

MASTER OF STATISTICS

STATISTICS IN BRIEF

A short handout for the MaSTAT students

GSEM, University of Geneva

CONTENTS

I Python guidance	7
1 Basic part	7
2 Advanced part	22
II R guidance	31
3 Basic part	31
4 Advanced part	60
5 Open data	64
6 Embedded Data Sets (in R and Python)	64
III Mathematics	66
7 Operations with vectors and matrices	66
8 Determinant and Trace	67
8.1 Definitions	67
8.2 Computation of the determinant	67
8.3 Properties	67
9 Eigenvectors and Eigenvalues	68
9.1 Definition	68
9.2 Quadratic forms and definite matrices	68
9.3 Theorem: Spectral decomposition	68
9.4 Theorem: Singular values decomposition	68
9.5 Differential calculus	69
10 Sequences and Series	70
10.1 Geometric series	70
11 Euler's number	70
12 Integrals	71
12.1 Antiderivatives of common functions	71
12.2 Integration by parts	71
12.3 Integration by substitution	73
IV Probability	75
13 Set Theory	75
13.1 Subset	75
13.2 Intersection	75
13.3 Union	76
13.4 Difference	76

13.5 Complement	76
13.6 Properties	76
14 Combinatorics	77
14.1 Factorial	77
14.2 Binomial Coefficient	77
15 Enumeration	77
15.1 simple arrangement	78
15.2 Arrangement with repetition	78
15.3 simple permutation	78
15.4 permutation with repetition	78
15.5 Simple combination	78
15.6 combination with repetition	78
16 Simple Probabilities	79
16.1 Conditional Probabilities	79
16.2 Independent Events	80
16.3 Theorem of total probabilities	80
16.4 Bayes' Theorem	80
17 Random Variable	80
17.1 Discrete Random Variable	81
17.2 Continuous Random Variable	82
17.3 Expected Values and Variance properties	83
18 Some Discrete Probability Distributions	84
18.1 Binomial Distribution	84
18.2 Poisson Distribution	84
19 Some Continuous Probability Distributions	84
19.1 Uniform Distribution	84
19.2 Exponential Distribution	85
19.3 Normal Distribution	86
20 Asymptotic theory	87
20.1 Convergence in Probability	87
20.2 Convergence in Distribution	87
20.3 Properties of the <i>plim</i> operator	87
21 Limit Theorems	88
21.1 Weak Law of Large Numbers	88
21.2 Central Limit Theorem	88
22 Concentration in Inequality	89
22.1 Markov's Inequality	89
22.2 Markov's Inequality	89
V Statictics	90
22.3 Quick summary	90
23 Random variables	90

23.1 Sampling schemes	91
23.2 Discrete random variables	92
24 Quantiles, percentiles...	94
25 Graphical representation	98
25.1 Q-Q plot	100
26 Estimation	104
26.1 Point estimation	104
26.2 Confidence intervals	105
26.2.1 Confidence interval for a population mean, σ known	106
26.2.2 Confidence interval for a population mean, σ unknown	106
27 Comparing estimators	106
27.1 Bias	107
27.2 Variance	108
27.3 MSE: Mean Squared Error	108
27.4 Efficiency	108
27.5 Consistency	109
28 Constructing estimators	111
28.1 Method of moments	111
28.2 Maximum likelihood	111
29 Type I and II error	117
30 Comparing 2 groups	118
31 Analysis of variance (ANOVA)	119
References	121

Acknowledgements

This guide is designed to help new graduate students in statistics to quickly review and familiarize themselves with the basic content of statistics.

The first version was jointly completed by Anna Venancio (anna.venancio@etu.unige.ch), Alice Scattolin (alice.scattolin@unige.ch), Chen Jiahao (Jiahao.Chen@etu.unige.ch), Eloi Favre (eloi.favre@etu.unige.ch), Laura Winkler (laura.winkler@etu.unige.ch), and Francisco Acosta (francisco.acosta@etu.unige.ch) in February 2023, under the supervision of Prof. Davide La Vecchia (davide.lavecchia@unige.ch). As time goes by, some of these items may no longer valid, please feel free to update and correct, and add useful topics.

Introduction

This book provides a short and to the point summary of the needed background knowledge and topics. If you feel like you are lacking some backgrorund knowledge, feel free to use [Moodle](#), the online student platform used at the University of Geneva. There, you can enroll into courses, even if you're not officially enrolled. You can download lecture materials and work on the practicals, to make sure you are up to speed and filled the blanks you've experienced.



Navigation

The exact locations of the topics that are being referenced will be inside these boxes.

Enjoy your statistics learning journey!

Part I.

Python guidance

1 BASIC PART

Python is a high-level scripting language that combines interpretability, compilation, interactivity, and object-oriented.

- **Python Download and Run**

Python3 can be downloaded from the official website of Python <https://www.python.org/>

You also can download Python documentation at <https://www.python.org/doc/link>

If you can successfully download and install python, then we have three methods to run python.

1.REPL (Read-Eval-Print Loop)

You can enter Python through the command line window and start writing code in the interactive interpreter. Generally, you can do this on Unix, DOS, or any other system that provides a command line or shell.

2.Command-line scripts

Python scripts can be executed on the command line by including the interpreter in your application, as follows:

```

1 $ python script.py      # Unix/Linux
2
3 #or
4
5 C:>python script.py    # Windows/DOS

```

Listing 1: Python example

3.IDE: Integrated Development Environment: PyCharm or Anaconda, Jupyter Notebook
 Integrated Development Environment (IDE) is an application program used to provide a program development environment, generally including tools such as code editors, compilers, debuggers, and graphical user interfaces. An integrated development software service suite that integrates code writing functions, analysis functions, compilation functions, and debugging functions. All software or software suites (groups) with this feature can be called integrated development environments. Such as Microsoft's Visual Studio series, Borland's C++ Builder, Delphi series, etc. The program can run independently,

or it can be used together with other programs. The integrated environment of python mainly includes PyCharm or Anaconda, Jupyter Notebook.

- **Identifier**

With Python, you name your own variables. Just follow these rules:

1. Characters must be letters of the alphabet or underscore
2. The rest of the identifier consists of letters, numbers and underscores
3. Identifiers are case sensitive
4. Words such as "False" and "None" have been used by python(Called "Reserved word"), so they cannot be used as identifiers. Python's standard library provides a keyword module that outputs all Reserved words of the current version:

```

1 import keyword
2
3 print(keyword.kwlist)
4
5 >>> ['False', 'None', 'True', 'and', 'as', 'assert', 'break', 'class', 'continue', 'def', 'del',
6      'elif', 'else', 'except', 'finally', 'for', 'from', 'global', 'if', 'import', 'in',
7      'is', 'lambda', 'nonlocal', 'not', 'or', 'pass', 'raise', 'return', 'try', 'while', 'with', 'yield']

```

Listing 2: Reserved words

- **Note**

Single-line comments in Python start with #, and multi-line comments are implemented using """.

```

1 # first comment
2
3 # Second comment
4
5 """
6 Third comment 1
7
8 Third comment 2
9 """
10
11 print ("Hello, Python!")
12 >>> Hello, Python!

```

Listing 3: Notes and Comments

- **Operator**

Python Operators in general are used to perform operations on values and variables. These are standard symbols used for the purpose of logical and arithmetic operations.

```

1 # Examples of Arithmetic Operator
2 a = 9

```

```

3 b = 4
4
5 # Addition of numbers
6 add = a + b
7
8 # Subtraction of numbers
9 sub = a - b
10
11 # Multiplication of number
12 mul = a * b
13
14 # Division(float) of number
15 div1 = a / b
16
17 # Division(floor) of number
18 div2 = a // b
19
20 # Modulo of both number
21 mod = a % b
22
23 # Power
24 p = a ** b
25
26 # print results
27 print(add)
28 >>>13
29 print(sub)
30 >>>5
31 print(mul)
32 >>>36
33 print(div1)
34 >>>2.25
35 print(div2)
36 >>>2
37 print(mod)
38 >>>1
39 print(p)
40 >>>656

```

Listing 4: Example of Arithmetic Operator

- **Data type**

Variables in Python do not need to be declared. Each variable must be assigned a value before use, and the variable will not be created until the variable is assigned. In Python, a variable is a variable, it doesn't have a type, what we mean by "type" is the type of object in memory that the variable refers to. The equal sign (=) is used to assign values to variables.

The left side of the equal sign (=) operator is a variable name, and the right side of the equal sign (=) operator is the value stored in the variable. E.g:

```

1 counter = 100          # integer
2 miles    = 1000.0       # float
3 name     = "unige"      # string
4
5 print (counter)
6 >>>100
7 print (miles)
8 >>>1000.0
9 print (name)
10 >>>unige

```

Listing 5: Examples of Basic Data Types

Python allows you to assign values to multiple variables at the same time. E.g.:

```

1 a = b = c = 1
2
3 aaa,bbb,ccc = 1, 2, "unige"

```

Listing 6: Multiple Variable Assignment

There are six standard data types in Python3: Number, String (string), List (list), Tuple, Set (collection), Dictionary, Among the six standard data types of Python3: Immutable data (3): Number (number), String (string), Tuple (tuple); Variable data (3): List (list), Dictionary (dictionary), Set (collection). We will discuss them later.

- **Number**

Python3 supports int, float, bool, complex (complex numbers). Like most languages, assignment and evaluation of numeric types is straightforward. The built-in type() function can be used to query what type of object a variable refers to.

```

1 a, b, c, d = 20, 5.5, True, 4+3j
2
3 print(type(a), type(b), type(c), type(d))
4
5 >>> <class 'int'> <class 'float'> <class 'bool'> <class 'complex'>

```

Listing 7: Function type() Example

- **String**

Strings in Python are enclosed in single quotes or double quotes while special characters are escaped with a backslash. The grammatical format of string interception is as follows: variable [head subscript: tail subscript]

```

1 str = 'unigestat'
2
3 print (str)
4 # output is string
5 >>> unigestat
6 print (str[0:-1])
7 # Output is all characters from the first to the second last
8 >>> unigesta

```

```

9 print (str[0])
10 # output is the first character of the string
11 >>> u
12 print (str[2:5])
13 # Output is the characters from the third to the fifth
14 >>> ige
15 print (str[2:])
16 # Output is all characters starting from the third
17 >>> igestat
18 print (str * 2)
19 # output is the string twice
20 >>> unigestatunigestat
21 print (str + "TEST")
22 # merge
23 >>> unigestatTEST

```

Listing 8: String Example 1

Python uses the backslash to escape special characters. If you don't want the backslash to be escaped, you can add an r in front of the string to represent the original string.

```

1 print('unige\stats')
2 >>>unige
3 >>>stats
4
5 print(r'unige\stats')
6 >>> unige\stats

```

Listing 9: String Example 2

- **List**

List is the most frequently used data type in Python. Lists can complete the data structure implementation of most collection classes. The types of elements in the list can be different, it supports numbers, strings can even contain lists (so-called nesting). A list is a comma-separated list of elements written between square brackets [].

Like strings, lists can also be indexed and truncated. After the list is truncated, a new list containing the required elements is returned. The syntax format of list interception is as follows:

variable [head subscript: tail subscript]

```

1 list1 = ['abcd', 786, 2.23, 'unige', 70.2]
2 tinylist = [123, 'unige']
3
4 print (list1)
5 # output is the complete list
6 >>> ['abcd', 786, 2.23, 'unige', 70.2]
7 print (list1[0])
8 # output is the first element of the list

```

```

9 >>> abcd
10 print (list1[1:3])
11 # output is from the second to the third element
12 >>> [786, 2.23]
13 print (list1[2:])
14 # output is all elements starting from the third element
15 >>> [2.23, 'unige', 70.2]
16 print (tinylist * 2)
17 # output is the list twice
18 >>> [123, 'unige', 123, 'unige']
19 print (list1 + tinylist)
20 # merge
21 >>> ['abcd', 786, 2.23, 'unige', 70.2, 123, 'unige']

```

Listing 10: List Example

The elements of a list can be changed:

```

1 a = [1, 2, 3, 4, 5, 6]
2 a[0] = 9
3 a[2:5] = [13, 14, 15]
4 print(a)
5 >>> [9, 2, 13, 14, 15, 6]
6 a[2:5] = [] # Set the corresponding element value to []
7 print(a)
8 >>> [9, 2, 6]

```

Listing 11: List Example 2

List has many built-in methods, such as append(), pop(), etc. Interested students can refer to relevant information.

- **Tuple**

A tuple is similar to list, except that the elements cannot be modified. The element types in the tuple can also be different:

```

1 tuple1 = ('abcd', 786, 2.23, 'unige', 70.2)
2 tinytuple = (123, 'unige')
3
4 print (tuple1)
5 # output is the complete tuple
6 >>> ('abcd', 786, 2.23, 'unige', 70.2)
7 print (tuple1[0])
8 # output is the first element of the tuple
9 >>> abcd
10 print (tuple1[1:3])
11 # output starts from the second element to the third element
12 >>> (786, 2.23)
13 print (tuple1[2:])
14 # output all elements starting from the third element
15 >>> (2.23, 'unige', 70.2)
16 print (tinytuple * 2)

```

```

17 # output is tuple twice
18 >>> (123, 'unige', 123, 'unige')
19 print (tuple1 + tinytuple)
20 # merge
21 >>> ('abcd', 786, 2.23, 'unige', 70.2, 123, 'unige')

```

Listing 12: Tuple Example 1

Similar to strings, tuples can be indexed and the subscript index starts from 0, and -1 is the position from the end. In fact, strings can be thought of as a special kind of tuple.

```

1 tup = (1, 2, 3, 4, 5, 6)
2 print(tup[0])
3 >>> 1
4 print(tup[1:5])
5 >>> (2, 3, 4, 5)
6 tup[0] = 11 # Operations that modify tuple elements are illegal
7 >>> Traceback (most recent call last):
8   File "<stdin>", line 1, in <module>
9 TypeError: 'tuple' object does not support item assignment

```

Listing 13: Tuple Example 2

- **Set**

A set is composed of one or several wholes of various shapes and sizes, and the things or objects that constitute a set are called elements or members. The basic functionality is to perform membership tests and remove duplicate elements. Sets can be created using curly braces {} or the set() function. Note: To create an empty set, you must use set() instead of {}, because {} is used to create an empty dictionary.

```

1 parame = {value01,value02,...}
2 # or
3 set(value)

```

Listing 14: Set Example 1

Here are two examples using set.

```

1 sites = {'Google', 'Taobao', 'Runoob', 'Facebook', 'Zhihu', 'Baidu'}
2
3 print(sites)
4 # Output collection, duplicate elements are automatically removed
5
6 >>> {'Zhihu', 'Baidu', 'Taobao', 'Runoob', 'Google', 'Facebook'}
7
8 # membership test
9 if 'Runoob' in sites:
10     print('Runoob is in the collection')
11 else:
12     print('Runoob is not in the set')

```

```

13
14 >>> 'Runoob is in the set'
15
16 # set can perform set operations
17 a = set('abracadabra')
18 b = set('alacazam')
19
20 print(a)
21 >>> {'b', 'c', 'a', 'r', 'd'}
22 print(a - b) # difference between a and b
23 >>> {'r', 'b', 'd'}
24 print(a | b) # union of a and b
25 >>> {'b', 'c', 'a', 'z', 'm', 'r', 'l', 'd'}
26 print(a & b) # intersection of a and b
27 >>> {'c', 'a'}
28 print(a ^ b) # elements that do not exist in a and b
29 >>> {'z', 'b', 'm', 'r', 'l', 'd'}

```

Listing 15: Set Example 2

- **Dictionary**

Dictionaries are another very useful built-in data type in Python. A list is an ordered collection of objects, and a dictionary is an unordered collection of objects. The difference between the two is that the elements in the dictionary are accessed by key, not by offset. A dictionary is a mapping type, and a dictionary is marked with , which is an unordered collection of key (key): value (value). Keys must use immutable types. Within the same dictionary, keys must be unique.

```

1 dict1 = {}
2 dict1['one'] = "1 - unige"
3 dict1[2] = "2 - stats"
4
5 tinydict = {'name': 'gsem', 'code':1, 'site': 'https://www.unige.ch/gsem'}
6
7
8 print (dict1['one']) # output the value whose key is 'one'
9 >>> 1 - unige
10 print (dict1[2]) # output the value whose key is 2
11 >>> 2 - stats
12 print (tinydict) # output the complete dictionary
13 >>> {'name': 'gsem', 'code': 1, 'site': 'https://www.unige.ch/gsem'}
14 print (tinydict.keys()) # output all keys
15 >>> dict_keys(['name', 'code', 'site'])
16 print (tinydict.values()) # output all values
17 >>> dict_values(['gsem', 1, 'https://www.unige.ch/gsem'])

```

Listing 16: Dictionary Example

- **Data Type Conversions**

Sometimes, we need to convert the built-in data type, data type conversion, in general, you only need to use the data type as the function name. Python data type conversion can

be divided into two types: Implicit type conversions - auto-completion and Explicit type conversion - need to use type function to convert.

1. Implicit type conversions

In implicit type conversion, Python automatically converts one data type to another without our intervention. In the following example, we operate on two different types of data, and the lower data type (integer) is converted to the higher data type (float) to avoid data loss.

```

1 num_int = 123
2 num_flo = 1.23
3
4 num_new = num_int + num_flo
5
6 print("datatype of num_int:",type(num_int))
7 >>> datatype of num_int : <class 'int'>
8
9 print("datatype of num_flo:",type(num_flo))
10 >>> datatype of num_flo : <class 'float'>
11
12 print("value of num_new:",num_new)
13 >>> value of num_new: 124.23
14
15 print("datatype of num_new:",type(num_new))
16 >>> datatype of num_new : <class 'float'>
```

Listing 17: Implicit Type Conversions 1

Code analysis:

In this example, we add two variables num-int and num-flo of different data types and store them in the variable num-new. Then we check the data types of the three variables. In the output, we see that num-int is an integer and num-flo is a float. Also, the new variable num-new is float because Python converts smaller data types to larger data types to avoid data loss. Let's look at another example, adding integer data and string data:

```

1 num_int = 123
2 num_str = "456"
3
4 print("Data type of num_int:",type(num_int))
5 print("Data type of num_str:",type(num_str))
6 print(num_int+num_str)
```

Listing 18: Implicit Type Conversions 2

output as follows:

```

1 >>> Data type of num_int: <class 'int'>
2 >>> Data type of num_str: <class 'str'>
3 >>> Traceback (most recent call last):
4   File "/Python-test/test.py", line 7, in <module>
```

```

5     print(num_int+num_str)
6 TypeError: unsupported operand type(s) for +: 'int' and 'str'

```

Listing 19: Implicit Type Conversions 2 Output

It can be seen from the output that the operation results of integer and string types will report an error, and `TypeError` will be output. Python cannot use implicit conversions in this case. However, Python provides a solution for these types of situations called explicit conversions.

2.Explicit type conversion

In explicit type conversion, the user converts the data type of the object to the desired data type. We use predefined functions like `int()`, `float()`, `str()` etc. to perform explicit type conversions.

`int()` converts to an integer:

```

1 x = int(1)
2 # x output is 1
3
4 y = int(2.8)
5 # y output is 2
6
7 z = int("3")
8 # z output is 3

```

Listing 20: Explicit Type Conversion 1

`float()` converts to a float:

```

1 x = float(1)
2 # x output is 1.0
3
4 y = float(2.8)
5 # y output is 2.8
6
7 z = float("3")
8 # z output is 3.0
9
10 w = float("4.2")
11 # w output is 4.2

```

Listing 21: Explicit Type Conversion 2

`str()` converts to a string type:

```

1 x = str("s1")
2 # The output of x is 's1'
3
4 y = str(2)
5 # y output is '2'

```

```

6
7 z = str(3.0)
8 # z output is '3.0'

```

Listing 22: Explicit Type Conversion 3

Operations on integer and string types can be done by those function.

```

1 num_int = 123
2 num_str = "456"
3
4 print("num_int data type is:",type(num_int))
5 print("Before type conversion, the data type of num_str is:",type(num_str))
6
7 num_str = int(num_str) # cast to integer
8 print("After type conversion, the data type of num_str is:",type(num_str))
9
10 num_sum = num_int + num_str
11
12 print("The result of adding num_int and num_str is:",num_sum)
13 print("sum data type is:",type(num_sum))

```

Listing 23: Explicit Type Conversion 4

output as follows:

```

1 The num_int data type is: <class 'int'>
2 Before type conversion, the data type of num_str is: <class 'str'>
3 After type conversion, the data type of num_str is: <class 'int'>
4 The result of adding num_int and num_str is: 579
5 The sum data type is: <class 'int'>

```

Listing 24: Explicit Type Conversion 4 Output

- **If-else Statement**

The general form of an if statement in Python is as follows:

```

1 if condition_1:
2     statement_block_1
3 elif condition_2:
4     statement_block_2
5 else:
6     statement_block_3

```

Listing 25: If-else Statement Example 1

If "condition 1" is True the "statement block 1" block statement will be executed. If "condition 1" is False, "condition 2" will be evaluated. If "condition 2" is True the "statement block 2" block statement will be executed. If "condition 2" is False, the "statement block 3" block statement will be executed.

```

1 dog_age = 15
2
3 if age <= 0:
4     print("You are kidding me!")
5 elif age == 1:
6     print("Equivalent to a 14-year-old.")
7 elif age == 2:
8     print("Equivalent to a 22-year-old.")
9 elif age > 2:
10    human = 22 + (age -2)*5
11    print("corresponding to human age: ", human)
12
13 >>> corresponding to human age: 87

```

Listing 26: If-else Statement Example 2

In nested if statements, you can place an if...elif...else structure inside another if...elif...else structure. The following is an example of nested if-else statements.

```

1 num = 6
2
3 if num%2==0:
4     if num%3==0:
5         print ("The number you entered is divisible by 2 and 3")
6     else:
7         print ("The number you entered is divisible by 2, but not by 3")
8 else:
9     if num%3==0:
10        print ("The number you entered is divisible by 3, but not by 2")
11    else:
12        print ("The number you entered is not divisible by 2 and 3")
13
14 >>> "The number you entered is divisible by 2 and 3"

```

Listing 27: If-else Statement Example 3

Python 3.10 adds the conditional judgment of match...case, no need to use a series of if-else to judge. The object after the match will be matched with the content after the case in turn. If the match is successful, the matched expression will be executed, otherwise it will be skipped directly. The underscore can match everything. The syntax format is as follows:

```

1 match subject:
2     case <pattern_1>:
3         <action_1>
4     case <pattern_2>:
5         <action_2>
6     case <pattern_3>:
7         <action_3>
8     case _:
9         <action_wildcard>

```

Listing 28: Match Case Statement Example 1

The following is an example of Match Case statements.

```

1 mystatus = 400
2 print(http_error(400))

3
4 def http_error(status):
5     match status:
6         case 400:
7             return "Bad request"
8         case 404:
9             return "Not found"
10        case 418:
11            return "I'm a teapot"
12        case _:
13            return "Something's wrong with the internet"
14
15 >>> Bad request

```

Listing 29: Match Case Statement Example 2

- **For-loop Statement**

A Python for loop can iterate over any iterable object, such as a list or a string. The general format of a for loop is as follows:

```

1 for <variable> in <sequence>:
2     <statements>
3 else:
4     <statements>

```

Listing 30: For-loop Statement Example 1

Python for loop example:

```

1 sites = ["Baidu", "Google", "Runoob", "Facebook"]
2 for site in sites:
3     print(site)
4
5 >>> Baidu
6     Google
7     Runoob
8     Taobao

```

Listing 31: For-loop Statement Example 2

Integer range values can be used with the range() function:

```

1 for number in range(1, 6):
2     print(number)
3
4 >>> 1

```

```

5   2
6   3
7   4
8   5

```

Listing 32: For-loop Statement Example 3

In Python, the for...else statement is used to execute a block of code after the loop ends. When the loop is executed (that is, all elements in the iterable are traversed), the code in the else clause will be executed. If a break statement is encountered during the loop, the loop will be interrupted, and the else clause will not be executed at this time.

```

1 for x in range(6):
2     print(x)
3 else:
4     print("Finally finished!")
5
6 >>> 0
7   1
8   2
9   3
10  4
11  5
12 Finally finished!

```

Listing 33: For-loop Statement Example 4

The break statement is used in the following for example, the break statement is used to jump out of the current loop body, and the else clause will not be executed:

```

1 sites = ["Baidu", "Google", "Taobao", "Runoob"]
2 for site in sites:
3     if site == "Taobao":
4         print("Taobao!")
5         break
6     print("loop" + site)
7 else:
8     print("no loop!")
9 print("finish loop!")
10
11 >>> loop Baidu
12   loop Google
13   Taobao!
14   finish loop!

```

Listing 34: For-loop Statement Example 5

- **While-loop Statement**

The general form of a while statement in Python:

```

1 while Judgment condition (condition):

```

2 Execute statements...

Listing 35: While-loop Statement Example 1

The following example uses while to calculate the sum from 1 to 100:

```

1 n = 100
2
3 sum = 0
4 counter = 1
5 while counter <= n:
6     sum = sum + counter
7     counter += 1
8
9 print("The sum of 1 to %d is: %d" % (n,sum))
10
11 >>> The sum of 1 to 100 is: 5050

```

Listing 36: While-loop Statement Example 2

Use break in while-loop Statement:

```

1 n = 5
2 while n > 0:
3     n -= 1
4     if n == 2:
5         break
6     print(n)
7 print('The end of the loop.')
8 >>> 4
9     3
10    The end of the loop.

```

Listing 37: While-loop Statement Example 3

Use continue in while-loop Statement:

```

1 n = 5
2 while n > 0:
3     n -= 1
4     if n == 2:
5         continue
6     print(n)
7 print('The end of the loop.')
8 >>> 4
9     3
10    1
11    0
12    The end of the loop.

```

Listing 38: While-loop Statement Example 4

- **Function**

Python defines functions using the def keyword, and the general format is as follows:

```
1 def function name (parameters):
2     body
```

Listing 39: Function Example 1

Let's use a function to output "Hello World!":

```
1 def hello() :
2     print("Hello World!")
3 hello()
4 >>> "Hello World!"
```

Listing 40: Function Example 2

More complex applications, with parameter variables in the function. Task: compare two numbers and return the larger number:

```
1 def max(a, b):
2     if a > b:
3         return a
4     else:
5         return b
6
7 a = 4
8 b = 5
9 print(max(a, b))
10 >>> 5
```

Listing 41: Function Example 3

Now you have mastered the basic python programming ability. Next, let me try some more advanced applications.

2 ADVANCED PART

- **Data Structure**

In this chapter, we mainly introduce the Python data structure based on the knowledge points learned earlier.

1. Advanced part of list

Lists in Python are mutable, which is the most important feature that distinguishes them from strings and tuples. In one sentence, lists can be modified, but strings and tuples cannot.

