     self.X_train, self.X_val, self.Y_train, self.
                          Y_val = train_test_split(train_val_X,
```

```
train_val_Y, test_size=validation_size/(1-
                    test_size), random_state=None)
               self.num_training = len(self.X_train.index)
               self.num_validation = len(self.X_val.index)
       print("Dataset \sqcup size: \sqcup training: \sqcup \{\}, \sqcup validation: \sqcup \{\}, \sqcup test:
            _{\sqcup}\{\}". \texttt{format}(\texttt{self.num\_training}, \ \texttt{self.num\_validation},
            self.num_test))
       self.scalar = MinMaxScaler()
       self.X_train = self.scalar.fit_transform(self.X_train)
       if validation_size == 0:
               pass
       else:
               self.X_val = self.scalar.transform(self.X_val)
       self.X_test = self.scalar.transform(self.X_test)
def design_model(self, hidden_layer_model: list, optimizer: str,
     learning_rate: float, loss: str):
       The input model describes the design of the model. It is
             a list of layers. Each layer is given by a
            dictionary describing layer type, number of nodes,
       self.model = tf.keras.Sequential()
       for ind in range(len(hidden_layer_model)):
               if ind == 0:
                       if hidden_layer_model[ind]["type"] == "
                            dense":
                               layer = tf.keras.layers.Dense(
                                   hidden_layer_model[ind]["node"
                                   ], activation =
                                   hidden_layer_model[ind]["
                                   activation"], input_shape = (
                                    self.num_input, ))
                       else:
                               pass
               else:
                       if hidden_layer_model[ind]["type"] == "
                            dense":
                               layer = tf.keras.layers.Dense(
                                   hidden_layer_model[ind]["node"
                                   ], activation =
                                   hidden_layer_model[ind]["
                                   activation"])
                       else:
                               pass
               self.model.add(layer)
       self.model.add(
               tf.keras.layers.Dense(self.num_output, activation
                    ='relu')
```

```
if optimizer == 'adam':
              optimizer = tf.keras.optimizers.Adam()
       self.model.compile(optimizer=optimizer, loss=loss,
           metrics=['mae'])
def train_model(self, epochs: int, batch_size: int):
       self.history = self.model.fit(self.X_train, self.Y_train
            , epochs=epochs, batch_size=batch_size,
           validation_split=0.2, verbose=2)
def evaluate_model(self):
       if self.num_validation == 0:
              print("Validation\_set\_is\_empty.")
       else:
              loss, mae = self.model.evaluate(self.X_val, self.
                  Y_val, verbose=2)
              print("Validation_set_test_result:_loss={},_mae
                   ={}".format(loss, mae))
def test_model(self):
       loss, mae = self.model.evaluate(self.X_test, self.Y_test
           , verbose=2)
       print("Test_set_test_result:_loss={},_mae={}".format(
           loss, mae))
```

The following code uses the above defined class to predict house price.

```
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import MinMaxScaler
import tensorflow as tf
import seaborn as sns
df_house_pricing = pd.read_csv("house_pricing.csv").dropna()
srm = SimpleRegressionModel()
srm.import_dataset(
   df_total=df_house_pricing,
   iplist=["bedrooms", "bathrooms", "sqft_living", "sqft_lot", "floors
        ", "sqft_above", "sqft_basement", "yr_built", "yr_renovated", '
        condition", "city"],
   oplist=["price"],
   validation_size=0, test_size=0.2)
hidden_layer_model = [
       {
              "type": "dense",
              "node": 128,
              "activation": "relu"
              "type": "dense",
              "node": 64,
```

```
"activation": "relu"
       },
               "type": "dense",
               "node": 64,
               "activation": "relu"
               "type": "dense",
               "node": 32,
               "activation": "relu"
       },
               "type": "dense",
               "node": 16,
               "activation": "relu"
       },
srm.design_model(hidden_layer_model=hidden_layer_model, optimizer='adam
    ', learning_rate=0.001, loss='mse')
srm.train_model(epochs=100, batch_size=50)
srm.evaluate_model()
srm.test_model()
```

The above example is self-explanatory. Some key components of the codes are:

- Use train_test_split() from sklearn.model_selection to split the data set into training set, validation set and testing set.
- Use MinMaxScaler() from sklearn.preprocessing to normalize the inputs to the AI model.
- Use tf.keras.Sequential() and tf.keras.layers to design the AI model structure.
- Use .compile() to configure optimization engine, loss function, validation matrix, etc., during the training.
- Use .fit() to train the model using the training set.
- Use .evaluate() to evaluate the AI model performance.
- Use .predict() to carry out prediction of input points.

Use .save("name") to save a model, and load_model() from tensorflow.keras.models to load a model. Some popular formats to store a model include H5 file.

14.2.3 Computer Vision

"nobreak

14.2.4 General Sequential Data Processing

"nobreak

14.2.5 Natural Language Processing

"nobreak

14.2.6 TensorFlow on Different Platforms

"nobreak

14.3 PyTorch

TensorFlow is another open-source software library for machine learning originally developed by Meta AI in 2016. It is now under the Linux foundation umbrella.

14.3.1 PyTorch Basics

"nobreak

14.3.2 Classification and Regression

"nobreak

14.3.3 Computer Vision

"nobreak

14.3.4 General Sequential Data Processing

"nobreak

14.3.5 Natural Language Processing

 ${\rm ``nobreak'}$

14.3.6 PyTorch on Different Platforms

Part IV Semantic Web

15

Semantic Web Basics

CONTENTS

15.1.3 Semantic Web Vision 126 15.1.4 Semantic Web Stack 126 15.1.5 Semantic Web Limitations and Challenges 130 15.2 Ontology 131 15.2.1 Philosophy Perspective 131 15.2.2 Semantic Web Perspective 131 15.2.3 Ontology Types and Categories 132 15.3 Logic 132 15.3.1 Syntax and Semantics 133 15.3.2 Logic Framework 134 15.3.3 Logical Expression 136 15.3.4 Logical Equivalence 137	15.1	Web o	f Data	123
15.1.3 Semantic Web Vision 126 15.1.4 Semantic Web Stack 126 15.1.5 Semantic Web Limitations and Challenges 130 15.2 Ontology 131 15.2.1 Philosophy Perspective 131 15.2.2 Semantic Web Perspective 131 15.2.3 Ontology Types and Categories 132 15.3 Logic 132 15.3.1 Syntax and Semantics 133 15.3.2 Logic Framework 134 15.3.3 Logical Expression 136 15.3.4 Logical Equivalence 137		15.1.1	Web 1.0 and 2.0	124
15.1.4 Semantic Web Stack 126 15.1.5 Semantic Web Limitations and Challenges 130 15.2 Ontology 131 15.2.1 Philosophy Perspective 131 15.2.2 Semantic Web Perspective 131 15.2.3 Ontology Types and Categories 132 15.3 Logic 132 15.3.1 Syntax and Semantics 133 15.3.2 Logic Framework 134 15.3.3 Logical Expression 136 15.3.4 Logical Equivalence 137		15.1.2	Web 3.0	125
15.1.5 Semantic Web Limitations and Challenges 130 15.2 Ontology 131 15.2.1 Philosophy Perspective 131 15.2.2 Semantic Web Perspective 131 15.2.3 Ontology Types and Categories 132 15.3 Logic 132 15.3.1 Syntax and Semantics 133 15.3.2 Logic Framework 134 15.3.3 Logical Expression 136 15.3.4 Logical Equivalence 137		15.1.3	Semantic Web Vision	126
15.2 Ontology 131 15.2.1 Philosophy Perspective 131 15.2.2 Semantic Web Perspective 131 15.2.3 Ontology Types and Categories 132 15.3 Logic 132 15.3.1 Syntax and Semantics 133 15.3.2 Logic Framework 134 15.3.3 Logical Expression 136 15.3.4 Logical Equivalence 137		15.1.4	Semantic Web Stack	126
15.2.1 Philosophy Perspective 131 15.2.2 Semantic Web Perspective 131 15.2.3 Ontology Types and Categories 132 15.3 Logic 132 15.3.1 Syntax and Semantics 133 15.3.2 Logic Framework 134 15.3.3 Logical Expression 136 15.3.4 Logical Equivalence 137		15.1.5	Semantic Web Limitations and Challenges	130
15.2.2 Semantic Web Perspective 131 15.2.3 Ontology Types and Categories 132 15.3 Logic 132 15.3.1 Syntax and Semantics 133 15.3.2 Logic Framework 134 15.3.3 Logical Expression 136 15.3.4 Logical Equivalence 137	15.2	Ontolo	ogy	131
15.2.3 Ontology Types and Categories 132 15.3 Logic 132 15.3.1 Syntax and Semantics 133 15.3.2 Logic Framework 134 15.3.3 Logical Expression 136 15.3.4 Logical Equivalence 137		15.2.1	Philosophy Perspective	131
15.3 Logic 132 15.3.1 Syntax and Semantics 133 15.3.2 Logic Framework 134 15.3.3 Logical Expression 136 15.3.4 Logical Equivalence 137		15.2.2	Semantic Web Perspective	131
15.3.1 Syntax and Semantics 133 15.3.2 Logic Framework 134 15.3.3 Logical Expression 136 15.3.4 Logical Equivalence 137		15.2.3	Ontology Types and Categories	132
15.3.2 Logic Framework 134 15.3.3 Logical Expression 136 15.3.4 Logical Equivalence 137	15.3	Logic		132
15.3.3 Logical Expression 136 15.3.4 Logical Equivalence 137		15.3.1	Syntax and Semantics	133
15.3.4 Logical Equivalence		15.3.2	Logic Framework	134
15.3.4 Logical Equivalence		15.3.3	Logical Expression	136
15.3.5 Logical Reasoning		15.3.4		137
		15.3.5	Logical Reasoning	138

Ontology is the philosophical study of the nature of being, existence, or reality. It concerns about "what is everything" and "how to define it" in the context of inductive and deductive reasoning. It discusses how we abstract and preserve knowledge for the generations to come.

Ontology inspires people in computer science about how we can store and exchange information efficiently using the internet. The solution is called the semantic web, an internet-based knowledge base schema. The internet powered by semantic web (together with other technologies) defines Web 3.0, a new-generation internet framework.

Notice that there are subtle differences between "internet" and "Internet". The lower case "internet" is a technology that bridges machines to form a network of any size, and the upper case "Internet" refers to the specific internet that links the entire world together. In this part of the notebook, however, they are used interchangeably.

15.1 Web of Data

The internet is a technology that enables information exchange among machines and humans. With the power of the internet, a user can obtain the information he needs from remote servers, databases, or knowledge bases.

In the early days, the use of internet required professional skills. Nowadays, everyone can access the internet using a graphical-interface browser that he can easily find on a computer or a mobile device. Public knowledge bases such as Wikipedia has made obtaining information much easier.

15.1.1 Web 1.0 and 2.0

Under the Web 1.0 framework which was popular in the early stage of internet development, information is stored on individual servers in a static manner. A user can browse the contents, but he cannot edit them. It is essentially a one-way data transmission. The "authority" such as a news company provides the information, and the users consume it. Examples of Web 1.0 implementations include news websites, static web gallery, etc.

One-way data transmission in Web 1.0 cannot meet the expectation from the users who want to share their information to other users on the internet. Thanks to the advent in information science and communication, under the Web 2.0 framework users can interact with the internet bidirectionally. A user can search and filter data from the servers and even upload his own data and share it with others. The Internet became far more powerful in terms of information exchange. The source of information does not necessarily come from the authority. The users are generators and consumers of information at the same time. Examples of Web 2.0 implementations include Blog, Twitter, YouTube and Weibo.

Web 2.0 is extremely popular and successful even to date. Yet, under the background of industry 4.0, big data and data-driven modeling of almost everything, there are still limitations and downsides to Web 2.0 that we would wish to improve.

One of the biggest challenges that we encounter with Web 2.0 is how to quickly locate the information we want among the vast amount of irrelevant data. Conventionally in Web 2.0, keyword-based searching engines are used. These searching engines do not comprehend the contextual knowledge of the contents, and as a result they may fail to return what a user truly expects. The user often needs to further manually pick up relevant information from the searching results. Nowadays searching engines are becoming smarter, and they can sometimes pre-filter the results. Nevertheless, it is still quite common that many returned results are useless to the user. Keyword-based engines also struggle with polysemous, synonyms and implicit information from pictures in the searching range, again due to the lack of understanding of the contents.

The problem here, however, is not caused by the searching engines alone. It is rather that the information stored on the internet does not come with its corresponding semantics in the first place. Today most of the information on the internet is stored in HTML format. HTML tells only the contents but not the meaning behind them. It is difficult for a machine to understand the meaning of the information from HTML corpus. This potentially makes it difficult for a machine to retrieve data efficiently. Even with the recent breakthrough in LLM enabling the machine to summarize articles, it is still unpractical for it to go through all the returned contents of a searching engine which sometimes contains hundreds of pages of information.

Model determines its functions. To solve the problem once for all, new data model to store and share the information need to be introduced.

15.1.2 Web 3.0

The goal of Web 3.0 is to allow more efficient information retrieval and sharing using the internet. Ideally, we would want the searching engine to truly understand the user's demand, and return only the most relevant information summarized in a nice manner. If there is no readily available response on the internet, the searching engine shall derive the response based on existing information, i.e., it should be capable of doing simple reasoning.

For this purpose, there are at least the following two development trends:

- Let the LLM-based AI remember all the knowledge. The LLM should be smart enough to precisely capture what the user wants and get back to him with useful and accurate information. This leads to chatbots and copilots.
- Create a powerful and flexible NoSQL database. Make the information
 in the database readable for both humans and machines, and somehow
 integrate the semantics of the contents into the database. This leads to
 semantic web.

In practice, it is most recommended to use both LLM-based AI and semantic web simultaneously. This is because both LLM-based AI and semantic web have drawbacks when used alone.

• LLM-based AI

- Time and computational load wise expensive to expand the knowledge base;
- Lack the ability of reasoning;
- A chance to produce misleading information;
- Cannot be merged and migrated easily because it is human or machine readable;

• Semantic web

- Steep learning curve to use its interface;
- Not good at summarizing and formatting the response in a humanfriendly manner.

Web 3.0 schema focuses on semantic web, which is essentially a database with specific data model that can build semantics into its framework. The "semantics" in this context specifically refers to the relations of objects and properties that can be used for machine-driven descriptive reasoning. See later sections for details. Semantic web is so important that it is often considered synonymous with Web 3.0, although the broader vision of Web 3.0 includes other features as well.

Another feature of Web 3.0 is distributed storage of information. Distributed storage has both advantages and disadvantages. On one hand, it provides higher resilience against data loss. On the other hand, it adds challenges to the searching and collecting of information for data stored in both semantic web structure or other database structure. Technologies like Inter-Planetary File System (IPFS), blockchain, and distributed databases are used to address these challenges.

This notebook discusses the technologies used in semantic web such as resource description framework (RDF) model, RDF schema (RDFS) and web ontology language (OWL). We will see how we can use semantic web to collect and organize information, and create a knowledge base for both humans and AI.

15.1.3 Semantic Web Vision

Semantic web database is "semantic" due to its data model. Semantic web data model not only stores objects and their properties, but also the relationships among objects and properties. The relationships, often represented by graphs, can be used for descriptive logic reasoning. The logic reasoning allows new information to be derived based on existing facts. In addition, semantic web is more scalable, reusable and reader-friendly for both humans and machines comparing with AI-based methods.

Many research institutes and organizations have tried building semantic webs of different scales. An example is DBPedia *dbpedia.org*. The goal of DBPedia is to create something like Wikipedia, but using semantic web.

15.1.4 Semantic Web Stack

A commonly seen semantic web stack looks like Fig. 15.1.

The layers and components in each associated layer are briefly introduced as follows.

• Web platform layer

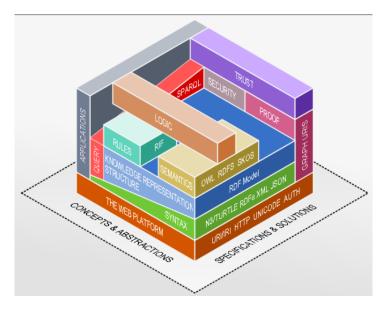


FIGURE 15.1 Semantic web stack [2].

- Uniform resource identifier (URI)/ internationalized resource identifier, often a string used to identify and trace an information resource;
- Communication protocols, such as HTTP/HTTPS;
- Text encoding methods, such as unicode, utf-8;
- Authentication methods;

• Syntax

- File formats to store information, such as and RDF/XML;
- Knowledge representation, semantics, and interfacing tools
 - RDF model by default, which uses triples to represent a graphical database;
 - RDFs and OWL which expand the capability of RDF model;
 - SPARQL which is the default language for semantic web query (SPARQL 1.0) and manipulation (SPARQL 1.1);

• Logic

- Logic statements and reasoning that semantic web supports for query;
- Rules that describes what query can be executed.

The knowledge and semantics are stored using the RDF/RDFS model. The RDF/RDFS is often further powered by OWL. RDF serves as the foundation with the basic structure. RDFS expands the capability of RDF by introducing new vocabularies such as "class", "subclass", "property". And finally OWL enables a set of tools to define complex ontologies and create rigorous and complex semantics. As an analogy, think of RDF as the paper and pencil, RDFS the 24-color crayon set, and OWL the painting skills. Together they make a sophisticated and richly structured picture.

It is worth mentioning that many syntax such as RDF/XML can be used as the markup languages to describe RDF/RDFS/OWL used in a semantic web. More details are introduced in later chapters.

SPARQL Protocol and RDF Query Language (SPARQL) is the recommended query language for querying and manipulating the semantic web. It allows a user to search, retrieve, and modify information stored in the RDF model. The syntax of SPARQL is designed to look similar with SQL, and many commands in SPARQL have counterparts in SQL, such as SELECT, WHERE, GROUP BY, and ORDER BY. However, the backend technologies of a semantic web engine (known as triplestore or RDF management system) and a relational database management system differ largely. For example, triplestore relies heavily on graph theory and graph pattern mapping when searching through the semantic web.

An example of semantic web is given below.

Example: Semantic Web on Animals

The following examples demonstrates the use of RDF/RDFS model and OWL to create a small semantic web on animals. The example is provided by *ChatGPT-4*.

RDF/RDFS

