

Python 学习笔记

Ysong

目录

1	Basic grammar	1
1.1	type	1
1.2	operator precedence	1
1.3	condition statements	2
1.4	loop	2
1.5	function	2
1.5.1	定义函数	2
1.5.2	调用函数	3
1.5.3	可变参数	3
1.5.4	静态类型函数	4
1.5.5	变量的作用域	4
1.5.6	High-order function	5
1.5.7	self-reference	6
1.5.8	decorator	6
1.5.9	function currying	7
1.5.10	匿名函数	7
1.6	input and output	7
1.6.1	打开关闭文件	7
1.6.2	读取写入文件	8
1.6.3	上下文管理器	9
1.6.4	处理.csv 文件	9
1.6.5	处理.json 文件	9
1.6.6	处理 zip 文件	10
2	Collections	10
2.1	iterator	10
2.1.1	itertools	11
2.2	list	11
2.3	tuple	12
2.4	set	12
2.5	dict	13
2.6	String	13
2.6.1	编码解码	13
2.6.2	常用操作	13
2.6.3	格式化字符串	14
2.7	collection	15
2.8	SQLite	16
3	Class	17
3.1	定义类模板	17
3.1.1	基本结构	17
3.1.2	访问限制	17
3.1.3	属性	17
3.1.4	类的特殊成员 (魔术方法)	18
3.2	类的使用	19
3.2.1	创建类的实例	19
3.2.2	类成员的使用	19
3.3	继承	20
3.4	有关类的内置函数	20

3.5	类的典型例子	20
3.5.1	文件管理类	20
4	Basic module	21
4.1	install and import module	21
4.1.1	install and manage modules	21
4.1.2	virtualenv	21
4.1.3	pipenv	21
4.1.4	anaconda	22
4.1.5	import modules	23
4.1.6	view the functions in modules	23
4.1.7	circular import	23
4.2	builtins module	23
4.3	sys module	23
4.4	os module	23
4.5	subprocess module	24
4.6	time module	24
4.7	threading module	25
4.7.1	threading	25
4.7.2	concurrent.futures	25
4.7.3	multiprocessing	25
4.8	logging module	26
4.9	re module	26
4.9.1	regular expression	26
4.9.2	match string	27
4.9.3	substitute string	28
4.9.4	split string	29
4.10	date module	29
4.10.1	datetime module	29
4.10.2	pytz module	30
4.10.3	calendar module	30
4.11	PIL module	30
4.12	random module	30
4.13	secrets module	31
4.14	cn2an module	31
4.15	Progress Bar	31
4.15.1	tqdm module	31
5	Maths & Statistics	31
5.1	math	31
5.2	numpy	32
5.2.1	定义矩阵	32
5.2.2	特殊函数	32
5.2.3	矩阵运算	32
5.3	Pandas	32
5.3.1	Series	32
5.3.2	DataFrame	33
5.3.3	读写数据	33
5.3.4	数据库交互	33
5.3.5	基本操作	34
5.3.6	统计操作	34

5.3.7	处理 datetime	35
6	visualization	36
6.1	matplotlib	36
6.1.1	画布预处理	36
6.1.2	子图表	36
6.1.3	数据获取	37
6.1.4	折线图	37
6.1.5	柱状图	38
6.1.6	饼图	38
6.1.7	散点图	38
6.1.8	箱型图	39
6.1.9	热力图	39
6.1.10	时间图像	39
6.1.11	极坐标图表	40
6.1.12	三维图表	40
6.1.13	动态实时图像	40
6.2	wordcloud	41
7	website	41
7.1	请求连接	41
7.1.1	urllib module	41
7.1.2	urllib3 module	42
7.1.3	requests module	42
7.2	处理 html 文件	42
7.3	Flask	43
7.3.1	create a website	43
7.3.2	template	43
7.3.3	template inheritance	45
7.3.4	bootstrap	45
7.3.5	External CSS	46
7.3.6	post data	47
7.3.7	form & field	47
7.3.8	flask-sqlalchemy	49
7.3.9	Bcrypt & LoginManager	51
7.3.10	FileField	52
7.3.11	Pagination	52
8	debug & tools	53
8.1	doctest	53
8.2	Try & Assert	54
8.3	Unit Testing	55
8.4	jupyter notebook	56
9	deep learning	56
9.1	pytorch 基本操作	56
9.1.1	数据操作	56
9.1.2	读写文件	57
9.1.3	线性代数	58
9.1.4	微分	58
9.1.5	概率	58
9.2	线性神经网络	58

9.2.1	线性回归	58
9.2.2	softmax 回归	59
9.2.3	多层向量机	60
9.3	获取数据集	60
9.3.1	线性回归数据集	60
9.3.2	Fashion-MNIST 图像分类数据集	61
9.3.3	数据加载切分	61

1 Basic grammar

1.1 type

Listing 1.1: variable type

```

1 >>> str = "python"
2 >>> dir(str) #获取相关用法
3 >>> type(str) #获取变量类型
4 >>> help(str) #获取 guide
5 >>> id(str) #获取内存地址
6 >>> isinstance(str, string) #若 str 为 string 类型, 则返回 True, 否则返回 False
7

```

Listing 1.2: basic functions

```

1 print('Hello World!')
2 print('Hello', 'World!', sep=', ', end='') #各字符串间用 ` , ` 分隔, 结尾没有换行符
3 print(round(3.75, 1)) #四舍五入保留一位小数
4 float_number = 3.3; int_number = int(float_number)
5

```

1.2 operator precedence

Table 1.1: Python 运算符优先级表

优先级	运算符
1	() 函数调用
2	[] 下标
3	. 属性访问
4	** 指数运算
5	+, - 正负号
6	~ 按位非
7	*, /, //, % 乘, 除, 地板除, 取余
8	+, - 加, 减
9	<<, >> 左移, 右移
10	& 按位与
11	按位或
12	^ 按位异或
13	<, <=, >, >=, ==, != 比较运算符
14	= 赋值运算符
15	in not in 逻辑包含/不包含
16	not 逻辑非
17	and 逻辑与
18	or 逻辑或

“==”checks for equality, “is” checks for identity(sharing the same memory).

Listing 1.3: is & ==

```

1 "Jack" == "Jack" #return True
2 "Jack" is "Jack" #return False
3

```

无论是 `x = x+1` 还是 `x += 1` 都会重新给 `x` 分配内存。

1.3 condition statements

Listing 1.4: condition statement

```
1 def absolute_value(x):
2     """Return the absolute value of x."""
3     if x < 0:
4         return -x
5     elif x == 0:
6         return 0
7     else:
8         return x
9
10 def absolute_value(x):
11     return -x if x < 0 else x
12
```

- False values in Python: False, 0, ' ', None
- True values in Python: Anything else

1.4 loop

Listing 1.5: iteration

```
1 i, total = 0, 0
2 while i < 3:
3     total += i
4     i += 1
5 else: #若没有 break 则执行
6     print(i)
7
8 for i in range(3):
9     total += i
10 else: #若没有 break 则执行
11     print(i)
12
```

1.5 function

函数用于处理需要多次重复运行的同一任务。python 的函数默认返回 None。

1.5.1 定义函数

Listing 1.6: definition of functions

```
1 def my_func(string1,string2):
2     print(string1)
3     print(string2)
4
5 def my_func1(string1='Hello',string2): #可以指定默认值
6     ''' 功能:打招呼
7         string1:默认hello
8         string2:对象名字 '''
9     print(string1)
10    print(string2)
11    return 1
12
13 def my_info(*args, **kwargs): #这里 *args 会创建一个元组, kwargs 会创建一个字典
14     print(args)
15     print(kwargs)
16
17 #可以使用 yield 函数返回, yield 允许函数在返回一个值的同时保存状态,并在下一次继续执行
```

```

18 #带 yield 的函数是一个 iterator
19 def count_up_to(n):
20     counter = 1
21     while counter <= n:
22         yield counter
23         counter += 1
24
25 #调用带有 yield 关键字的函数时，它会返回一个生成器对象
26 gene = count_up_to(5)
27 print(next(gene))
28 for number in gene:
29     print(number)
30 #当生成器函数不再遇到 yield 语句时，生成会终止
31

```

1.5.2 调用函数

以下是几种正确的调用方法

Listing 1.7: use of function

```

1 my_func('Hello', 'world')
2 str1='Hello';str2='world'
3 my_func(str1,str2)
4 my_func(string1=str1,string2=str2) #写明形式参数名称为关键字参数
5 my_func(str1,string2=str2) #允许前面的形参不写名称，后面的写明
6 my_func(string2=str2,string1=str1) #写明形式参数名时可以交换顺序
7 info_value = ['Math', 'Art']; dic_value = {'name': 'John', 'age': 22}
8 my_info('Math', 'Art', name='John', age=22)
9 my_info(*info_value, **dic_value)
10