Here are the methods for lists in Python:

Methods	Description
list.append(x)	Adds an element, equivalent to <code>a[len(a):] = [x]</code>
list.extend(L)	Extends the list, equivalent to <code>a[len(a):] = L</code>
list.insert(i, x)	Inserts an element at the specified position.
list.remove(x)	Removes the first element in the list with value x.
list.pop([i])	Removes the element from the specified position.
list.clear()	Removes all items from the list, equals <code>del a[:]</code>
list.index(x)	Returns the index of the first element in the list whose value is x
list.count(x)	Returns the number of times x occurs in the list
list.sort()	Sorts the elements in the list
list.reverse()	Reverses the elements in the list
list.copy()	Returns a shallow copy of the list, equal to <code>a[:]</code>

The following example demonstrates most of the list methods:

```

1 a = [66.25, 333, 333, 1, 1234.5]
2 print(a.count(333), a.count(66.25), a.count('x'))
3 >>> 2 1 0
4 a.insert(2, -1)
5 a.append(333)
6 print(a)
7 >>> [66.25, 333, -1, 333, 1, 1234.5, 333]
8 print(a.index(333))
9 >>>1
10 a.remove(333)
11 print(a)
12 >>>[66.25, -1, 333, 1, 1234.5, 333]
13 a.reverse()
14 print(a)
15 >>>[333, 1234.5, 1, 333, -1, 66.25]
16 a.sort()
17 print(a)
18 >>>[-1, 1, 66.25, 333, 333, 1234.5]

```

Listing 42: Data Structure 1

The list method makes it convenient to use the list as a stack. The stack is a specific data structure, and the first element entered is the last one released (last in, first out). An element can be added to the top of the stack with the `append()` method.

Use the `pop()` method without specifying an index to release an element from the top of the stack. E.g:

```

1 stack = [3, 4, 5]
2 stack.append(6)
3 stack.append(7)
4 stack

```

```

5 >>>[3, 4, 5, 6, 7]
6 stack.pop()
7 >>>7
8 stack
9 >>>[3, 4, 5, 6]
10 stack.pop()
11 >>>6
12 stack.pop()
13 >>>5
14 stack
15 >>>[3, 4]

```

Listing 43: Data Structure 2

2.Advanced part of tuples and sequences

As follows, tuples are always output with parentheses to facilitate the correct expression of nested structures. They may or may not be entered with or without parentheses, but they are usually required (if the tuple is part of a larger expression).

```

1 t = 12345, 54321, 'hello!'
2 t[0]
3 >>> 12345
4 t
5 >>> (12345, 54321, 'hello!')
6 # Tuples may be nested:
7 ... u = t, (1, 2, 3, 4, 5)
8 >>> u
9 >>> ((12345, 54321, 'hello!'), (1, 2, 3, 4, 5))

```

Listing 44: Data Structure 3

3.Advanced part of Set

A set is an unordered collection of unique elements. Basic functionality includes relationship testing and elimination of duplicate elements.

Collections can be created with braces (). Note: If you want to create an empty set, you must use set() instead of ; the latter creates an empty dictionary.

```

1 basket = {'apple', 'orange', 'apple', 'pear', 'orange', 'banana'}
2
3 print(basket) # delete duplicates
4 >>> {'orange', 'banana', 'pear', 'apple'}
5
6 'orange' in basket # detect members
7 >>> True
8
9 'crabgrass' in basket
10 >>> False
11
12 # The following demonstrates the operation of two collections
13 a = set('abracadabra')
14 b = set('alacazam')

```

```

15 a # unique letter in a
16 >>> {'a', 'r', 'b', 'c', 'd'}
17
18 a - b # letter in a but not in b
19 >>> {'r', 'd', 'b'}
20
21 a | b # letter in a or b
22 >>> {'a', 'c', 'r', 'd', 'b', 'm', 'z', 'l'}
23
24 a & b # letters in both a and b
25 >>> {'a', 'c'}
26
27 a ^ b # letter in a or b, but not both a and b
28 >>> {'r', 'd', 'b', 'm', 'z', 'l'}

```

Listing 45: Data Structure 4

4.traversals skills traversal skills

When traversing through a dictionary, the key and the corresponding value can be interpreted at the same time using the items() method:

```

1 basket = {'apple', 'orange', 'apple', 'pear', 'orange', 'banana'}
2 knights = {'gallahad': 'the pure', 'robin': 'the brave'}
3 for k, v in knights.items():
4     print(k, v)
5
6 >>> gallahad the pure
       robin the brave
7

```

Listing 46: Data Structure 5

When traversing in a sequence, the index position and corresponding value can be obtained at the same time using the enumerate() function:

```

1 for i, v in enumerate(['tic', 'tac', 'toe']):
2     print(i, v)
3
4 >>> 0 tic
5     1 tac
6     2 toe

```

Listing 47: Data Structure 6

Iterating over two or more sequences simultaneously can be combined using zip():

```

1 questions = ['name', 'quest', 'favorite color']
2 answers = ['lancelot', 'the holy grail', 'blue']
3 for q, a in zip(questions, answers):
4     print('What is your {0}? It is {1}.'.format(q, a))
5
6 >>> What is your name? It is lancelot.

```

```

7 What is your quest? It is the holy grail.
8 What is your favorite color? It is blue.

```

Listing 48: Data Structure 7

To traverse a sequence in reverse, first specify the sequence, then call the reversed() function:

```

1 for i in reversed(range(1, 10, 2)):
2     print(i)
3 >>> 9
4     7
5     5
6     3
7     1

```

Listing 49: Data Structure 8

To iterate over a sequence in order, use the sorted() function to return a sorted sequence without modifying the original value:

```

1 basket = ['apple', 'orange', 'apple', 'pear', 'orange', 'banana']
2 for f in sorted(set(basket)):
3     print(f)
4
5 >>> apple
6     banana
7     orange
8     pear

```

Listing 50: Data Structure 9

• Errors and Exceptions

As a Python beginner, when you first learn Python programming, you often see some error messages, which we did not mention before, and we will introduce them in this chapter. There are two types of errors in Python that are easily recognizable: syntax errors and exceptions. Python assert (assertion) is used to judge an expression and trigger an exception when the expression condition is false.

Python grammatical errors or parsing errors are often encountered by beginners, as shown in the following example:

```

1 while True print('Hello world')
2
3 >>> File "<stdin>", line 1, in ?
4     while True print('Hello world')
5         ^
6 SyntaxError: invalid syntax

```

Listing 51: Errors and Exceptions 1

In this example, the function print() is checked for an error because it is missing a colon : in front of it.The parser points out the wrong line and marks a small arrow at the location of the first error found.

Even if the syntax of a Python program is correct, errors may occur when running it. Errors detected at runtime are called exceptions.Most exceptions will not be handled by the program, and are displayed here in the form of error messages:

```

1 10 * (1/0)
2 # 0 cannot be used as a divisor, triggering an exception
3 >>> Traceback (most recent call last):
4   File "<stdin>", line 1, in ?
5 ZeroDivisionError: division by zero
6
7 4 + spam*3
8 # spam is undefined, trigger exception
9 >>> Traceback (most recent call last):
10  File "<stdin>", line 1, in ?
11 NameError: name 'spam' is not defined
12
13 '2' + 2
14 # int cannot be added to str, an exception is triggered
15 >>> Traceback (most recent call last):
16  File "<stdin>", line 1, in <module>
17 TypeError: can only concatenate str (not "int") to str

```

Listing 52: Errors and Exceptions 2

Exceptions occur in different types, which are all printed as part of the message: the types in the example are ZeroDivisionError, NameError and TypeError.

The front part of the error message shows the context where the exception occurred, and the specific information is displayed in the form of a call stack.We can be try to show exception. But I won't discuss too much here, and those who are interested can refer to relevant information.

- **Random**

The Python random module is mainly used to generate random numbers. The random module implements pseudorandom number generators for various distributions. To use the random function you must first import.

Next we use the random() method to return a random number in the half-open interval [0,1), including 0 but excluding 1.

```

1 # import the random package
2 import random
3
4 # Generate random numbers
5 print(random.random())
6

```

```
7 >>> 0.4784904215869241
```

Listing 53: Random 1

The `seed()` method changes the seed of the random number generator, and this function can be called before calling other random module functions.

```
1 import random
2
3 random.seed()
4 print ("Generate a random number using the default seed:", random.random())
5 >>> Generate random numbers with default seed: 0.7908102856355441
6
7 print ("Generate a random number using the default seed:", random.random())
8 >>> Generate random numbers with default seed: 0.81038961519195
9
10 random.seed(10)
11 print("Generate a random number using the integer 10 seed:", random.random())
12 >>> Generate random numbers with default seed: 0.81038961519195
13
14 random.seed(10)
15 print("Generate a random number using the integer 10 seed:", random.random())
16 >>> Generate random numbers with default seed: 0.81038961519195
17
18 random.seed("hello",2)
19 print ("Generate a random number using a string seed: ", random.random())
20 >>> Generate random numbers with default seed: 0.81038961519195
```

Listing 54: Random 2

We have no way to provide all the random variables that random generators want, which obey a specific distribution. But we've provided an explanation of the method below so you can find it when you need it.

Methods	Description
seed()	Initialize the random number generator
getrandbits(k)	Return a non-negative Python integer with k random bits
randrange()	Return a randomly selected element from range(a, b, c)
randint(a, b)	Return a random integer N such that a <= N <= b.
choice(seq)	Returns a element from the non-empty sequence seq.
shuffle(x[, random])	Randomly shuffle the sequence x
sample(N, k)	Returns k-length list of unique elements from the N
random()	Returns the random number in the range [0.0, 1.0).
uniform()	Returns the random number in the range [0.0, 1.0).
betavariate(alpha, beta)	Returns Beta distribution
expovariate(lambd)	Returns Exponent distribution
gammavariate()	Returns Gamma distribution
gauss(u, sigma)	Returns normal distribution(faster than the below one)
lognormvariate(u, sigma)	Returns Lognormal distribution.
normalvariate(u, sigma)	Returns normal distribution.
vonmisesvariate(u, sigma)	Returns von Mises distribution.
paretovariate(u, sigma)	Returns Pareto distribution
weibullvariate(u, sigma)	Returns Weibull distribution.

- **Pip**

We know that there are many libraries available in python. We can use pip to install and manage these libraries. In a word, pip is a Python package management tool that provides functions for finding, downloading, installing, and uninstalling Python packages.

To check if pip is installed you can use the following command on the terminal:

```
1 pip --version
```

Listing 55: Pip Example 1

Download the installation package using the following command: pip install some-package-name For example, we install the numpy package:

```
1 pip install numpy
```

Listing 56: Pip Example 2

We can also easily remove packages with the following command: pip uninstall some-package-name

For example we remove the numpy package:

```
1 pip uninstall numpy
```

Listing 57: Pip Example 3

If we want to see the packages we have installed, we can use the following command:

```
1 pip list
```

Listing 58: Pip Example 4

The above is the content of all python guides. We only list some of the frequently used content. As time goes by, some of them may no longer apply, welcome to update and correct.

Part II.

R guidance

3 BASIC PART

R language is a mathematical programming language designed for mathematical researchers, mainly used for statistical analysis, graphics, and data mining.

Both the R language and the C language are the research results of Bell Labs, but they have different focus areas. The R language is an interpreted language for mathematical researchers, while the C language is designed for computer software engineers.

The R language is an interpreted language (different from the compiled language of the C language), and its execution speed is much slower than that of the C language, which is not conducive to optimization. But it provides a richer data structure operation at the grammatical level and can output text and graphic information very conveniently, so it is widely used in mathematics, especially in the field of statistics.

R language official website: <https://cran.r-project.org/>.

List of official mirror sites: <https://cran.r-project.org/mirrors.html>

If you are slow to download files from r's official website, you can try to download from a mirror website.

- **R installation**

The development environment of R language itself has a graphical development environment, which is different from many other engineering languages, so the development environment is best installed on an operating system designed for desktop personal computers (such as Windows, macOS or Ubuntu desktop version, etc.).

First, we need to download the installation package for the R locale:

Windows:

Official address: <https://cloud.r-project.org/bin/windows/base/>

Linux:

Official address: <https://cloud.r-project.org/bin/linux/>

macOS:

Official address: <https://cloud.r-project.org/bin/macosx/>

- **Grammar**

The "Hello, World!" program code in R language is as follows:

```
1 myString <- "Hello, World!"
2
3 print( myString )
```

Listing 59: R Example

The above example assigns the string "Hello, World!" to the myString variable, and then uses the print() function to output.

Note: R language assignments use the left arrow <- symbol, but some newer versions also support the equals sign =.

- **Note**

Comments are mainly used to analyze a piece of code, which can make it easier for readers to understand. Comments in programming languages will be ignored by the compiler and will not affect the execution of the code.

Comments in general programming languages are divided into single-line comments and multi-line comments, but R language only supports single-line comments, and the comment symbol is #.

In fact, if there are multiple lines of comments, we only need to add the # sign to each line.

```
1 # This is my first programming code
2 myString <- "Hello, World!"
3
4 print( myString )
```

Listing 60: Note Example

- **Operations**

The assignment symbol is a formal advantage and an operational disadvantage of the R language: the form is more suitable for mathematicians. After all, not all mathematicians are accustomed to using = as an assignment symbol.

Operationally speaking, neither the < symbol nor the - symbol are very easy to type characters, which will make many programmers uncomfortable. Therefore, newer versions of the R language also support = as an assignment operator: The following examples demonstrate simple mathematical operations:

```
1 > 1 + 2 * 3
2 [1] 7
```

```

3 > (1 + 2) * 3
4 [1] 9
5 > 3 / 4
6 [1] 0.75
7 > 3.4 - 1.2
8 [1] 2.2
9 > 1 - 4 * 0.5^3
10 [1] 0.5
11 > 8 / 3 %% 2
12 [1] 8
13 > 8 / 4 %% 2
14 [1] Inf
15 > 3 %% 2^2
16 [1] 3
17 > 10 / 3 %/% 2
18 [1] 10

```

Listing 61: Operations Example

- **Data types**

Data types refer to a broad system for declaring variables or functions of different types. The type of a variable determines how much space the variable is stored in, and how the stored bit pattern is interpreted. There are three basic data types in R language: numeric, logical, and text.

There are two main types of numeric constants: scientific notation and general form.

The logical type is often called Boolean in many other programming languages, and the only constant values are TRUE and FALSE. Logical vectors are mainly used for logical operations on vectors, for example:

```

1 > c(11, 12, 13) > 12
2 [1] FALSE FALSE TRUE

```

Listing 62: Data types Example 1

The which function is a very common logical vector processing function, which can be used to filter the subscripts of the data we need:

```

1 > c(11, 12, 13) > 12
2 [1] FALSE FALSE TRUE

```

Listing 63: Data types Example 2

For example, we need to filter data greater than or equal to 60 and less than 70 from a linear table:

```

1 > vector = c(10, 40, 78, 64, 53, 62, 69, 70)
2 > print(vector[which(vector >= 60 & vector < 70)])

```

```
3 [1] 64 62 69
```

Listing 64: Data types Example 3

Similar functions are all and any:

```
1 > all(c(TRUE, TRUE, TRUE))
2 [1] TRUE
3 > all(c(TRUE, TRUE, FALSE))
4 [1] FALSE
5 > any(c(TRUE, FALSE, FALSE))
6 [1] TRUE
7 > any(c(FALSE, FALSE, FALSE))
8 [1] FALSE
```

Listing 65: Data types Example 4

all() is used to check whether a logical vector is all TRUE, and any() is used to check whether a vector contains TRUE.

The most intuitive data type is the text type. Text is a string that often appears in other languages, and constants are enclosed in double quotes. In R, literal constants can be enclosed in single or double quotes, for example:

```
1 > 'UNIGE' == "UNIGE"
2 [1] TRUE
```

Listing 66: Data types Example 5

The string data type itself is not complicated. Here we focus on introducing the operation functions of strings:

```
1 > toupper("unige") # convert to uppercase
2 [1] "UNIGE"
3 > tolower("UNIGE") # convert to lowercase
4 [1] "unige"
5 > substr("123456789", 1, 5) # Intercept the string, from 1 to 5
6 [1] "12345"
7 > substring("1234567890", 5) # Intercept the string, from 5 to the end
8 [1] "567890"
9 > as.numeric("12") # convert string to number
10 [1] 12
11 > as.character(12.34) # convert a number to a string
12 [1] "12.34"
13 > strsplit("2019;10;1", ";") # delimiter to split the string
14 [[1]]
15 [1] "2019" "10" "1"
16 > gsub("/", "-", "2019/10/1") # replace string
17 [1] "2019-10-1"
```

Listing 67: Data types Example 6

If we divide by object type, there are the following 6 types: vector, list, matrix, array, factor, data.frame.

1. Vector

Vectors are an indispensable tool in mathematical operations - our most common vector is a two-dimensional vector, which is bound to be used in a planar coordinate system. From the perspective of data structure, a vector is a linear table, which can be regarded as an array. The existence of vectors as a type in R language makes it easier to manipulate vectors:

```
1 > a = c(3, 4)
2 > b = c(5, 0)
3 > a + b
4 [1] 8 4
```

Listing 68: Data types Example 7

`c()` is a function that creates vectors. Here two 2D vectors are added to get a new 2D vector (8, 4). If a two-dimensional vector and a three-dimensional vector are operated, the mathematical meaning will be lost. Although the operation will not stop, it will be warned. I suggest that you get used to putting an end to this kind of situation.

Each element in the vector can be taken out individually by subscript:

```
1 > a = c(10, 20, 30, 40, 50)
2 > a[2]
3 [1] 20
```

Listing 69: Data types Example 8

Note: The "subscript" in the R language does not represent the offset, but represents the number, that is to say, it starts from 1.

R can also easily take out a part of the vector:

```
1 > a[1:4] # Take items 1 to 4, including items 1 and 4
2 [1] 10 20 30 40
3 > a[c(1, 3, 5)] # get items 1, 3, 5
4 [1] 10 30 50
5 > a[c(-1, -5)] # remove items 1 and 5
6 [1] 20 30 40
```

Listing 70: Data types Example 9

Vectors support scalar calculations:

```

1 > c(1.1, 1.2, 1.3) - 0.5
2 [1] 0.6 0.7 0.8
3 > a = c(1,2)
4 > a^2
5 [1] 1 4

```

Listing 71: Data types Example 10

"Vector" has some commonly used linear table processing functions: vector sorting/vector statistics/vector generation

Example of vector sorting:

```

1 > a = c(1, 3, 5, 2, 4, 6)
2 > sort(a)
3 [1] 1 2 3 4 5 6
4 > rev(a)
5 [1] 6 4 2 5 3 1
6 > order(a)
7 [1] 1 4 2 5 3 6
8 > a[order(a)]
9 [1] 1 2 3 4 5 6

```

Listing 72: Data types Example 11

Note: The order() function returns a vector of subscripts after the vector is sorted.

Examples of vector statistics:

```

1 > sum(1:5)
2 [1] 15
3 > sd(1:5)
4 [1] 1.581139
5 > range(1:5)
6 [1] 1 5

```

Listing 73: Data types Example 12

Vectors can be generated using the c() function, or a continuous sequence using the min:max operator. If you want to generate an arithmetic sequence with gaps, you can use the seq function. Examples of vector generation:

```

1 > seq(1, 9, 2)
2 [1] 1 3 5 7 9
3
4 > seq(0, 1, length.out=3)
5 [1] 0.0 0.5 1.0
6
7 > rep(0, 5)
8 [1] 0 0 0 0 0

```

Listing 74: Data types Example 13

2.Matrix

The R language provides a matrix type for the study of linear algebra. This data structure is very similar to a two-dimensional array in other languages, but R provides language-level matrix operation support. Elements in a matrix can be numbers, symbols, or mathematical expressions. A matrix in R language can be created using the `matrix()` function.

```

1 # byrow is TRUE elements are arranged by row
2 M <- matrix(c(3:14), nrow = 4, byrow = TRUE)
3 print(M)
4
5 # Ebyrow is FALSE elements are arranged in columns
6 N <- matrix(c(3:14), nrow = 4, byrow = FALSE)
7 print(N)
8
9 # define row and column names
10 rownames = c("row1", "row2", "row3", "row4")
11 colnames = c("col1", "col2", "col3")
12
13 P <- matrix(c(3:14), nrow = 4, byrow = TRUE, dimnames = list(rownames, colnames))
14 print(P)
```

Listing 75: Data types Example 14

The output is:

```

1 [,1] [,2] [,3]
2 [1,] 3 4 5
3 [2,] 6 7 8
4 [3,] 9 10 11
5 [4,] 12 13 14
6      [,1] [,2] [,3]
7 [1,] 3 7 11
8 [2,] 4 8 12
9 [3,] 5 9 13
10 [4,] 6 10 14
11      col1 col2 col3
12 row1 3 4 5
13 row2 6 7 8
14 row3 9 10 11
15 row4 12 13 14
```

Listing 76: Output of Data types Example 14

The R language matrix provides the `t()` function, which can exchange the rows and columns of the matrix. For example, if there is a matrix with m rows and n columns, use the `t()` function to convert it into a matrix with n rows and m columns.

```
1 # Create a matrix with 2 rows and 3 columns
```

```

2 M = matrix( c(2,6,5,1,10,4), nrow = 2, ncol = 3, byrow = TRUE)
3 print(M)
4      [,1] [,2] [,3]
5 [1,] 2 6 5
6 [2,] 1 10 4
7 # convert to a matrix with 3 rows and 2 columns
8 print(t(M))

```

Listing 77: Data types Example 15

The output is:

```

1 [,1] [,2] [,3]
2 [1,] 2 6 5
3 [2,] 1 10 4
4 [1] "-----Conversion-----"
5      [,1] [,2]
6 [1,] 2 1
7 [2,] 6 10
8 [3,] 5 4

```

Listing 78: Output of Data types Example 15

If you want to get the matrix elements, you can use the column index and row index of the element, similar to the coordinate form.

```

1 # define row and column names
2 rownames = c("row1", "row2", "row3", "row4")
3 colnames = c("col1", "col2", "col3")
4
5 # create matrix
6 P <- matrix(c(3:14), nrow = 4, byrow = TRUE, dimnames = list(rownames, colnames))
7 print(P)
8 # Get the elements of the first row and third column
9 print(P[1,3])
10
11 # Get the element of the fourth row and second column
12 print(P[4,2])
13
14 # get the second row
15 print(P[2,])
16
17 # Get the third column
18 print(P[,3])

```

Listing 79: Data types Example 16

The output is:

```

1 col1 col2 col3
2 row1 3 4 5
3 row2 6 7 8
4 row3 9 10 11
5 row4 12 13 14
6 [1] 5
7 [1] 13
8 col1 col2 col3
9      6 7 8
10 row1 row2 row3 row4
11      5 8 11 14

```

Listing 80: Output of Data types Example 16

Example of Matrix Calculations

```

1 # Create a matrix with 2 rows and 3 columns
2 matrix1 <- matrix(c(7, 9, -1, 4, 2, 3), nrow = 2)
3 print(matrix1)
4
5 matrix2 <- matrix(c(6, 1, 0, 9, 3, 2), nrow = 2)
6 print(matrix2)
7
8 # add two matrices
9 result <- matrix1 + matrix2
10 cat("Addition result: ", "\n")
11 print(result)
12
13 # Subtract two matrices
14 result <- matrix1 - matrix2
15 cat("Subtraction result: ", "\n")
16 print(result)
17
18 # Multiply two matrices
19 result <- matrix1 * matrix2
20 cat("multiplication result: ", "\n")
21 print(result)
22
23 # Divide two matrices
24 result <- matrix1 / matrix2
25 cat("Division result: ", "\n")
26 print(result)

```

Listing 81: Data types Example 17

The output is:

```

1 [,1] [,2] [,3]
2 [1,] 7 -1 2
3 [2,] 9 4 3
4      [,1] [,2] [,3]
5 [1,] 6 0 3
6 [2,] 1 9 2
7 Addition result:
8      [,1] [,2] [,3]
9 [1,] 13 -1 5
10 [2,] 10 13 5
11 Subtract result:
12      [,1] [,2] [,3]
13 [1,] 1 -1 -1
14 [2,] 8 -5 1
15
16 Multiplication result:
17      [,1] [,2] [,3]
18 [1,] 42 0 6
19 [2,] 9 36 6
20 Divide result:
21      [,1] [,2] [,3]
22 [1,] 1.166667-Inf 0.6666667
23 [2,] 9.000000 0.4444444 1.5000000

```

Listing 82: Output of Data types Example 17

3.List

A list is a collection of objects in the R language that can be used to hold different types of data, which can be numbers, strings, vectors, another list, etc., and of course can also contain matrices and functions. The R language creates lists using the `list()` function.