```
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
@prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#> .
@prefix ex: <http://example.org/> .

ex:Animal rdf:type rdfs:Class .

ex:Mammal rdf:type rdfs:Class ;
  rdfs:subClassOf ex:Animal .

ex:Reptile rdf:type rdfs:Class ;
  rdfs:subClassOf ex:Animal .

ex:hasLegs rdf:type rdf:Property ;
  rdfs:domain ex:Animal ;
  rdfs:range rdfs:Literal .
```

```
ex:Dog rdf:type rdfs:Class ;
   rdfs:subClassOf ex:Mammal .

ex:Lizard rdf:type rdfs:Class ;
   rdfs:subClassOf ex:Reptile .

ex:Max rdf:type ex:Dog ;
   ex:hasLegs 4 .

ex:Lizzy rdf:type ex:Lizard ;
   ex:hasLegs 4 .
```

In this example, RDF/RDFS is used to:

- Define classes (Animal, Mammal, Reptile, Dog, and Lizard);
- Define a property (hasLegs);
- Set the domain and range of the property;
- Create subclass relationships (Mammal and Reptile are subclasses of Animal; Dog is a subclass of Mammal; Lizard is a subclass of Reptile);
- Define individuals (Max and Lizzy) and their properties.

OWL

```
@prefix rdf: <a href="mailto://www.w3.org/1999/02/22-rdf-syntax-ns"> .
@prefix rdfs: <a href="mailto://www.w3.org/2000/01/rdf-schema"> .
@prefix owl: <a href="mailto://www.w3.org/2002/07/owl#> .
@prefix ex: <a href="mailto://example.org/"> .
ex:hasParent rdf:type owl:0bjectProperty;
    rdfs:domain ex:Animal;
    rdfs:range ex:Animal .

ex:hasChild rdf:type owl:0bjectProperty;
    owl:inverseOf ex:hasParent .

ex:isWarmBlooded rdf:type owl:Class;
    rdfs:subClassOf ex:Animal .

ex:Mammal rdfs:subClassOf ex:isWarmBlooded .
ex:Reptile owl:disjointWith ex:isWarmBlooded .

ex:Max ex:hasParent ex:Buddy .
```

```
ex:Buddy rdf:type ex:Dog ;
ex:hasLegs 4 .
```

In this example, OWL is used to:

- Define object properties (hasParent and hasChild);
- Specify an inverse relationship between properties (hasParent and hasChild);
- Define a new class (isWarmBlooded) and set it as a superclass of Mammal;
- Specify a disjoint relationship between Reptile and isWarmBlooded;
- Define a new individual (Buddy) and his properties;
- Specify a relationship between individuals (Max and Buddy).

In this demonstration example, RDF and RDFS provided the basic structure and hierarchy for the knowledge base, while OWL added more expressivity by defining complex relationships, additional semantics, and constraints.

The above semantic web can be queried by SPARQL. An example is given below.