```

以下是几种错误的调用方法

Listing 1.8: mistakes in using functions

```

1 my_func(string1=str1,str2) #前面写明形参名称，后面必须写
2 my_func(string2=str2,str1)
3

```

1.5.3 可变参数

可变参数是 Python 特有的设计，可变参数也被称为不定长参数，即传入函数中的实际参数可以是 0 个、1 个或多个。

Listing 1.9: variable parameters

```

1 def greet(*names): #这里会将输入接受并放在一个元组中
2     print("Hello",end=' ')
3     for item in names:
4         print(item,end=' ')
5
6     if __name__=='__main__':
7         greet('Tom','Jerry')
8
9     param=['Tom','Jerry'] #也可以调用列表
10    greet(*param) #这里不能写形参名
11

```

另一种方法会将输入以“形参名: 变量名”的形式存为一个字典。

Listing 1.10: other variable parameters

```

1 def greet(**names): #输入转为字典
2     print('Hello')

```



```

3         for key,value in names.items():
4             if value=='male':
5                 print('Mr',end=' ')
6             else:
7                 print('Miss',end=' ')
8             print(key)
9
10        if __name__=='__main__':
11            greet(Tom='male',Jerry='female') #直接调用
12
13        dict={'Tom':'male','Jerry':'female'} #也可以用字典调用
14        greet(**dict) #这里也不能用关键字参数
15

```

1.5.4 静态类型函数

Python 的函数默认是动态类型的，可以接受所有类型的输入。当然也可以将函数定为静态类型的 (和 C++ 类似)。

Listing 1.11: static functions

```

1 def square(number:int|float)->int|float:
2     return num**2
3

```

1.5.5 变量的作用域

- 在函数内的定义的变量是局部变量，外部不能调用。
- 函数外部定义的变量是全局变量，可以在任意位置调用该变量。
- 函数内定义的变量可以通过 `global` 关键字定为全局变量。
- `environment` 就是 `frame` 的顺序，当调用某变量时，会从里向外逐层寻找该变量。
- 在一个函数中，一个变量名只能始终为一个全局变量，或者始终是一个局部变量。
- 在内部函数中的变量名前加上 `nonlocal`，则内部函数的该变量与外部函数的变量是同一个。

Listing 1.12: local and global

```

1 message='Hello'
2 def my_func():
3     global message
4     cnt = 0
5     def in_func():
6         nonlocal cnt
7         message='World'
8         cnt++
9     print(message) #输出 'World'
10
11 #下面这段代码会报错
12 x = 'global x'
13
14 def test():
15     print(x)
16     x = 'local x'
17     print(x)
18
19 test()
20 #解决这个报错的方法有两种，即函数内声明全局或不使用全局
21 #可以在第一个 print(x) 前面加上 global x
22 #也可以删去第一个 print(x)
23

```

1.5.6 High-order function

High-order function 是一类返回函数的函数，用于处理同类型的任务。

Listing 1.13: High-order function

```
1 def make_adder(n):
2     """Return a function that takes K and return K + N.
3
4     >>> add_three = make_adder(3)
5     >>> add_three(4)
6     """
7     def adder(k):
8         return k + n
9     return adder
10
11 add_three = make_adder(3)
12 print(make_adder(100)(4))
13
```

High-order function 在音频领域经常被使用，用于制作波形。

Listing 1.14: mario

```
1 from wave import open
2 from struct import Struct
3 from math import floor
4
5 frame_rate = 11025
6
7 def encode(x):
8     i = int((1 << 14) * x)
9     return Struct('h').pack(i)
10
11 def play(sampler, name='song.wav', seconds=2):
12     out = open(name, 'wb')
13     out.setnchannels(1)
14     out.setsampwidth(2)
15     out.setframerate(frame_rate)
16     t = 0
17     while t < seconds * frame_rate:
18         sample = sampler(t)
19         out.writeframes(encode(sample))
20         t += 1
21     out.close()
22
23 def tri(frequency, amplitude=0.3):
24     period = frame_rate // frequency
25     def sampler(t):
26         saw_wave = t / period - floor(t / period + 0.5)
27         tri_wave = 2 * abs(2 * saw_wave) - 1
28         return amplitude * tri_wave
29     return sampler
30
31 c_freq, e_freq, g_freq = 261.63, 329.63, 392.00
32
33 def both(f, g):
34     return lambda t: f(t) + g(t)
35
36 def note(f, start, end, fade=0.01):
37     def sampler(t):
38         seconds = t / frame_rate
39         if seconds < start:
40             return 0
41         elif seconds > end:
42             return 0
43         elif seconds > end - fade:
44             return (end - seconds) / fade * f(t)
45         elif seconds < start + fade:
46             return (seconds - start) / fade * f(t)
47         else:
```

```

48         return f(t)
49     return sampler
50
51     c, e, g = tri(c_freq), tri(e_freq), tri(g_freq)
52     g_low = tri(g_freq / 2)
53     z = 0
54     song = note(e, z, z + 1/8)
55     z += 1/8
56     song = both(song, note(e, z, z + 1/8))
57     z += 1/4
58     song = both(song, note(e, z, z + 1/8))
59     z += 1/4
60     song = both(song, note(c, z, z + 1/8))
61     z += 1/8
62     song = both(song, note(e, z, z + 1/8))
63     z += 1/4
64     song = both(song, note(g, z, z + 1/8))
65     z += 1/2
66     song = both(song, note(g_low, z, z + 1/8))
67     z += 1/2
68
69     play(song)
70

```

1.5.7 self-reference

self-reference function 会返回自身或者与自身相关的函数，这样可以递归的执行不定长度的任务。

Listing 1.15: self-reference

```

1  def print_all(x):
2      print(x)
3      return print_all
4  print_all(1)(2)(3) #这里会全部输出来
5

```

1.5.8 decorator

decorator 是一种用于修改函数或方法行为的工具，它接受另一个函数作为输入，并返回一个新的函数。decorator 可以叠加多层

Listing 1.16: decorator

```

1  def decorator_function(original_function):
2      def wrapper_function():
3          print('wrapper executed before {}'.format(original_function.__name__))
4          return original_function()
5      return wrapper_function
6
7  @decorator_function #将装饰器应用到 display 函数上，使 display 自动被装饰器的逻辑包裹
8  def display():
9      print('display function ran')
10
11  #display = decorator_function(display) 装饰器的效果和这一句相同
12
13  display() #上面的两句都会打印
14
15  #若待装饰的 function 是有参数的，则需要给装饰器的返回值也添加上参数
16  def decorator_function(original_function):
17      def wrapper_function(*args, **kwargs):
18          print('wrapper executed before {}'.format(original_function.__name__))
19          return original_function(*args, **kwargs)
20      return wrapper_function
21
22  #decorator 也可以带有别的参数

```

```
23 def prefix_decorator(prefix):
24     def decorator_function(original_function):
25         def wrapper_function():
26             print(prefix, 'wrapper executed before {}'.format(original_function.__name__))
27             return original_function()
28         return wrapper_function
29     return decorator_function
30
31 @prefix_decorator('TESTING:')
32 def display():
33     print('display function ran')
34
```

1.5.9 function currying

function currying 将一个接受多个参数的函数分解为一系列每个只接受一个参数的函数，即一个函数链。

Listing 1.17: function currying

```
1 def curry2(f):
2     def g(x):
3         def h(y):
4             return f(x, y)
5         return h
6     return g
7
8 from operator import add
9 m = curry2(add)
10 add_three = m(3)
11 print(add_three(2))
12
```

1.5.10 匿名函数

Python 的匿名函数也就是 lambda 表达式。匿名函数可以是另一个函数的参数，如 sort 函数的 key 变量。

- 只有 def 得到的函数会有一个本征名，在 shell 中输入变量名即可看到这一差异。
- def 函数可以有多步运行过程，而 lambda 表达式只能有一个表达式。

语法为 `result=lambda [arg1[arg2,...,argn]]:expression`

Listing 1.18: lambda function

```
1 result = lambda r:math.pi*r*r
2 area = (lambda r: math.pi * r * r)(3)
3 print(result(5))
4
5 def area(r):
6     return r * r * math.pi
7 result = area
8
```

1.6 input and output

1.6.1 打开关闭文件

Listing 1.19: open and close files

```

1 file = open(<filename>[,mode]) #打开文件
2 print(file.name); print(file.mode) #输出文件名和打开方式
3 file.close() #关闭文件
4 print(file.closed) #检查文件是否已经关闭, 已关闭会返回 True
5
6 #用 with 语句打开文件可以确保文件被正确关闭, 且可以避免文件打开错误引起的异常
7 with open(<filename>[,mode]) as file:
8