In the following example, we create a list containing strings, vectors and numbers:

```

1 list_data <- list("facebook", "google", c(11,22,33), 123, 51.23, 119.1)
2 print(list_data)

```

Listing 83: Data types Example 18

The output is:

```

1 [[1]]
2 [1] "facebook"
3
4 [[2]]
5 [1] "google"
6
7 [[3]]
8 [1] 11 22 33
9
10 [[4]]
11 [1] 123
12
13 [[5]]
14 [1] 51.23
15
16 [[6]]
17 [1] 119.1

```

Listing 84: Output of Data types Example 18

We can use the `names()` function to name the elements of the list:

```

1 # list contains vectors, matrices, lists
2 list_data <- list(c("Google","facebook","Tesla"), matrix(c(1,2,3,4,5,6), nrow = 2),
3   list("facebook",12.3))
4
5 # Set names for list elements
6 names(list_data) <- c("Sites", "Numbers", "Lists")
7
8 # show list
9 print(list_data[1])
10
11 # access the third element of the list
12 print(list_data[3])
13
14 # access the first vector element
15 print(list_data$Numbers)

```

Listing 85: Data types Example 19

The output is:

```

1 $Sites
2 [1] "Google", "facebook", "Tesla"
3
4 $Numbers
5      [,1] [,2] [,3]
6 [1,] 1 3 5
7 [2,] 2 4 6
8
9 $Lists
10 $Lists[[1]]
11 [1] "facebook"
12
13 $Lists[[2]]
14 [1] 12.3

```

Listing 86: Output of Data types Example 19

The elements in the list can be accessed using the index. If we use the names() function to name, we can also use the corresponding name to access:

```

1 # list contains vectors, matrices, lists
2 list_data <- list(c("Google", "facebook", "Tesla"), matrix(c(1,2,3,4,5,6), nrow = 2),
3                   list("facebook", 12.3))
4
5 # Set names for list elements
6 names(list_data) <- c("Sites", "Numbers", "Lists")
7
8 # show list
9 print(list_data[1])
10
11 # access the third element of the list
12 print(list_data[3])
13
14 # access the first vector element
15 print(list_data$Numbers)

```

Listing 87: Data types Example 20

The output is:

```

1 $Sites
2 [1] "Google", "facebook", "Tesla"
3
4 $Lists
5 $Lists[[1]]
6 [1] "facebook"
7
8 $Lists[[2]]
9 [1] 12.3
10
11
12 [,1] [,2] [,3]
13 [1,] 1 3 5
14 [2,] 2 4 6

```

Listing 88: Output of Data types Example 20

We can add, delete, and update operations on the list, as shown in the following example:

```

1 # list contains vectors, matrices, lists
2 list_data <- list(c("Google", "facebook", "Tesla"), matrix(c(1,2,3,4,5,6), nrow = 2),
3   list("facebook", 12.3))
4
5 # Set names for list elements
6 names(list_data) <- c("Sites", "Numbers", "Lists")
7
8 # add elements
9 list_data[4] <- "new element"
10 print(list_data[4])
11
12 # delete element
13 list_data[4] <- NULL
14
15 # Output is NULL after deletion
16 print(list_data[4])
17
18 # update element
19 list_data[3] <- "I replaced the third element"
20 print(list_data[3])

```

Listing 89: Data types Example 21

The output is:

```

1 [[1]]
2 [1] "New element"
3
4 $<NA>
5 NULL
6
7 $Lists
8 [1] "I replaced the third element"
```

Listing 90: Output of Data types Example 21

We can combine multiple lists into one list using the `c()` function:

```

1 # Create two lists
2 list1 <- list(1,2,3)
3 list2 <- list("Google","facebook","Tesla")
4
5 # Merge list
6 merged.list <- c(list1,list2)
7
8 # display the merged list
9 print(merged.list)
```

Listing 91: Data types Example 22

The output is:

```

1 [[1]]
2 [1] 1
3
4 [[2]]
5 [1] 2
6
7 [[3]]
8 [1] 3
9
10 [[4]]
11 [1] "Google"
12
13 [[5]]
14 [1] "facebook"
15
16 [[6]]
17 [1] "Tesla"
```

Listing 92: Output of Data types Example 22

To convert a list into a vector, you can use the `unlist()` function to convert the list into a vector, which is convenient for us to perform arithmetic operations::

```

1 # create list
2 list1 <- list(1:5)
3 print(list1)
4
5 list2 <- list(10:14)
6 print(list2)
7
8 # convert to vector
9 v1 <- unlist(list1)
10 v2 <- unlist(list2)
11
12 print(v1)
13 print(v2)
14
15 # add two vectors
16 result <- v1+v2
17 print(result)

```

Listing 93: Data types Example 23

The output is:

```

1 [[1]]
2 [1] 1 2 3 4 5
3
4 [[1]]
5 [1] 10 11 12 13 14
6
7 [1] 1 2 3 4 5
8 [1] 10 11 12 13 14
9 [1] 11 13 15 17 19

```

Listing 94: Output of Data types Example 23

4.Array

Arrays are also objects of the R language, and the R language can create one-dimensional or multi-dimensional arrays. The R language array is a collection of the same type. The matrix matrix we learned earlier is actually a two-dimensional array. The syntax of the array() function is as follows:

```
1 array(data = NA, dim = length(data), dimnames = NULL)
```

Listing 95: Data types Example 24

In the following example, we create a two-dimensional array with 3 rows and 3 columns:

```

1 # Create two vectors of different lengths
2 vector1 <- c(5,9,3)
3 vector2 <- c(10,11,12,13,14,15)
4
5 # create array
6 result <- array(c(vector1,vector2),dim = c(3,3,2))
7 print(result)

```

Listing 96: Data types Example 25

The output is:

```

1 , , 1
2
3 [,1] [,2] [,3]
4 [1,] 5 10 13
5 [2,] 9 11 14
6 [3,] 3 12 15
7
8 , , 2
9
10 [,1] [,2] [,3]
11 [1,] 5 10 13
12 [2,] 9 11 14
13 [3,] 3 12 15

```

Listing 97: Output of Data types Example 25

Use the `dimnames` parameter to set the names of the individual dimensions:

```

1 # Create two vectors of different lengths
2 vector1 <- c(5,9,3)
3 vector2 <- c(10,11,12,13,14,15)
4 column.names <- c("COL1","COL2","COL3")
5 row.names <- c("ROW1","ROW2","ROW3")
6 matrix.names <- c("Matrix1","Matrix2")
7
8 # Create an array and set the names of each dimension
9 result <- array(c(vector1,vector2),dim = c(3,3,2),dimnames = list(row.names,column.names,
10   matrix.names))
10 print(result)

```

Listing 98: Data types Example 26

The output is:

```

1 , , Matrix1
2
3     COL1 COL2 COL3
4 ROW1 5 10 13
5 ROW2 9 11 14
6 ROW3 3 12 15
7
8 , , Matrix2
9
10    COL1 COL2 COL3
11 ROW1 5 10 13
12 ROW2 9 11 14
13 ROW3 3 12 15

```

Listing 99: Output of Data types Example 26

If you want to get an array element, you can use the column index and row index of the element, similar to the coordinate form.

```

1 # Create two vectors of different lengths
2 vector1 <- c(5,9,3)
3 vector2 <- c(10,11,12,13,14,15)
4 column.names <- c("COL1", "COL2", "COL3")
5 row.names <- c("ROW1", "ROW2", "ROW3")
6 matrix.names <- c("Matrix1", "Matrix2")
7
8 # create array
9 result <- array(c(vector1,vector2),dim = c(3,3,2),dimnames = list(row.names, column.names,
   matrix.names))
10
11 # Display the elements of the third row in the second matrix of the array
12 print(result[3,,2])
13
14 # Display the elements in the first row and third column of the first matrix of the array
15 print(result[1,3,1])
16
17 # output the second matrix
18 print(result[,,2])

```

Listing 100: Data types Example 27

The output is:

```

1 COL1 COL2 COL3
2     3 12 15
3 [1] 13
4     COL1 COL2 COL3
5 ROW1 5 10 13
6 ROW2 9 11 14
7 ROW3 3 12 15

```

Listing 101: Output of Data types Example 27

Since arrays are composed of matrices of multiple dimensions, we can access array elements by accessing the elements of the matrix.

```

1 # Create two vectors of different lengths
2 vector1 <- c(5,9,3)
3 vector2 <- c(10,11,12,13,14,15)
4
5 # create array
6 array1 <- array(c(vector1,vector2),dim = c(3,3,2))
7
8 # Create two vectors of different lengths
9 vector3 <- c(9,1,0)
10 vector4 <- c(6,0,11,3,14,1,2,6,9)
11 array2 <- array(c(vector1,vector2),dim = c(3,3,2))
12
13 # create matrix from array
14 matrix1 <- array1[,,2]
15 matrix2 <- array2[,,2]
16
17 # matrix addition
18 result <- matrix1+matrix2
19 print(result)
```

Listing 102: Data types Example 28

The output is:

```

1 [,1] [,2] [,3]
2 [1,] 10 20 26
3 [2,] 18 22 28
4 [3,] 6 24 30
```

Listing 103: Output of Data types Example 28

In addition, we can use the `apply()` element to perform cross-dimensional calculations on array elements

```

1 # Create two vectors of different lengths
2 vector1 <- c(5,9,3)
3 vector2 <- c(10,11,12,13,14,15)
4
5 # create array
6 new.array <- array(c(vector1,vector2),dim = c(3,3,2))
7 print(new.array)
8
9 # Calculate the sum of the numbers in the first row of all matrices in the array
10 result <- apply(new.array, c(1), sum)
11 print(result)
```

Listing 104: Data types Example 29

The output is:

```

1 , , 1
2
3     [,1] [,2] [,3]
4 [1,]    5   10   13
5 [2,]    9   11   14
6 [3,]    3   12   15
7
8 , , 2
9
10    [,1] [,2] [,3]
11 [1,]    5   10   13
12 [2,]    9   11   14
13 [3,]    3   12   15
14
15 [1] 56 68 60

```

Listing 105: Output of Data types Example 29

5.Factor

Factors are used to store different types of data. For example, the gender of a person can be divided into two categories: male and female, and age can be divided into minors and adults. The R language creates factors using the factor() function, with a vector as an input parameter. factor() function syntax format:

```

1 factor(x = character(), levels, labels = levels,
2       exclude = NA, ordered = is.ordered(x), nmax = NA)

```

Listing 106: Data types Example 30

The following example converts a character vector to a factor:

```

1 x <- c("Male", "Female", "Male", "Male", "Female")
2 sex <- factor(x)
3 print(sex)
4 print(is.factor(sex))

```

Listing 107: Data types Example 31

The output is:

```

1 [1] Male Female Male Male Female
2 Levels: Male Female
3 [1] TRUE

```

Listing 108: Output of Data types Example 31

The following example sets the factor level to c('male','female'):

```

1 x <- c("Male", "Female", "Male", "Male", "Female", levels=c('Male','Female'))
2 sex <- factor(x)
3 print(sex)
4 print(is.factor(sex))

```

Listing 109: Data types Example 32

The output is:

```

1 levels1 levels2
2 Male Female Male Male Female Male Female
3 Levels: Male Female
4 [1] TRUE

```

Listing 110: Output of Data types Example 32

Next, we use the labels parameter to add labels to each factor level. The character order of the labels parameter must be consistent with the character order of the levels parameter, for example:

```

1 sex=factor(c('f','m','f','f','m'),levels=c('f','m'),labels=c('female','male'),ordered=TRUE)
2 print(sex)

```

Listing 111: Data types Example 33

The output is:

```

1 [1] female male female female female male
2 Levels: female < male

```

Listing 112: Output of Data types Example 33

6.Data.frame

The R language data frame is created using the `data.frame()` function, and the syntax format is as follows:

```
1 data.frame(..., row.names = NULL, check.rows = FALSE,
2             check.names = TRUE, fix.empty.names = TRUE,
3             stringsAsFactors = default.stringsAsFactors())
```

Listing 113: Data types Example 34

The following creates a simple data frame, including name, job number, monthly salary:

```
1 table = data.frame(
2   name = c("alex", "venti"),
3   Job number = c("001", "002"),
4   Monthly salary = c(1000, 2000)
5 )
6
7 print(table) # view table data
```

Listing 114: Output of Data types Example 34

The output is:

```
1 Name Job Number Monthly Salary
2 1 alex 001 1000
3 2 venti 002 2000
```

Listing 115: Data types Example 35

The data structure of the data frame can be displayed through the `str()` function:

```
1 table = data.frame(
2   name = c("alex", "venti"),
3   Job number = c("001", "002"),
4   Monthly salary = c(1000, 2000)
5 )
6 # Get the data structure
7 str(table)
```

Listing 116: Data types Example 35

The output is:

```
1 'data.frame': 2 obs. of 3 variables:
2   $ Name: chr "alex", "venti"
3   $ Job number: chr "001" "002"
4   $ Monthly salary: num 1000 2000
```

Listing 117: Output of Data types Example 35

`summary()` can display the summary information of the data frame:

```

1 table = data.frame(
2   name = c("alex", "venti"),
3   Job number = c("001", "002"),
4   Monthly salary = c(1000, 2000)
5
6 )
7 # show summary
8 print(summary(table))

```

Listing 118: Data types Example 36

The output is:

```

1 Name Job Number Monthly Salary
2 Length: 2 Length: 2 Min. : 1000
3 Class :character Class :character 1st Qu.:1250
4 Mode :character Mode :character Median :1500
5                                     Mean :1500
6                                     3rd Qu.:1750
7                                     Max. :2000

```

Listing 119: Output of Data types Example 36

We can also extract specified columns:

```

1 table = data.frame(
2   name = c("alex", "venti"),
3   Job number = c("001", "002"),
4   Monthly salary = c(1000, 2000)
5 )
6 # extract the specified column
7 result <- data.frame(table$name, table$monthly salary)
8 print(result)

```

Listing 120: Data types Example 37

The output is:

```

1 table.name table.monthly salary
2 1 alex 1000
3 2 venti 2000

```

Listing 121: Output of Data types Example 37

We can read the data of a certain column of a specified row in a form similar to coordinates. Below we read the data of columns 1 and 2 of rows 2 and 3:

```

1 table = data.frame(
2   Name = c("alex", "venti", "huh"),
3   Job number = c("001", "002", "003"),
4   Monthly Salary = c(1000, 2000, 3000)
5 )
6 # Read the data in columns 1 and 2 of rows 2 and 3:
7 result <- table[c(2,3),c(1,2)]
8 print(result)

```

Listing 122: Data types Example 38

The output is:

```

1 name job number
2 2 venti 002
3 3 huh 003

```

Listing 123: Output of Data types Example 38

We can extend the existing data frame. In the following example, we add department columns:

```

1 table = data.frame(
2   Name = c("alex", "venti", "huh"),
3   Job number = c("001", "002", "003"),
4   Monthly Salary = c(1000, 2000, 3000)
5 )
6 # Add department column
7 table$Department <- c("Operation", "Technology", "Editor")
8
9 print(table)

```

Listing 124: Data types Example 39

The output is:

```

1 Name Job Number Monthly Salary Department
2 1 alex 001 1000 Operations
3 2 venti 002 2000 technology
4 3 huh 003 3000 edit

```

Listing 125: Output of Data types Example 39

We can use the rbind() function to combine multiple vectors into a data frame:

```

1 table = data.frame(
2   Name = c("alex", "venti","huh"),
3   Job number = c("001","002","003"),
4   Monthly Salary = c(1000, 2000, 3000)
5 )
6 newtable = data.frame(
7   name = c("s1", "s2"),
8   Job number = c("101","102"),
9   Monthly salary = c(5000, 7000)
10 )
11 # merge two dataframes
12 result <- rbind(table,newtable)
13 print(result)

```

Listing 126: Data types Example 40

The output is:

```

1 Name Job Number Monthly Salary
2 1 alex 001 1000
3 2 venti 002 2000
4 3 huh 003 3000
5 4 s1 101 5000
6 5 s2 102 7000

```

Listing 127: Output of Data types Example 40

We can use the cbind() function to combine multiple vectors into a data frame:

```

1 # create vector
2 sites <- c("Google","facebook","Tesla")
3 likes <- c(222,111,123)
4 url <- c("www.google.com","www.facebook.com","www.Tesla.com")
5
6 # Combine the vectors into a data frame
7 addresses <- cbind(sites,likes,url)
8
9 # view the data frame
10 print(addresses)

```

Listing 128: Data types Example 41

The output is:

```

1 sites like url
2 [1,] "Google" "222" "www.google.com"
3 [2,] "facebook" "111" "www.facebook.com"
4 [3,] "Tesla" "123" "www.Tesla.com"

```

Listing 129: Output of Data types Example 41

- **If-else Statement**

An if statement consists of a Boolean expression followed by one or more statements. The syntax format is as follows:

```
1 if(boolean_expression) {
2     // statements that will be executed if the boolean expression is true
3 }
```

Listing 130: If-else Statement 1

Execute the code inside if the Boolean expression is true, and do not execute it if it is false.

```
1 x <- 50L
2 if(is.integer(x)) {
3     print("X is an integer")
4 }
5
6 [1] "X is an integer"
```

Listing 131: If-else Statement 2

An if statement can be followed by an optional else statement, which executes if the Boolean expression is false. The syntax format is as follows:

```
1 if(boolean_expression) {
2     // statement that will be executed if the boolean expression is true
3 } else {
4     // statement that will be executed if the boolean expression is false
5 }
```

Listing 132: If-else Statement 3

The code inside the if block is executed if the Boolean expression boolean expression is true. If the Boolean expression is false, the code inside the else block is executed.

```
1 x <- c("Google", "facebook", "Tesla")
2
3 if("facebook" %in% x) {
4     print("contains facebook")
5 } else {
6     print("does not contain facebook")
7 }
8
9 [1] "contains facebook"
```

Listing 133: If-else Statement 4

`summary()` can display the summary information of the data frame:

```

1 x <- c("google", "facebook", "taobao")
2
3 if("weibo" %in% x) {
4     print("The first if contains weibo")
5 } else if ("facebook" %in% x) {
6     print("The second if contains facebook")
7 } else {
8     print("not found")
9 }
10 # show summary
11 print(summary(table))
12
13 [1] "The second if contains facebook"

```

Listing 134: If-else Statement 5

- **Loop Statement**

The loop types provided by R language are: repeat loop, while loop, for loop. The loop control statements provided by R language are: break statement, Next statement.

The repeat loop will execute the code until the conditional statement is true before exiting the loop, and the break statement is used to exit. The syntax format is as follows:

```

1 repeat {
2     // code
3     if(condition) {
4         break
5     }
6 }

```

Listing 135: Loop Statement 1

The following example exits the loop when the variable cnt, which is the count variable, is 5:

```

1 v <- c("Google", "facebook")
2 cnt <- 2
3
4 repeat {
5     print(v)
6     cnt <- cnt+1
7
8     if(cnt > 5) {
9         break
10    }
11 }

```

Listing 136: Loop Statement 2

The output is:

```

1 [1] "Google" "facebook"
2 [1] "Google" "facebook"
3 [1] "Google" "facebook"
4 [1] "Google" "facebook"

```

Listing 137: Loop Statement 3

A while loop statement in R repeatedly executes a target statement as long as a given condition is true. The syntax format is as follows:

```

1 while(condition)
2 {
3     statement(s);
4 }

```

Listing 138: Loop Statement 4

The following example outputs the content in the while statement block when the variable cnt is less than 7, and cnt is the counting variable:

```

1 v <- c("Google", "Runoob")
2 cnt <- 2
3
4 while (cnt < 7) {
5     print(v)
6     cnt = cnt + 1
7 }

```

Listing 139: Loop Statement 5

The output is:

```

1 [1] "Google" "facebook"
2 [1] "Google" "facebook"
3 [1] "Google" "facebook"
4 [1] "Google" "facebook"
5 [1] "Google" "facebook"

```

Listing 140: Loop Statement 6

The for loop statement in the R programming language can repeatedly execute the specified statement, and the number of repetitions can be controlled in the for statement. The syntax format is as follows:

```

1 for (value in vector) {
2     statements
3 }

```

Listing 141: Loop Statement 7

The following example outputs 26 letters for the first four letters:

```

1 v <- LETTERS[1:4]
2 for (i in v) {
3   print(i)
4 }
```

Listing 142: Loop Statement 8

The output is:

```

1 [1] "A"
2 [1] "B"
3 [1] "C"
4 [1] "D"
```

Listing 143: Loop Statement 9

The break statement in R language is inserted in the loop body to exit the current loop or statement and start the script to execute the following statement. If you use nested loops, the break statement will stop the execution of the innermost loop and start the execution of the outer loop statement.

The following example uses break to exit the loop when the variable cnt is 5, where cnt is the count variable:

```

1 v <- c("Google", "facebook")
2 cnt <- 2
3
4 repeat {
5   print(v)
6   cnt <- cnt+1
7
8   if(cnt > 5) {
9     break
10 }
11 }
```

Listing 144: Loop Statement 10

The output is:

```

1 [1] "Google" "facebook"
2 [1] "Google" "facebook"
3 [1] "Google" "facebook"
4 [1] "Google" "facebook"
```

Listing 145: Loop Statement 11

The next statement is used to skip the current loop and start the next loop (similar to Python's continue). The following example outputs the first 6 letters of the 26 letters, and skips the current cycle when the letter is D, and proceeds to the next cycle:

```

1 v <- LETTERS[1:6]
2 for (i in v) {
3
4   if (i == "D") { # D will not output, skip this loop and enter the next one
5     next
6   }
7   print(i)
8 }
```

Listing 146: Loop Statement 12

The output is:

```

1 [1] "A"
2 [1] "B"
3 [1] "C"
4 [1] "E"
5 [1] "F"
```

Listing 147: Loop Statement 13

- **Function**

A function is a group of statements that together perform a task. The R language itself provides many built-in functions, and of course we can also create our own functions. The function definition in R language uses the function keyword, and the general form is as follows:

```

1 function_name <- function(arg_1, arg_2, ...) {
2   // function body
3 }
```

Listing 148: Function 1

The following demonstrates how to define two functions:

```

1 # Define a function for counting a series to the squared value
2 new.function <- function(a) {
3   for(i in 1:a) {
4     b <- i^2
5     print(b)
6   }
7 }
```

Listing 149: Function 2

Next we can call the function:

```

1 new.function <- function(a) {
2   for(i in 1:a) {
```

```

3     b <- i^2
4     print(b)
5   }
6 }
7
8 # Call the function and pass parameters
9 new.function(6)

```

Listing 150: Function 3

The output is:

```

1 [1] 1
2 [1] 4
3 [1] 9
4 [1] 16
5 [1] 25
6 [1] 36

```

Listing 151: Function 4

Now you have mastered the basic R programming ability. Next, let me try some more advanced applications.

4 ADVANCED PART

- **ggplot2**

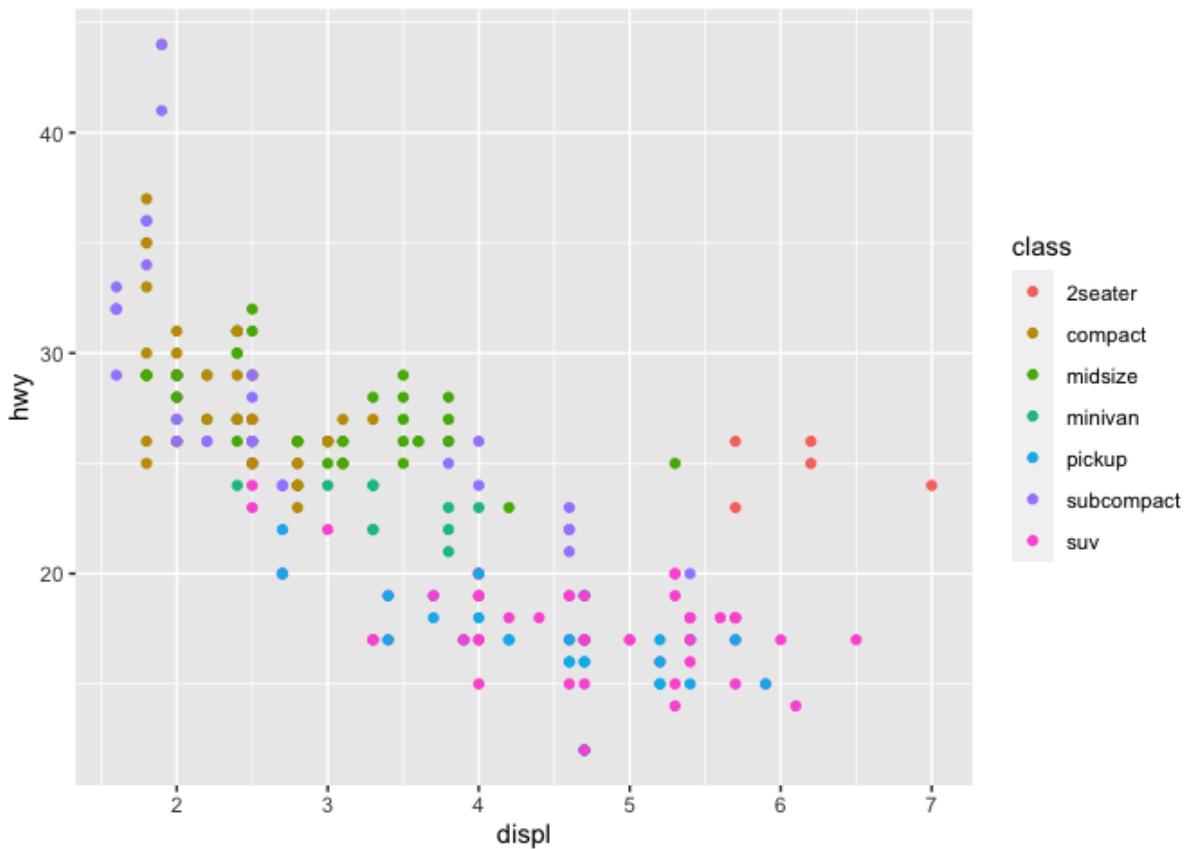
ggplot2 is a system for declaratively creating graphics.

```

1 library(ggplot2)
2
3 ggplot(mpg, aes(displ, hwy, color = class)) +
4   geom_point()

```

Listing 152: ggplot2 example 1



ggplot2 also has more visualization capabilities. Here we do not introduce them one by one. However, we have provided a cheat sheet to help you look it up if you need to.

- **Package**

A package is a collection of R functions, example data, precompiled code including R programs, annotated documentation, examples, test data, etc. R language-related packages are generally stored in the "library" directory under the installation directory. By default, some commonly used packages have been brought with the R language installation. Of course, we can also customize and add some packages to be used later.

We can use the following function to view the installation directory of R packages:

```
1 > .libPaths()
2 [1] "/Library/Frameworks/R.framework/Versions/4.0/Resources/library"
3 >
```

Listing 153: Package 1

We can use the following function to view installed packages:

```
1 library()
```

Listing 154: Package 2

We can use the following function to view the packages loaded by the compilation environment:

```
1 > search()
2 [1] ".GlobalEnv" "package:stats" "package:graphics"
3 [4] "package:grDevices" "package:utils" "package:datasets"
4 [7] "package:methods" "Autoloads" "package:base"
```

Listing 155: Package 3

To install new packages, you can use the `install.packages()` function in the following format:

```
1 install.packages("The name of the package to be installed")
```

Listing 156: Package 4

The above is the content of all R guides. We only list some of the frequently used content. As time goes by, some of them may no longer apply, welcome to update and correct.

Programming and Data Ressources - a first step

We wanted to collect some resources for you to give you some help if you do not know where to start to look and to find some useful programming resources on the web.

First full introduction to Python:

- <https://youtu.be/rfscVS0vtbw> (freeCodeCamp)

- <https://youtu.be/kqtD5dpn9C8> (1 hour course intro in Python)
- <https://youtu.be/8mAITcNt710> (Harvard CS50 – Full Computer Science University Course, includes partly C, Python, HTML, SQL,... generally for the computer science mindset)

However, the most useful thing to do when you are stuck with an error is to look on the internet on websites such as stackoverflow, github, reddit and many more, where probably your question has been asked already (at least the questions when you start programming).

5 OPEN DATA

Also we provide you with some websites for data sets, just to give you a start.