```
PREFIX ex: <a href="http://example.org/">
SELECT ?mammal
WHERE {
    ?mammal rdf:type/rdfs:subClassOf* ex:Mammal .
}
```

There is a learning curve for SPARQL. Commercialized semantic web applications such as WolframAlpha (wolframalpha.com) often provide a "search bar" with some natural language processing capability which triggers SPARQL query according to the user's input. It is possible to use a more powerful LLM-based chatbot to assist SPARQL query as well.

15.1.5 Semantic Web Limitations and Challenges

Though semantic web is powerful, building semantic web can be challenging and requires a lot of careful design and human labor. As of 2023, the majority of web sites and applications have not yet embraced the full potential of the semantic web, some of which only partially adopt semantic web concepts or technologies.

However, things may change due to the recent advent in internet-of-things (IoT, as defined in Industry 4.0) and LLM-based copilots. IoT devices gen-

erates large amount of data, which makes the base of building large-scale knowledge. Copilots are useful with converting data from other formats into RDF model.

15.2 Ontology

Ontology has very rich meanings from both philosophy and semantic web perspectives.

15.2.1 Philosophy Perspective

Knowledge is the overlapping part of ground truth and human beliefs, i.e., it is the truth that humans know of being the truth. Ontology is the methodology of storing and communicating knowledge.

In the context of philosophy, ontology discusses the meaning of objects being "existing", how objects can be categorized, and how objects relate to each other. Ontology mainly discusses the following questions:

- What is existence? What does it mean for something to exist or not exist?
- How many different types of "existences" are there, and what are their natures?
- What is the nature of abstract entities like numbers, properties, and relations?
- How do different entities relate to and interact with each other?
- Can something exist independently of our perception or thought?

The discussion of these questions can be traced back to 300 BC or even earlier. Aristotle defined a system to structure and reason knowledge. The famous Aristotelian logic "major premise + minor premise \rightarrow conclusion" is a systematic way of reasoning new knowledge. Aristotelian logic is an important tool of ontology, and it inspires the proposition of the semantic web.

15.2.2 Semantic Web Perspective

In the context of the Semantic Web, ontology is a formal, machine-readable representation of the domain knowledge in a specific area. It serves as a shared vocabulary for describing and reasoning the knowledge within that domain. Notice that unlike natural language which can be ambiguous, semantic web shall be described clearly, precisely and consistently.

In practice, RDF/RDFS and OWL can be used to express the ontology. The followings vocabularies are defined in RDF/RDFS and OWL.

- Class. A class is an abstraction of objects sharing some similarities.
- Properties. A property defines a feature of a class. Triples are often used to describe a property. A triple follows the form of "subject + property + object".
- Relations. A relation describes how classes relates to each other. A relation is often described using the triple where both subject and object are classes.
- Constraints. A constraint describes the rules enforced on a property.
- Instance. An instance is a realization of a class.

15.2.3 Ontology Types and Categories

There are different layers of ontology. The higher the layer, the more general the knowledge. The lower the layer, the more specific the knowledge within a particular domain, application or task. An example is given in Fig. 15.2. In practice, there might not be a clear boundary between two adjacent layers.

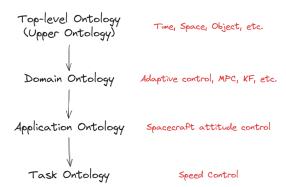


FIGURE 15.2

A demonstration of ontology level.

Lower-layer ontology can inherit and modify knowledge from the upperlayer ontology. Other ways of leveling ontology include using the expressiveness. The "light-weight" ontology is informal, less semantic, and supports less comprehensive logic. The "heavy-weight" ontology, on the other hand, is formal, more semantic, and supports more comprehensive logic and reasoning up to first-order logic. The vocabulary also grows with the ontology complexity.

15.3 Logic

Logic and logical expression form a research area by themselves and it is not possible to include all the details in the scope of this notebook. Only a brief scratch is given.

Humans are good at deriving new knowledge from existing knowledge. The systematic way of doing so is logic. The term "formal logic" rigorously defines logic reasoning methods and procedures to automate logic inference by machines.

15.3.1 Syntax and Semantics

The meaning behind a sentence is hidden in the sequence of words it uses. As an analogy, the semantics is the water, and the sentence is the cup that holds the water. Likewise, in logic, the syntax that express the logic holds the semantics of the logic. Depending on who or what tool is used to process the syntax, different levels of semantics can be interpreted.

Consider the following example. This is a piece of code written in Python syntax.