```

Table 1.2: Methods of Opening a File

符号	说明	注意
r	只读模式, 指针放在开头	
rb	二进制只读模式	
r+	可以读取, 也可以从文件的开头覆盖	文件必须存在
rb+	二进制读写模式	
w	只写模式	
wb	二进制只写模式	若文件存在, 则覆盖, 否则
w+	只读模式打开文件并清空	创建新文件
wb+	二进制只读模式并清空文件	
a	追加模式	
ab	二进制追加模式	若文件存在, 则指针放在末尾, 否则创建文件
a+	追加且可读 (注意指针位置)	
ab+	二进制追加且可读	

1.6.2 读取写入文件

Listing 1.20: read and write files

```

1 #读取文件
2 with open('test1.txt', 'r', encoding='utf-8') as file:
3     f_contents = file.read(); f_contents = file.read(100) #读取文件前 100Byte 的内容
4     f_list = file.readlines() #按行切分并组成列表
5     f_list = file.readline() #只读一行, 结尾带有换行符
6     for line in file:
7         print(line, end='') #遍历每一行
8     file.tell() #返回 cursor 在文件中的位置
9     file.seek(7) #移动 cursor 到第 7Byte 的位置
10 #写入文件
11 with open('test2.txt', 'w') as f:
12     f.write('Test')
13     f.seek(0); f.write('Hello World!') #这里的内容会覆盖原有的内容
14
15 #读写非文本文件, 可以用二进制读写, 下面以创建一个图片副本为例
16 with open('puppy.jpg', 'rb') as rf:
17     with open('puppy.jpg', 'wb') as wf:
18         chunk_size = 4096
19         rf_chunk = rf.read(chunk_size)
20         while len(rf_chunk) > 0:
21             wf.write(rf_chunk)
22             rf_chunk = rf.read(chunk_size)
23

```

read() 和 write() 函数还有以下可选参数:

- offset: 移动字符个数
- whence: 指定从什么位置开始计算, 0 表示开头, 1 表示当前位置, 2 表示文件末尾 (只能在二进制模式下使用)

1.6.3 上下文管理器

上下文管理器就是可以在 with 语句中被调用的特殊函数，可以自动管理资源的分配和释放。利用 contextmanager decorator 可以将普通的函数转换成上下文管理器。

Listing 1.21: contextmanager

```
1 import contextlib
2
3 @contextlib.contextmanager
4 def open_file(file, mode): #自定义打开文件的函数
5     try:
6         f = open(file, mode) #yield 之前的语句将在进入 with 语块时执行
7         yield f
8     finally:
9         f.close() #yield 之后的语句将在退出 with 语块时执行
10
11 with open_file('test.txt', 'w'):
12     f.write('Hello, world!')
13
14 import os
15
16 @contextlib.contextmanager
17 def change_dir(destination): #自定义切换目录的函数
18     try:
19         cwd = os.getcwd()
20         os.chdir(destination)
21         yield
22     finally:
23         os.chdir(cwd)
24
```

1.6.4 处理.csv 文件

csv 是一种表格类型的文件，其每行各单元之间有分隔符，默认分隔符是‘,’。若单元格内部也有域分隔符相同的字符，该单元格会带引号。

Listing 1.22: read and write .csv files

```
1 import csv
2
3 with open('test.csv', 'r') as csv_file:
4     csv_reader = csv.reader(csv_file) #csv_reader 是一个迭代器
5     next(csv_reader) #跳过第一行
6     for line in csv_reader:
7         print(line) #这里每一行是一个 List
8
9     dic_reader = csv.DictReader(csv_file) #以字典形式读取
10    #若行数为 n，则得到一个长度为 n-1 的迭代器，第 1 行是 key，其余 n-1 行为 value
11
12 with open('new_names.csv', 'w') as new_file:
13     csv_writer = csv.writer(new_file, delimiter='\t') #这里设定了分隔符
14     csv_writer.writerow(list(range(5)))
15
16     dic_writer = csv.DictWriter(new_file, fieldnames=['number', 'name', 'score'])
17     #以字典形式写入，fieldnames 为 column 名，即第一行
18     csv_writer.writeheader() #写入第一行的 fieldnames
19     csv_writer.writerow({'number': 1, 'name': 'Jack', 'score': 98})
20
```

1.6.5 处理.json 文件

json 库可以用于处理.json 文件 (JavaScript Object Notation)。json 文件中的 object 会被处理称 dict，array 会被处理生 list，其他类型的变量也会对应处理。

Listing 1.23: json module

```

1 import json
2
3 data = json.loads(json_string) #将 json 文件的数据类型转换为 python 中的数据类型
4 new_json_string = json.dumps(data, indent=2, sort_key=true) #将 python 语句转换称 json 语句
5
6 with open('test.json') as f:
7     data = json.load(f) #读取 json 文件的内容
8 with open('new_test.json', 'w') as wf:
9     json.dump(data, f) #写入 json 文件
10 #可选:indent 设置缩进为 2, sort_key 设置排序所有键
11

```

1.6.6 处理 zip 文件

Listing 1.24: zip file

```

1 import zipfile
2
3 my_zip = zipfile.ZipFile('files.zip', 'w') #创建 zip file
4 my_zip.write('test.txt') #将已有的文件写入 zip file
5 my_zip.write('thumbnail.png')
6 my_zip.close()
7
8 with zipfile.ZipFile('files.zip', 'w') as my_zip: #使用 file manager
9     my_zip.write('test.txt')
10    my_zip.write('thumbnail.png')
11 with zipfile.ZipFile('files.zip', 'w', compression=zipfile.ZIP_DEFLATED) as my_zip: #压缩文件
12
13 with zipfile.ZipFile('files.zip', 'r') as my_zip: #压缩文件
14     print(my_zip.namelist()) #获取 zip file 中的文件表
15     my_zip.extractall('files') #提取文件到 files 文件夹下
16     my_zip.extract('thumbnail.png') #提取一个文件
17
18 import shutil
19
20 shutil.make_archive('another', 'zip', 'files') #打包 files 文件夹到 another.zip
21 shutil.unpack_archive('another.zip', 'another') #解压 another.zip 到 another 文件夹
22 #压缩文件可选 zip,tar,gztar,bztar,xztar
23

```

2 Collections

2.1 iterator

iterator 和 generator(带 yield 的函数) 有些类似, 区别在于内存的使用, iterator 会生成所有的值并存在内存里, generator 则按需生成值, 节省内存。例如, 当需要读取一个大文件中的前几行时, iterator 可以避免吧整个文件都读取进来。

Listing 2.1: iterator

```

1 i_nums = iter([1, 2, 3, 4, 5]) #从 iterable 生成 iterator
2 ite = map(lambda x: x * x, range(10), [i for i in range(20)])
3 ite = map(pow, range(10), itertools.repeat(2))
4 ite = filter(lambda x: x % 2 == false, range(10))
5 #这两个函数里第二个及更后面的参数可以是 iterator,range,tuple,list,string 等可以迭代的内容
6 #map 会对遍历过程中的每个元素进行 function 操作, filter 只会留下 condition 为 True 的元素
7 ite = zip('abcd', range(4), [i*i for i in range(4)])
8 #zip 返回的迭代器是各可循环容器的元素组成的元组, 遍历直到有容器被遍历完
9 ite = zip_longest('abcd', range(4), [i*i for i in range(4)])
10 #这里会一直遍历到所有 iterable 都遍历完, 用 None 补充没有的部分
11 print(next(result)); print(next(result)) #逐项查找
12 for value in ite:
13     print(value) #遍历迭代器

```

```

14 #注意, 迭代器在遍历后会变为空, 这一点与 list, range, tuple 不同。
15 #迭代器转为 list 等其他变量类型时会将其遍历一遍, 因此同样会变空。
16