- <https://data.worldbank.org/> (World Bank Data)
- <https://www.census.gov/data.html> (US Census Data, general source)
- <https://www.kaggle.com/datasets> (great variety of data sets)
- <https://archive.ics.uci.edu/ml/index.php> (Machine Learning Data)
- <https://trends.google.com/trends/explore> (General trends, almost every topic)
- <https://healthdata.gov/> (Health Data)
- <https://www.ncei.noaa.gov/products> (large provider of climate and weather information)
- https://www.earthdata.nasa.gov/?_fsi=BqJ6IiI5 (NASA)
- <https://registry.opendata.aws/> (Amazon Web Service, Public Datasets)
- <https://www.aeaweb.org/resources/data> (Economic Data, mainly US)
- <https://github.com/awesomedata/awesome-public-datasets> (all topics)

Furthermore, every country (at least in Europe and North America) has a Federal Statistical Office, where you can get access to many kinds of data and services. Here is the Swiss one:
<https://www.bfs.admin.ch/bfs/en/home/services/ogd/portal.html>

6 EMBEDDED DATA SETS (IN R AND PYTHON)

For R:

Sometimes the data sets are embedded in a package, so you have to install the package first.

```
install.packages("name_of_package") #do not forget the strings

data("name_of_data") #to load the data into your environment

e.g.

data("mtcars") #Loading

head(mtcars, 6) #Print the first 6 rows
```

In the package *datasets*, you find already a big selection of available in-built data sets:
<https://stat.ethz.ch/R-manual/R-devel/library/datasets/html/00Index.html>

Some very popular data sets in R that you will encounter in your studies with probability = 1:

- mtcars Motor Trend Car Road Tests
- Iris Edgar Anderson's Iris Data
- Titanic Survival of passengers on the Titanic
- JohnsonJohnson Quarterly Earnings per Johnson Johnson Share

For Python:

A curated collection of datasets for data analysis machine learning is provided by open-datasets. *OpenDatasets* is a Python library for downloading datasets from online sources like Kaggle and Google Drive using a simple Python command.

<https://pypi.org/project/opendatasets/> (website that gives information on how to install and import data e.g. from Kaggle and some further links for data sets)

Part III.

Mathematics

Linear Algebra

7 OPERATIONS WITH VECTORS AND MATRICES

Given two matrices $\mathbf{A} \in \mathbb{R}^{p \times q}$ and $\mathbf{B} \in \mathbb{R}^{r \times s}$, the following operations can be defined.

- **Sum.** If $p = r$ and $q = s$, $\mathbf{A} + \mathbf{B} = \mathbf{C} \in \mathbb{R}^{p \times q}$, where each element of \mathbf{C} is the sum of the corresponding elements in \mathbf{A} and \mathbf{B} ,

$$c_{ij} = a_{ij} + b_{ij}.$$

- **Transpose.** The transpose of $\mathbf{A} \in \mathbb{R}^{p \times q}$ is the matrix $\mathbf{A}^\top \in \mathbb{R}^{q \times p}$, where the i -th column of \mathbf{A} corresponds to the i -th row of \mathbf{A}^\top , i.e. $(\mathbf{A}^\top)_{ji} = a_{ij}$
- **Scalar Multiplication.** Given a scalar α , the matrix $\alpha\mathbf{A}$ is obtained by multiplying each entry of \mathbf{A} by α , i.e. $(\alpha\mathbf{A})_{ij} = \alpha \cdot a_{ij}$.
- **Matrix Multiplication.** The matrix multiplication \mathbf{AB} is possible only if $q = r$. $\mathbf{AB} = \mathbf{C} \in \mathbb{R}^{p \times s}$, where $c_{ij} = \sum_{k=1}^q a_{ik}b_{kj}$.
- **Power to an integer.** If $p = q$ and $n \in \mathbb{N}$, the n -th power of the matrix \mathbf{A} is $\mathbf{A}^n = \mathbf{A} \times \cdots \times \mathbf{A}$, where the matrix \mathbf{A} is multiplied by itself n times.
- **Kronecker product.** The Kronecker product of two matrices \mathbf{A} and \mathbf{B} is $\mathbf{A} \otimes \mathbf{B} \in \mathbb{R}^{(pr) \times (qs)}$ and is such that

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \dots & a_{1m}\mathbf{B} \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} & \dots & \dots \\ \dots & \dots & \dots & \dots \\ a_{n1}\mathbf{B} & \dots & \dots & a_{nm}\mathbf{B} \end{bmatrix}.$$

- **Inverse.** The inverse of the matrix \mathbf{A} is denoted by \mathbf{A}^{-1} and is such that

$$\mathbf{A}^{-1}\mathbf{A} = \mathbf{AA}^{-1} = \mathbf{I},$$

where \mathbf{I} is the identity matrix. Only if \mathbf{A} is quadratic, the inverse exists, but a square matrix does not necessarily have an inverse.

Note. Note that these operations apply also for vectors, considering that a column vector $v \in \mathbb{R}^p$ corresponds to a matrix of dimension $p \times 1$.

8 DETERMINANT AND TRACE

8.1 Definitions

The **determinant** of a squared matrix $\mathbf{A} \in \mathbb{R}^{p \times p}$ is a function which associates the matrix A to a scalar. It is normally denoted as $\det(\cdot)$.

The **trace** of a matrix $\mathbf{A} \in \mathbb{R}^{p \times p}$ is defined as $\text{Tr}(\mathbf{A}) = \sum_{1 \leq i \leq p} a_{ii}$.

8.2 Computation of the determinant

- If $p = 1$, $\det(\mathbf{A}) = \mathbf{A}$;
- if $p = 2$, $\det(\mathbf{A}) = a_{11}a_{22} - a_{12}a_{21}$;
- if $p > 2$, the Leibniz formula for the determinant is applied.

8.3 Properties

The determinant has the following properties.

- $\det(\mathbf{A}) = \det(\mathbf{A}^\top)$,
- $\det(\mathbf{A}^{-1}) = 1 / \det(\mathbf{A})$ if \mathbf{A} is invertible,
- $\det(\mathbf{AB}) = \det(\mathbf{A}) \det(\mathbf{B})$,
- $\det(\alpha\mathbf{A}) = \alpha \det(\mathbf{A})$,
- $\det(\mathbf{A}) = \prod_{i=1}^p \lambda_i$, where λ_i is the i -th eigenvalue of the matrix \mathbf{A} (see Section 9).

The trace has the following properties.

- $\text{Tr}(\mathbf{A} + \mathbf{B}) = \text{Tr}(\mathbf{A}) + \text{Tr}(\mathbf{B})$,
- $\text{Tr}(\mathbf{AB}) = \text{Tr}(\mathbf{BA})$,
- $\text{Tr}(\mathbf{B}^{-1}\mathbf{AB}) = \text{Tr}(\mathbf{A})$,
- $\text{Tr}(\mathbf{A}) = \sum_{i=1}^p \lambda_i$, where λ_i is the i -th eigenvalue of the matrix \mathbf{A} (see Section 9).

Note. $\text{Tr}(\mathbf{ABC}) = \text{Tr}(\mathbf{CAA})$, but $\text{Tr}(\mathbf{ABC}) \neq \text{Tr}(\mathbf{BAC})$.

9 EIGENVECTORS AND EIGENVALUES

9.1 Definition

Let \mathbf{A} be a $p \times p$ matrix and let $\lambda \in \mathbb{R}$. If λ is such that $\mathbf{v} \in \mathbb{R}^p$, $v \neq \mathbf{0}$ and

$$\mathbf{Av} = \lambda \mathbf{v},$$

then λ is called an eigenvalue of \mathbf{A} and \mathbf{v} is called an eigenvector of \mathbf{A} associated to the eigenvalue λ .

9.2 Quadratic forms and definite matrices

Let $\mathbf{A} \in \mathbb{R}^{p \times p}$ be a symmetric matrix and $\mathbf{x} \in \mathbb{R}^p$, the **quadratic form** of \mathbf{A} is defined as $\mathbf{x}^\top \mathbf{Ax}$.

Let λ_i be the i -th eigenvalue of \mathbf{A} , then

- if $\mathbf{x}^\top \mathbf{Ax} > 0 \forall \mathbf{x} \in \mathbb{R}^p$, \mathbf{A} is called *definite positive*, and $\lambda_i > 0 \forall i = 1, \dots, p$;
- if $\mathbf{x}^\top \mathbf{Ax} \geq 0 \forall \mathbf{x} \in \mathbb{R}^p$, \mathbf{A} is called *semi-definite positive*, and $\lambda_i \geq 0 \forall i = 1, \dots, p$;
- if $\mathbf{x}^\top \mathbf{Ax} < 0 \forall \mathbf{x} \in \mathbb{R}^p$, \mathbf{A} is called *definite negative*, and $\lambda_i < 0 \forall i = 1, \dots, p$;
- if $\mathbf{x}^\top \mathbf{Ax} \leq 0 \forall \mathbf{x} \in \mathbb{R}^p$, \mathbf{A} is called *semi-definite negative*, and $\lambda_i \leq 0 \forall i = 1, \dots, p$.

9.3 Theorem: Spectral decomposition

Theorem 9.1. Let \mathbf{A} be a real symmetric matrix. Then, \mathbf{A} can be expressed as

$$\mathbf{A} = \mathbf{E}\Lambda\mathbf{E}^\top = \sum_{i=1}^p \lambda_i \mathbf{e}_i \mathbf{e}_i^\top,$$

where

- \mathbf{E} is a $p \times p$ matrix and the i -th column of \mathbf{E} , \mathbf{e}_i , is the i -th eigenvector of \mathbf{A} ,
- Λ is a $p \times p$ diagonal matrix and each diagonal element λ_i corresponds to the i -th eigenvalue of \mathbf{A} (with possibly identical eigenvalues).

9.4 Theorem: Singular values decomposition

Theorem 9.2. Let \mathbf{A} be a $n \times p$ matrix, then \mathbf{A} can be expressed as

$$\mathbf{A} = \mathbf{UDV}^\top,$$

where

- \mathbf{U} is a $n \times n$ matrix such that $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_n$;
- \mathbf{D} is a diagonal $n \times p$ matrix and the i -th diagonal element is called the i -th singular value of \mathbf{A} ;
- \mathbf{V} is a $p \times p$ matrix such that $\mathbf{V}^\top \mathbf{V} = \mathbf{I}_p$.

9.5 Differential calculus

In this section, we report the main techniques to compute the derivatives taken with respect to a vector $\mathbf{x} \in \mathbb{R}^p$. The notation \mathbf{x}' indicates the transpose of the vector \mathbf{x} .

Consider the scalar function $y = f(x_1, x_2, \dots, x_p) = f(\mathbf{x})$. The **gradient** ∇_y , or vector of first derivatives, is defined as

$$\nabla_y = \frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_n} \end{bmatrix}$$

The matrix of second derivatives is called the **Hessian matrix**, H_y . It is defined as

$$\begin{aligned} H_y &= \frac{\partial^2 y}{\partial \mathbf{x} \partial \mathbf{x}'} = \frac{\partial (\partial y_i / \partial \mathbf{x})}{\partial \mathbf{x}'} \\ &= \begin{bmatrix} \frac{\partial^2 y}{\partial x_1 \partial x_1} & \cdots & \frac{\partial^2 y}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 y}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 y}{\partial x_n \partial x_n} \end{bmatrix} \end{aligned}$$

Consider a vector $\mathbf{y} \in \mathbb{R}^m$. The derivative of \mathbf{y} with respect to \mathbf{x}' is a matrix of dimension $m \times p$, called the **Jacobian matrix** of \mathbf{y} with respect to \mathbf{x}' .

$$\begin{aligned} \frac{\partial \mathbf{y}}{\partial \mathbf{x}'} &= [\partial y_i / \partial x_j]_{ij} \\ &= \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} = [\partial \mathbf{y}' / \partial \mathbf{x}]' \end{aligned}$$

Let's now focus on linear functions.

Let \mathbf{c} be a vector in \mathbb{R}^p and $z = \mathbf{c}'\mathbf{x}$, then

$$\frac{\partial z}{\partial \mathbf{x}} = \frac{\partial \mathbf{c}'\mathbf{x}}{\partial \mathbf{x}} = \frac{\partial \mathbf{x}'\mathbf{c}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial z}{\partial x_1} \\ \vdots \\ \frac{\partial z}{\partial x_p} \end{bmatrix} = \begin{bmatrix} c_1 \\ \vdots \\ c_p \end{bmatrix} = \mathbf{c}.$$

Let \mathbf{A} be a matrix in $\mathbb{R}^{p \times p}$ and $\mathbf{z} = \mathbf{A}'\mathbf{x}$, then

$$\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \frac{\partial \mathbf{A}'\mathbf{x}}{\partial \mathbf{x}} = (\mathbf{a}_1 \quad \mathbf{a}_2 \quad \cdots \quad \mathbf{a}_p) = \mathbf{A},$$

where \mathbf{a}_j is the j^{th} column of \mathbf{A} .

Let's now focus on the quadratic function $z = \mathbf{x}'\mathbf{A}\mathbf{x}$. Then

$$\begin{aligned} \frac{\partial z}{\partial \mathbf{x}} &= \frac{\partial \mathbf{x}'\mathbf{A}\mathbf{x}}{\partial \mathbf{x}} = \mathbf{A}'\mathbf{x} + \mathbf{A}\mathbf{x} = (\mathbf{A}' + \mathbf{A})\mathbf{x}, \\ \frac{\partial^2 z}{\partial \mathbf{x} \partial \mathbf{x}'} &= \frac{\partial^2 (\mathbf{x}'\mathbf{A}\mathbf{x})}{\partial \mathbf{x} \partial \mathbf{x}'} = \mathbf{A} + \mathbf{A}'. \end{aligned}$$

Note. If A is symmetric, $\frac{\partial \mathbf{x}' A \mathbf{x}}{\partial \mathbf{x}} = 2A\mathbf{x}$ and $\frac{\partial^2 \mathbf{x}' A \mathbf{x}}{\partial \mathbf{x} \partial \mathbf{x}'} = 2A$.

Furthermore,

$$\frac{\partial \mathbf{x}' B \mathbf{y}}{\partial B} = \mathbf{x}\mathbf{y}'$$

$$\frac{\partial \text{tr}(A)}{\partial A} = I$$

$$\frac{\partial |A|}{\partial A} = |A| (A')^{-1}$$

$$\frac{\partial \ln |A|}{\partial A} = (A')^{-1}$$

$$\frac{\partial AB}{\partial \mathbf{x}} = A(\partial B / \partial \mathbf{x}) + (\partial A / \partial x)B$$

$$\frac{\partial A^{-1}}{\partial x} = -A^{-1}(\partial A / \partial x)A^{-1}$$

$$\frac{\partial \text{tr}(A'A)}{\partial A} = \text{tr}(AA') / \partial A = 2A$$

$$\frac{\partial \text{tr}(A'XA)}{\partial A} = (X + X')A, \quad X \text{ symmetric} \Rightarrow \frac{\partial \text{tr}(A'XA)}{\partial A} = 2XA$$

$$\frac{\partial \text{tr}(AXA')}{\partial A} = A(X + X'), \quad A \text{ symmetric} \Rightarrow \frac{\partial \text{tr}(AXA')}{\partial A} = 2AX$$

$$\frac{\partial \text{tr}(XAB)}{\partial A} = X'B', \quad \text{with } A \in \mathbb{R}^{M \times N}, X \in \mathbb{R}^{P \times M} \text{ and } B \in \mathbb{R}^{N \times P}$$

$$\frac{\partial \text{tr}(XA'B)}{\partial A} = BX, \quad \text{with } A \in \mathbb{R}^{M \times N}, X \in \mathbb{R}^{P \times N} \text{ and } B \in \mathbb{R}^{M \times P}$$

$$\frac{\partial \text{tr}(XAA'X')}{\partial A} = 2X'XA$$

$$\frac{\partial \text{vec}(AX)}{\partial A} = (I_T \otimes \text{vec } I_M)X', \quad \text{with } A \in \mathbb{R}^{M \times N} \text{ and } X \in \mathbb{R}^{N \times T}$$

$$\frac{\partial \text{vec}(XA)}{\partial A} = I_N \otimes \text{vec } X', \quad \text{with } A \in \mathbb{R}^{M \times N} \text{ and } X \in \mathbb{R}^{K \times M}.$$

10 SEQUENCES AND SERIES

10.1 Geometric series

Intuition. A geometric sequence is a sequence that is formed by multiplying the previous element by a ratio.

For $r \neq 1$, the *finite geometric series* is defined as:

$$\sum_{j=0}^M r^j = \frac{r^{M+1} - 1}{r - 1};$$

if $|r| < 1$, the *infinite geometric series* is decreasing and converges to

$$\sum_{j=0}^{\infty} r^j = \frac{1}{1 - r}.$$

11 EULER'S NUMBER

Euler's number, better known as the number e , is approximately equal to 2.71828 and can be characterized in different ways, such as

1. Limit. $\lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n = e;$

2. Maclaurin Series. $\sum_{n=0}^{\infty} \frac{1}{n!} = e$;

The **Euler's formula** states that

$$e^{i\theta} = \cos \theta + i \sin \theta.$$

12 INTEGRALS

In this section the most common antiderivatives are recalled, as well as main techniques of integration: integration by parts and integration by substitution.

Intuitively, the definite integral of a scalar function represents the area under its curve between two points.



Definition of integral

If you need a refresher on the definition of Riemann's integral, definite integral and indefinite integral, you can check out [Mathematics I, Section 2](#).

12.1 Antiderivatives of common functions

Table 1 is a sum up of common antiderivatives. Note that $k \in \mathbb{N}$, while $a, b \in \mathbb{R} \setminus \{0\}$.

12.2 Integration by parts

The formula of integration by parts is

$$\int f(x) \cdot g(x) dx = F(x)g(x) - \int F(x)g'(x) dx, \quad (1)$$

where $F(x)$ is the antiderivative of $f(x)$.

Example 12.1. Consider the integral $\int \log(ax) dx$.

To solve this integral, one can apply the integration by parts method and notice that it can be rewritten as

$$\int 1 \cdot \log(ax) dx.$$

Considering the notation the formula 1, in this specific example we have $f(x) = 1$ and $g(x) = \log(ax)$. Hence, $F(x) = x$ and $g'(x) = a \frac{1}{ax} = \frac{1}{x}$.

By following the formula 1, we obtain

$$\begin{aligned} \int 1 \cdot \log(ax) dx &= x \log(ax) - \int x \frac{1}{x} dx \\ &= x \log(ax) - \int 1 dx = \\ &= x \log(ax) - x + c. \end{aligned}$$

Table 1: Most common Antiderivatives

Integral	Antiderivative	Conditions
$\int x^k dx$	$\frac{1}{k+1}x^{k+1} + c$	$k \neq -1$
$\int (f(x))^k f'(x) dx$	$\frac{1}{k+1}(f(x))^{k+1} + c$	$k \neq -1$
$\int \frac{1}{x} dx$	$\ln x + c$	
$\int \frac{1}{f(x)} f'(x) dx$	$\ln f(x) + c$	
$\int e^{ax} dx$	$\frac{1}{a}e^{ax} + c$	
$\int e^{f(x)} f'(x) dx$	$e^{f(x)} + c$	
$\int a^x dx$	$\frac{a^x}{\ln a} + c$	
$\int \cos(ax) dx$	$\frac{1}{a} \sin(ax) + c$	
$\int \sin(ax) dx$	$-\frac{1}{a} \cos(ax) + c$	
$\int \frac{1}{\cos^2(ax)} dx$	$\frac{1}{a} \tan(ax) + c$	
$\int \frac{1}{\sin^2(ax)} dx$	$-\frac{1}{a \tan(ax)} + c$	
$\int \frac{1}{ax^2 + b} dx$	$-\frac{1}{\sqrt{ab}} \arctan\left(\frac{\sqrt{a}}{\sqrt{b}}x\right) + c$	
$\int \frac{1}{\sqrt{b - ax^2}} dx$	$\frac{1}{\sqrt{a}} \arcsin\left(\frac{\sqrt{a}}{\sqrt{b}}x\right) + c$	

12.3 Integration by substitution

The formula of integration by substitution comes from the chain rule, and is the following

$$\int_a^b f(g(x))g'(x)dx = \int_{g(a)}^{g(b)} f(u)du. \quad (2)$$

Example 12.2. Let's solve the integral $\int_a^b \left(\frac{x-l}{s}\right)^2 dx$.

This integral can be solved by substitution and defying

$$u = g(x) = \frac{x-l}{s},$$

from which we can compute

- $du = d\left(\frac{x-l}{s}\right) = dx/s \rightarrow dx = s du,$
- $g(a) = \frac{a-l}{s},$
- $g(b) = \frac{b-l}{s}.$

Substituting these quantities in the original integral leads to

$$\begin{aligned} \int_a^b \left(\frac{x-l}{s}\right)^2 dx &= \int_{\frac{a-l}{s}}^{\frac{b-l}{s}} (u^2)s du = \\ &= s \int_{\frac{a-l}{s}}^{\frac{b-l}{s}} u^2 du = \\ &= s \left[\frac{u^3}{3} \right]_{\frac{a-l}{s}}^{\frac{b-l}{s}} = \\ &= \frac{s}{2} \left[\left(\frac{b-l}{s}\right)^3 - \left(\frac{a-l}{s}\right)^3 \right]. \end{aligned}$$

Help with checking calculus and plotting

To solve specific math problems or to find information on mathematical subjects and topics e.g. arithmetic, algebra, calculus, differential equations, statistics, Wolfram Alpha can be useful:

<https://www.wolframalpha.com/>

There is a lot to play around to visualise and plot equations, solve equations.



For statistics, Wolfram Alpha is able to compute all manner of descriptive and inferential statistical properties and to produce regression analyses and equation fitting. However, we think for learning new concepts and hand-written practice, you should rely on the materials given in class. Additionally, you can use R or Python to check your results or have an additional understanding.

To have a visual representation or intuition, you can also check GeoGebra <https://www.geogebra.org/>.

Specifically helpful for checking derivatives and integral computations:

<https://www.integralrechner.de/>

<https://www.ableitungsrechner.net/>

Part IV.

Probability

In this chapter, the main focus will be presenting a general and basic introduction to probabilities, especially some that will be used throughout the master courses. The goal is to present the main components of probability theory as well as useful properties for the master. However, it will not be as explanatory as a general class on the subject. This is why we advice you to either follow probability 1 or probability 2 given your own probability's knowledge.

13 SET THEORY

Let us start with sets of elements **A**, **B** and **C** belonging to the sample space **S**. It can also be called Omega, written Ω .

13.1 Subset

We say that **A** belongs to **B** if $\forall x \mid x \in A \Rightarrow x \in B$.

It is noted $A \subset B$.

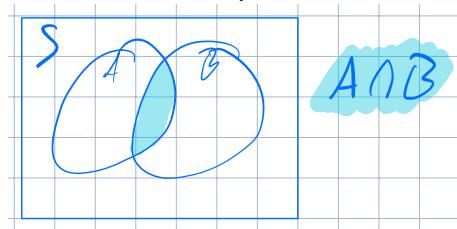
- If $A \subset B$ and $B \subset A$, then $A = B$.
- If $A \subset B$ and $B \subset C$, then $A \subset C$.

13.2 Intersection

We write the intersection between two sets:

$$A \cap B = \{x \mid x \in A \& x \in B\}.$$

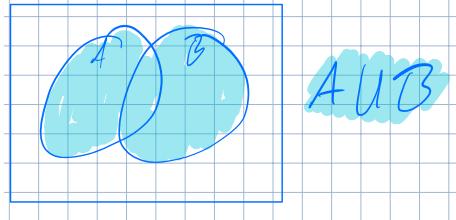
If **A** and **B** are disjoint, then $A \cap B = \emptyset$.



13.3 Union

We write the union between two sets:

$$A \cup B = \{x \mid x \in A \text{ or } x \in B\}.$$



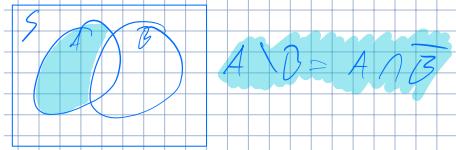
$$A \cup B$$

13.4 Difference

We write the difference between two sets:

$$A \setminus B = \{x \mid x \in A \text{ & } x \notin B\}.$$

Note that $A \setminus B = A \cap \bar{B}$.

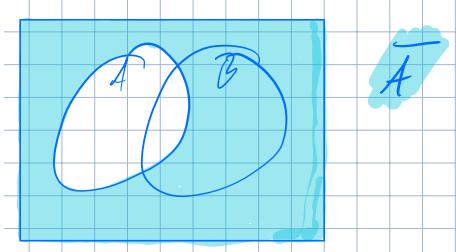


$$A \setminus B = A \cap \bar{B}$$

13.5 Complement

We say that a set \bar{A} is the complement of A . It is defined as follow: $\bar{A} = \{x \mid x \in S \text{ & } x \notin A\}$.

Note that $A \cup \bar{A} = S$.



$$\bar{A}$$

13.6 Properties

- Commutative Laws

$$A \cup B = B \cup A,$$

$$A \cap B = B \cap A.$$

- Associative Laws

$$(A \cup B) \cup C = A \cup (B \cup C),$$

$$(A \cap B) \cup C = A \cap (B \cap C).$$

- Distributive laws

$$A \cup (B \cap C) = (A \cup B) \cap (A \cup C),$$

$$A \cap (B \cup C) = (A \cap B) \cup (A \cap C).$$
- Morgan's Laws

$$\overline{A \cup B} = \bar{A} \cap \bar{B},$$

$$\overline{A \cap B} = \bar{A} \cup \bar{B}.$$



Set Theory

If you need a refresher on this topic, you can check out [Probability I lectures](#), Lecture 1 and Probability Theory: a concise course. [8]

14 COMBINATORICS

Combinatorics is the field of mathematics which is concerned by counting on finite subset. This can be particularly useful when doing basic probability.

14.1 Factorial

The factorial is the product of all natural numbers until a given n . It is noted: $n! = 1 * 2 * \dots * (n - 1) * n = \prod_{i=1}^n i$.

Note that $n! = n * (n - 1)!$ and $0! = 1$.

14.2 Binomial Coefficient

$$\binom{n}{k} = \frac{n!}{k!(n-k)!} , n \& k \in \mathbb{N} \text{ and } 0 \leq k \leq n.$$

- $\binom{n}{0} = \binom{n}{n} = 1,$
- $\binom{n+1}{k+1} = \binom{n}{k} + \binom{n}{k+1},$
- $\binom{n}{n-k} = \binom{n}{k}$

15 ENUMERATION

This section will explain how to count elements on a set by different way.