```
ax = []
for i in range(20):
    if i<=1:
        ax.append(i)
    else:
        ax.append(ax[i-1]+ax[i-2])</pre>
```

Three interpreters are studied, the Python interpreter, the AI model, and the human cognition. Different interpreters draw different semantics from the same syntax.

The Python interpreter is essentially purely syntax-driven. When it executes a script, it follows the provided instructions without any broader comprehension or anticipation of the syntax or the results. In the context of this example, the interpreter does not recognize that it is generating the famous Fibonacci sequence. It simply follows the rules of the script, calculating and outputting each number in the series as instructed.

On the other hand, an AI model like GPT-3 does possess a level of semantic understanding, but it is quite different from that of a human. Through training on a diverse range of data corpus, it has learned to associate the Python script for generating a Fibonacci series with the mathematical concept of the Fibonacci sequence. This understanding is not innate but rather a result of recognizing patterns in the data it was trained on. When asked about the script, it can provide a summary or explanation based on its learned associations, including the name of the series, and the purpose of each line of the code. It follows the corpus it has been trained on.

Finally, there's the human level of semantics, which is by far the most advanced and intuitive. Not only can a human understand the Python script and the concept of the Fibonacci series, they can also infer its broader mathematical properties, such as its relation to the golden ratio, and its general behavior. Humans can understand that the series will converge to the golden ratio not only when starting with 0 and 1, but with any two initial positive integers. This level of understanding is a combination of learned knowledge, pattern recognition, and the ability to extrapolate or generalize from existing information.

One of the goals of semantic web is to enable a machine to understand the semantics as much as possible. The following are some commonly seen terminologies regarding semantics:

- Intentional Semantics: the semantics that the information producer intends to share.
- Formal Semantics: the semantics contained in the formally defined language.
- Procedural Semantics: the semantics in programming script.
- Model-theoretic Semantics: the semantics understood from natural language by a formally defined model.

In the context of logic, model-theoretic semantics are the main interest of study.

15.3.2 Logic Framework

Different logic frameworks have different levels of complexity, hence different capabilities of expressing semantics. Commonly seen logic frameworks are briefly introduced as follows.

Propositional Logic

Propositional logic is the fundamental of logic reasoning. In propositional logic, knowledge is represented by either simple facts which are known as propositions, or facts connected with "AND", "OR", "NOT", "IF ... THEN ...", "IF AND ONLY IF" which are known as compound propositions.

First-Order Logic

First-order logic (FOL) is the most commonly used logic as of today. In FOL, quantifiers and logic connectives are introduced, including universal quantification \forall , existential quantification \exists , conjunction \land , disjunction \lor and negation \neg . A formal way of representing logic expressions are defined. For example, using the following to represent "all humans are mortal"

$$\forall x \, (\operatorname{Human}(x) \to \operatorname{Mortal}(x))$$

where CLASS(x) is equivalent of a proposition "x belongs to CLASS", which can be either true or false. The above proposition says "for any item, if that item belongs to human, then that item belongs to mortal".

Description Logic

Description logic (DL) is a subset of the first-order logic. Deriving DL from FOL is not under the scope of this notebook. Some key features of DL are:

- Define classes and subclasses
- Define properties (roles) and associated domains and ranges
- Define restrictions on classes and properties
- Support universal and existential quantifiers

Semantic web uses OWL to implement DL on computers. More details are given in Section 17.

Undecidability in FOL

We would surely want to introduce highly expressive and complex formalisms to the existing RDF/RDFS model. This motivates OWL, which enables DL in semantic web. More details of OWL is introduced later in Section 17. Notice that DL instead of FOL is widely used in semantic web due to its undecidability. DL is a subset of FOL and it is simpler.

A logic is decidable if there is an algorithm that can determine whether any given statement in the logic is true or false (valid or invalid). If no such algorithm exists for a logic, it is undecidable.

It is undecidable whether a FOL formula is provable (or true under all possible interpretations). An FOL reasoning algorithm may not terminate in finite time. The proof of this theorem can be found elsewhere and is not given in this notebook.

While FOL inference is not decidable, DL inference, on the other hand, is decidable. That is one of the reasons DL is preferred over FOL in semantic web.

Attributive language with complement (ALC) is one of the ways to describe a simple DL, and it forms an important subset of OWL. Understanding ALC is helpful with learning OWL in later sections. A brief introduction of ALC is given below.

"Concept" (corresponding with class in RDF/RDFS) and "role" (corresponding with property in RDF/RDFS) are introduced in ALC. Top concept (root class) and bottom concepts (leaf classes) are defined. Each property is associated with a range, which is a basically a set of values that the property can take.

Constructors are used to describe the concept, role and range used in

ALC. Conjunction, disjunction, negation, existential and universal quantifiers are supported. For example, $\forall R.C$ means that all roles "R" (a concept can have multiple roles of the same name) must have value taken from concept "C". Similarly, $\exists R.C$ means that there is at least one role "R" whose value is taken from concept "C". Concept relations are defined. Commonly used concept relation constructors are inclusion (to describe sub class), equality (to assign class), union, intersection, complement, etc.

ALC uses terminological knowledge and assertional logic statements to define classes. Examples are given below. Using terminological knowledge we can define a teacher as

Teacher \equiv Person \land \exists HasStudent.Student \land \exists Teaches.Class

which translates to "a Teacher is a Person, and it has at least one role Has-Student whose value is from Student, and has at lease one role Teaches whose value is Class", where "Person", "Student", "Class" are concepts and "Has-Student", "Teaches" are roles. There are also teachers who teach only tutorials but not classes. To include them into the Teacher class, consider using

Teacher
$$\equiv$$
 Person \land \exists HasStudent.Student \land (\exists Teaches.Class \lor \exists Teaches.Tutorial)

With terminological knowledge, ALC can enforce restrictions on concepts and rules flexibly. Using assertional knowledge, on the other hand, allows defining instances and subclasses as follows.

Teacher (Peter) Teaches (Peter, Linear Algebra)

All the above ALC can be translated into OWL then implemented in the semantic web. More details are given in later sections.

15.3.3 Logical Expression

A logic is defined by

$$L := (S, \models)$$

where S is the set that contains all the statements of of our interests, and \models the entailment relation

$$\models = \models^1 \bigcup \models^2$$

where

$$\models^{1} = \{(\Phi, \phi) | \Phi \subseteq S, \phi \in S, \Phi \to \phi\}$$
 (15.1)

$$\models^{1} = \{(\Phi, \phi) | \Phi \subseteq S, \phi \in S, \Phi \to \phi\}$$

$$\models^{2} = \{(\Phi, \Psi) | \Phi, \Psi \subseteq S, \forall \psi \in \Psi, \Phi \models^{1} \psi\}$$
(15.1)
$$(15.2)$$

In (15.1), ϕ is known as the logical consequence of Φ . If two logical assertions satisfy $\Phi \models \Psi$ and $\Psi \models \Phi$, then they are logically equivalent $\Phi \equiv \Psi$.

A logic statement is meaningful only if its interpretation I and formula F are clearly articulated. I and F are two very important terms in logic expression. Interpretation is a formal construct that defines the meaning of a symbol in a formal language. For example, in the ontology of people, an interpretation can map an instance symbol, such as "Alice", to an actual person in the real world, and class symbol "Person", to the concept of a human being, etc. Formula, on the other hand, is a statement formulated by string of symbols from a formal language. It is an assertion trying to represent some fact. For example, a formula can be "Alice hasSibling Bob".

The true or false of the formula depend not only by the formula itself, but also by the interpretation. In the earlier example "Alice hasSibling Bob", if "Alice" and "Bob" are indeed mapped to two siblings, and the relation "hasSibling" is as it literally represents, then it is true. Here is another example. Consider formula "10 is greater than 5". Intuitively, this is a universal truth. But from the interpretation and formula perspective, this also depends on the interpretation. If "10" and "5" are interpreted as numerical quantities and "is greater than" as the standard numerical greater-than relation, then it is true. However, if "10" and "5" are interpreted as amounts of debt and "is greater than" as "richer" (with less debt being richer), then it would be false under this interpretation.

The interpretations that make the formula true form the model of the formula denoted by I(F) or $I \models F$, which read as "I models F" or "I satisfies F".

With the above, we can express FOL using logic expression. Here is an example

$$\forall X : \mathrm{Child}(X) \to \mathrm{lovesIcecream}(X)$$

which says "for all elements denoted by X, if X interpreted as Child holds true, then X interpreted as lovesIcrcream must also be true".

Multiple formulas can be grouped into a theory (T), which can be used interchangeably with F. A theory can be treated as a knowledge base.

15.3.4 Logical Equivalence

Logical equivalence has already been introduced in earlier section. Recall (15.2) where we say if $\Phi \models \Psi$ and $\Psi \models \Phi$, Φ and Ψ are logically equivalent denoted by $\Phi \equiv \Psi$. Here Φ and Ψ are sets of statements.

Consider a simplified case where Φ and Ψ each contains only one statement. Let the formula of the statements of Φ and Ψ be F and G respectively. The notations applied to Φ and Ψ applies to F and G similarly. For example, $F \models G$ means that under the same interpretation if F is true, G must also be true. If $F \models G$ and $G \models F$, F and G are logically equivalent denoted by $F \equiv G$.

There is a huge table of logical equivalence formulas. The proof of the formulas is not given here. Commonly seen equivalence is given in Table 15.1.

TABLE 15.1

Numerical calculations.			
Statement	Equivalent	Comment	
$\neg \neg p$	p	Double negation law	
$(p \wedge q)$	$(q \wedge p)$	Commutative laws	
$(p \lor q)$	$(q \lor p)$	Commutative laws	
$(p \wedge (q \wedge r))$	$((p \land q) \land r)$	Association laws	
$(p \lor (q \lor r))$	$((p \lor q) \lor r)$	Association laws	
$(p \land (q \lor r))$	$((p \land q) \lor (p \land r))$	Distributive laws	
$(p \lor (q \land r))$	$((p \lor q) \land (p \lor r))$	Distributive laws	
$(p \to q)$	$(\neg p \lor q)$	Conditional statements	
$(p \to q)$	$(\neg q \to \neg p)$	Conditional Statements	
$(p \leftrightarrow q)$	$((p \to q) \land (q \to p))$	Conditional Statement	
$(\neg(p \land q))$	$(\neg p \lor \neg q)$	De Morgan's laws	
$(\neg(p \lor q))$	$(\neg p \land \neg q)$	De Morgan's laws	

The interpretation of conjunction \land , disjunction \lor and negation \neg are "AND", "OR", "NOT" respectively.

Canonical forms have been defined for logical expressions. There are at least 6 different canonical forms for a logical expression, namely conjunctive normal form, disjunctive normal form, prenex normal form, skolem normal form, negation normal form and clausal normal form. Canonical forms are not necessarily the easiest and most intuitive forms of an expression (in fact it is quite the opposite), but somethings they are simple for further analysis and process, such as finding contradictions.

Notice that when deriving certain canonical forms from the original logical expression, information may get lost and they are not necessarily always equivalent, but they usually are.

15.3.5 Logical Reasoning

Logical reasoning is about proving the true or false of a formula F given a theory T. Notice that it is often not easy, and the proof may not be unique.

Before talking about how knowledge is retrieved from the knowledge base using DL inference and reasoning, we need to discuss what to return if knowledge is not found. When open world assumption (OWA) is made, the knowledge base considers itself as underdevelopment, and it is open-minded to unknowns. While when closed world assumption (CWA) is made, the knowledge base considers itself as developed, and unknown means nonexistence.

Consider an example where it is asked "are Alice's children all males?", and in the database there are two children of Alice, both of which are male. Under OWA, the inference would return "maybe", as it does now know whether there are more children of Alice. While in CWA, the inference would return "yes", as it checks both children to be male, and it assumes they would be all the

children Alice has. However, when Alice does have female child registered, both OWA and CWA inferences would return "no", because the female child contradicts the assertion regardless of the completeness of the database.

Logic reasoning using the law of contradiction is one of the most commonly used approach. To use the law of contradiction to prove $T \models F$, in the high level just find contradictions in $\{\neg F, T\}$. The detailed procedure may differ depending on the logic framework being used.

A commonly used way of looking for contradictions in $\{\neg F, T\}$ is to resolve the formula into something more structured, for example to one of its canonical forms. Commonly used canonical forms for contradiction check are clausal normal form and disjunctive normal form.

16

Resource Description Framework

CONTENTS

16.1	Unifor	m Resource Identifier	141
16.2	Resour	ce Description Framework (RDF)	143
	16.2.1	Triple Representation	143
	16.2.2	Multi-valued Relation and Blank Node	144
	16.2.3	Lists	146
	16.2.4	Reification	147
	16.2.5	Converting RDB to RDF	147
	16.2.6	Conclusion	148
16.3	RDF S	Schema (RDFS)	148
	16.3.1	RDFS Motivation	149
	16.3.2	RDF Versus RDF/RDFS via an Example	150
	16.3.3	RDFS Expanded Class and Properties	151
	16.3.4	Semantics inside RDF/RDFS	152
16.4	SPARO	QL Protocol and RDF Query Language (SPARQL)	152
	16.4.1	SPARQL for Basic Query	153
	16.4.2	SPARQL for Advanced Operations	155
	16.4.3	Default Graph and Named Graph	157
	16.4.4	SPARQL Programming Returns	158
	16.4.5	Underlying Data Structure of Triplestores	158

Semantic web uses URI to identify an existing interpretation, RDF/RDFS and OWL to store semantics, and SPARQL to query and manipulate the database. RDF and RDFS are discussed in this chapter.

16.1 Uniform Resource Identifier

Human use symbols to represent objects in reality. The concept associated with a symbol is used to help define and interpret objects. An example of interpreting "Apple" is given in Fig. 16.1.

URI points to the interpretation of the symbols used in the semantic web. In practice, the resource is either the definition of the class/object/property,

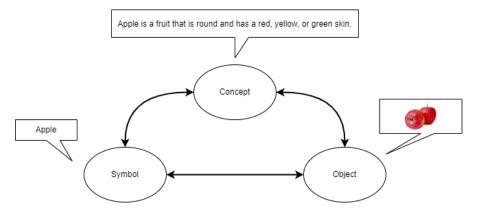


FIGURE 16.1 Semiotic triangle of Apple.

or the higher-level ontology (maybe also in the format of semantic web). For example, consider building a semantic web to describe an apple farm. Consider using dbpedia.org/page/Apple from DBpedia as the URI of interpretation of apple, which already gives a clear and rich definition of apple such as abstract, Wikipedia link, emoji, sugar level, etc. This saves the time of building the interpretation of apple from scratch.

URI shall contain at least two pieces of information: the address (locator), and the identity (name). URI is often ASCII encoded, but it is possible to extend to unicode. The generic syntax is given below. If it looks similar with Uniform Resource Locator (URL), that is because URL is considered as a type of URI.

scheme: [//authority]path[?query][#fragment]

where

- scheme specifies the protocol used to access the resource. Common schemes include http, https, ftp, mailto, file, and data. The scheme is followed by a colon:.
- authority specifies the user information, host, and port, separated by an at sign @ and a colon:, respectively. The authority is preceded by a double forward slash //.
- path identifies the specific resource within the context of the scheme and authority. It is a sequence of segments separated by forward slashes /.
- query provides additional information that the resource can use for processing. It is a series of key-value pairs separated by an ampersand &. The query starts with a question mark?.

• fragment identifies a specific part or section within the resource. It is typically used with HTML documents to indicate a specific anchor or location within the page. The fragment starts with a hash sign #.

More details can be found at

https://www.rfc-editor.org/rfc/rfc3986.html#section-3.1

which happens to be a good example of an URI in https scheme.

16.2 Resource Description Framework (RDF)

RDF is the backbone data model of the semantic web. It stores data in the form of triples. Notice that RDF alone is not powerful enough to record DL. In practice, RDF usually works with RDFS which expand its vocabulary to include class, subclass, etc., and OWL which further expand its vocabulary and introduce DL features, to together form a comprehensive semantic web. RDF is the fundamental of RDFS and OWL.

Notice that RDF is not a syntax by itself. A markup language is needed to host the RDF framework. Commonly used markup languages are listed below. They can be translated from one to the other.

- RDF/XML: has good compatibility with old machines.
- Terse RDF Triple Language (Turtle): easy to use, human-readable.
- JSON for Linked Data (JSON-LD): popular in web applications and APIs where JSON-based format is required.
- N-Triples: simple, machine-readable format for data exchange between tools.

Unless otherwise mentioned, Turtle is used in this notebook.

16.2.1 Triple Representation

RDF uses "subject-predicate(property)-object" triple to represent knowledge. A triple is corresponding with an edge in a directed graph, which is often used to visualize RDF.

For example, consider representation of knowledge "Einstein was born in Ulm". In RDF, "Einstein" is the subject, and "Ulm" the object. The predicate is "has birthPlace", where "birthPlace" is a property assigned to "Einstein". The graph representation is given by Fig. 16.2. The RDF in Turtle is given as follows.



FIGURE 16.2

Graph representation of triple for knowledge "Einstein was born in Ulm".

```
@prefix dbo: <http://dbpedia.org/ontology/> .
@prefix dbr: <http://dbpedia.org/resource/> .
dbr:Albert_Einstein dbo:birthPlace dbr:Ulm .
```

where Albert_Einstein and Ulm are defined in dbpedia.org/resource/ as name and place, and birthPlace in dbpedia.org/ontology/ as a property. To breakdown the Turtle in more details:

- Qprefix keyword is used to define prefixes for namespaces, making it easier to write URIs.
- dbr:Albert_Einstein is the subject, representing Albert Einstein as a resource in the DBpedia namespace.
- dbo:birthPlace is the predicate, representing the "birthPlace" property from the DBpedia ontology.
- dbr: Ulm is the object, representing the city of Ulm as a resource in the DBpedia namespace.
- . indicates the end of a statement.

We use

```
<object> <predicate> <subject> .
```

to claim a statement, and

```
<object> <predicate1> <subject1>; <predicate2> <subject2>; <predict3> <
    subject3> .
```

to assign multiple predicates to a subject to avoid repeating the same object in statements.

16.2.2 Multi-valued Relation and Blank Node

It is possible to use multi-valued relations and blank nodes to enforce combining of information. Consider an example given in Fig. 16.3, where the graph is used to demonstrate a lecture taking place at a specific room at given time slot. What if the lecture takes place twice a week, each time at a different location? With the help of multi-valued relations and blank nodes, this can be demonstrated clearly as shown in Fig. 16.4.

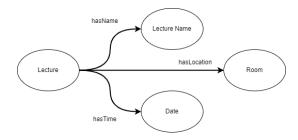


FIGURE 16.3

An example that demonstrate when and where a lecture takes place.

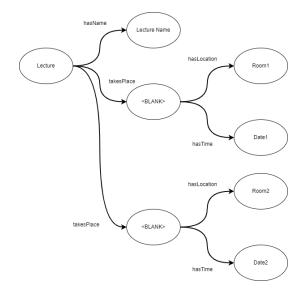


FIGURE 16.4

An example that demonstrate a lecture taking place at multiple locations and time slots using multi-valued relations and blank nodes.

Turtle and other markup languages provides syntax for creating blank nodes that looks like the following.

To refer to a blank node, it is also possible to give the blank node a name. The syntax looks like the following.