```

2.1.1 itertools

itertools 可以处理 iterator。

Listing 2.2: itertools

```

1 counter = itertools.count(5) #返回一个从 5 到无穷的 iterator
2 counter = itertools.cycle([1, 2, 3]) #返回列表循环的 iterator
3 counter = itertools.repeat(2, times=3) #返回 2 重复 3 次的 iterator, times 默认无穷
4 counter = itertools.starmap(pow, [(0, 2), (1, 2), (2, 2)]) #与 map 类似
5 counter = itertools.combinations("abcd", 2) #两两无序组合为 tuple, 组成 iterator
6 counter = itertools.permutations("abcd", 2) #两两有序组合为 tuple, 组成 iterator
7 counter = itertools.combinations_with_replacement("abcd", 4) #允许 tuple 中元素重复
8 counter = itertools.permutations_with_replacement("abcd", 4) #允许 tuple 中元素重复
9 combined = itertools.chain("abcd", [1, 2, 3]) #拼接 iterables, 类似 zip
10 result = itertools.islice(range(10), 1, 5, 2) #从 index=1 迭代到 index=4, step=2
11 result = itertools.compress("abcd", [True, True, False, True]) #abd 组成的 iterator
12 result = itertools.filterfalse(lambda x: x%2 == 0, range(10)) #判断与 filter 相反
13
14 students = [
15     {
16         "name": "Jack",
17         "city": "Shanghai",
18         "group": 1,
19     },
20     {
21         "name": "Amy",
22         "city": "Beijing",
23         "group": 2,
24     },
25     {
26         "name": "Jerry",
27         "city": "Hangzhou",
28         "group": 1,
29     },
30     {
31         "name": "Jane",
32         "city": "Wuhan",
33         "group": 2,
34     }
35 ]
36 students.sort(key=lambda x: x["group"])
37 student_group = itertools.groupby(students, lambda x: x["group"])
38 for key, group in student_group: #这是一个二级 iterator
39     print(key) #group 名, 即 List 中 group 的值
40     for person in group:
41         print(person)
42

```

2.2 list

Listing 2.3: basic operations of list

```

1 #创建列表
2 empty_list = []; empty_list = list()
3 student1=["Hermione","Harry","Ron"]
4 student2=["Draco","Padma"]
5 student_score=list(range(10)) #[range(10)] 会创建一个只有 range(10) 一个元素的列表
6 student_score=[(letter, num) for letter in 'abcd' for num in range(4)]
7 score = [(letter, num) for letter, num in zip('abcd', range(4))] #注意这两个是不同的
8 student_score = list(map(lambda n: n*n, range(10)))
9 #对列表元素进行操作

```



```

10 student="Potter"
11 student1.append(student) #在列表的末尾追加
12 student1.insert(2, student) #在指定 index 处插入
13 student2.extend(student1) #列表的拼接, 这样可以比 + 更快
14 del student[-1] #删除最后一个元素
15 student_popped=student1.pop(1) #删除指定 index 处的元素, 默认删除最后一个
16 student_popped=student1.remove('Padma') #寻找第一个指定元素并删除, 找不到会报错
17 student_reversed=student1.reverse() #反转元素顺序
18 seq_Ron=student.index("Ron") #查找第一个匹配的元素并输出, 找不到会报错
19 #列表拼接、扩展、排序、遍历
20 student1=student1+student2
21 student3=sorted(student1); student1.sort() #sorted 是一个 function, sort 是一个 method
22 student3=sorted(student1, key=str.lower(), reverse=True) #不区分大小写倒序排序
23 #sorted 函数会对每个元素进行 key 的操作, 按照返回值的顺序, 对原数据进行排序
24 student1.sort(reverse=True) #True stands for down
25 if "Hermione" in student1:
26     print("Hermione is in student1")
27 for a in student1: #遍历列表
28     for index, item in enumerate(student1):
29         for index, student in enumerate(student1, start=1): #这里 index 可以从 1 开始编号
30             students = ', '.join(student1) #将 List 中的元素合并成一个字符串, 以逗号为分隔符
31 #列表统计
32 number=student1.count("Padma") #统计元素出现的数量
33 ind=student2.index("Padma") #找到元素首次出现的位置
34 min_one, max_one = min(student_score), max(student_score)
35 total=sum(student_score) #列表求和
36

```

2.3 tuple

元组是不可变的序列, 但可以重新赋值

Listing 2.4: basic operation of tuple

```

1 #元组创建
2 empty_tuple = (); empty_tuple = tuple() #创建空元组
3 a=tuple(range(10,20,3)); a=(); a=(1,2,3)
4 a_tup = sorted(a) #tuple 没有 sort method
5

```

2.4 set

Python 的 set 是无序集合, 这与 C++ 不同。

Listing 2.5: basic operation of set

```

1 students=[
2     {"name":"Hermione","house":"Gryffindor"},
3     {"name":"Harry","house":"Gryffindor"},
4     {"name":"Ron","house":"Gryffindor"},
5     {"name":"Draco","house":"Slytherin"},
6     {"name":"Padma","house":"Ravenclaw"},
7 ]
8 houses=set() #创建空集合, 注意 houses= 会创建一个字典
9 houses.pop() #删除第一个元素, 注意顺序是随机的
10 houses.clear() #清空集合
11 for student in students:
12     houses.add(student["house"])
13 for house in sorted(houses):
14     print(house)
15 if "Slytherin" in houses:
16     print("Slytherin is in houses.")
17

```

2.5 dict

Listing 2.6: basic operation of dictionary

```

1  #创建字典
2  students1={
3      "Harry":"Gryffindor",
4      "Ron":"Gryffindor",
5      "Draco":"Slytherin",
6      "Padma":"Ravenclaw",
7  }
8  name=["Harry","Ron","Draco","Padma"]
9  house=["Gryffindor","Gryffindor","Slytherin","Ravenclaw"]
10 student1=dict(zip(name,sign))
11 student1=dict(n: s for n, s in zip(name,sign))
12 student1=dict(((("Harry","Gryffindor"),("Ron","Gryffindor"),
13                ("Draco","Slytherin"),("Padma","Ravenclaw")))) #这三句的效果相同
14 student1=dict(Harry="Gryffindor",Ron="Gryffindor",Draco="Slytherin",Padma="Ravenclaw")
15 student_empty=dict.fromkeys(name) #创建值为空的字典
16
17 #访问, 排序, 删除, 加入, 遍历
18 ron_house=student1["Ron"] #访问字典, 找不到会报错
19 ron_house=student1.get("Ron") #访问字典, 找不到会返回 None, 可以通过这种方法避免报错
20 david_house=student1.get("David", "Not Found") #访问字典, 找不到返回 Not Found
21 student1["Hermione"]="Gryffindor" #加入元素
22 student1.update({"Hermione":"Gryffindor"}) #合并两字典
23
24 s_house = sorted(student1) #根据 key 来排序
25 ron_house=student1.pop("Ron") #删除元素并返回值
26 del student1["Ron"] #删除元素
27
28 for key,value in student1.items() #遍历字典, 注意 item 是函数, 要加 ()
29 for key in student1.keys(); for key in student1 #遍历键数组
30 for value in student1.values() #遍历值数组
31

```

2.6 String

2.6.1 编码解码

Listing 2.7: create a string

```

1  raw_string = r'\tTab' #这里的\t 不会被替换
2  name="Harry"; text='Hello world!'
3  byte=name.encode('GBK') #编码, 返回编码的 16 进制值
4  name1=byte.decode('GBK') #解码
5

```

2.6.2 常用操作

Listing 2.8: basic operations of string

```

1  length=len(name) #返回字符串长度
2  size=len(name.encode()) #返回内存大小
3  text.split(' ') #切分字符串
4  strnew=string.join(house) #合并字符串
5  sub='ab'
6  str_replace=sub.replace('a', 'c') #字符串替换
7  num=string.count(sub) #检索子字符串出现次数
8  start=string.find(sub) #子字符串首次出现的索引, 若没有出现, 返回-1
9  str_lower=str.lower() #转为小写, 不会改变原字符串
10 str_upper=str.upper() #转为大写
11 str_strip=str.strip(['@']) #删去字符串左右两侧的 ' ', \t, \r, \n 等 (默认), 以及 '@' (可选)
12 str_lstrip=str.lstrip(['@']) #删去字符串左侧的特定字符
13 str_rstrip=str.rstrip(['@']) #删去字符串右侧的特定字符

```

```

14 str_filled = str.ljust(10, fillchar=' ') #字符串左对齐, 右侧填充至 10 个字符
15 str_filled = str.rjust(10, fillchar=' ') #字符串右对齐
16 str_filled = str.center(10, fillchar='*') #字符串中间对齐
17

```

2.6.3 格式化字符串

格式化字符串也就是先制定一个模板，并在模板中以占位符为变量留下位置。一种实现方式是与 C 语言相似的占位符，更适合 python 的做法是使用 format 函数或者 f 字符串。

Listing 2.9: formatting string1

```

1 '%[-][+][0][m][.n]格式字符'%exp #这是模板
2 template = '编号: %09d\t 公司名称: %s\t官网: http://www.%s.com'
3 print(template%(7, '百度', 'baidu'))
4

```

可选项如下

- -: 左对齐，正数前面无负号，负数前面有负号
- +: 右对齐，正数前面加正号，负数前面加负号
- 0: 右对齐，正数前面无负号，负数前面加负号，用 0 补足位数
- m: 占有宽度
- .n: 小数位数

Table 2.1: placeholder

格式字符	说明	格式字符	说明
%s	字符串	%c	单个字符
%x	十六进制整数	%o	八进制整数
%f %F	浮点数	%d %i	十进制整数
%%	字符%	%e	指数 (基底为 e)

format 函数使用模板

Listing 2.10: format string

```

1 '{[index][:[fill][align][sign][#][width][.precision][type]]}'.format(item1, item2)
2