15.1 simple arrangement

Among n distinct elements, we choose k distinct elements. By classifying them in a particular order, we obtain an arrangement simple:

$$A_k^n = \frac{n!}{(n-k)!}.$$

15.2 Arrangement with repetition

Among n distinct elements, we choose k distinct or not elements. By classifying them in a particular order, we obtain an arrangement avec répétition:

$$\overline{A}_k^n = n^k.$$

15.3 simple permutation

By classifying in a particular order n distinct elements, we obtain a simple permutation:
 $P_n = n!$

15.4 permutation with repetition

By classifying in a particular order n elements of which n_1 are of type 1, n_2 of type 2,..., n_p of type p , we obtain a permutation with repetition:

$$\overline{P}_{(n_1, \dots, n_p)} = \frac{n!}{\prod_{i=1}^p n_i!}.$$

Note that $\sum_{i=1}^p n_i = n$

15.5 Simple combination

Among n distinct elements, we choose k distinct elements without ordering them in any order:

$$C_k^n = \binom{n}{k} = \frac{n!}{k!(n-k)!}.$$

15.6 combination with repetition

Among n distinct elements, we choose k distinct or not elements without ordering them in any order:

$$\overline{C}_k^n = \binom{n+k-1}{k} = \frac{(n+k-1)!}{k!(n-1)!}.$$



Combinatorics and Enumeration

If you need a refresher on this topic, you can check out the Formulaires et Tables CRM [3]. This book is available at the Uni mail Library.

16 SIMPLE PROBABILITIES

Let us start with some basic notion as already defined in the Set Theory section 13. Define Ω the universe; set of all possible outcomes associated with a given random event. Define \mathbf{A} , \mathbf{B} and \mathbf{C} , which are different events belonging to \mathbf{U} . The notation $P(\mathbf{A}) = p$ defines the probability of observing the event \mathbf{A} with associated probability p .

You can find more information in the following book: A First Course in Probability. [7]

- $P(\mathbf{U}) = 1$: The probability of the universe is certain,
- $P(\emptyset) = 0$: The probability of the empty set is null,
- $0 \leq P(\mathbf{A}) \leq 1$: The probability of an event is always between $[0; 1]$,
- $P(\bar{\mathbf{A}}) = 1 - P(\mathbf{A})$,
- If $\mathbf{A} \subset \mathbf{B}$, then $P(\mathbf{A}) \leq P(\mathbf{B})$,
- $P(\mathbf{A} \cup \mathbf{B}) = P(\mathbf{A}) + P(\mathbf{B}) - P(\mathbf{A} \cap \mathbf{B})$: The probability of the union is the sum of the probabilities minus the probability of the intersection,
- If \mathbf{A} and \mathbf{B} are disjoints, then $P(\mathbf{A} \cup \mathbf{B}) = P(\mathbf{A}) + P(\mathbf{B})$,
- $P(\bar{\mathbf{A}} \cup \bar{\mathbf{B}}) = 1 - P(\mathbf{A} \cap \mathbf{B})$,
- $P(\bar{\mathbf{A}} \cap \bar{\mathbf{B}}) = 1 - P(\mathbf{A} \cup \mathbf{B})$,
- $P(\mathbf{A} \cap \bar{\mathbf{B}}) = P(\mathbf{A}) - P(\mathbf{A} \cap \mathbf{B})$.

16.1 Conditional Probabilities

We denote $P(\mathbf{B}|\mathbf{A})$ the probability of the event \mathbf{B} given the event \mathbf{A} , if $P(\mathbf{A}) \neq 0$

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

Furthermore, we can write the probability of the intersection of two events as follow:

$$\Rightarrow P(\mathbf{A} \cap \mathbf{B}) = P(\mathbf{A}) \cdot P(\mathbf{B}|\mathbf{A}) = P(\mathbf{B}) \cdot P(\mathbf{A}|\mathbf{B})$$

16.2 Independent Events

If $P(\mathbf{A} \cap \mathbf{B}) = P(\mathbf{A}) \cdot P(\mathbf{B})$, then \mathbf{A} and \mathbf{B} are called independent events.

Furthermore, conditioning two independent events yields the unconditional probability.

$$\Rightarrow P(\mathbf{B}|\mathbf{A}) = \frac{P(\mathbf{A} \cap \mathbf{B})}{P(\mathbf{A})} = \frac{P(\mathbf{A}) \cdot P(\mathbf{B})}{P(\mathbf{A})} = P(\mathbf{B}).$$

16.3 Theorem of total probabilities

If $\bigcup_{i=1}^n \mathbf{B}_i = \mathbf{U}$, $\forall i, j \quad \mathbf{B}_i \cap \mathbf{B}_j = \emptyset$ and $\forall i \quad P(\mathbf{B}_i) \neq 0$, then

$$P(\mathbf{A}) = P(\mathbf{A}|\mathbf{B}_1) \cdot P(\mathbf{B}_1) + \dots + P(\mathbf{A}|\mathbf{B}_n) \cdot P(\mathbf{B}_n) = \sum_{i=1}^n P(\mathbf{A}|\mathbf{B}_i) \cdot P(\mathbf{B}_i).$$

16.4 Bayes' Theorem

If $\bigcup_{i=1}^n \mathbf{B}_i = \mathbf{U}$, $\forall i, j \quad \mathbf{B}_i \cap \mathbf{B}_j = \emptyset$ and $\forall i \quad P(\mathbf{B}_i) \neq 0$, then

$$P(\mathbf{B}_i|\mathbf{A}) = \frac{P(\mathbf{A}|\mathbf{B}_i) \cdot P(\mathbf{B}_i)}{\sum_{j=1}^n P(\mathbf{A}|\mathbf{B}_j) \cdot P(\mathbf{B}_j)}, \text{ if } P(\mathbf{A}) \neq 0$$



Combinatorics and Enumeration

If you need a refresher on this topic, you can check out [Probability I lectures](#), Lecture 2.

17 RANDOM VARIABLE

To start this chapter, let us define a random variable X as a function that can takes different outcomes/realizations x_1, x_2, \dots with different probabilities associated p_1, p_2, \dots . Such function can be called discrete random variable if the set of different outcomes is finite or countable, or continuous if the set is infinite/uncountable.

Example 1: Discrete Random Variable

For example, if we roll a dice the set of all possible outcomes is as followed: $\{x_1 = 1, x_2 = 2, x_3 = 3, x_4 = 4, x_5 = 5, x_6 = 6\}$. The associated probabilities are all the same: $\{p = \frac{1}{6}\}$. As we have a finite number of outcomes, this function is called discrete random variables.

Example 2: Continuous Random Variable

On the other hand, if we define a function X being the waiting time in minutes before being served at the restaurant. Assuming it won't take the server longer than 75 minutes to serve, the set of all possible outcomes is $]0; 75] = \{x : 0 < x \leq 75\}$. As there is an uncountable number of outcomes for this function, we say that it follows a continuous random variable.

Furthermore, the exact probability of observing one outcome in a continuous random variable

is always 0: $P(X = x) = 0$. We rather compute the probability that the server will arrive before a certain time: $P(X \leq x)$

To formalize these concepts, we say that to characterize a random variable, one needs a cumulative distribution function function (**CDF**) and a probability mass function (**PMF**) for discrete variable or a probability density function (**PDF**) for a continuous one.

17.1 Discrete Random Variable

The PMF is defined as follow:

$$P_a = P(a) = P(X = a).$$

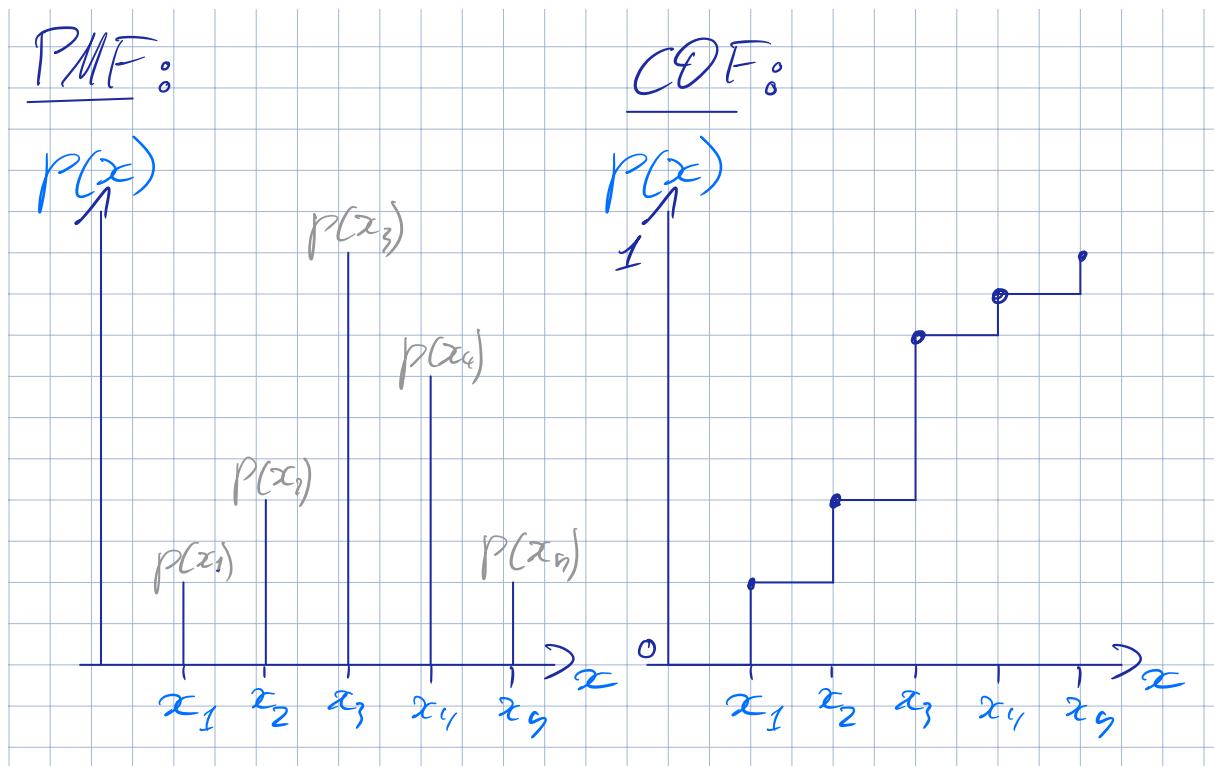
All outcomes a have their according probability $P_a \geq 0$ and as there is a countable number n of outcome, one can sum their probability to obtain the universe.

$$\sum_{i=1}^n P(X_i) = 1 = P(\Omega).$$

The CDF is a step function in the case of discrete random variable. It is defined as:

$$F_x(a) = P(X \leq a) = \sum_{x \leq a} p_x.$$

- $E(X) = \sum_{i=1}^n x_i \cdot p_i$
- $Var(X) = \sum_{i=1}^n p_i \cdot (x_i - E(X))^2$
 $\iff \sum_{i=1}^n p_i \cdot x_i^2 - E(X)^2$

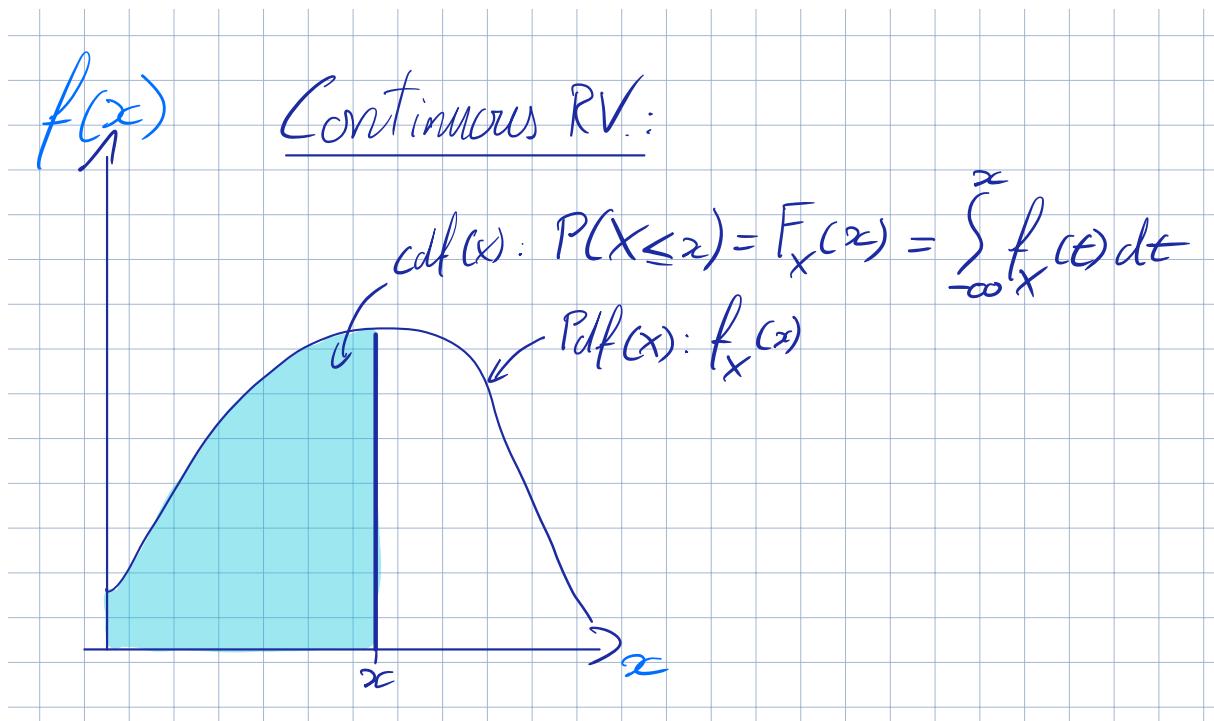


17.2 Continuous Random Variable

To define the PDF, we need that: $f_X(x) \geq 0 \quad \forall x$ and $\int_{-\infty}^{\infty} f_X(x)dx = 1$.

As X can take an uncountable number of outcome, we cannot sum them all to define the CDF. Rather we take the area under the density function, namely: $\int_{-\infty}^{\infty} f_X(x)dx = F_X(x)$ defines the CDF.

- $F_X(x) = \int_{-\infty}^x f_X(t)dt = P(X \leq x)$
- $P(a \leq X \leq b) = \int_a^b f_X(x)dx$
- $E(X) = \int_{-\infty}^{\infty} x \cdot f_X(x)dx$
- $Var(X) = \int_{-\infty}^{\infty} f_X(x) \cdot (x - E(x))^2 dx = \int_{-\infty}^{\infty} x^2 \cdot f_X(x)dx - E(X)^2 = E(X^2) - E(X)^2$.



17.3 Expected Values and Variance properties

The expected value is as its name suggests the anticipated amount a random variable will take on average. It is denoted $E(X)$. The variance is a dispersion value to the mean. Consider X and Y two random variables, a any scalar or real number, then

- $E(X + Y) = E(X) + E(Y)$
- $E(a) = a$
- $E(aX) = aE(x)$
- $E(X + a) = E(X) + a$
- If X and Y are independents, then $E(X \cdot Y) = E(X) \cdot E(Y)$
- $Var(X) = E[(X - E(X))^2] = E(X^2) - E(X)^2$
- $Var(aX) = a^2 Var(X)$
- $Var(a) = 0$
- $Var(X + a) = Var(X)$
- $Var(X + Y) = Var(X) + Var(Y) - 2 \cdot Cov(X, Y)$ note that if X and Y are independent, their covariance equals 0.

- $\sqrt{Var(X)} = \sigma$ which is called the standard deviation
- $Cov(X, Y) = E[(X - E(X)) \cdot (Y - E(Y))]$

18 SOME DISCRETE PROBABILITY DISTRIBUTIONS

Here is a short list of representative distributions that you may see during the MASTAT courses. A more comprehensive list is either available in the Probability 1 slides or in A First course in Probability. [7]

18.1 Binomial Distribution

Denoted $X \sim Bin(n; p)$, with n being the number of trials and p the number of success.

- $f_X(x) = \binom{x}{p} \cdot p^x \cdot (1-p)^{n-x}$
- $E(X) = np$
- $Var(X) = np \cdot (1-p)$

Note that if $n = 1$, $X \sim Bernouilli(p)$

18.2 Poisson Distribution

Denoted $X \sim \mathcal{P}(\lambda)$, with parameter λ .

- $f_X(x) = \frac{\lambda^x e^{-\lambda}}{x!}$
- $F_X(x) = e^{-\lambda} \sum_{i=0}^x \frac{\lambda^i}{i!}$
- $E(X) = \lambda$
- $Var(X) = \lambda$

19 SOME CONTINUOUS PROBABILITY DISTRIBUTIONS

19.1 Uniform Distribution

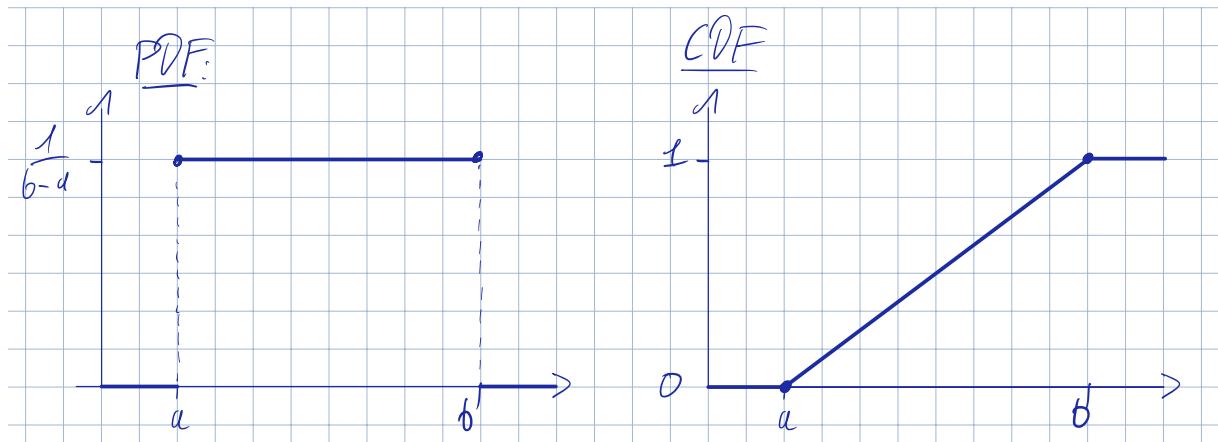
Denoted $X \sim \mathcal{U}(a; b)$, with bounds $[a; b]$

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

- $F_X(x) = \begin{cases} 0 & \text{if } x < a \\ \frac{x-a}{b-a} & \text{if } a \leq x \leq b \\ 1 & \text{if } x > b \end{cases}$

- $E(X) = \frac{a+b}{2}$

- $Var(X) = \frac{(b-a)^2}{12}$



19.2 Exponential Distribution

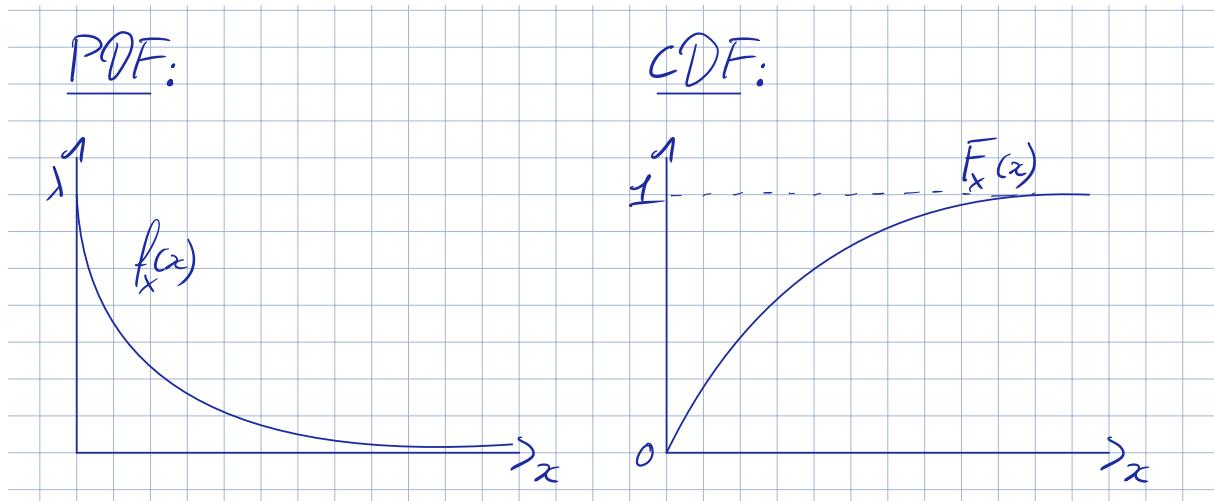
Denoted $X \sim \mathcal{E}(\lambda)$, with parameter $\lambda > 0$

- $f_X(x) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{otherwise} \end{cases}$

- $F_X(x) = \begin{cases} 1 - e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{otherwise} \end{cases}$

- $E(X) = \frac{1}{\lambda}$

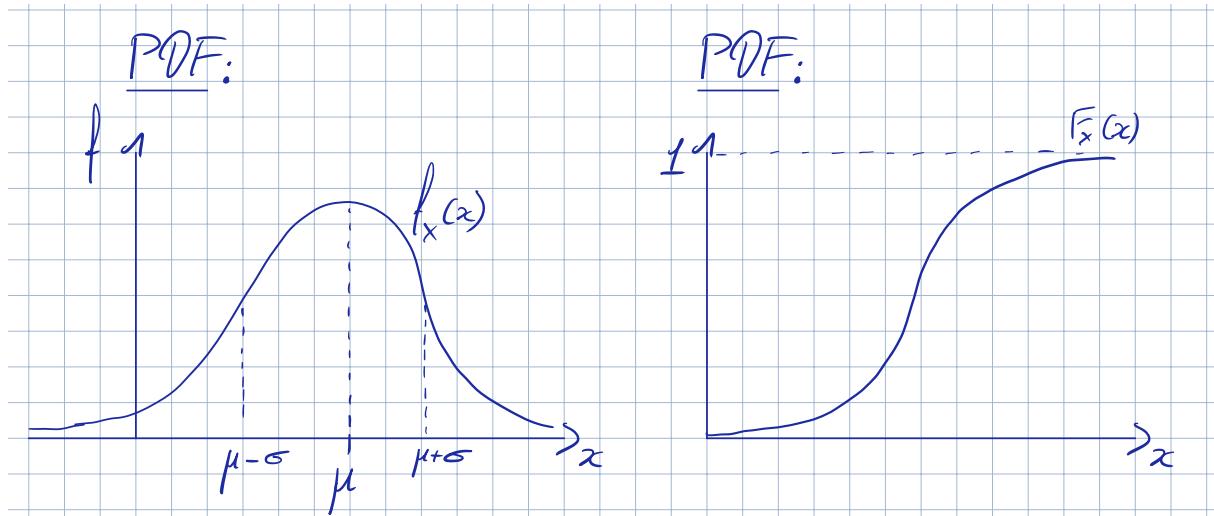
- $Var(X) = \frac{1}{\lambda^2}$



19.3 Normal Distribution

Denoted $X \sim \mathcal{N}(\mu, \sigma^2)$, with parameter μ as its mean and σ^2 as its variance.

- $f_X(x) = \frac{1}{\sqrt{2\sigma^2\pi}} e^{-\frac{1}{2}(\frac{(x-\mu)}{\sigma})^2}$
- $F_X(x) = \int_{-\infty}^x f_X(t)dt$
- $E(X) = \mu$
- $Var(X) = \sigma^2$



Another useful normal distribution is the so-called log-normal distribution, written $X \sim \ln\mathcal{N}(\mu, \sigma^2)$.

$$E(X) = e^{\mu + \frac{\sigma^2}{2}}$$

$$\text{and } Var(X) = (e^{2\mu + \sigma^2} - 1)e^{2\mu + \sigma^2}.$$

20 ASYMPTOTIC THEORY

This section mainly focuses on the basic of the asymptotic theory in probability, meanings what happens when $n \rightarrow \infty$. We will only speak about the two main elements of asymptotic theory, namely convergence in probability and convergence in distribution as you will be seeing other specificities later during the MASTAT.

20.1 *Convergence in Probability*

We say that T_n converges in probability to a constant c if and only if $\lim_{n \rightarrow \infty} P(|T_n - c| > \epsilon) = 0, \forall \epsilon > 0$.

We denote that convergence in probability as $\text{plim}_{n \rightarrow \infty} T_n = c$ or $T_n \xrightarrow{p} c$.

T_n can also converge to a random variable T . This is the case if and only if $\lim_{n \rightarrow \infty} P(|T_n - T| > \epsilon) = 0, \forall \epsilon > 0$.

20.2 *Convergence in Distribution*

Now consider T_n with $F_n(t)$ as its CDF. If $\lim_{n \rightarrow \infty} |F_n(t) - F(t)| = 0$ at every point of continuity $F(t)$, then $T_n \xrightarrow{d} T \sim F(t)$.

Furthermore, $F(t)$ is called the limiting distribution of T_n .

20.3 *Properties of the plim operator*

Let $T_n \xrightarrow{p} c$ and $S_n \xrightarrow{p} d$ be two random sequences. Then:

- $T_n + S_n \xrightarrow{p} c + d$
- $T_n * S_n \xrightarrow{p} c * d$
- $T_n / S_n \xrightarrow{p} c / d$
- $\text{plim } g(T_n) = g(\text{plim } T_n)$ for any $g(\cdot)$ continuous

Let $T_n \xrightarrow{d} T$ and $S_n \xrightarrow{p} c$ be two random sequences. Then:

- $S_n * T_n \xrightarrow{d} c * T$ meaning that the limiting distribution of $S_n * T_n$ is cT
- If $g(\cdot)$ is a continuous function, then $g(T_n) \xrightarrow{d} g(T)$
- If $S_n \xrightarrow{d} S$ and $\text{plim } (T_n - S_n) = 0$, then T_n and S_n have the same limiting distribution



Asymptotic Theory

If you need more information on Asymptotic Theory, we strongly recommend to follow [Econometrics - Masters](#) course as it focuses a lot on that.

21 LIMIT THEOREMS

Here is a brief presentation of the Weak Law of Large Number (WLLN) and the Central Limit Theorems (CLT), which are key theorems in Probabilities and Statistics. You will be using those a lot throughout the master, so it is best to be acquainted with them.

21.1 Weak Law of Large Numbers

The WLLN describes what happens when an experiment is conducted repeatedly n numbers of time. It states that the sample mean of the results will converge to the expected value. Let X_1, \dots, X_n be a collection of i.i.d. random variables from a distribution F with mean μ and $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ their sample average, then:

$$\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| > \epsilon) = 0 \text{ or simply } \bar{X}_n \xrightarrow{P} \mu.$$

21.2 Central Limit Theorem

The CLT states that for i.i.d random variables, their standardized mean converges to a normal distribution even though they were not normally distributed.

Let X_1, \dots, X_n a collection of i.i.d random variable with $E[X] = \mu$ and $Var[X] = \sigma^2 < \infty$. Let $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ be their sample mean, then:

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} \mathcal{N}(0, \sigma^2).$$



Limit Theorems

You will be seeing those theorems in a different classes throughout your courses at the MASTAT such as [The statistical analysis of time series](#).

22 CONCENTRATION IN INEQUALITY

To end this probability chapter, we will present two famous inequalities in probability, namely the Markov's inequality and Chebyshev's inequality.