```
<subject1> <predicate1> _:<blank-node-name> .
_:<blank-node-name> <predicate2> <object2> .
_:<blank-node-name> <predicate3> <object3> .
```

where _: <black-node-name> is used to declare a blank node and assign it a name.

16.2.3 Lists

Lists help to make the code clean and readable. There are two types of lists, the container (open list, extendable) and the collections (closed list, fixed). Container is helpful to handle the situation given in Fig. 16.5. Notice that

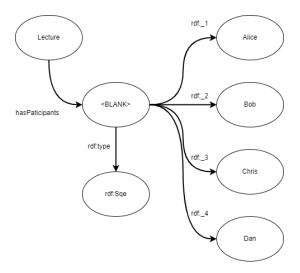


FIGURE 16.5

An example of a container.

the blank node has a type rdf:Seq. This tells that the items stored in the container follows an ordered set. There are other container types, such as rdf:Bag (unordered set) and rdf:Alt (alternatives of elements; only one element is relevant for the application).

The collection, on the other hand, defines a closed list as shown by Fig. 16.6. In Turtle, this can be done by using nested [] iteratively. A short cut

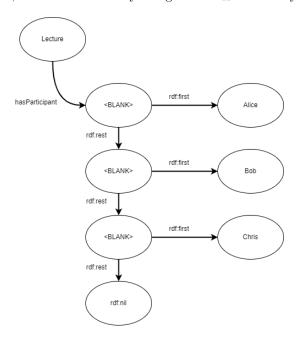


FIGURE 16.6

An example of a collection.

is to use (), with the items in the collection listed down in the bracket as follows.

<subject> <predicate> (<object1> <object2> <object3>) .

16.2.4 Reification

RDF permits interleaving of statements, i.e., to make a statement about another statement. For example, consider "Alice says that Bob ate the cake". This example contains nested statement, where "Bob ate the cake" is a statement, and "Alice says ..." is a statement on that statement. RDF reification follows Fig. 16.7.

16.2.5 Converting RDB to RDF

For small-scale relational database, it is possible to convert it to semantic web RDF manually. For example, consider an RDB for all the books in a study. There are three tables in the database, books, authors and publishers, respectively. Each table contains a few dozens of entries. The RDB can be

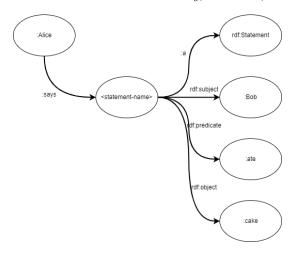


FIGURE 16.7

An example of RDF reification.

converted to RDF as shown in Fig. 16.8. It can be realized very easily, simply by defining everything as nodes and stack triples for all relationships.

However, when comes to a large and complicated database, converting it to RDF manually would consume too much labor. Several ways to systematically do the conversion have been proposed. Details are not covered here. See [6] for more details.

16.2.6 Conclusion

In summary, RDF, different from a plain XML or JSON which can also store independent classes, defines not only the objects (nodes) themselves but also the relationships among objects. This is essentially how RDF differs from a conventional key-value based NoSQL. RDF can be easily scaled up as nodes with the same name naturally merge together.

The underlying mechanism behind RDF, such as how machine stores graphical database including nodes and relationships, and how it enables query using SPARQL, is out of the scope of this notebook. In short, there are several ways to do that, for example by transforming the graph links into multiple small tables, etc. Different approaches may have their pros and cons. Details are neglected here.

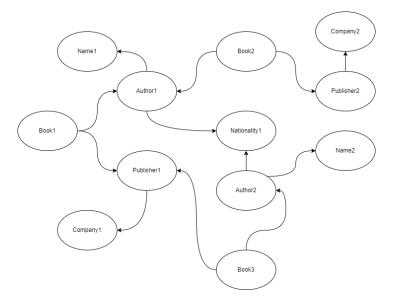


FIGURE 16.8

Semantic web of a few books, and their authors and publishers.

16.3 RDF Schema (RDFS)

RDFS enforces schema to the RDF model. By using RDF/RDFS, a more consistent and semantic RDF model can be achieved compared with using RDF along.

16.3.1 RDFS Motivation

RDF is flexible. It is so flexible that sometimes it becomes difficult to maintain consistency of the RDF model, let alone performing sophisticated reasoning from it. RDFS, also known as RDF vocabulary description language, enforce schema to the RDF model by adding more "meta information" which builds more connections among the nodes.

RDFS expands the vocabulary of RDF. It introduces the concepts of "class" and "subclass" to RDF. It provides built-in predefined classes such as rdfs:Literal, rdfs:Resource, rdfs:Datatype, etc., and enforce the nodes to be linked to these classes. RDF already defines rdf:Property. RDFS further expands the properties and relations. All above makes the RDF/RDFS modeling more consistent and semantic than using RDF alone.

In the deeper insights, RDFS helps to add "ontology" to the RDF model by introducing the schema. It essentially integrate common understanding and domain knowledge to the information. For example, by creating a property ":hasSpouse" whose domain and range person, it demonstrates the common knowledge that a person can be married to another person.

16.3.2 RDF Versus RDF/RDFS via an Example

The following example compares RDF and RDF/RDFS implementations on the same context. Consider statement "banana is yellow, apple is red, and orange is orange color". In a RDF implementation, this would look like the green-colored elements in Fig. 16.9. It is a disconnected graph. The semantic web failed to bridge the triples together and realize that all of them are discussing colors of fruits.

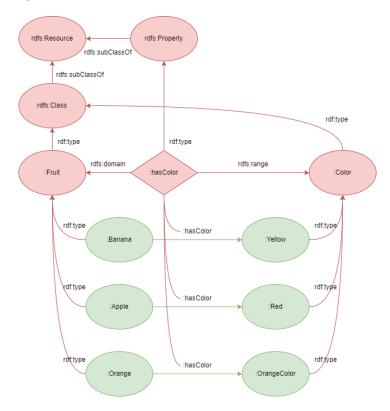


FIGURE 16.9

Semantic web of fruits and their colors, with RDF implementation in green RDF/RDFS in red.

In a RDF/RDFS implementation, class/subclass and property are enforced in the RDF model, and each property must be associated with clearly de-

fined domain and range. All classes eventually root to rdfs:Resource. This is shown by the green + red color in Fig. 16.9. The corresponding Turtle is

```
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
@prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#> .
@prefix example: <http://example.org/> .
# Define classes (optional; they are implicit)
rdfs:Class rdfs:subClassOf rdfs:Resource .
rdf:Property rdfs:subClassOf rdfs:Resource .
example:Fruit rdf:type rdfs:Class;
rdfs:subClassOf rdfs:Resource .
example:Color rdf:type rdfs:Class;
rdfs:subClassOf rdfs:Resource .
# Define properties
example:hasColor rdf:type rdf:Property;
rdfs:domain example:Fruit ;
rdfs:range example:Color .
# Define fruits and colors
example:Banana rdf:type example:Fruit .
example:Apple rdf:type example:Fruit .
example:Orange rdf:type example:Fruit .
example:Yellow rdf:type example:Color .
example:Red rdf:type example:Color
example:Green rdf:type example:Color .
# Define relationships between fruits and colors
example:Banana example:hasColor example:Yellow .
example:Apple example:hasColor example:Green .
example:Orange example:hasColor example:Orange
```

Notice that to build the RDF model shown by the green-colored elements in Fig. 16.9 (without RDFS), only the last 3 lines would be required. It can be seen from this example that RDF/RDFS uses more "complicated" structures to enforce schema of the model. In this example, "a fruit has color of a color" is enforced. Notice that predicate rdf:type can be used interchangeably with Turtle keyword a. They both declares a instance of a class.

16.3.3 RDFS Expanded Class and Properties

RDFS expands the vocabulary of RDF by introducing classes and properties hierarchy as well as domain and range for a property. As a brief summary, the following list shows some of the introduced concepts by RDFS.

• rdfs:Resource

- rdfs:Class
- rdf:Property (Notice that "Property" is already defined in RDF framework; RDFS enhanced its capability by introducing new mechanisms.)
- rdfs:subClassOf
- rdfs:subPropertyOf
- rdfs:domain
- rdfs:range

RDFS introduces additional powerful properties for a class. They help to make the RDF model more human-readable. The following is a short list of some of the new properties.

- rdfs:seeAlso points to a resource where a detailed explanation of the node can be found.
- rdfs:isDefinedBy defines the relation of a resource to its definition.
- rdfs:comment points to text comment.
- rdfs:label assign a more human-readable name to the node.

More details of the latest version of RDF/RDFS defined classes and properties are given by W3C and can be found at w3.org/TR/rdf-schema/.

16.3.4 Semantics inside RDF/RDFS

The semantics in RDF/RDFS are given by both the triples in the RDF model as well as the class and property hierarchies introduced by RDFS. Here are some examples of different types of hierarchies, and the semantics behind them.

- Class inheritance. If apple is a subclass of fruit, and fruit a subclass of plant, then apple must also be a subclass of a plant.
- Property domain and range. Let "hasColor" predicate be associated with domain "fruit" and range "color". If an object "hasColor", then the object must be a fruit, and corresponding subject in the triple must be a color.
- Property inheritance. Consider two predicates, "isMotherOf" and "isParentOf". Both predicates have the same domain and range of "person". Predicate "isMotherOf" is a sub property of "isParentOf". In this case "A is the mother of B" leads to "A is the parent of B".

The semantics allows some extent of reasoning, allowing "hidden information" to be derived.

16.4 SPARQL Protocol and RDF Query Language (SPARQL)

SPARQL is a widely used query and manipulation language for semantic web. It is SQL-like in the syntax, but its underlying mechanisms differ largely from how a relational database management system run SQL.

In the lately released SPARQL 1.1 standard, more operations such as advanced query and interfering are included, making it more powerful and capable than what SPARQL 1.0 standard described. SPARQL 1.1 is already supported by many triplestores.

Notice that SPARQL is not only a language, but also a protocol layer. The input and return have specific formats which are also defined by the SPARQL standard. Introducing SPARQL from a protocol perspective is not the focus of this chapter, hence it is not covered in details.

More details of SPARQL 1.1 can be found at w3.org/TR/sparql11-query/. SPARQL is not the only language for semantic web query and manipulation. It is good at general tasks. However, when comes to specific tasks such as expressive querying, data validation and advanced reasoning, there might be better choices.

16.4.1 SPARQL for Basic Query

SPARQL offers flexible ways of querying data. There are different commands to trigger a query, each fulfilling a different purpose. For example,

- SELECT returns a tabular just like SQL.
- CONSTRUCT returns a new RDF graph based on the query result.
- ASK returns true and false of whether a query has a solution.
- DESCRIBE returns the schematic of a resource; this is useful when the structure of RDF data in the data source is unclear.

The basic syntax of SPARQL looks similar with SQL as shown below.

where ?variableName denotes a variable, and without ?, a constant. Recall the fruit example used in Section 16.3.2. The following is an example of querying that semantic web.