```

- index: 可选，表示该位置填入的是 format 的那个参数 (0-base)
- fill: 可选，填充字符
- align: 可选，< 左对齐，> 右对齐，^ 内容居中
- sign: 可选，+，-，用法与占位符选项相同
- #: 可选，进制前缀显示，也可以选 ',' 表示添千位分隔符
- width: 可选，字段宽度
- .precision: 可选，小数位数

- type: 可选, 指定类型

Table 2.2: type options

格式字符	说明	格式字符	说明
S	字符串类型 (默认)	b	十进制整数转为二进制
D	十进制整数	o	十进制整数转为八进制
C	按 ASCII 码转为字符	x X	十进制整数转为十六进制
e E	转为科学计数法	f F	转为浮点数, 默认 6 位小数
g G	自动切换	%	转为百分数

下面给出一些例子

Listing 2.11: examples of format string

```

1 number, name, url = 7, "百度", "baidu"
2 template='编号:{0>9s}\t公司名称:{s}\t官网:http://www.{s}.com'
3 context=template.format(number, name, url)
4 context_f=f'编号:{number}\t公司名称:{name}\t官网:http://www.{url}.com'
5
6 f_template='编号:{f_number}\t公司名称:{f_name}\t官网:http://www.{f_url}.com'
7 f_context = f_template.format(f_number = number, f_name = name, f_url = url)
8
9 context='编号:{2}\t公司名称:{0}\t官网:http://www.{1}.com'.format(name, url, number)
10
11 #下面假设有一个 company 变量, 他有三个成员, 分别是 name、url 和 number
12 context='编号:{0.number}\t公司名称:{0.name}\t官网:http://www.{0.url}.com'.format(company)
13
14 company = {'name': 'baidu', 'number': 7, 'url': 'baidu'}
15 context = '编号:{number}\t公司名称:{name}\t官网:http://www.{url}.com'.format(**company)
16 context = 'This function ran with args: {}, and kwargs: {}'.format(args, kwargs)
17 #这里 {} 中的内容可以是表达式或其他类型的数据 (如 List, dict)
18

```

2.7 collection

collection 包含了 dict, list, str 等 Container, 同时提供了一些额外的数据结构, 相当于 C++ 标准库。

Listing 2.12: namedtuple

```

1 import collection
2
3 #namedtuple 可以通过给 tuple 中每个变量加名称增加代码的可读性。
4 Color = collection.namedtuple('Color', ['red', 'green', 'blue'])
5 color = Color(55, 155, 255); print(color[0], color.red)
6
7 #deque 双端队列, 可以从两端添加和删除元素
8 collection.deque([1, 2, 3])
9
10 #Counter 计数器, 计算可迭代对象中元素的频率
11 c = collection.Counter(['a', 'b', 'c', 'a', 'b', 'b'])
12 c.update(['a', 'b']) #继续计数
13 c.most_common(5) #输出最多的 5 项, tuple 组成的 list
14
15 #OrderedDict 有序字典, 记住元素插入顺序
16 collection.OrderedDict([('a', 1), ('b', 2)])
17
18 #defaultdict 可以为字典中没有的元素添加默认值, 默认值类型为 int
19 dd = collection.defaultdict(int); print(dd['key']) #输出 0
20
21 #heapq 是优先级队列, 用最小堆实现
22 heap = collection.heapq.heapify([10, 17, 50, 7, 30, 24, 27, 45, 15, 5, 36, 21])

```

```

23 collection.heappush(heap, 13) #插入元素, heap 会以满二叉堆形式存储
24 print(heapq.heappop(heap)) #弹出堆顶元素并返回
25
26 c = Counter('hello world!'); c = Counter(['cat', 'dog'])
27 d = Counter({'red': 3, 'blue': 2}); d = Counter(cat=4, dogs=8)
28 print(c['cat'], c.items())
29 ite = c.elements() #返回一个迭代器
30 print(c.most_common(3)) #出现最多的三个
31 print(c.total()) #总数
32 c.subtract(d) #两个 counter 相减
33

```

2.8 SQLite

sqlite3 是 python 标准库中的数据库模块。sqlite3 中数据类型有 5 种, 分别是 NULL, INTEGER, REAL(float), TEXT, BLOB。

Listing 2.13: SQLite

```

1 import sqlite3
2
3 #创建并连接数据库
4 conn = sqlite3.connect(':memory:') #将数据存在 RAM 中, 每次关闭都会清空
5 conn = sqlite3.connect('my_database.db') #将数据存在文件中
6
7 c = conn.cursor() #创建指针, 数据库中的内容将由 cursor 完成
8
9 c.execute("""CREATE TABLE IF NOT EXISTS students( #创建名为students的表, 并创建3个column
10     first_name text, #不能创建已经存在的column否则会报错
11     last_name text,
12     id integer
13 )""") #IF NOT EXISTS 是可选的
14
15 #向表中添加数据
16 c.execute("INSERT INTO students VALUES ('Amy', 'Elisabeth', 11)")
17 c.execute("INSERT INTO students VALUES ('{}', '{}', {})".format(student1.first_name, student1.last_name, student1.id))
18 c.execute("INSERT INTO students VALUES (?, ?, ?)", #可以用一个 tuple 简化操作
19     (student1.first_name, student1.last_name, student1.id))
20 c.execute("""INSERT INTO students VALUES
21     (:first_name, :last_name, :pay)""", #可以用 dict 使操作更清晰
22     {'first_name': student1.first_name,
23      'last_name': student1.last_name,
24      'id': student1.id})
25
26 #从表中读取数据, * 表示提取全部
27 c.execute("SELECT * FROM students WHERE last_name='Elisabeth'")
28 c.execute("SELECT * FROM students WHERE last_name=?", ('Elisabeth',)) #1 个元素也要加 `,'
29 c.execute("SELECT * FROM students WHERE last_name=:last", {'last': 'Elisabeth'})
30 get_data = c.fetchone() #从 select 得到的内容中读取 1 行
31 get_data = c.fetchmany(5) #从 select 得到的内容中读取 5 行, 返回一个列表
32 get_data = c.fetchall() #从 select 得到的内容中读取所有行
33
34 conn.commit() #需要提交 cursor 的操作
35 conn.close() #关闭数据库
36
37 #with 语句可以让连接, 提交等操作更加简便, 下面给出两个例子
38 conn = sqlite3.connect('my_database.db') #将数据存在文件中
39 with conn:
40     c.execute("INSERT INTO students VALUES ('Amy', 'Elisabeth', 11)")
41     #离开 with 代码块时自动执行 conn.commit()
42
43 with sqlite3.connect('my_database.db') as conn:
44     c = conn.cursor()
45     c.execute("INSERT INTO students VALUES ('Amy', 'Elisabeth', 11)")
46     #离开 with 代码块时自动执行 conn.commit() 和 conn.close()
47

```

3 Class

3.1 定义类模板

3.1.1 基本结构

Listing 3.1: define a class

```

1  class User:
2      '''这是一个学生用户类''' #类的说明
3      #以下是静态成员变量, Python 中成为类的属性
4      student_user = 'Student User'
5      student_number = 0
6
7      def __init__(self, name, number): #构造函数
8          self.__name=name #这里名字设置为私有成员, 这里变量是非静态成员, python 称之为实例属性
9          self.number = number;
10         User.student_number += 1
11         #这里的 name, number 是非静态成员变量, Python 中成为实例的属性
12
13     def __str__(self): #输出运算的重载
14         return f'{self.number} {self.__name}'
15
16     @property
17     def name(self):
18         return self.__name
19
20     @name.setter
21     def name(self, name):
22         self.__name = name
23
24     def introduce(self): #self 类似 *this 指针
25         print("I'm ", self.__name, ", my number is ", self.number, sep='')
26
27     @staticmethod
28     def greet(input_student): #静态成员函数
29         print(f"Hello, my name's {self.__name}!")
30
31     @classmethod
32     def get_info(cls): #类方法, 其第一个变量是类, 常用于类型转换, 类属性操作等
33         print(f"There are {cls.student_number} {cls.student_user} in total.")
34

```

3.1.2 访问限制

Python 没有对属性和方法的访问权限进行限制。为了保证类内部某些属性不被外部访问, 可以在属性或方法名前(或前后)加上双下划线。

- `__foo__`: 首尾双下划线表示定义特殊方法, 一般是系统定义名字。
- `__foo`: 双下划线表示私有成员, 只允许所在的类调用。
- `_foo`: 单下划线表示保护成员, 即 C++ 中的 `protected` 成员。