22.1 *Markov's Inequality*

let X be a non-negative random variable and $a > 0$, then:

$$P(X > a) \leq \frac{E[X]}{a}.$$

22.2 *Markov's Inequality*

Let X be a random variable having a finite non-zero variance σ^2 and finite mean μ . Let $k \in \mathbb{R}_+^*$, then:

$$P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}.$$

Part V.

Statistics

Introduction

As you are all master of statistics students, you must already have some notions of statistics. We would like to point out the following. It is always advisable to stay critical about the results of the functions, tables, data. Before proceeding with the analysis, always ask yourselves the following questions:

- Does the data look realistic? Are there any observations that don't make sense, for example age being -700 for a human?
- Always double check the most basic of the calculations.
- Always make sure the sample size is big enough. In the masters, you will learn about the asymptotic analysis, the asymptotic properties of the estimators as well as small-sample analysis.
- Does the data collection entail an estimation bias? For example, asking exclusively people at the gym how many hours per week they sport won't produce any reliable results for the whole population of your city!

Just for fun, try and uncover misleading statistics in the media!

All of that being said, let's proceed to the quick summary of the material.

22.3 Quick summary

This section regroups the basic notions and introduction into Statistics.

Population is the whole (often unobserved) group of people.

Sample is a smaller, representative group of the population.

Data are the actual measurements of the study.

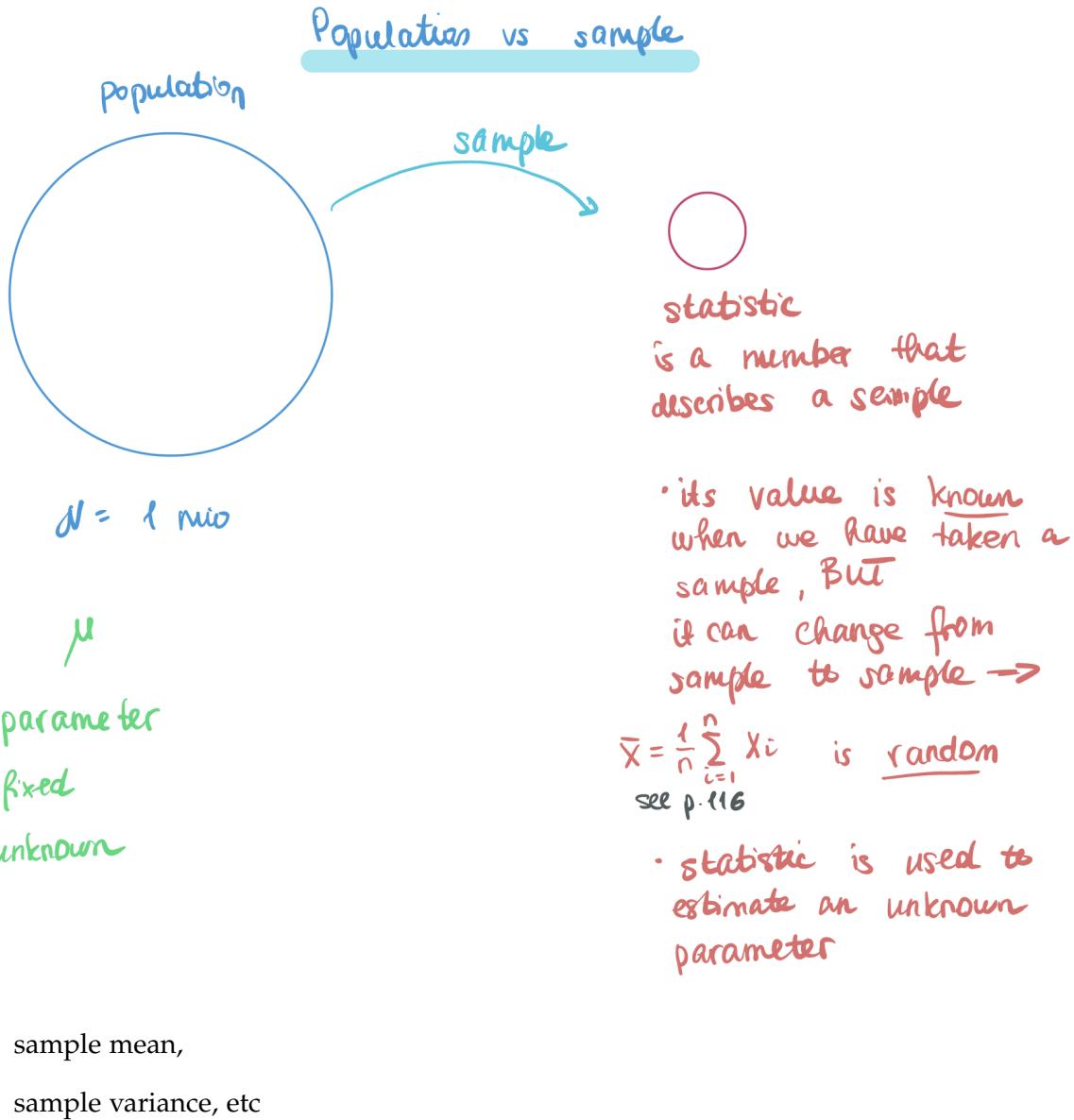
23 RANDOM VARIABLES

Talking about population and sample, we can define the following variables as fixed:

- population mean,
- population size,
- population variance, etc.

On the other hand, when we take a (random) sample from a population, we start dealing with random variables. The statistics summarize a sample, for example:

Figure 1: Population vs sample.



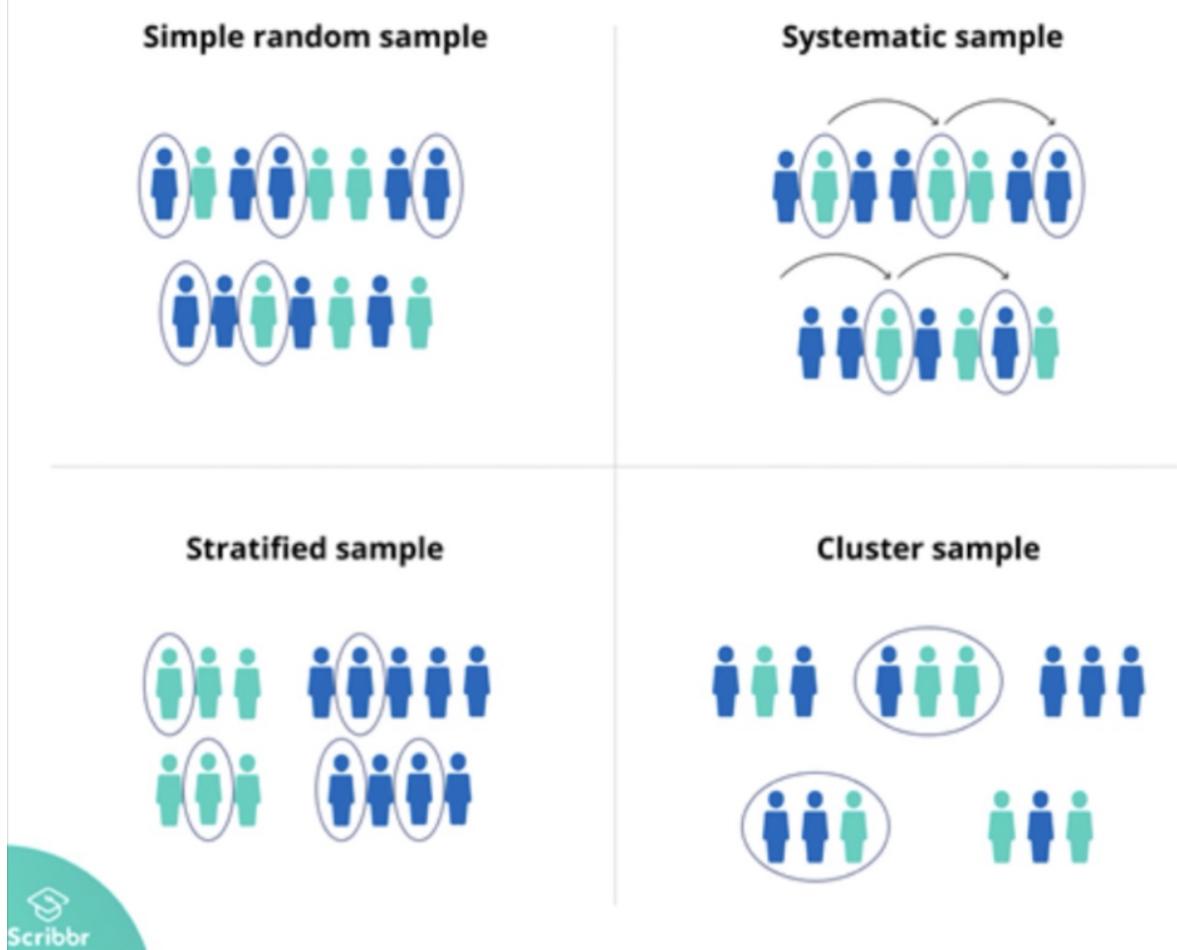
Why are they random? Because their value changes from sample to sample.

23.1 Sampling schemes

How do we sample then? There are different techniques, one of the most basic being a **simple random sampling (SRS)**. Then, every possible sample of a given size has every chance to be chosen.

There is also a **stratified random sampling**, where the idea is to divide the population into *strata*, or groups of individuals similar in some way.

Figure 2: Sampling schemes.



Sampling schemes

If you need a refresher on this topic, you can check out [Statistics lectures](#), Preliminary Work Week 4 and Moore et al. book, pages 401 and on.

23.2 Discrete random variables

A discrete random variable is one that can take values from an ordered list. [5].

Endogenous variables are variables that are changed or determined by its relationship with other variables in the model.

Exogenous variables are the independent variables, sometimes called the outside forces.



Population and sample

If you need a refresher on this topic, you can check out [Statistics lectures](#), Preliminary Work Week 1.

The **mean** is the arithmetic mean of the sum of all the observed values. Consider the sum of all the observations divided by the number of the observations.

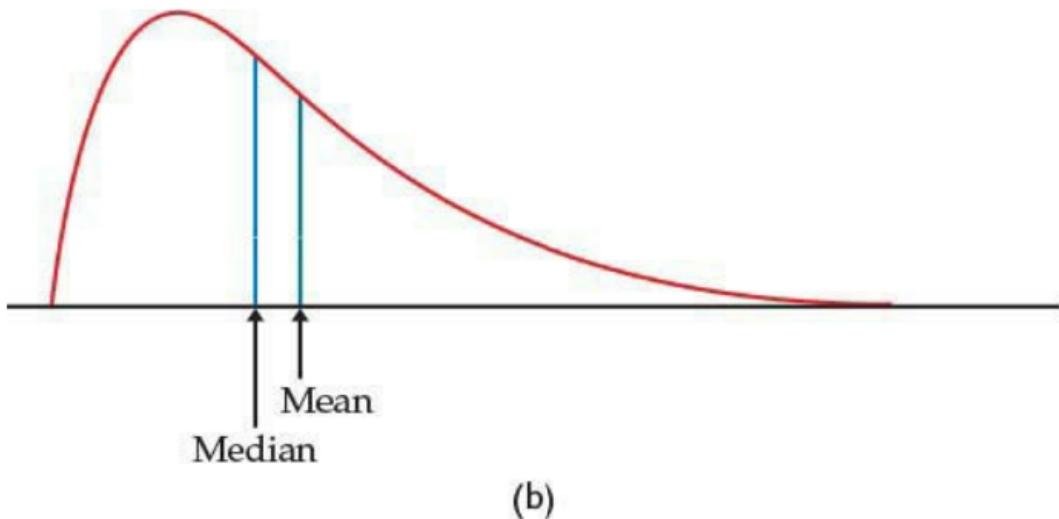
Median is another way to measure the center of the data. It truly represents the "middle" of the dataset. It is obtained by first sorting the data and then finding the middle value by using the rank:

$$r = p \cdot (n - 1) + 1, \quad (3)$$

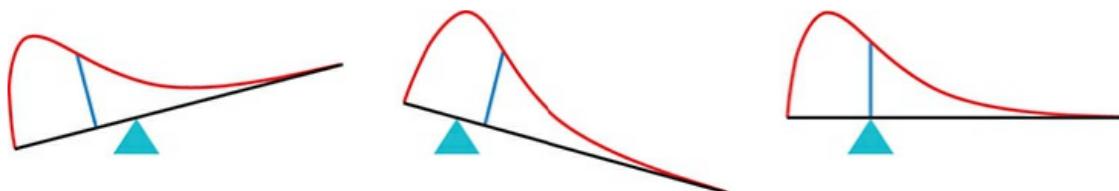
where r is the rank, p the probability, and n the number of observations.

These images give a great intuition behind the mean and median and come from a fairly easy to read book you can turn to if you feel like you need to know more about statistics [5]:

Figure 3: The median and the mean.

**FIGURE 1.21**

(a) A symmetric Normal density curve with its mean and median marked. (b) A right-skewed density curve with its mean and median marked.

**FIGURE 1.22**

The mean of a density curve is the point at which it would balance.



Measures of centrality

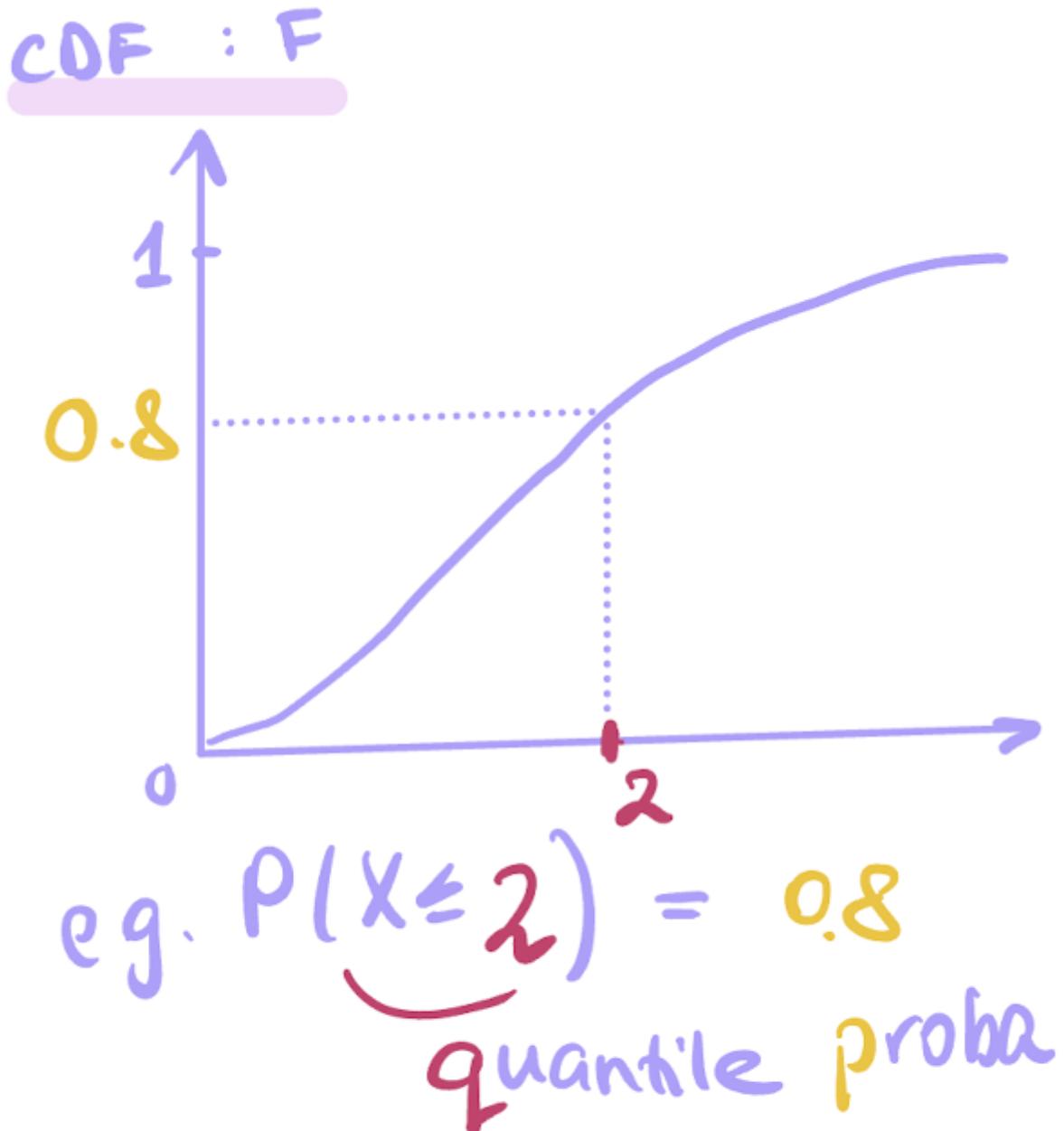
If you feel you need more information about this topic, you can take a look at the [Statistics lectures](#), Preliminary Work Week 2.

Quite a lot of confusion can be caused by the misunderstanding of the quantiles and percentiles. **Quantiles** actually divide the distribution into areas with a certain probability. For example, the **median** always is the 0.5 quantile: per definition, the median divides the probability distribution

into areas of equal probability. With each quantile, there is an associated **probability**, or area under the cumulative distribution function curve.

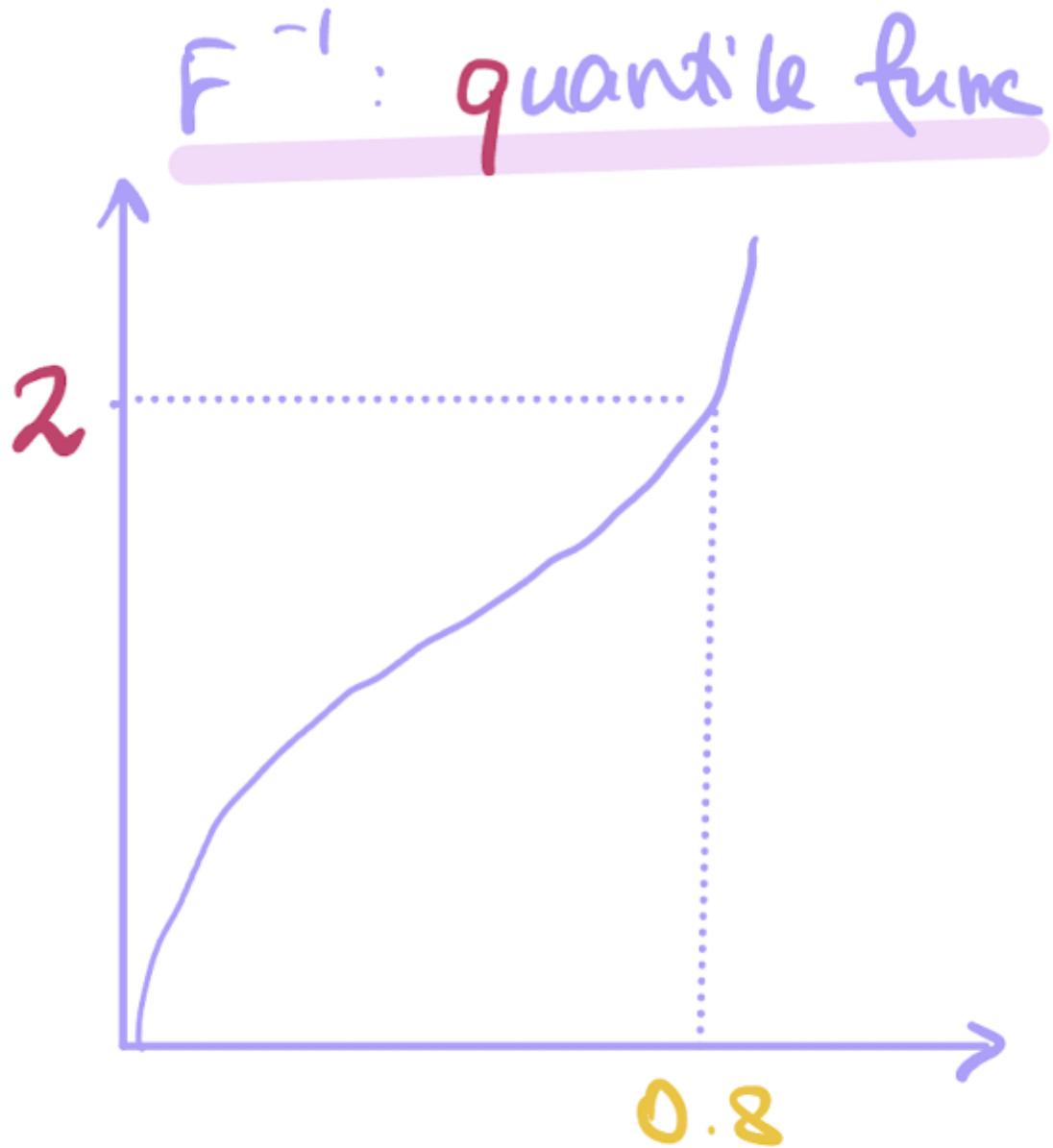
Here is a quick visual summary of the concept for the continuous random variables.

Figure 4: CDF of a continuous random variable.



$$F(q) = 0.8$$

Figure 5: The quantile function.



$$F^{-1}(0.8) = Q(0.8) = 2$$

When working with R, it is helpful to remember which command to use when you want to find out the quantiles, percentiles or simulate from a certain distribution.

Figure 6: Useful functions in R.

In R:	P values of cdf	$F(q) = P(X \leq q)$	uniF
	q quantiles	$F^{-1}(p) = Q(p)$	norm
	d evaluate pdf/pmf	$f(x) \text{ or } P(X=x)$	binom
	r randomly generate from distributions	simulation	pois

+

exp	chisq
t	F



Quantiles

To find out how to compute the quantiles, median, and construct graphs based on those quantities, check out [Statistics lectures](#), Preliminary Work Week 2.

25 GRAPHICAL REPRESENTATION

When working with data, the first thing you always want to do is represent the data graphically. Here is why:

- it lets you simplify data sets,
- shows the basic tendencies,
- helps you spot the outliers,
- and so on.

You can represent data graphically in many ways. The [R Cheatsheet](#) for the ggplot2 package lets you always keep a clear summary of what graph you can use for which kind of data (one, two, or more variables, continuous or discrete data)[23](#). Make sure you at least understand exactly how to build and read the following plots:

- histogram,
- boxplot,
- mosaic plot,

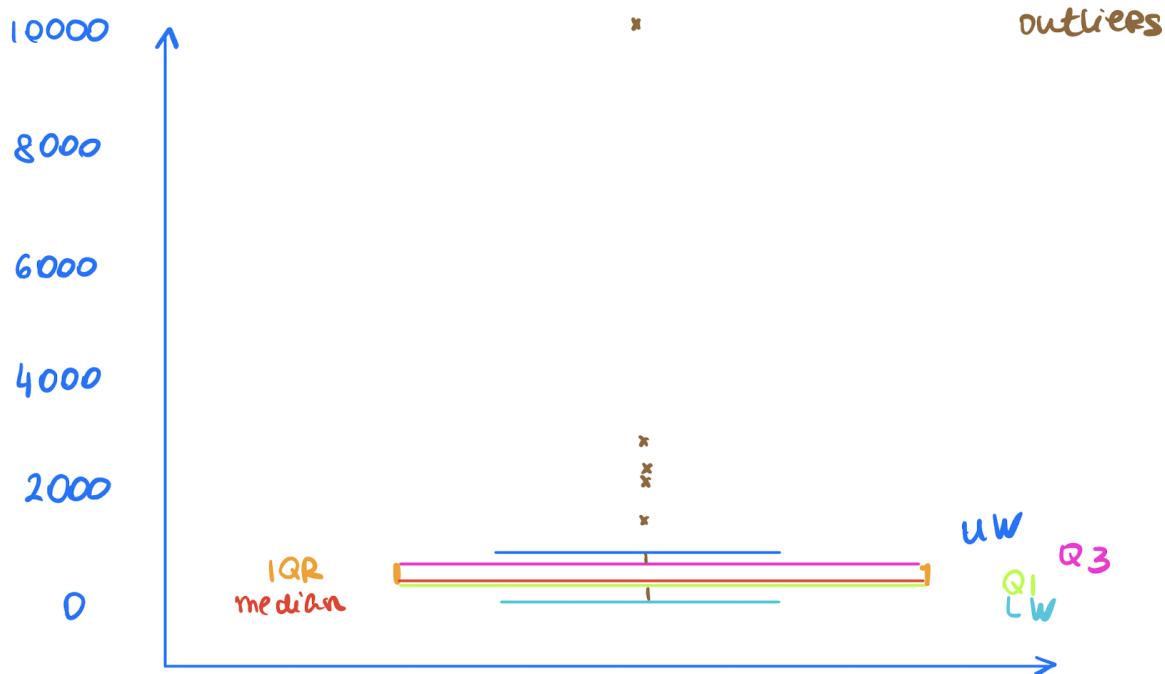
- density plot,
- QQ-plot.

In order to understand the working of R, it is sometimes useful to recreate the plots by hand. For example, Figure 7 represents a boxplot computed by hand using the ranks and the quantiles of the data. The same graph can be created using the `boxplot()` command.

The steps are the following.

1. Sort the data from smallest to largest observation.
2. Compute the Q_1 , Q_2 and Q_3 using formula specified in equation 3, where Q_1 represents the 0.25 quantile, Q_2 0.5 quantile and Q_3 0.75 quantile.
3. Compute the IQR, Interquantile Range: $IQR = Q_3 - Q_1$.
4. Compute the LB, Lower Bound: $LB = Q_1 - 1.5 \times IQR$.
5. Compute the UB, Upper Bound: $UB = Q_3 + 1.5 \times IQR$.
6. To find the UW (Upper Whisker) and LW (Lower Whisker), find the data point closest to the UB and LB, respectively, but still within the bounds.
7. All data outside the LB and UB are the outliers.

Figure 7: Boxplot drawn by hand.