```
PREFIX ex: <http://www.example.com/>
PREFIX rdf: <a href="http://www.w3.org/1999/02/22-rdf-syntax-ns#">http://www.w3.org/1999/02/22-rdf-syntax-ns#</a>
ASK WHERE {
       ?fruit rdf:type ex:Fruit .
       ?fruit ex:hasColor ex:Yellow .
} # return yes
SELECT ?fruit
WHERE {
        ?fruit ex:hasColor ex:Yellow .
} # return a table with one element, ex:Banana
CONSTRUCT {?fruit ex:hasColor ?color}
WHERE {
        ?fruit rdf:type ex:Fruit .
       ?color rdf:type ex:Color .
       ?fruit ex:hasColor ?color .
       FILTER (?color = ex:Yellow || ?color = ex:Red)
} # returns 2 triples, banana has color yellow and apple has color red
DESCRIBE ?fruit
WHERE {
       ?fruit ex:hasColor ex:Yellow .
} # return triples related to ex:Banana
```

The "FILTER" keyword can be used to quantitatively filter a numerical or string-like variable. An example is given below. It is worth mentioning that only triples clauses need to end with a period ".". The filter condition lead by FILTER () does not require the period.

Commonly used keywords and operators in FILTER clause include && (and), | | (or), ! (not), as well as string functions STR, LANG, CONTAINS, STRSTARTS, STRENDS, STRLEN, SUBSTR, REGEX (regular expression matching), etc., and numeric functions +, -, *, /, >, >=, <, <=, ABS, ROUND, CEIL, FLOOR, etc., and comparison operators =, !=.

There are also other clauses such as OPTIONAL and UNION. When OPTIONAL clause is applied on a property as part of the WHERE clause, it allows both entities with the correct property values as well as entities without the property to pass the filter. An example is given below. RDF:

```
@prefix foaf: <http://xmlns.com/foaf/0.1/> .
```

```
:alice
a foaf:Person;
foaf:name "Alice";
foaf:mbox <mailto:alice@example.com> .
:bob
a foaf:Person;
foaf:name "Bob" .
```

SPARQL:

The above returns both names and Alice's email. Though Bob does not have a property of foaf:mbox, his name is still included in the query result.

The UNION clause allows combining the returns of two queries together, given that the returns follow the same structure. An example is given below.

which would return both people and company names in one go.

16.4.2 SPARQL for Advanced Operations

SPARQL 1.0 provides basic query functions. In the lately released SPARQL 1.1 standard, advanced query and triple manipulation are supported. They are introduced in this section.

Advanced Query

In the SELECT clause, it allows simple calculations on top of the returned result, and name it as a new column. For example,

returns ?y*1.1 instead of ?y, and furthermore rename it as ?z. SPARQL 1.1 enables aggregate functions. For example,

counts the total number of fruits in the database. Notice that when using aggregate functions, new variable name (in this example, ?numOfFruit) must be assigned. Otherwise, there will be an syntax error.

SPARQL 1.1 supports nested query. An example is given below.

Triple Manipulation

SPARQL provides operations to insert, edit, and delete elements in a semantic web as follows.

• INSERT inserts triples into a graph.

• DELETE deletes triples from a graph.

An example of creating a semantic web that contains fruits and their colors is given below.

```
PREFIX ex: <http://www.example.com/>
PREFIX rdfs: <a href="http://www.w3.org/2000/01/rdf-schema">http://www.w3.org/2000/01/rdf-schema">http://www.w3.org/2000/01/rdf-schema</a>
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
INSERT DATA {
        # Define resources
        ex:Banana rdf:type ex:Fruit .
        ex:Apple rdf:type ex:Fruit .
        ex:Pear rdf:type ex:Fruit .
        ex:Yellow rdf:type ex:Color .
        ex:Red rdf:type ex:Color .
        ex:Green rdf:type ex:Color .
        # Define properties
        ex:hasColor rdf:type rdf:Property ;
        rdfs:domain ex:Fruit;
        rdfs:range ex:Color .
        # Link resources with properties
        ex:Banana ex:hasColor ex:Yellow .
        ex:Apple ex:hasColor ex:Red .
        ex:Pear ex:hasColor ex:Green .
```

16.4.3 Default Graph and Named Graph

There can be multiple semantic webs in a triplestore, each semantic web corresponding with a graph name. When querying, data can be retrieved across multiple semantic webs. If graph name is not specified in the SPARQL command in a insert command, the default graph of the database is used.

An example of inserting and querying a specific graph is given below.

```
GRAPH <http://example.org/mygraph> {
          ?book <http://purl.org/dc/elements/1.1/title> ?title .
}
```

The semantic web name (graph name) is given in the form of a URI, and can be used as a reference resources in other semantic webs.

16.4.4 SPARQL Programming Returns

SPARQL is an HTML based protocol. The returned result of a SPARQL query in the protocol layer can be specified in the HTML header. There are several return types defined in the standard. When SELECT or ASK are used, the return types can be:

- XML
- JSON
- TSV (table-separated values, similar to CSV)

When CONSTRUCT or DESCRIBE are used, the return is an RDF that can be in

- RDF/XML
- Turtle
- N-Triples
- JSON-LD

and a few more. The decoding of the return can be done using variety of tools and packages, and they are not discussed here.

16.4.5 Underlying Data Structure of Triplestores

RDF is a data model or framework that stores data in the form of triples. The benefit of doing so is to enable DL reasoning and hence build semantics into the data

When comes to the underlying data structure that actually processes the data in the computer memory and runs the RDF functions (such as SPARQL query based on graph pattern matching) efficiently, each triplestore may have its preferences. Some triplestores may develop dedicated graph database engines to store data and process queries. Others may utilize existing SQL or NoSQL database structures and engines and build an RDF interface on top of them.

The mechanisms of these different engines and data structures are not discussed in details here. Only a brief introduction is given as follows.

Table-Based Data Structure

An intuitive way of storing triples is to use a structured table that contains three columns: subject, predicate (property), and object. To save some space, all subjects, predicates and objects can be encoded into integers, where there are additional translation tables that can be used to decode the data. RDF should support multiple graphs. To enable multiple graphs, an additional column indicating graph ID needs to be added to all the tables mentioned above.

In this case, a SPARQL query needs to be converted into SQL (introducing self-join) then performed on the table. One of the major problems of this data structure is that it is too computationally expensive when the query is complicated.

To reduce self-join in the query, table size needs to be made small. The big structured tables need to be split into multiple small tables via aggregation. Due to the nature of triples, there will be many "NULL" in the tables. Implementing multi-value property, etc., can also be complicated. Overall, it is just too difficult to design the optimal RDB schematics, let along creating a database of that schematics.

To systematically and automatically split big tables into small tables without warring too much about the schematics, one idea is to destruct the table to the very ground level. Each property becomes a small table that contains its subjects and objects. The problem of this approach, however, is that it becomes time consuming to loop over all the tables when the RDF model is large. The performance of the architecture will be bad, and things such as inserting will become expensive.

NoSQL-Based Data Structure

Structured table and RDB approaches are rarely used in commercialized triplestores due to the expensive computational cost. The most widely used approaches nowadays are NoSQL based.

Two commonly seen NoSQL data structures are vertically partitioned tables and hexastores. These structures help to either reduce or speed up joins, making the computation more efficient than the RDB based methods.

Dedicated Graph Database

Some triplestores use dedicated graph database engines designed to natively store and process graph-structured data. These engines are optimized for the kinds of operations that are common in RDF and SPARQL, such as complex graph traversals and pattern matching.

17

Web Ontology Language

CONTENTS

17.1	RDF/RDFS Limitations	161
	OWL Vision	
17.3	OWL Basic Syntax	164
	OWL Advanced Syntax	
17.5	Semantic Web with Rules	169

By adding schema to RDF, RDFS already boosts the capability and adds ontology to the RDF model. As demonstrated in the previous chapter, RDF/RDFS is able to process simple logic reasoning. However, RDF/RDFS alone cannot handle DL. OWL is, therefore, proposed to enable DL in the semantic web.

Notice that it is OWL the abbreviation of web ontology language, not WOL, for historical contexts reasons.

17.1 RDF/RDFS Limitations

Several examples are given to demonstrate the limitations of RDF/RDFS.

RDFS brings schema to the RDF model by introducing class and property hierarchy. For example, in the schematic design of the model we can create "animal eats food", where "eat" is a property with domain "animal" and range "food". We can then add instances to animal and food classes. This is demonstrated in Fig. 17.1. The SPARQL to create such an RDF model is given below.

```
PREFIX ex: <a href="http://www.example.com/">http://www.w3.org/2000/01/rdf-schema#>
PREFIX rdf: <a href="http://www.w3.org/1999/02/22-rdf-syntax-ns#">http://www.w3.org/1999/02/22-rdf-syntax-ns#>
INSERT DATA {
    # Define resources
    ex:Human rdf:type ex:Animal .
    ex:Cow rdf:type ex:Animal .
```

```
ex:Vegetable rdfs:subClassOf ex:Food .
ex:Meat rdfs:subClassOf ex:Food .
ex:Cabbage rdf:type ex:Vegetable .
ex:Ribeye rdf:type ex:Meat .