3.1.3 属性

Python 中, 数据成员被称为属性。可以通过 `@property` 将一个方法转换为属性, 转换后可以直接通过方法名调用该方法, 无需添加 `()`。这样做可以简化代码, 也为属性添加安全保护, 即添加了 `@property` 的属性是只读的 (这是由于 `return` 时经过了复制传递)。

Listing 3.2: property of class

```

1  class User:
2      ''' 用户类 ''' #类的说明

```

```

3     @property
4     def name(self):
5         return self.__name #这个 name 属性被设定为只读的
6

```

@name.setter 可以在给属性 name 赋值时运行

Listing 3.3: setter of class

```

1     @name.setter #这里 name 是属性名称
2     def name(self, name):
3         if len(name) < 20:
4             self.__name=name
5

```

@name.deleter 属性可以定义一个在删除 name 属性时执行的函数。

Listing 3.4: deleter

```

1     @name.deleter
2     def name(self):
3         print('Delete Name!')
4         self.__name = None
5     del student1.__name
6

```

3.1.4 类的特殊成员 (魔术方法)

Listing 3.5: Special members for class

```

1     print(User.__class__) #类的名称
2     print(User.__bases__) #类的基类
3     print(User.__dict__) #类的数据成员及对应的值，输出一个字典
4
5     def __new__(cls, *args, **kwargs): #创建类的实例
6         print(f"Run new with: cls={cls}, args={args}, kwargs={kwargs}")
7         return super().__new__(cls) #调用父类的 __new__ 方法，并传入当前类作为参数
8
9     def __del__(self): #对象删除时调用的方法
10
11     def __init__(self, name, number): #接收并初始化实例
12         self.__name=name #这里名字设置为私有成员，这里的变量是非静态成员，python 中称之为实例属性
13         self.number = number;
14         User.student_number += 1
15         #这里的 name, number 是非静态成员变量，Python 中成为实例的属性
16         #Python 中，一个类函数只能有一个 __init__ 函数，若需要多个构造函数，可以用 classmethod 实现
17         #在实例创建时，会先调用 __new__ 在调用 __init__
18
19     def __str__(self): #print() 和 str() 的重载
20         return f'{self.number} {self.__name}'
21     def __repr__(self): #返回一个 Python 合法的字符串，使之可以用 eval 重建
22         return(f"User({self.__name}, {self.number})")
23         #The goal of __repr__ is to be unambiguous
24         #The goal of __str__ is to be readable
25     student1 = User("Amy", 11); repr_str = repr(student1)
26     copy_student1 = eval(repr_str)
27     def __format__(self, format_spec): #重载 str.format() 和 f-string
28
29     def __bool__(self): #bool() 的重载
30
31     def __add__(self, other): #重载 +
32     def __iadd__(self, other): #重载 +=
33     def __sub__(self, other): #重载 -
34     def __isub__(self, other): #重载 -=
35     def __mul__(self, other): #重载 *
36     def __imul__(self, other): #重载 *=
37     def __truediv__(self, other): #重载 /
38     def __itruediv__(self, other): #重载 /=

```

```

39 def __floordiv__(self, other): #重载//
40 def __mod__(self, other): #重载%
41 def __pow__(self, other): #重载
42 def __eq__(self, other): #重载 ==
43 def __ne__(self, other): #重载 !=
44 def __lt__(self, other): #重载 <
45 def __le__(self, other): #重载 <=
46 def __gt__(self, other): #重载 >
47 def __ge__(self, other): #重载 >=
48
49 def __len__(self): #重载 len()
50 def __getitem__(self, key): #重载 [], 读取值
51 def __setitem__(self, key, val): #重载 []=, 修改值
52 def __delitem__(self, key): #重载 del [], 删除值
53 def __iter__(self): #返回迭代器, 重载 iter()
54 def __next__(self): #返回迭代器的下一个值, 重载 next()
55
56 def __call__(self, x): #允许实例像 function 一样被调用
57     return self.__name + x
58 #既然 class 可以像 function 一样被调用, 那 class 也可以成为decorator
59 class decorator_class:
60     def __init__(self, original_function):
61         self.original_function = original_function
62
63     def __call__(self, *args, **kwargs):
64         print('call method executed before {}'.format(self.original_function.__name__))
65         return self.original_function(*args, **kwargs)
66 @decorator_class
67 def display():
68     print('display function ran')
69
70 def __enter__(self): #进入 with 语句时自动调用
71 def __exit__(self, exc_type, exc_value, traceback):
72 #离开 with 语句时自动调用 (正常和非正常退出都会调用)
73 #三个参数分别是异常类型、异常值和异常回溯信息, 返回 True 会一直异常, 否则异常会继续传播
74

```

3.2 类的使用

3.2.1 创建类的实例

Listing 3.6: Create an instance of class

```

1 student1 = User('Amy', 11)
2

```

3.2.2 类成员的使用

Listing 3.7: Using class members

```

1 print(student.student_number); print(User.student_number) #都输出 1
2 print(student.number); student1.number = 1 #这里 number 无法修改
3
4 student1.introduce() #输出 I'm Amy, my number is 11
5 User.introduce(student1) #效果与上一句相同
6
7 User.greet(student1) #输出 Hello, my name's Amy!
8 student1.greet(student1) #效果与上一句相同
9
10 User.get_info() #输出 There are 1 Student User in total.
11 student1.get_info() #效果与上一句相同
12 print(student1) #这里会调用 __str__ 函数, 输出 11 Amy
13

```


3.3 继承

Python 的继承默认是公有继承。继承的内容有类的属性，实例方法，类方法，静态方法，特殊方法，私有成员不会被继承。类方法在继承时，会以子类作为第一个参数。

Python 中的所有类都继承自 `builtins.object` 类。

Listing 3.8: Class inheritance

```

1 class Student(User)
2     def __init__(self, grade, name, number):
3         super().__init__(name, number) #调用基类的构造函数
4         self.__grade = grade
5

```

子类可以重写父类中的方法，这与 C++，Java 没什么区别，只是 Python 没有虚函数机制。

3.4 有关类的内置函数

由于 python 中的类在访问和添加成员的操作上限制较少，可以调用一些内置函数来检查这些类的合理性。

Listing 3.9: built-in functions about class

```

1 isinstance(object, classinfo) #检查 object 的类别是不是 classinfo
2 isinstance(student1, User) #返回 True
3
4 issubclass(cls, classinfo) #检查 cls 类是不是 classinfo 的子类
5 issubclass(Student, User) #返回 True
6
7 hasattr(object, name) #检查 object 是否有 name 属性
8 hasattr(student1, "number") #返回 True
9
10 callable(object) #检查 object 是否可以像函数一样被调用
11 callable(student1.name) #返回 False
12 callable(student1.greet) #返回 True
13

```

3.5 类的典型例子

3.5.1 文件管理类

文件管理类可以实现打开创建文件，读写文件，文件异常处理等工作

Listing 3.10: File management class

```

1 class Open_File():
2
3     def __init__(self, filepath, mode):
4         self.filepath = filepath
5         self.mode = mode
6
7     def __enter__(self):
8         print(f"Entering the file: {self.filepath}")
9         self.file = open(self.filepath, self.mode)
10        return self.file
11
12    def __exit__(self, exc_type, exc_val, traceback):
13        print(f"Exiting the file: {self.filepath}")
14        if exc_type:
15            print(f"An exception occurred: {exc_type}, {exc_value}")
16        self.file.close()
17
18    with Open_File('test.txt', 'w') as f:
19        f.write('Testing')
20

```

4 Basic module

4.1 install and import module

4.1.1 install and manage modules

可以通过 pip 安装和管理所需的 module。

Listing 4.1: download modules

```

1 pip install pandas #安装 module
2 pip uninstall pandas #卸载安装的 module
3 pip list #列出所有已安装的 module, 可选-o 查看是否是最新版本
4 pip install -U pandas #更新 module
5 pip freeze --local | grep -v '^-\e' | cut -d = -f 1 | xargs -n1 pip install -U
6 pip search pandas #搜索指定的 module
7 pip freeze > requirements.txt #将项目所需的 module 整理成文档
8 pip install -r requirements.txt #将文件中的 module 全部安装下来
9

```

4.1.2 virtualenv

virtualenv 是指项目运行的 module 环境，这个环境可以不包含 global 环境中的所有 module。

Listing 4.2: setup virtualenv

```

1 pip install virtualenv
2 mkdir Environments
3 cd Environments
4 virtualenv project1_env -p /usr/bin/python3.10 #创建 virtualenv 环境, python 版本可选
5 source project1_env/bin/activate #activate virtualenv
6 pip install numpy #这里下载所需的 module 即可
7 pip install -r requirements.txt #由文件下载所有 module
8 pip freeze --local > requirements.txt #将 virtualenv 中的 module 整理成文档
9 deactivate #退出 virtualenv
10 rm -rf project1_env #删除已创建的 virtualenv 环境
11