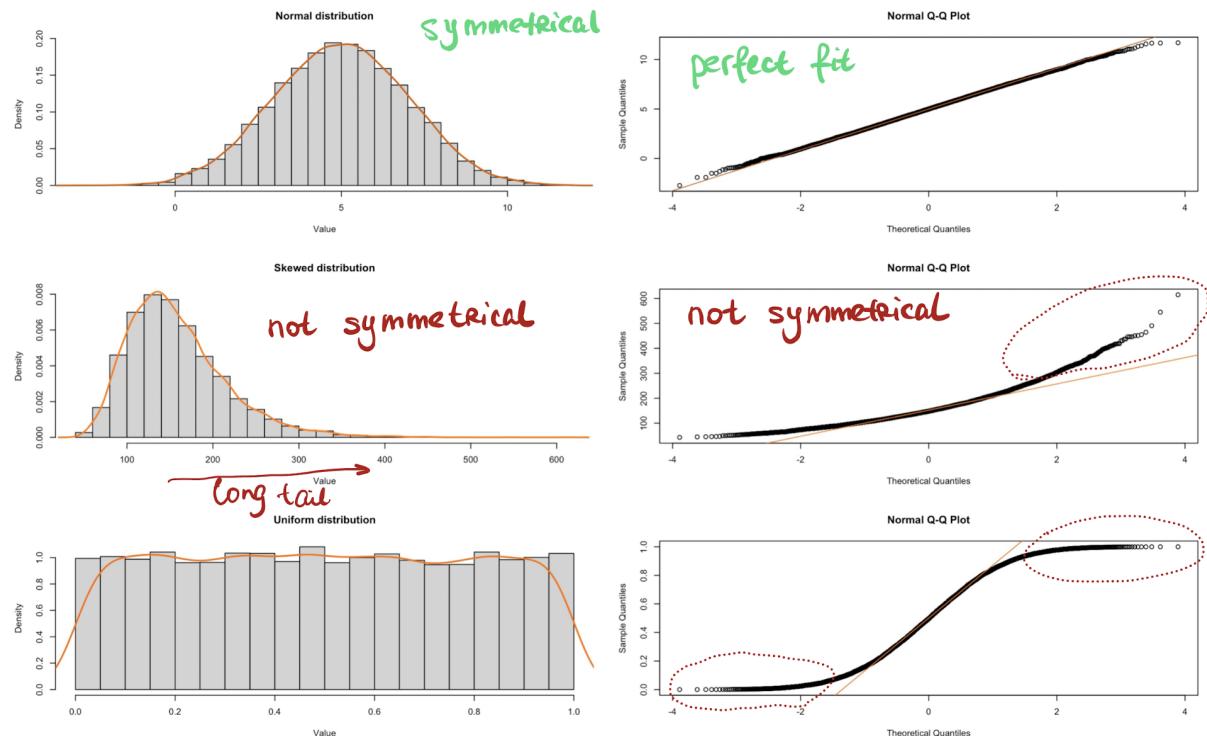
Graphical representation

For more information about the histograms, boxplots and contingency tables see [Statistics lectures](#), Preliminary Work Weeks 2 and 3.

25.1 Q-Q plot

Q-Q plot helps you compare two probability distributions one against another. The following plots compare the simulated distributions to the theoretical quantiles of the normal distribution with the given parameters. The plots will give you an intuition behind the concept.

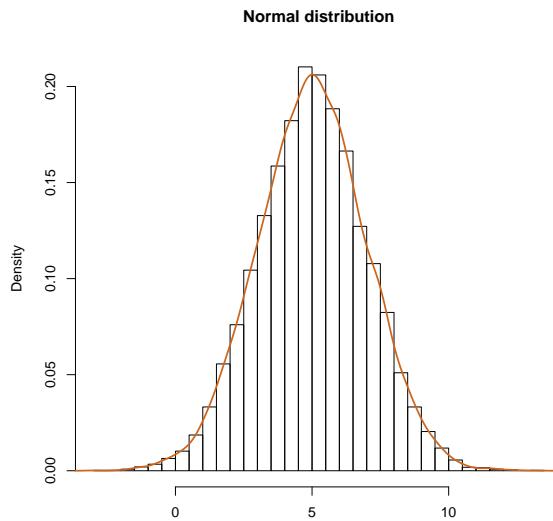
Figure 8: Distribution and QQ-plot.



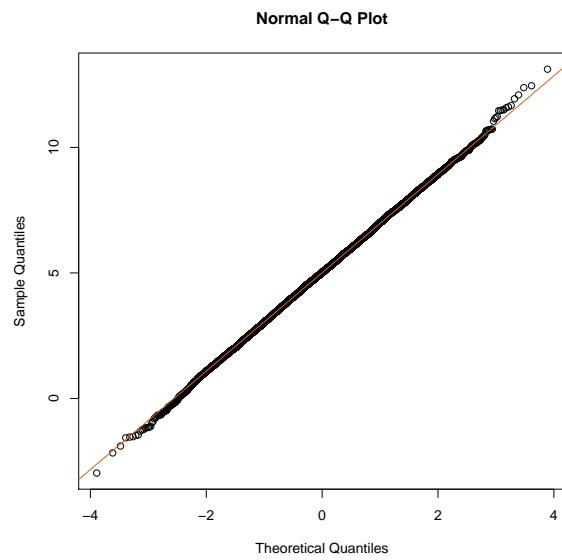
You can recreate the plots using the following R codes. The codes contain the function `density`, which computes the kernel density estimates. You can find more information about it [here](#).

Here we simulate from the normal distribution with mean 5 and standard deviation of 2.

```
# Normal distribution
norm.distr.sim=rnorm(10000, mean=5,
sd=2)
hist(norm.distr.sim,breaks = 30,
probability = T,
#histogram plot
xlab="",main = "Normal distribution")
# add a density plot
lines(density(norm.distr.sim),
lwd = 2, # thickness of line
col = "chocolate3")
```



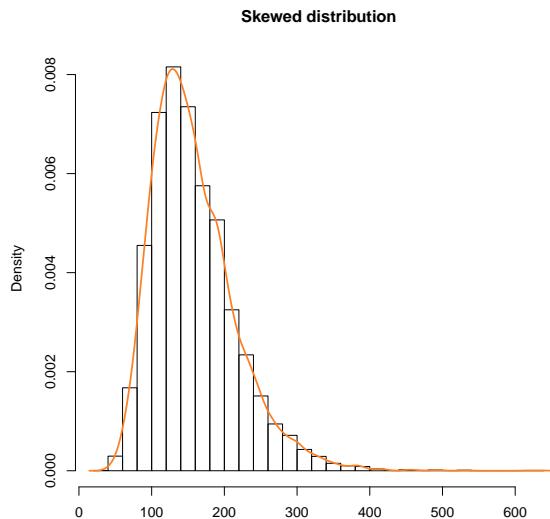
```
# QQ plot of the normal distribution
qqnorm(norm.distr.sim)
qqline(norm.distr.sim, col="chocolate3")
```



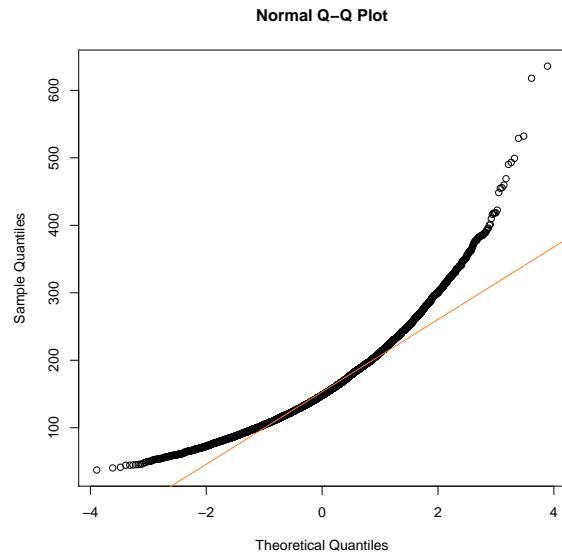
Here we simulate from a skewed (lognormal) distribution.

```
# Simulate from a skewed distribution
skewed.distr <- rlnorm(10000, 5, 0.35)

#histogram plot
hist(skewed.distr, breaks = 30,
probability = T, xlab="", 
main = "Skewed distribution")
# add a density plot
lines(density(skewed.distr),
lwd = 2, # thickness of line
col = "chocolate1")
```



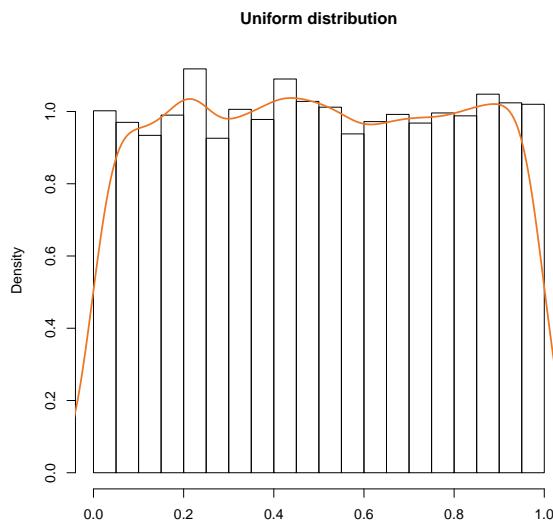
```
# QQ plot of the skewed distribution
qqnorm(skewed.distr)
qqline(skewed.distr, col="chocolate1")
```



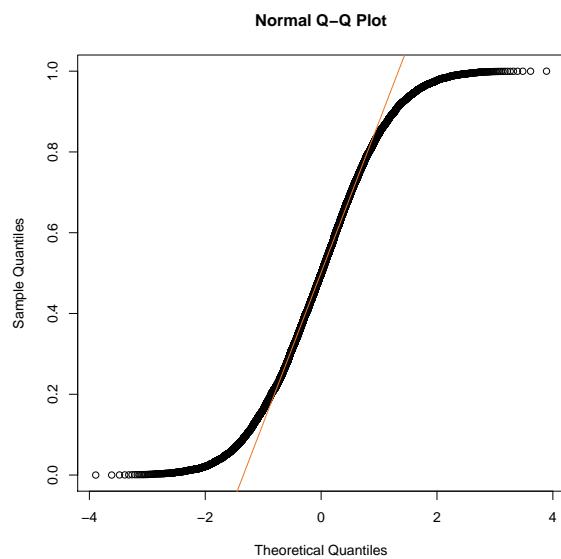
Here we simulate from a uniform distribution.

```
# Simulate from uniform
skewed.distr2 <- runif(10000)

#histogram plot
hist(skewed.distr2, breaks = 30,
probability = T, xlab="", 
main = "Uniform distribution")
# add a density plot
lines(density(skewed.distr2),
lwd = 2, # thickness of line
col = "chocolate2")
```



```
# QQ plot of the skewed distribution
qqnorm(skewed.distr2)
qqline(skewed.distr2, col="chocolate2")
```



Q-Q plot

For more information about Q-Q plot, see [Statistics lectures](#), Preliminary Work Week 3.

26 ESTIMATION

Often, the population parameters are unknown and you'll need to estimate them through sampling.

26.1 Point estimation

What is an estimator, really? Remember that it's a function of the sampled variables, X_1, \dots, X_n and is thus a *random* variable. Estimator is a gentleman and always wears a **hat**, for example, θ is a population parameter, while $\hat{\theta}$ denotes an estimator of the univariate parameter. Also, it is important to distinguish between the random variables X_1, \dots, X_n , because each sample is random and thus the values change from one time to another, while small letters x_1, \dots, x_n denote the realizations and are thus known values. This is very important to distinguish. Here is an example in R, where we can simulate from a standard normal distribution many times - here, 3, just for the purpose of the illustration, and each time the values are different, while, as long as we store the simulation in a variable, it becomes a certain *realisation* of a random variable.

In the left-side section three observations of the random variable X are observed three times. In the section on the left, the recorded observations are represented three times. The values do not change as these observations x_1, x_2 and x_3 are not random anymore.

```
#simulation from standard normal
# X1,X2,X3 are all different:
rnorm(3, mean=0, sd=1)

## [1]  0.8890266  2.4347347 -0.7719102

rnorm(3, mean=0, sd=1)

## [1]  0.4152474 -1.1185879 -1.5135935

rnorm(3, mean=0, sd=1)

## [1]  0.1491298  1.6576958 -0.3087861
```

	# realisations x1,x2,x3:
x=rnorm(3, mean=0, sd=1)	## [1] -1.03292099 0.01447077 0.72651270
x	## [1] -1.03292099 0.01447077 0.72651270
x	## [1] -1.03292099 0.01447077 0.72651270
x	## [1] -1.03292099 0.01447077 0.72651270

The two estimators you'll see most often are the mean and the variance. Here we give the formulas without any further explanation.

Mean

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

A mean has a very natural intuitive meaning: it's an arithmetic mean of the sum of all the observations divided by the number of the observations. By noting \bar{X} the sample mean and

μ the population mean and under some assumptions, we have that asymptotically (in other words, when the number of observations tends to infinity), the sample mean will converge to the population mean. Let $\bar{X} = \sum_{i=1}^n X_i$. We then have that $E(\bar{X}) = \mu$ and $\text{Var}(\bar{X}) = \sigma^2/n$. Applying Chebyshev's inequality to \bar{X} , we get

$$\Pr\left(\frac{|\bar{X} - \mu|}{\sigma/\sqrt{n}} \geq x\right) \leq \frac{1}{x^2},$$

so

$$\Pr\left(|\bar{X} - \mu| \geq \frac{x\sigma}{\sqrt{n}}\right) \leq \frac{1}{x^2}.$$

This holds if X is any positive number. Hence it holds when $x = \varepsilon\sqrt{n}/\sigma$, where ε is any positive number. That gives then

$$\Pr(|\bar{X} - \mu| \geq \varepsilon) \leq \frac{\sigma^2}{\varepsilon^2 n} \rightarrow 0 \text{ as } n \rightarrow \infty,$$

or, in other words

$$\bar{X} \xrightarrow{p} \mu$$

, or \bar{X} converges in probability to μ .

Variance

The variance is a measure of variability, the squared differences from the mean. It tells us the degree at which the data is spread. We can compute it by

$$\text{Var} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

26.2 Confidence intervals

When you compute an estimator, an important part is to know by how much that estimator is expected to vary. Confidence intervals give important information about the accuracy of an estimate.

To put it into a context, imagine you draw a sample of 100 observations 100 times. Each time you draw a sample of 100, you compute the mean. Then, you plot these means using a histogram. Then, a 95% confidence interval will cover 95 of these means.

How to compute the confidence intervals?

You have 2 situations, namely, when σ is known and when σ is unknown. This distinguish will entail the important theoretical results and properties of the estimators. Putting it grossly, a confidence interval can be seen as

$$\text{estimate} \pm \text{margin of error}.$$

26.2.1 Confidence interval for a population mean, σ known

$$\left[\bar{X} \pm z_{(1-\frac{\alpha}{2})} \frac{\sigma}{\sqrt{n}} \right], \quad (4)$$

where $z = (x - \mu)/\sigma$ and $z_{(1-\frac{\alpha}{2})}$ is the $1 - \frac{\alpha}{2}$ quantile of the standard normal distribution. In order to find out the value, you can either use a z-table or refer to statistical programs, such as R and remember the recap table in section 6.

26.2.2 Confidence interval for a population mean, σ unknown

When σ is unknown, you'll have to estimate it and thus will be dealing with more uncertainty, which will result in a wider interval, given all other parameters are constant.

$$\left[\bar{X} \pm t_{n-1,1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}} \right], \quad (5)$$

where $t_{n-1,1-\frac{\alpha}{2}}$ is the $1 - \frac{\alpha}{2}$ quantile of the t-distribution with $n - 1$ degrees of freedom.



Estimation and confidence intervals

For more information about point estimation, confidence intervals, formulas for confidence intervals for proportion, σ^2 , difference of means and the theoretical properties of the estimators, see [Statistics lectures](#), Preliminary Work Weeks 5 and 6, as well as Moore et al. book, pages 356 and on.

27 COMPARING ESTIMATORS

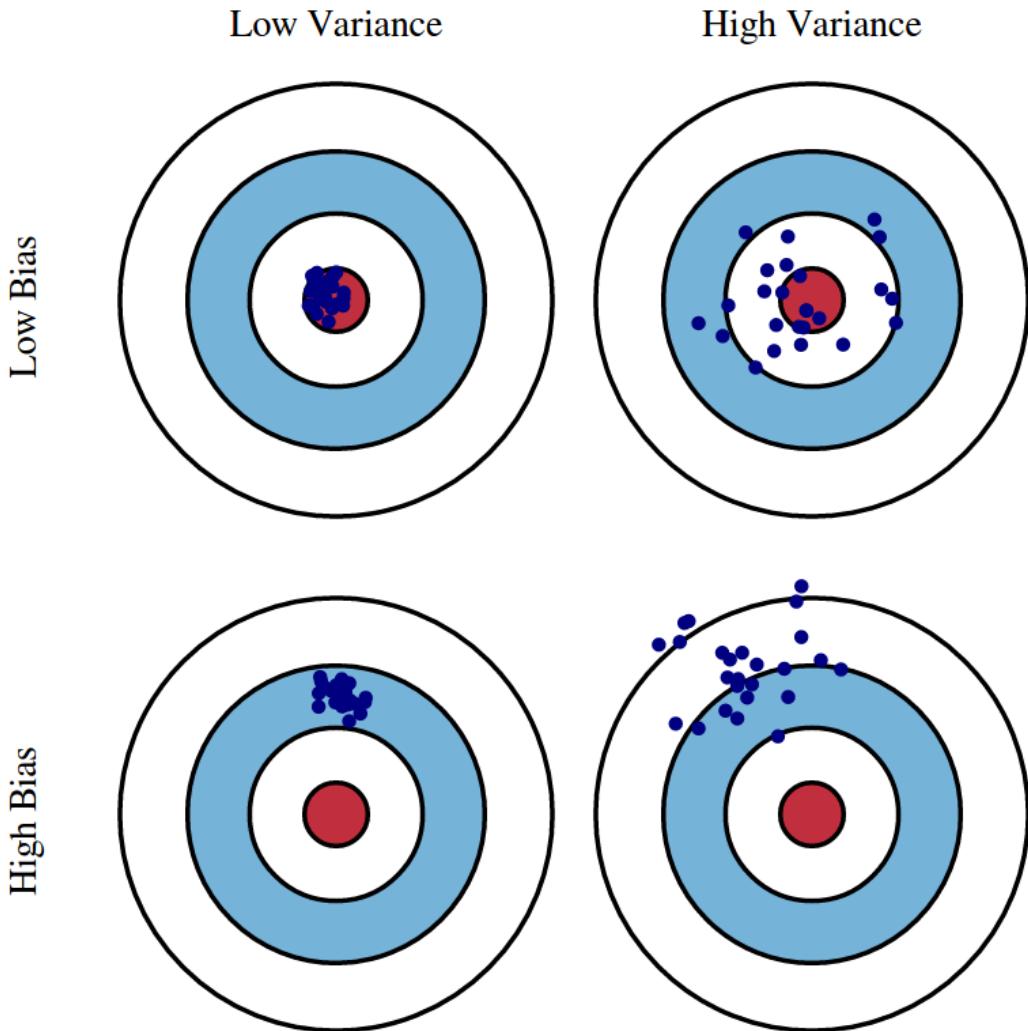
When comparing the estimators, two important notions are those of **accuracy** and **precision**. To get an intuition behind the concept, the following graphical representation might help: Figure 9.

Accuracy can be defined as the closeness or nearness of the measurements to the true or actual value of the quantity being measured.[\[1\]](#) The term accuracy is also used to as a ratio between the correct predictions and the total number of predictions in the classification problems.

The variance of the observations can be used as a measure of the accuracy. You can think about it in the following manner. If an estimator has a very big variance, the chances the true value is someone in there get higher, as the variance get bigger. Thus, the estimation becomes less and less accurate.

Precision is the term mostly associated with bias. Although, some sources say that bias is a systematic error whereas precision is a random error [\[1\]](#).

Figure 9: Bias and variance.



Let us go through both of them step by step.

27.1 Bias

When you draw a sample and compute an estimator, it is important to know how close to the true value that estimator is. When thinking of bias, you should directly think of the **expectation**. An **unbiased** estimator is the one that is equal to the true parameter.

When thinking about bias, you should immediately think of expectation. Bias is defined as follows:

$$\text{bias} (T_n, \theta) = E (T_n) - \theta, \quad (6)$$

where $T_n := T(X_1, \dots, X_n)$ is an estimator of X_1, \dots, X_n and θ the true parameter.

If $\text{bias}(T_n, \theta) = 0$, we say that the estimator is unbiased. $\text{Bias}(T_n, \theta) < 0$ indicates underestimation of θ . $\text{Bias}(T_n, \theta) > 0$ indicates overestimation of θ .

27.2 Variance

The formal definition of the estimator's variance is the following.

$$\text{Var}(T_n) = E((T_n - E(T_n))^2) \quad (7)$$

In other words, variance expresses variability around the expectation of the estimator T_n .

27.3 MSE: Mean Squared Error

It goes without saying the the best estimator is the one that has both high accuracy and precision in figure 9. Let's formalize this.

Mean squared error is defined as the sum of bias squared plus variance.

$$\text{MSE}(T_n, \theta) = \text{Var}(T_n) + \text{bias}((T_n, \theta))^2, \quad (8)$$

27.4 Efficiency

When comparing two estimators, one can compare the MSE's for both estimators. The following is the definition of efficiency.

The efficiency of an estimator $T := T(X_1, \dots, X_n)$ of θ with respect to an alternative estimator $U := U(X_1, \dots, X_n)$ of θ is defined by:

$$\text{eff}(T, U) = \frac{\text{MSE}(U, \theta)}{\text{MSE}(T, \theta)} \quad (9)$$

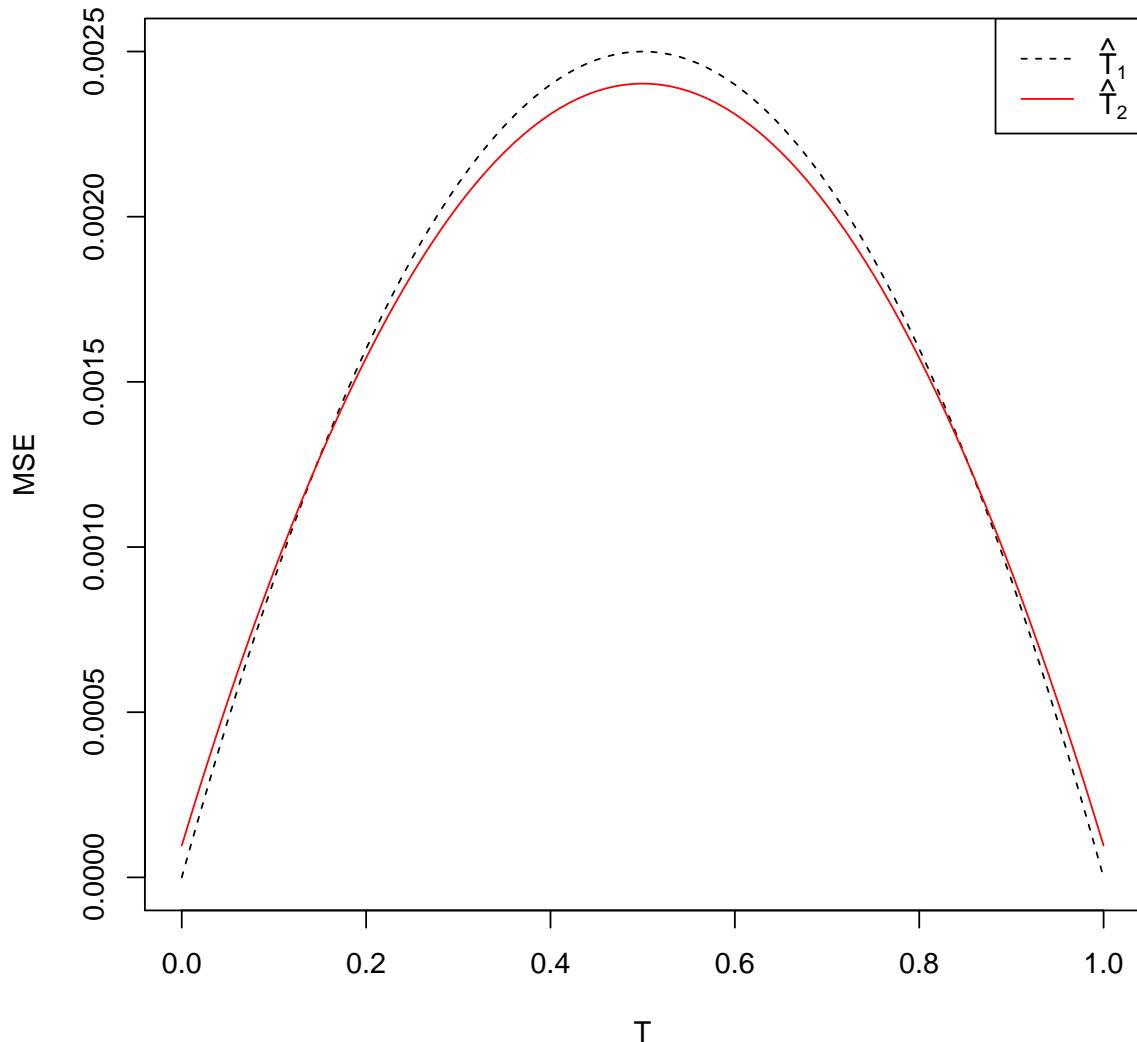
We are looking for an estimator with the **lowest** error, thus making:

- if $\text{eff}(T, U) < 1$, estimator U should be preferred,
- if $\text{eff}(T, U) > 1$, estimator T should be preferred.

This can be demonstrated graphically in the following way. MSE represents a curve as a function of parameter T . We can observe that the curve separates the MSE in two regions. Choosing the right estimator thus depends on many factors.

- when $T \in [0, 0.15]$ or $T \in (0.85, 1]$: $\text{MSE}(\hat{T}^{(2)}, T) > \text{MSE}(\hat{T}^{(1)}, T)$,
- when $T \in (0.15, 0.85)$: $\text{MSE}(\hat{T}^{(1)}, T) > \text{MSE}(\hat{T}^{(2)}, T)$.

Illustration for the MSE. Here T represents the proportion of the observed values.