# Define properties
ex:eat rdf:type rdf:Property ;
rdfs:domain ex:Animal ;
rdfs:range ex:Food .

# Define triples
ex:Human ex:eat ex:Meat .
ex:Human ex:eat ex:Cabbage .
ex:Cow ex:eat ex:Cabbage .
```

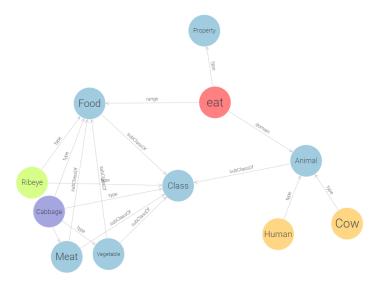


FIGURE 17.1

An RDF model that demonstrates "animal eats food".

From Fig. 17.1, under OWA when querying what animals eat what food, the following would be returned:

- Human eats rib-eye.
- Human eats cabbage.
- Cow eats rib-eye (this is not what we expect).
- Cow eats cabbage.

Since it is not explicitly denied that cow cannot eat meat, under OWA, semantic web thinks that there is the possibility of cows eating meat. RDF/RDFS cannot exclude meat from the menu of a cow, or to say it in general, RD-F/RDFS cannot add restrictions to a property.

A walk around is to define two properties, "eatMeat" and "eatVegetable", each with their associated ranges. Let human have both properties while cow have only "eatVegetable" property. However, this complicates the connection of symbols and destroys the schema of the model, essentially making null the effort of RDFS.

Similarly, consider an example where "hasVisualAcuity" is used to describe human's clarity of vision. A person usually have two visual acuity values for the two eyes. In RDF, that means a person should have two and only two "hasVisualAcuity" property. This cannot be enforced by RDF/RDFS.

Consider another example where "man is a subclass of human" and "woman is a subclass of human". With only RDF/RDFS, however, there is no way to add restrictions that says "a person cannot be man and woman at the same time", i.e., to disjoint man class and woman class.

RDF/RDFS does not support class combinations. For example, for existing classes "Car", "Motorcycle", "Bicycle", "Ship", "Plane", RDF/RDFS cannot create a new class "AllTranportationTool" to automatically include every instance from the above classes. In general, with RDF/RDFS alone, it is not possible to define a class which is a union/intersect/complement of other classes.

In semantic web stack Fig. 15.1, one layer above RDF/RDFS, OWL is proposed to solve the above problems.

17.2 OWL Vision

OWL essentially introduces DL inference to the semantic web, making it more expressive and capable of complicated reasoning. With that in mind, OWL provides additional features to enforce and enhance the schema of an RDF model, thus addressing the problems mentioned in Section 17.1. It allows adding relations and property constraints to the RDF model. The following summarizes the main features introduced by OWL.

- Allow the disjoint of sub classes (an instant cannot be man and woman at the same time).
- Allow enforcing the number of attributes of each property type (one can have only one social security number, but can have many mail addresses)
- Allow enforcing the range of the values an attribute can take (the age of a person must be a positive integer).

- Offer more expressive class definitions including union, intersection, complement, etc.
- Enlarge the vocabulary sets on top of RDF/RDFS.

As of this writing, OWL 2 is the latest version of OWL recommended by W3C. It makes the following assumptions:

- OWA. In this context, the knowledge base is considered open, and absence of information must not be valued as negative.
- An instance can have non unique names, i.e., multiple syntax can be mapped to the same instance via interpretation.

In practice, OWL comes in different flavors, including OWL Lite, OWL DL and OWL Full, just to name a few. OWL 2 supports different choices of syntax, including RDF-syntax, XML-syntax, functional-syntax, Manchester-syntax, and Turtle.

17.3 OWL Basic Syntax

Unless otherwise mentioned, Turtle is used when introducing OWL syntax. Classes, individuals, and properties are already introduced in RDF/RDFS. OWL also defines these concepts, each with enriched features. It is worth mentioning that OWL defines two special classes, owl:Thing and owl:Nothing, corresponding with the top class (universal set) and bottom class (empty set) in DL, respectively.

Class and Subclass

Defining a class from the root owl: Class using OWL is given below.

:Fruit a owl:Class;

where a can be replaced by rdf:type. To define an individual via class membership, simply use

:Apple a :Fruit .

which is equivalent of saying Fruit(Apple) in DL. It is also possible to define an individual without a named class as follows, which indicates that the named class will be added in a later stage. In this case, owl:NamedIndividual is used.

$: \verb"ANewThing" a owl: \verb"NamedIndividual" .$

Subclass can be defined similarly as follows.

```
:Fruit a owl:Class;
    rdfs:subClassOf :Food .
:Meat a owl:Class;
    rdfs:subClassOf :Food .
:Food a owl:Class .
```

Notice that it was impossible to enforce class disjoint using RDF/RDFS alone. This is made possible in OWL as follows.

```
:Fruit owl:disjointWith :Meat .
```

Or alternatively,

```
[] a owl:AllDisjointClasses ;
    owl:members
    ( :Fruit
    :Meat
    :SoftDrink
    :Beer ) .
```

which is a shortcut when disjoint is claimed on multiple classes. Similarly, class equivalence can be defined using owl:equivalentWith keyword.

Property

There are two types of properties defined in OWL, namely the object property which links to another resource URL (another object), and datatype property which links to a literal. Examples are given below.

```
:hasColor a owl:ObjectProperty ;
    rdfs:domain :Fruit ;
    rdfs:range :Color .
:hasShelfLife a owl:DatatypeProperty ;
    rdfs:domain :Fruit ;
    rdfs:range xsd:integer .
```

where notice that XML schema definition (XSD) defines many data types and sub data types, including xsd:integer, xsd:string, xsd:decimal, xsd:boolean, xsd:date (a calendar date), xsd:time (time in a day), xsd:dateTime, xsd:double (64-bit floating number), xsd:float (32-bit floating number), xsd:anyURI, xsd:language, etc.

17.4 OWL Advanced Syntax

Equivalent and Different Individuals

In RDF/RDFS, different individuals are distinguished by their names, and each individual can have one and only one name (it is technically possible

for an individual to have no name, but it is not often a good practice). In OWL, it is possible to either map multiple names to the same individual, or to emphasize that two names are pointing to different individuals (this is sometimes necessary due to the OWA). Examples are given below.

```
:Computer a owl:Class .
:HarryPc a :Computer .
:PotterPc a :Computer ;
owl:sameAs :HarryPc .
:DumbledorePc a :Computer ;
owl:differentFrom :HarryPc .
```

To state that individuals are different, alternatively use the following

```
[] a owl:AllDifferent;
    owl:distinctMembers
    ( :HarryPc,
    :DumbledorePc,
    :HermionePc,
    :HagridPc ) .
```

Closed Class

It is possible to define an enumerate class, whose instances are taken from individuals. An example is given below.

which indicates that any individuals to be defined under class "Weekdays" must be one of the 5 individuals "Monday", "Tuesday", etc. Do not confuse owl:oneOf with owl:unionOf, as the former defines a class as an enumerate of individuals, and the later defines a class as an enumerate of subclasses.

Class Constructors

Intersection, union and complement are the commonly seen class (set) constructors. They can be defined as follows.

```
:MediumWine
:SweetWine
)
].
:Plant a owl:Class;
rdfs:subClassOf [
owl:complementOf :Animal
].
```

Property Restrictions

It has been introduced earlier that the type of the property can be specified by using either owl:ObjectProperty (specify range as a class) or owl:DatatypeProperty (specify range as a datatype such as string, decimal, integer, etc). In the context of OWL, more restrictions can be enforced to properties, including the number of appearance of each property, and the values each property can take. Details are given below.

The syntax of adding different property restrictions may differ slightly. A general syntax is given by

```
:<Class> rdfs:subClassOf [
   rdf:type owl:Restriction ;
   owl:onProperty :<PropertyName> ;
   owl:<RestrictionType> <RestrictionDetails>
] .
```

where the blank node is used for property restrictions.

For example, to indicate that the value of a property must be taken from a class, use

```
:<Class> owl:restriction [
          owl:onProperty :<PropertyName> ;
          owl:allValuesFrom :<RestrictionClass>
] .
```

And to indicate that among all the property values of a property, at least one of them must take values from a class, replace owl:allValuesFrom with owl:someValuesFrom.

To indicate that a property must take specific value, use owl:hasValue at the restriction type, and put the value as <RestrictionDetails>, whether an individual or a datatype value.

To restrict the number of a property, use owl:maxCardinality, owl:minCardinality or owl:cardinality, and put the maximum, minimum or fixed number as <RestrictionDetails>.

It is possible to define a class from property restrictions, essentially saying "the class is defined as a collection of individuals that satisfy the following property restrictions". To do that, use rdfs:subClassOf with property restrictions as given in the example below.

where :Me is a predefined individual that maps to myself. To put it in words, a book is considered an instance of :MyBook if it has a :hasOwner property with the value :Me.

Property Hierarchy

Properties can have hierarchy like classes. Properties hierarchy can be realized using rdfs:subPropertyOf in RDF/RDFS framework. OWL further enriches that idea, by introducing new concepts such as owl:inverseOf. Examples below are used to demonstrate property hierarchy.

To define property hierarchy, use

```
:isOwnerOf a owl:ObjectProerty ;
    rdfs:subPropertyOf :isMasterOf ;
    rdfs:domain :Person ;
    rdfs:range :Book .
:isOwnedBy a owl:ObjectProperty ;
    owl:inverseOf :isOwnerOf ;
    rdfs:domain :Book ;
    rdfs:range :Person .
```

Notice that when defining a inverse property, it is not necessary, but still good practice, to point out the domain and range. This is for better readability and easier error checking.

OWL further defines the following features of a property.

- owl:TransitiveProperty: if Property(a,b) and Property(b,c), then Property(a,c). An example is "isAncestorOf".
- owl:SymmetricProperty: if Property(a,b), then Property(b,a). An example is "isColleagueWith"; there is also owl:AsymmetricProperty, which indicates that if Property(a,b), it is not possible that Property(b,a).
- owl:FunctionalProperty: if Property(a,b) and Property(a,c), then b=c. An example is hasMother.
- owl:InverseFunctionalProperty:ifProperty(a,c) and Property(b,c), then a=b. Examples are isMotherOf, hasId (assuming ID is unique for each individual).

Restrictions owl:disjointWith and owl:AllDisjointClasses can be

used to claim disjoint of classes, it is possible to declare disjunctive properties using owl:propertyDisjointWith and owl:AllDisjointProperties as follows. Disjunctive properties cannot take the same value. For example,

```
:hasParent a owl:ObjectProperty ;
    rdfs:domain :Person ;
    rdfs:range :Person .
[] a owl:AllDisjointProperties ;
    owl:members ( :hasParent :hasChild :hasSibling ) .
```

In OWL, top and bottom classes are defined, corresponding to the universal set and empty set, respectively. Similar ideas apply to properties. Object and datatype properties each has its top and bottom properties.

It is possible to declare a negative assertion as a property to state that a relation does not hold for two individuals. An example is given below.

```
[] a owl:negativePropertyAssertion ;
    owl:sourceIndividual :SherlockHolmes ;
    owl:assertionProperty :isMurderOf ;
    owl:targetIndividual :Watson .
```

which claims that it is not true that "SherlockHolmes" has the property "is-MurderOf" with the range "Watson". Notice that negative assertion is supported only at individual-to-individual level. It is currently not possible to claim that Sherlock Holmes is not a murder of any victims, i.e, :Watson cannot be replaced by a class such as :Person. In FOL, this would have been possible. But currently OWL does not adopt FOL for a good reason (FOL can be undecidable, and too computationally expensive).

General Role Inclusion

If we define "B is the father of A" and "C is the brother of B", then naturally, "C is the uncle of A" can be derived. This role chain can be realized via OWL general role inclusion. However, this adds uncertainty to the reasoning, and may cause the model to be undecidable. This is one of the reasons why different tiers of OWL have been proposed. More features usually mean more computations and higher risks of undecidable results, and the user needs to decide which tier to use.

And do notice that due to the Gödel's incompleteness theorems, there is no sophisticated enough system or knowledge base that can use finite input to derive all the knowledge in a complete (every statement can be proved true or false) and consistent (no contradiction) way.

Using role chain, we can define isUncleOf as follows.