```

4.1.3 pipenv

管理环境是很有必要的，若所有项目都在 base 环境中运行，package 的更新可能导致项目不可用。

pipenv 结合了 pip 和 virtualenv，可以用来管理代码运行环境。pipenv 由 packages(生产环境运行必须有的包)，dev-packages(开发环境所需的包)

Pipfile 列出了基本信息，Pipfile.lock 是决定性的，列出了具体信息。

Listing 4.3: pipenv

```

1 pipenv install requests #向 package 中添加 module, 环境保存在 Pipfile 中
2 pipenv install pytest --dev #向 dev 中添加 module
3 pipenv uninstall requests #删除 module
4
5 pipenv install #安装 Pipfile 中的所有 module
6 pipenv install --ignore-pipfile #安装 Pipfile.lock 中的所有 module
7 pipenv install -r requirement.txt #利用文件加载环境
8 pipenv graph #列出环境内容, 包括依赖关系
9 pipenv lock -r > requirements.txt #整理环境内容
10 pipenv lock #更新 pip.lock
11
12 pipenv --python 3.6 #修改 python 版本, 直接修改 Pipfile 不会修改环境
13 pipenv --rm #删除环境, 不会删除 Pipfile
14 pipenv install #由 Pipfile 重建环境
15 pipenv --venv #查看虚拟环境目录
16 pipenv check #检查是否能更新, Pipfile 和环境是否对应

```

```

17 pipenv shell #激活环境
18 pipenv run python #在环境里运行指令 (打开 python)
19 exit #退出环境
20
21 touch .env #这个文件可以设置独属于这个环境的环境变量
22
23

```

Listing 4.4: .env

```

1 SECRET_KEY="MySuperSecretKey" #局部环境变量
2

```

4.1.4 anaconda

anaconda 也可以用来管理 module 和环境, 其效果相当于 pip+virtualenv。anaconda 的优势在于可以安装不属于 python 的 module, 且可提供图形化管理功能。

注意, anaconda 和 pipenv 有冲突, 不能一起使用。

Listing 4.5: anaconda

```

1 conda create --name my_app python=2.7 flask sqlalchemy
2 #创建一个名为 my_app 的项目, 项目环境中带有 flask 和 sqlalchemy, python 版本为 2.7
3 conda activate my_app #激活环境, 默认环境为 base
4 conda deactivate #退出环境
5 conda env list #列出所有已创建的环境
6 conda remove --name my_app --all #删除已有的环境
7 conda env export > environment.yaml #导出 environment 的内容
8 conda env export create -f environment.yaml #通过文件创建环境
9

```

有时我们需要记录环境与对应项目的目录, 这是我们可以环境目录 (可以通过 `conda env list` 查看) `etc/conda/activate.d/env_vars.sh`, `etc/conda/deactivate.d/env_vars.sh`, 这两个文件分别会在 `activate` 和 `deactivate` 时自动运行。

Listing 4.6: activate.d/env_vars.sh

```

1 #!/bin/sh
2 export DATABASE_URI="postgresql://user:pass@db_server:5432/test_db"
3 #just add anything you need
4

```

Listing 4.7: deactivate.d/env_vars.sh

```

1 #!/bin/sh
2 unset DATABASE_URI
3

```

可以通过 `~/.bashrc` 中加入以下函数以自动启用文件夹中的 `environment.yaml`。

Listing 4.8: conda_auto_env

```

1 function conda_auto_env() {
2     if [ -e "environment.yaml" ]; then
3         ENV_NAME=$(head -n 1 environment.yaml | cut -f2 -d " ")
4         #Check if you are already in the environment
5         if [[ $CONDA_PREFIX != *$ENV_NAME* ]]; then
6             #Try to activate environment
7             conda activate $ENV_NAME &>/dev/null
8         fi
9     fi
10 }
11
12 #export PROMPT_COMMAND="conda_auto_env;$PROMPT_COMMAND"
13 #PROMPT_COMMAND 每次按下 enter 时都会运行

```

```

14 #若不需要，可以注释掉最后一行，并使用 conda_auto_env 手动调用环境
15

```

4.1.5 import modules

Listing 4.9: import module

```

1 import math
2 import math as ma
3 from math import sqrt
4 from math import *
5

```

4.1.6 view the functions in modules

Listing 4.10: view the functions in builtin module

```

1 import builtins
2 print(dir(builtins)) #这个查看方法对所有 module 都有效
3

```

4.1.7 circular import

circular import 是一个常见的错误。当 main.py import B.py 的其中部分时，B.py 整个文件都会被运行，若 B.py 同样 import main，那么 A.py 也会从头开始运行整个文件，而这时 main.py 会找不到其 import B.py 中的部分，这时就触发了 circular import。

4.2 builtins module

最基础的 module，包含了 print len range abs 等常用函数，TypeError ValueError KeyError IndexError 等常见异常，True False None 等内置常量，int str list dict set 等内置类型，不需要导入即可使用。

4.3 sys module

有关系统操作的 module

Listing 4.11: sys module

```

1 print(sys.path)
2 #系统路径列表，可以通过 append 追加 module 所在地址
3 #也可以追加在 ~/.bashrc 的 PYTHONPATH 里面（作为环境变量）
4 #如果需要 import 的 module 是一个文件夹，需要保证文件夹中有一个 '__init__.py' 文件
5 # '__init__.py' 文件会在 import 时运行，该文件可以为空，也可以导入文件夹中的子模块
6 print(sys.executable) #输出 python 的文件位置
7 print(sys.version) #输出 python 版本，用来检验编译器
8

```

4.4 os module

os module 可以实现操作系统相关的功能。

Listing 4.12: os module

```

1 #文件管理
2 os.getcwd() #输出工作区目录
3 os.chdir('/mnt/d/Desktop') #移动到指定目录

```

```

4 os.listdir() #输出当前目录下的所有文件及文件夹
5 os.mkdir('os_test') #在当前目录下新建, 注意只能新建一层
6 os.makedirs('os_test/my_profile') #新建目录, 允许多层
7 os.rmdir('os_test'); os.removedirs('os_test/myprofile') #删除目录
8 os.rename('test.txt', 'demo.txt') #重命名文件
9 print(os.stat('demo.txt')) #查看文件属性, 如文件大小 (st_size), 最后一次修改时间 (st_mtime)
10 for dirpath, dirnames, filenames in os.walk('/mnt/d/Desktop/python') #递归的遍历目录下所有文件
11     if f.endswith('.py'): #筛出所有 python 脚本文件
12         #dirpath 表示当前所在的目录, dirname 表示当前目录下的文件夹, filenames 表示当前目录下的文件名

13
14 os.environ.get('HOME') #返回 home 的地址, 这里还可以查询其他环境变量, 如 PATH
15 os.path.join(os.environ.get('HOME'), 'test.txt') #拼接目录, 直接字符串相加容易遗漏或多加slash
16 os.path.exists('/mnt/d/test.txt') #返回一个bool值, 即目录是否存在
17 os.path.basename('/mnt/d/test.txt') #返回`test.txt`这里无论路径是否存在都不会报错
18 os.path.dirname('/mnt/d/test.txt') #返回 `/mnt/d`
19 os.path.split('/mnt/d/test.txt') #返回 ('/mnt/d', 'test.txt')
20 os.path.splitext('/mnt/d/test.txt') #返回 ('/mnt/d/test', '.txt')
21

```

4.5 subprocess module

subprocess 可以用于在终端中执行任务。

Listing 4.13: subprocess module

```

1 subprocess.run('ls') #在 terminal 中运行 ls, 若运行失败不会 abort
2 subprocess.run('ls', check=True) #在 terminal 中运行 ls, 若运行失败会 abort
3 subprocess.run('ls', stderr=subprocess.DEVNULL) #在 terminal 中运行 ls, 忽略所有 error
4
5 p1 = subprocess.run(['ls', '-la']) #捕获运行的结果
6 p1.args; p1.returncode; p1.stderr #指令, 返回值 (0 表示没有出错, 1 表示出错), 标准错误流
7
8 p1 = subprocess.run(['ls', '-la'], capture_output=True) #捕获运行的输出 (字节流)
9 print(p1.stdout.decode()) #返回标准输出
10 p1 = subprocess.run(['ls', '-la'], capture_output=True, text=True) #捕获运行的输出 (字符串)
11 print(p1.stdout); #返回标准输出
12
13 with open('output.txt', 'w') as f:
14     p1 = subprocess.run(['ls', '-la'], stdout=f, text=True) #输出到文件中
15
16 p1 = subprocess.run(['cat', 'test.txt'], capture_output=True, text=True)
17 p2 = subprocess.run(['grep', '-n', 'test'], capture_output=True, text=True, input=p1.stdout)
18 #传递输出和输入值, 逐步完成 process
19