27.5 Consistency

A consistent estimator "approaches" the true parameter θ as the sample size increases. Put more formally, an estimator $T(X_1, \dots, X_n)$ is said to be consistent if $\text{Plim}_{n \rightarrow \infty} T(X_1, \dots, X_n) = \theta$, that is, if $\forall \epsilon > 0$ ²

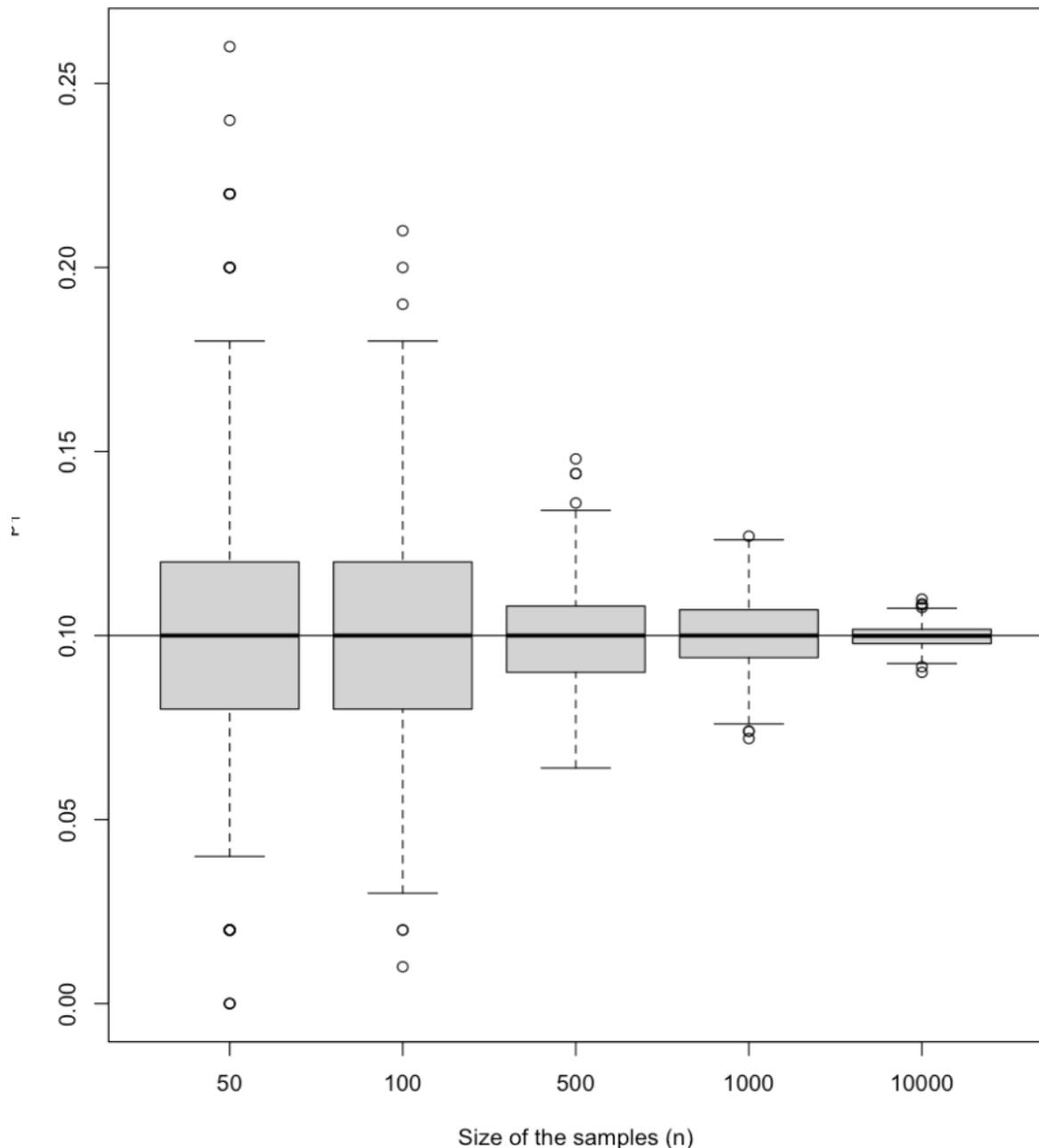
$$P(|T(X_1, \dots, X_n) - \theta| > \epsilon) \xrightarrow{n \rightarrow \infty} 0.$$

² See also Probability section for more information. This will come in useful during the first Time Series lecture.

Let's again illustrate this with an example. In figure 10, we are looking at simulations of a random variable following a Bernoulli distribution. On the x -axis different n are represented. The true parameter p is represented by the black horizontal line.

Figure 10: Bias and consistency.

Boxplot of simulated \hat{p}_1 for different sample size



The boxplots show that \hat{p}_1 is an unbiased and consistent estimator. It is unbiased because $E(\hat{p}_1) = p$ for any n ; it is consistent since the sampling distribution of \hat{p}_1 is becoming more concentrated at p as the sample size increases.



Bias, variance, MSE and efficiency

For more information about choosing the right estimator, see [Statistics lectures](#), Preliminary Work Week 7, [Probability and Statistical Learning lectures](#), Section 7, Subsections 3 and 4, Heumann and Shalabh [4] book chapter 9, Moore and McCabe [5] pages 210-212 for bias and variance, Rice [6] book p. 127 for MSE.

28 CONSTRUCTING ESTIMATORS

When constructing estimators, two main approaches exist. You will use both throughout the whole masters, make sure you are very familiar with the concepts as well as procedures to derive these estimators.

28.1 Method of moments

The method of moments works by matching the population moments to the empirical moments. The method is very intuitive and simple to implement. The procedure is as follows.

Method of moments implementation.

Step 1: compute the empirical moment of order k for the population

Step 2: replace the empirical moment with sample moment

The population moment of order k is defined as

$$\mu_k = E(X^k),$$

and the empirical sample moment as

$$\hat{\mu}_k = \frac{1}{n} \sum_{i=1}^n X_i^k.$$

For example, $X \sim Ber(p)$ (X following Bernoulli distribution), we know $E(X) = p$, we thus replace $E(X)$ by \bar{X} and get $\bar{X} = p \Rightarrow \hat{p}_{MM} = \bar{X}$.

28.2 Maximum likelihood

The idea behind the maximum likelihood is to find a pdf that is **most plausible** based on the observations.

For the intuition, you can watch [this video](#) - by the way, the whole channel is great to gain intuition behind many statistical concepts.

The method is very intuitive and allows you to get the full exploit of the information at hand. The likelihood function can be used to summarize data. Whereas when working with distribution, we focus on the observations given the same function parameter, when working with the likelihood function, one focuses more in the probability of the sample we actually observed under various possible values of θ , the function parameter [2].

Let $f(\mathbf{x} | \theta)$ denote the joint pdf or pmf of the sample $\mathbf{X} = (X_1, \dots, X_n)$. Then, given that $\mathbf{X} = \mathbf{x}$ is observed, the function of θ defined by

$$L(\theta | \mathbf{x}) = f(\mathbf{x} | \theta)$$

is called the likelihood function. Note the difference in the order of notation for θ and x .

While the likelihood function will be absolutely a necessary concept throughout the whole master's program, if you are particularly interested in the mathematical statistics and the deeper understanding of the purely statistical concepts, we recommend you to follow the [Advanced Statistical Inference](#) course, that will give you very deep insights and concepts you will probably not see anywhere else. You can also read [2] or [9].

Maximum likelihood computation.

Step 1: determine what pdf X follows

Step 2: determine the pdf for X_1, \dots, X_n . Generally speaking, the likelihood is the product of pdf's of X_1, \dots, X_n . If X are iid (independent identically distributed), then the likelihood function can be written as

$$\begin{aligned} L(\theta | x_1, \dots, x_n) &= \\ &= \begin{cases} \prod_{i=1}^n f_\theta(x_i) = f_\theta(x_1) \cdots f_\theta(x_n) & X_i \text{ continuos} \\ \prod_{i=1}^n P_\theta(X_i = x_i) = P_\theta(X_1 = x_1) \cdots P_\theta(X_n = x_n) & X_i \text{ discrete} \end{cases} \end{aligned}$$

Step 3: take the log of the likelihood function. The reason we take the log is because it's easier to work with sums rather than products.

$$\begin{aligned} \ell(\theta | x_1, \dots, x_n) &= \log(L(\theta | x_1, \dots, x_n)) = \\ &= \begin{cases} \sum_{i=1}^n \log(f_\theta(x_i)) & X_i \text{ continuos} \\ \sum_{i=1}^n \log(P_\theta(X_i = x_i)) & X_i \text{ discrete} \end{cases} \end{aligned}$$

Step 4: take the FOC, first order condition and solve it for θ .

$$\frac{\partial \ell(\theta | x_1, \dots, x_n)}{\partial \theta} = 0 = \begin{cases} \sum_{i=1}^n \frac{\partial \log f_\theta(x_i)}{\partial \theta_i} & X_i \text{ continuous} \\ \sum_{i=1}^n \frac{\partial \log P_\theta(X_i = x_i)}{\partial \theta} & X_i \text{ discrete} \end{cases}$$

The function $\frac{\partial \ell(\theta | x_1, \dots, x_n)}{\partial \theta}$ is called the **score** function.

Step 5: take the second derivative to make sure it's a maximum.

$$\frac{\partial^2 \ell(\theta | x_1, \dots, x_n)}{\partial \theta^2} \stackrel{?}{=} 0$$

A quick reminder from calculus: when a function's slope (first derivative) is zero at x , and the second derivative at x is:

1. less than 0, it is a local maximum
2. greater than 0, it is a local minimum
3. equal to 0, then the test fails (there may be other ways of finding out though).

Thus, for a maximum likelihood function we want the second derivative to be less than zero:

$$\frac{\partial^2 \ell(\theta | x_1, \dots, x_n)}{\partial \theta^2} \stackrel{?}{<} 0$$

The result of the equation in *Step 4* is the estimator $\hat{\theta}_{MLE}$ that maximizes the likelihood function with respect to θ , thus finding an estimator with the most *alike* the theoretical pdf.

The maximum likelihood estimators have a lot of nice properties and the technique is general. The results are the following:

1. $\hat{\theta}_{ML}$ can be biased for θ .
2. $\hat{\theta}_{ML}$ is consistent for θ .

There is a general distributional asymptotic result:

$$\sqrt{n} (\hat{\theta}_{ML} - \theta) \sim_{n \rightarrow \infty} N(0, J^{-1}(\theta)),$$

where $J(\theta) = E \left(\left(\frac{\partial \log(f(\theta | x_1, \dots, x_n))}{\partial \theta} \right)^2 \right)$ is called the **Fisher information**.

This implies that $\hat{\theta}_{ML} \sim_{approx} N(\theta, \frac{1}{n} J^{-1}(\theta))$ because $E(\hat{\theta}_{ML}) \simeq \theta$ and

$$\text{Var}(\hat{\theta}_{ML}) = \frac{1}{n} \text{Var}(\sqrt{n} (\hat{\theta}_{ML} - \theta)) \simeq \frac{1}{n} J^{-1}(\theta).$$

Method of moments and Maximum Likelihood



For more information about method of moments, maximum likelihood, examples of computations, asymptotics and the confidence intervals based on those estimators, see [Statistics lectures](#), Preliminary Work Week 8, [Probability and Statistical Learning lectures](#), Section 7, Subsection 1, Section 9, Heumann and Shalabh [4] book chapter 9.3, [Advanced Statistical Inference](#) course (advanced) and Casella and Berger [2] book chapter 6 (advanced).

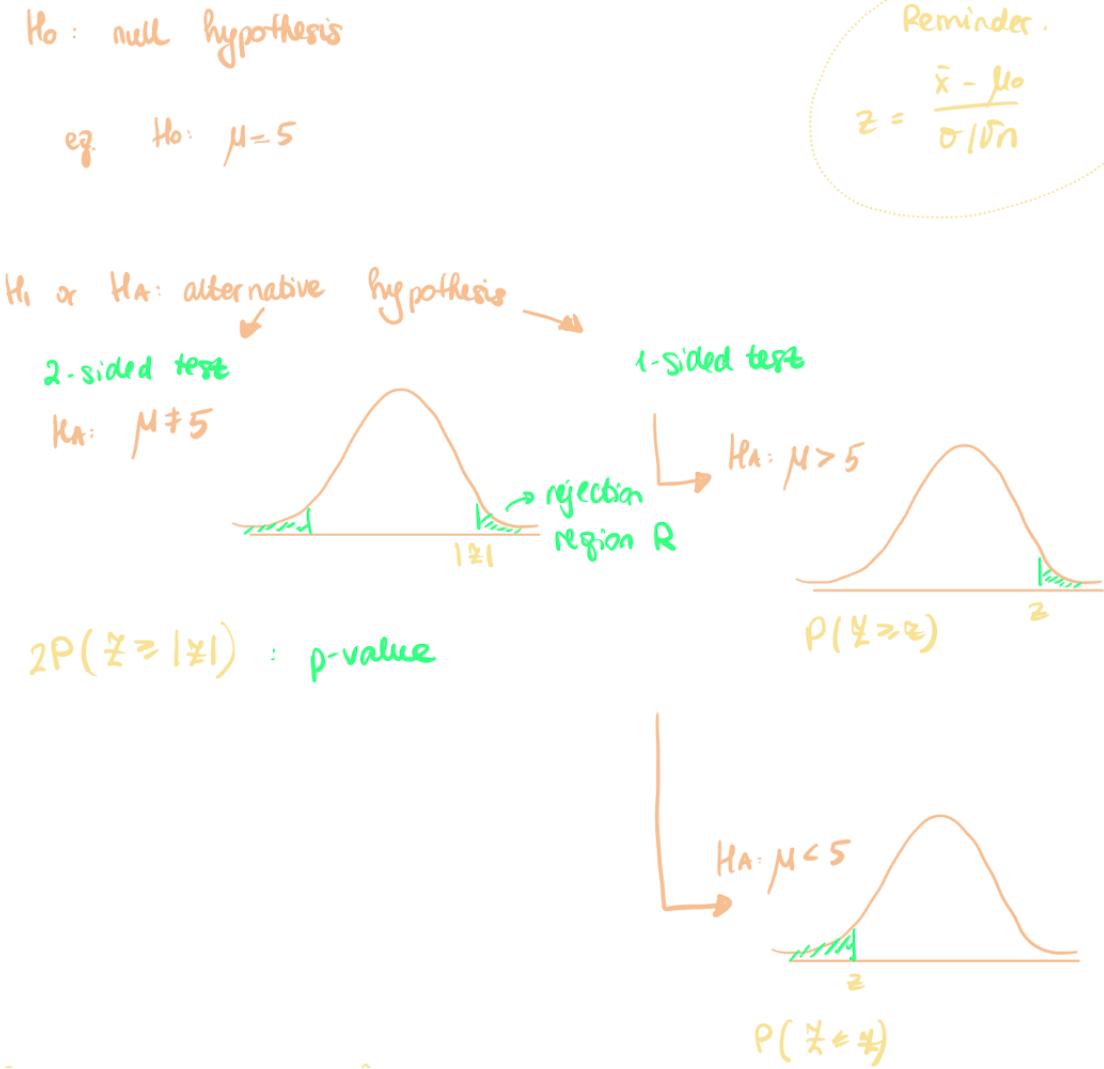
Hypothesis testing

When making decisions, a formal procedure needs to be implemented to compare the observed data with a hypothesis. This is exactly the idea behind hypothesis testing. It contains several "ingredients":

1. Random variable and its distribution,
2. Confidence level,
3. Test statistic.

Hypothesis testing can be summarized as in figure 11.

Figure 11: Hypothesis testing.



Now some definitions and tips:

- Null hypothesis, denoted H_0 , is the assumed or specific value you want to test for,
- Alternative hypothesis is what you believe the value to be,
- The confidence level α is always set **before** the experiment, the most common values being $\alpha = 0.05, \alpha = 0.01, \alpha = 0.1$,
- As to the test statistic, you can take mean, proportion, difference of means or another statistic of your choice.

To determine whether to reject or not reject the H_0 , you can use one of the following methods presented in figure 12. They all yield the same result.

Figure 12: Three approaches to hypothesis testing.

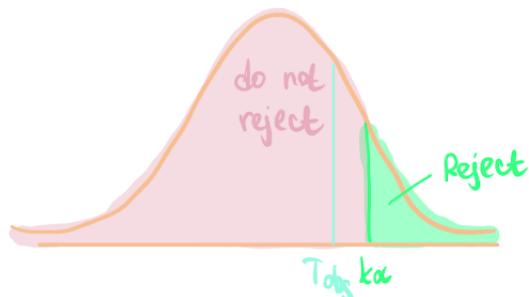
3 approaches that give the same result (think about why!)

- 1) Compare T_{obs} and $k\alpha$ (critical value)
- 2) Determine whether $T_{obs} \in R$ (rejection region)
- 3) Compare the p-value to α

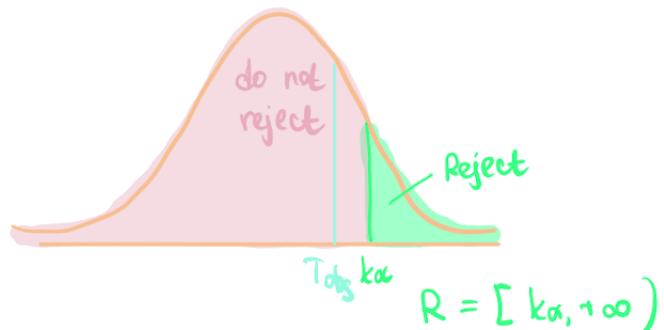
1) $T_{obs} < k\alpha$

Question:

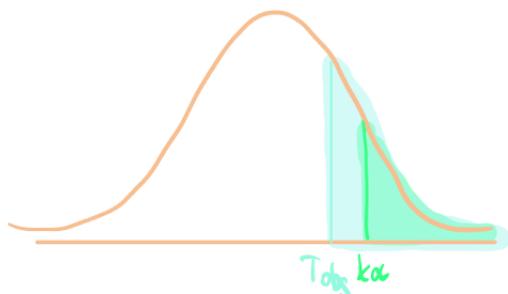
What would be
your conclusion?



2) $T_{obs} \notin R$



3) $P(Z \geq k\alpha) > \alpha$



observed area > rejection region area

- In method 1, the observed value of the statistic, T_{obs} , is compared to the critical value for the specified level of α , k_α .
- In method 2, rejection region R is being constructed from k_α : $R = [k_\alpha, \infty)$.
- In method 3, $p - value$ (observed area) is computed and compared to the rejection region area.

For the mean of the observations, assuming the normal distribution, known σ and n , $H_0 : \mu = \mu_0$ and $H_A : \mu > \mu_0$, the $p - value$ under H_0 can be computed as follows:

$$P_{H_0}(\bar{X} > t.\text{obs}) = P_{H_0}\left(\frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} > \frac{t.\text{obs} - \mu_0}{\sigma/\sqrt{n}}\right) \quad (10)$$

Then, H_0 is rejected when $p - value \leq \alpha$. When $p - value > \alpha$, we do not have enough evidence to reject H_0 .

29 TYPE I AND II ERROR

When dealing with hypothesis testing, different types of errors have to be dealt with. They can be summarized as follows.

	H_0 true	H_0 not true
Reject H_0	Type I error	Correct decision
Don't reject H_0	Correct decision	Type II error

We define probabilities associated with each situation

$$\alpha = P(\text{reject } H_0 \mid H_0 \text{ true}) = P_{H_0}(\text{reject } H_0)$$

and

$$\beta = P(\text{do not reject } H_0 \mid H_0 \text{ not true}) = P_{H_A}(\text{do not reject } H_0).$$

The **power** is defined as $power = 1 - \beta = 1 - P_{H_A}(\text{do not reject } H_0) = P_{H_A}(\text{reject } H_0)$.
In order to increase the power of a test, following can be done:

- increase n ,
- increase the distance between the two hypotheses ($(\mu - \mu_0) \uparrow$),
- increase α .

Figure ?? illustrates the power of a test.

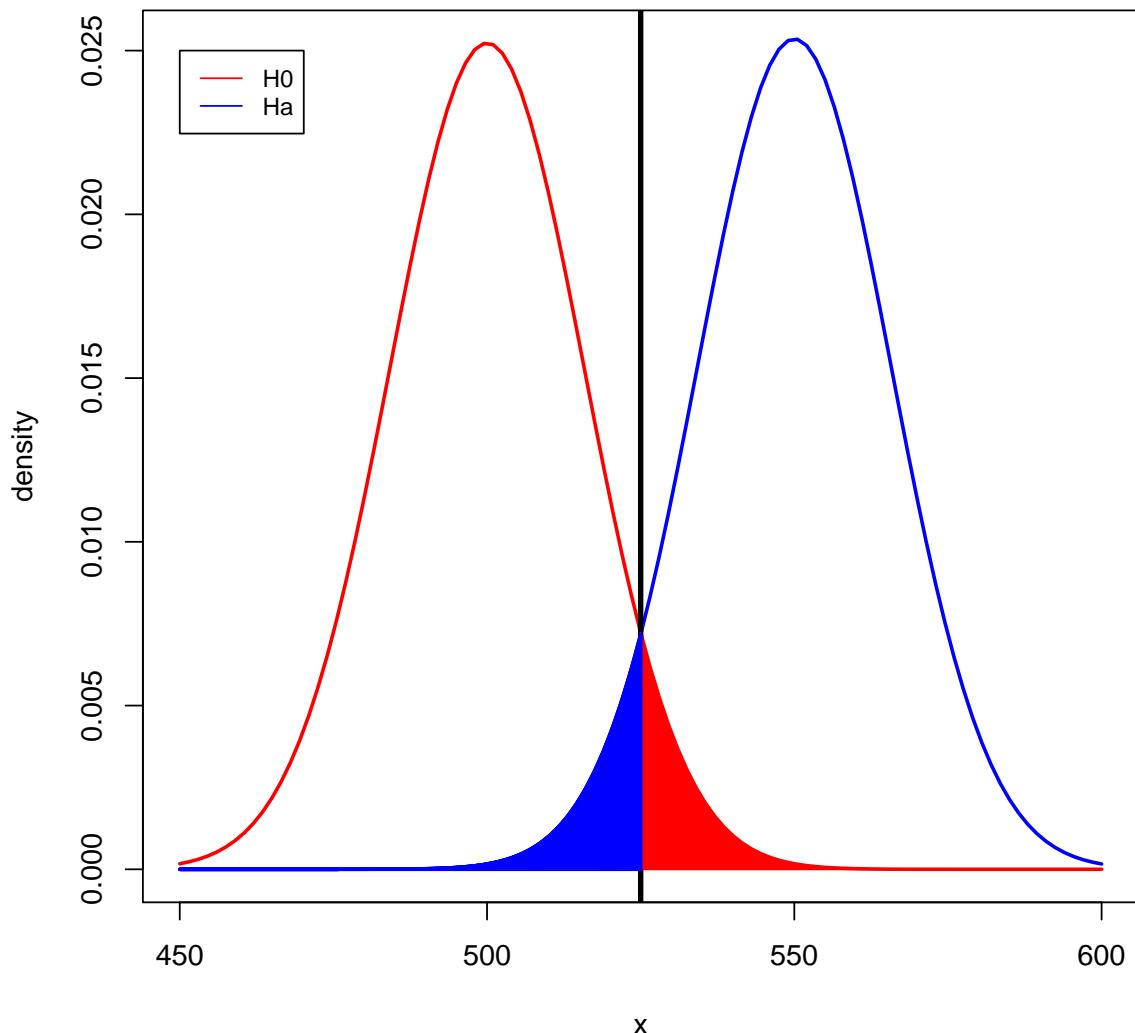


Figure 13: Power of the test.

30 COMPARING 2 GROUPS

When testing for location (mean), two groups can be compared. The idea stays the same, but the formula has to be adjusted.

Let X_1, \dots, X_{n_X} and Y_1, \dots, Y_{n_Y} two independent samples, with $E(X_i) = \mu_1$ and $E(Y_i) = \mu_2$. To test $H_0 : \mu_1 = \mu_2$ (or equivalently $H_0 : \mu_1 - \mu_2 = 0$), use

$$T = \frac{\bar{X} - \bar{Y}}{\sqrt{\frac{s_X^2}{n_X} + \frac{s_Y^2}{n_Y}}} \sim_{\text{approx}, H_0} N(0, 1),$$

where $s_X^2 = \frac{1}{n-1} \sum_{i=1}^{n_X} (X_i - \bar{X})^2$ and $s_Y^2 = \frac{1}{n-1} \sum_{j=1}^{n_Y} (Y_j - \bar{Y})^2$.

T can then be compared to k_α using Method 1 mentioned in section 28.2.

31 ANALYSIS OF VARIANCE (ANOVA)

When comparing more than 2 groups, it is useful to do analysis of variance test (ANOVA). A following example can be helpful. Here, we compare means of several (I) groups (samples), having J measures per group.

Model (one-way layout):

$$\begin{array}{lll} Y_{ij} & = \mu_i + \varepsilon_{ij} & = \underset{\text{j-measure}}{\mu} + \underset{\text{i-group}}{\alpha_i} + \underset{\text{random error}}{\varepsilon_{ij}}, \\ & & \end{array} \quad (11)$$

for $i = 1, \dots, I$ et $j = 1, \dots, J$.

We assume $\varepsilon_{ij} \sim \mathcal{N}(0, \sigma^2)$ and require $\sum_{i=1}^I \alpha_i = 0$ (or $\alpha_1 = 0$,

par ex.).

We test: $H_0 : \alpha_1 = \alpha_2 = \dots = \alpha_I = 0$ (which implies

$$\mu_1 = \mu_2 = \dots = \mu_I = \mu.$$

The analysis of variance is based on the following decomposition of Y_{ij} :

$$Y_{ij} = Y_{..} + (Y_{i.} - Y_{..}) + (Y_{ij} - Y_{i.})$$

with $Y_{i.} = \frac{1}{J} \sum_j Y_{ij}$ the group i mean, and $Y_{..} = \frac{1}{IJ} \sum_{i,j} Y_{ij}$ the global mean. We therefore have the following

$$\underbrace{\sum_{i,j} (Y_{ij} - Y_{..})^2}_{SS_T} = \underbrace{\sum_i (Y_{i.} - Y_{..})^2}_{SS_B} + \underbrace{\sum_{i,j} (Y_{ij} - Y_{i.})^2}_{SS_W}, \quad (12)$$

where SS_T is the total sum of squares,

SS_B is the sum of squares between groups and

SS_W is the sum of squares within groups.



Hypothesis testing

For more information about hypothesis testing, see [Statistics lectures](#), Preliminary Work Weeks 9-12, [Probability and Statistical Learning lectures](#), Week 11, Subsection 2, Heumann and Shalabh [4] book chapter 10, Moore and McCabe [5] pages 660 and on, Rice [6] book chapter 9.

REFERENCES

- [1] Ayyub, B. M. and McCuen, R. H. (2016). *Probability, statistics, and reliability for engineers and scientists*. CRC press.
- [2] Casella, G. and Berger, R. L. (2021). *Statistical inference*. Cengage Learning.
- [3] Commissions romandes de mathématique, d. p. e. d. c. (2006). *Formulaires et tables*. Editions du Tricorne.
- [4] Heumann, C. and Shalabh, M. S. (2016). *Introduction to statistics and data analysis*. Springer.
- [5] Moore, D. S. and McCabe, G. P. (1989). *Introduction to the Practice of Statistics*. WH Freeman/Times Books/Henry Holt & Co.
- [6] Rice, J. A. (2006). *Mathematical statistics and data analysis*. Cengage Learning.
- [7] Ross, S. (2014). *A First course in Probability*. Pearson New International Edition.
- [8] Rozanov, Y. (1977). *Probability Theory: A Concise Course*. Dover Publications.
- [9] Shao, J. (2003). *Mathematical statistics*. Springer Science & Business Media.