```
:isUncleOf a owl:ObjectProperty ;
    owl:PropertyChainAxiom ( :hasFather :hasBrother ) .
```

It is not recommended of use role chain in a semantic web.

17.5 Semantic Web with Rules

Rules (rule-based systems, also known as rule-based engines) are used to express logic beyond DL, i.e., beyond what RDF/RDFS and OWL can describe. From this perspective, one can think of rules-based semantic web as a further enhancement beyond RDF/RDFS and OWL.

The basic syntax of a rule is simply

IF A THEN B

where A and B are premises and conclusion respectively. There are many variations, including logical rules (FOL rules, for example), procedural rules, and logic programming rules.

A FOL rule often looks like the following:

$$A \to B$$

or

$$A_1 \wedge \ldots \wedge A_n \to B$$

They can be equivalently converted to

$$\neg A \lor B$$

and

$$\neg A_1 \lor \ldots \lor \neg A_n \lor B$$

respectively.

Notice that implementing FOL rules (by using ALC just like we did for DL) in general would cause undecidability. However, it is possible to implement a subset of FOL rules while keeping the semantic web decidable. There are syntaxes to do that such as Declarative Logic Programming Language (also known as Datalog). By implementing the (decidable subset of) FOL rules, the expressiveness of the semantic web can be boosted.

There are semantic web triplestores that support the implementation of rules. Rules add expressiveness and meantime complexity and additional computational load to the semantic web.

Semantic Web Rule Language (SWRL) is based on the combination of parts of OWL and Datalog. The motivation of SWRL is to implement Datalog rules to the OWL ontology. Rule Interchange Format (RIF) is a collection of languages, rules, datatypes and frameworks recommended by W3C that tries to standardize rule-based semantic web description.

18

Semantic Web Practice

CONTENTS

18.1	Ontological Engineering				
18.2	Ontology Design	172			
	18.2.1 Ontology Management Activity	172			
	18.2.2 Ontology Design Basics	173			
18.3	Linked Data Engineering				
18.4					
18.5	Triplestore				
18.6	Example: Semantic Web for Home Assets	174			
	18.6.1 Define Classes Hierarchy	175			
	18.6.2 Define Properties Hierarchy	175			
	18.6.3 Add OWL	175			
	18.6.4 Data Retrieval Examples	175			
18.7		175			

This chapter studies the design, development and deployment of ontology and semantic web model. Both the methodologies and the tools are concerned.

18.1 Ontological Engineering

Different ontologies (class, property, hierarchy, interpretation, logic, etc.) can be used to describe the identical knowledge. This is known as the "problem of semantic gap". It is challenging to find the optimal and consistent way of representing the knowledge using ontology and semantic web.

Ontological engineering studies the systematic ways to design the ontology and the semantic web for a specific domain, application or task (recall Fig. 15.2). It concerns with at least the following:

- Ontology design: it studies the methods to systematically design, develop, and develop ontology models.
- Ontology mapping: it studies the methods to efficiently compare different ontology models.

- Ontology merging: it studies the methods to efficiently combine ontology models.
- Ontology learning: it studies the automatic learning of new knowledge by an existing ontology model, when new data sets are provided.

18.2 Ontology Design

Ontology design describes all activities necessary for the construction of an ontology model. It helps with consistent, efficient, and distributed (collaborative) ontology model design.

18.2.1 Ontology Management Activity

Design, develop and deploy ontology model for a sophisticated system can be time and manpower consuming. It is important to manage each step during the entire procedure to ensure healthy development of the system. This includes

- 1. Scheduling: identify tasks and problems to be solved by the semantic web; plan and arrange resources, time, manpower and money ahead.
- 2. Control: guarantee correct execution of tasks and problems to be solved.
- 3. Quality assurance: guarantee all steps are done correctly, including using the correct software, and everything is documented in details.

In pre-development stage, environment study needs to be carried out. This is mainly to identify what software to use to host the semantic web, and what interface/API shall the semantic have, and what applications would call the API to talk to the semantic web. A feasibility study is also necessary. For example, we need to consider whether it makes sense to develop a semantic web for the application, and whether the semantic web design is realizable.

In development stage, domain expert needs to come in to structure domain knowledge in a conceptual model. Knowledge engineer or data scientist then formalize the conceptual model into a computable (formal) model, then into ontology representation language.

Finally, in the post-development phase, pipeline needs to be designed to maintain, update and scale up and down the ontology model. In the case where the ontology model needs to be migrated into different platforms, used by unplanned applications, or merged with other models, necessary changes need to be made to the model. In the case when the knowledge in a model needs to be exported, knowledge recycle needs to be supported.

There are many ontology support activities. These activities need to be

carried out during the different stages of ontology development. Below is a list of these activities summarized in bullet points.

- Knowledge acquisition. Interview experts in the field and learn domain knowledge from them. This is referred as ontology learning. This is often done manually by the knowledge engineer or data scientist in the beginning stage. It is also possible to develop tools to automatically gather information and transfer it into ontology models, and even merge it with existing models.
- Technical evaluation. A domain expert checks occasionally to make sure the ontology model is correct.
- Integration and merging. This refers to the case where a big (scaled-up) ontology model can be built from a small existing ontology model.
- Alignment. Where there are multiple ontology models describing the same physical thing, alignment needs to be made to make sure knowledge consistency.
- Documentation and configuration management.

18.2.2 Ontology Design Basics

To design a complicated ontology model, many methods and techniques need to be used comprehensively. It is too complicated to be covered here.

This section introduces the basics of ontology design. It covers a subset of the aforementioned methods and techniques in a very brief manner.

18.3 Linked Data Engineering

"nobreak

18.4 Semantic Search

"nobreak

18.5 Triplestore

Triplestore is the database engine of the RDF model.

When it comes to relational database, there are many choices of DBMS such as Microsoft SQL Server, Oracle DBMS, MySQL and MariaDB. Similarly, there are many choices of triplestores for semantic web. A list of widely used triplestores is given below, just to name a few.

- GraphDB: a commercialized enterprise-tier semantic graph database management system compliant with W3C standards. It is famous for its performance and inference capabilities. It also provides free-tier for learning and for small projects, with limited capability.
- Apache Jena: an open-source Java framework for building semantic web and linked data applications. It has RDF APIs that can read and process RDF and SPARQL written in XML, Turble, JSON-LD and N-Triples.
- Virtuoso: a multi-model DBMS for both RDB and NoSQL databases such as RDF. It is famous for its scalability and standards compliance.
- AllegroGraph: a closed source triplestore which is designed to store RDF triples. It also operates as a document store designed for storing, retrieving and managing document-oriented information, in JSON-LD format. AllegroGraph is currently in use in commercial projects and a US Department of Defense project. It offers both free-tier license and enterprise-tier license.

More triplestores can be found at W3C, where a list of triplestores is maintained [1, 3]. As of this writing, there are 50 triplestores registered at W3C.

For the demonstrations given in this notebook, GraphDB is used, unless otherwise mentioned. A full instruction on using GraphDB, including applying for access and installation of the software, can be found at graphdb.ontotext.com.

Notice that RDF/RDFS/OWL/SPARQL APIs can be enabled using packages or libraries in many programming languages. For example, in Python, there are several packages available for semantic web operations, including rdflib, Owlready2, SPARQLWrapper, etc. Some of these packages have "lite" triplestore engine built-in which provides some SPARQL features for in-memory operations. They do not have all the features of a triplestore. However, they can connect to a triplestore API, in which case they serve as Python-to-triplestore interfaces.

18.6 Example: Semantic Web for Home Assets

As an example, we are creating a semantic web for a the assets in a household. Here "assets" refer to the electrical products, furniture, LEGOs, and other non-consumable, relatively static elements.

We will start with defining classes to divide everything into large groups, including electrical product, furniture, toy, etc. Under each class, sub-classes are defined, such as TV, computer, game console under electrical product, bed, chair, sofa, lamp under furniture, and LEGO under toy. Lastly, we will define instances under each sub-class. For example, bedroom TV and living room TV under TV, Nintendo SWITCH under game console, living room sofa under sofa, study computer, living room computer, TV attached computer under computer, etc.

We will then use RDFS to enforce schema to the RDF model as follows. Consider electrical product class for example. All elements in this class shall have a property called "hasBrand", which maps them to a pre-defined "electricalBrand" class, inside which are commonly seen electrical brands such as Boche, Siemens, Nintendo, Sony, Google, etc. The electrical product shall also have a "hasPrice" property, "warrantyExpiresAt" property, etc. The similar concept applies to all other produces including furniture, etc.

Finally, use OWL to setup some limits of the properties. For example, for electrical products, the price is usually between 0 to 5000 dollars, etc.

The result should be something like Fig. 18.1, after visualization.

18.6.1 Define Classes Hierarchy

"nobreak

18.6.2 Define Properties Hierarchy

"nobreak

18.6.3 Add OWL

"nobreak

18.6.4 Data Retrieval Examples

"nobreak

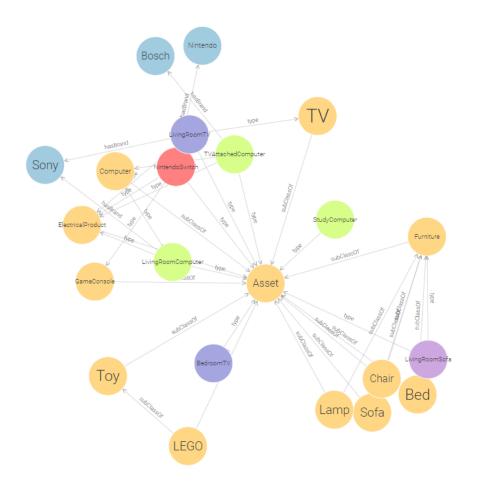


FIGURE 18.1

An example of an RDF model in GraphDB that describes house assets. This is only a demonstration graph and the information inside is artificial and not true.

TABLE 18.1

Commonly used name spaces in RDF models. URI is neglected since they can be easily found online.

Namespace	Description
rdf	RDF syntax.
rdfs	RDFS syntax.
xsd	XML syntax.
foaf	Friend-of-a-friend. It describes people, their activities and re-
	lations to other people and object.

18.7 Reference: Commonly Used Namespace

Commonly used built-in name spaces are summarized in Table 18.1. To search URI for a name space, use prefix.cc.

Bibliography

- [1] Category:triplestore.
- [2] The common layered semantic web technology stack.
- [3] Largetriplestores.
- [4] The R project for statistical computing.
- [5] Rstudio.
- [6] Franck Michel, Johan Montagnat, and Catherine Faron Zucker. A survey of RDB to RDF translation approaches and tools. PhD thesis, I3S, 2014.