```

4.6 time module

time module 可以用于计时

Listing 4.14: time module

```

1 time.sleep(1) #将进程挂起 1s
2 t = time.time() #返回 linux 时间戳, 单位都是秒
3 t1 = time.perf_counter(); t2 = time.perf_counter() #高精度时间计数器
4 t3 = time.process_time(); t4 = time.process_time() #高精度的 CPU 时间
5 print(t2 - t1) #计算时间差, 会包含 sleep() 的时间
6 print(t4 - t3) #计算时间差, 不会包含 sleep() 的时间
7

```

Listing 4.15: Timer Class

```

1 class Timer: #@save
2     """记录多次运行时间"""
3     def __init__(self):
4         self.times = []
5         self.start()

```

```

6
7     def start(self):
8         """启动计时器"""
9         self.tik = time.time()
10
11    def stop(self):
12        """停止计时器并将时间记录在列表中"""
13        self.times.append(time.time() - self.tik)
14        return self.times[-1]
15
16    def avg(self):
17        """返回平均时间"""
18        return sum(self.times) / len(self.times)
19
20    def sum(self):
21        """返回时间总和"""
22        return sum(self.times)
23
24    def cumsum(self):
25        """返回累计时间"""
26        return np.array(self.times).cumsum().tolist()
27

```

4.7 threading module

4.7.1 threading

threading module 用于处理多线程。

Listing 4.16: threading module

```

1 my_function(sec):
2     time.sleep(sec)
3     return "accomplished"
4 t1 = threading.Thread(target=my_function, args=(1,)) #创建一个线程
5 t2 = threading.Thread(target=my_function, args=(2,)) #创建一个线程
6 t1.start(); t2.start(); #启动线程, 同时运行余下部分 (3 个线程)
7 t1.join(); t2.join(); #待 t1, t2 线程结束后再继续运行余下部分 (2 个线程)
8

```

4.7.2 concurrent.futures

Listing 4.17: concurrent.futures module

```

1 import concurrent.futures
2
3 with concurrent.futures.ThreadPoolExecutor() as executor:
4     f1 = executor.submit(my_function, 1)
5     f2 = executor.submit(my_function, 2) #这两个线程是并行的
6     print(f1.result()); print(f2.result())
7
8     results = [executor.submit(my_function, i) for i in range(10)]
9     for f in concurrent.futures.as_completed(results):
10         print(f.result()) #根据完成先后打印
11
12     results = executor.map(my_function, [i for i in range(5)])
13     for result in results:
14         print(result) #根据提交先后打印
15

```

4.7.3 multiprocessing

Listing 4.18: multiprocessing module

```

1 processes = []
2 for _ in range(i):
3     p = multiprocessing.Process(target=my_function, args=(i,))
4     p.start()
5     processes.append(p)
6 for process in processes:
7     process.join()
8

```

4.8 logging module

logging module 是 python 标准库中用于记录日志的模块，它可以在程序运行时输出各种级别的日志。

Listing 4.19: logging module

```

1 #basicConfig 用于设置日志的基本信息，这会设置 root logger，影响所有的 logger
2 logging.basicConfig(filename='app.log', filemode='w',
3                     level=logging.DEBUG, #默认是 WARNING
4                     format='%(asctime)s - %(name)s - %(levelname)s - %(message)s',
5                     datefmt='%Y-%m-%d %H:%M:%S'
6                     #format 的更多选项见 LogRecord attributes
7                     )
8 #注意，多次设置 root logger 只有第一个生效，这在导入文件时要特别注意
9
10 #可以使用 handlers 将日志同时输出到多个文件或输出到控制台
11 logging.basicConfig(
12     level=logging.INFO,
13     format='%(asctime)s - %(levelname)s - %(message)s',
14     handlers=[
15         logging.FileHandler("logfile.log"), # 文件处理器
16         logging.StreamHandler() # 控制台处理器
17     ]
18 )
19
20 #为了让每个文件有单独的 logging 文件和等级，可以在每个文件里单独设置 logger
21 logger = logging.getLogger(__name__) #一般以 __name__ 作为 logger 的命名
22 logger.setLevel(logging.INFO) #设置 logger 的 level
23 file_handler = logging.FileHandler('test.log') #设置日志文件和格式
24 file_handler.setFormatter(logging.Formatter('%(levelname)s: %(name)s: %(message)s'))
25 logger.addHandler(file_handler) #将 logger 的内容输出到文件中
26 string_handler = logging.StreamHandler()
27 string_handler.setFormatter(logging.Formatter('%(levelname)s: %(name)s: %(message)s'))
28 logger.addHandler(string_handler) #将 logger 的内容输出到终端
29
30 #以下信息将会以上面设定的形式存到 app.log 文件中，只有 level 以上级别会显示出来
31 logging.debug("这是调试信息")
32 logging.info("这是一般信息")
33 logging.warning("这是警告信息")
34 logging.error("这是错误信息")
35 logging.critical("这是严重错误信息")
36
37 #若为每个文件单独设置了 logger，则可以这样调用
38 logger.info("这是一般信息")
39
40 #error 等级的错误可以添加报错信息，即 traceback
41 logging.exception("这是一条带有 Traceback 的错误信息")
42

```

4.9 re module

4.9.1 regular expression

正则表达式由元字符，限定符，选择字符，排除字符等组成。

元字符中的 ‘.’ ‘\$’ ‘^’, 限定符中的 ‘?’ ‘+’ ‘*’ 以及 ‘\’ 若需要匹配, 则应该使用转义字符。

Table 4.1: meta-character

元字符	匹配对象
.	换行符以外的任意字符 $[\backslash n \backslash r]$
$\backslash w$	字母、数字、下划线或汉字
$\backslash W$	非字母、数字、下划线、汉字
$\backslash s$	任意空白字符 $([\backslash f \backslash n \backslash r \backslash t \backslash v])$
$\backslash S$	任意非空白字符
$\backslash d$	数字字符
$\backslash D$	非数字字符
$\backslash b$	单词的边界 (开始或结束)
$\backslash B$	非单词的边界
\wedge	字符串的开始
$\$$	匹配字符串的结束

Table 4.2: qualifier

限定符	含义
?	匹配 0-1 次
+	匹配至少 1 次, 会贪婪匹配
*	匹配 0 次或多次, 会贪婪匹配
$+? *?$	匹配规则同上, 但是非贪婪
$\{n\}$	匹配 n 次
$\{n,\}$	匹配至少 n 次
$\{n,m\}$	匹配至少 n 次, 至多 m 次
	匹配左右其一

- 方括号表示字符集合, 如 $[aeiou]$ 匹配元音字符, $[a-z]$ 匹配小写字符。 $[\text{u4e00}-\text{u9fa5}]$ 匹配汉字。
- 方括号内的 \wedge 表示排除字符, 如 $[\wedge aeiou]$ 匹配非元音字母的所有字符。
- 方括号中的无法使用限定符, 限定符都判定为原来的字符。
- 小括号可以改变限定符的作用范围, 如 $(\text{thir}|\text{four})\text{th}$ 相当于 $\text{thir}|\text{fourth}$, $(\backslash.[0-9]\{1,3\})\{3\}$ 会将括号中的内容重复三次。
- 小括号还可以基于匹配模式从字符串中提取子字符串, 组成一个元组。

Python 中使用正则表达式需要将部分\转义, 由于需要转义的\可能很多, 可以使用模式字符串, 即在字符串前加上 r 或 R, 例如: `r'\bm\w*\b'`。

4.9.2 match string

Python 中的 re 模块可以通过正则表达式处理字符串。

下面是常用的匹配函数, 这里匹配的字符串区间没有重叠的部分。

Listing 4.20: string matching

```
import re #导入模块
```



```

3 pattern=r'mr_\w+' #模式字符串
4 string='MR_SHOP mr_shop' #待匹配字符串
5
6 #match 方法可以从字符串开始处开始匹配, 若匹配失败, 返回 None
7 match=re.match(pattern,string,re.I) #不区分大小写匹配字符串, 返回一个 match 对象
8 print(match.start());print(match.end()) #输出匹配字符串起始位置
9 print(match.span()) #输出匹配字符串起止位置元组
10 print(match.string) #输出匹配前的字符串, 即 string
11 print(match.group()) #输出匹配的数据
12
13 #search 方法可以搜索字符串中第一个可以匹配的子串, 返回一个 match 对象
14 match=re.search(pattern,string,re.I) #不区分大小写匹配字符串, 返回一个 match 对象
15
16 #findall 和 finditer 方法可以搜索字符串中所有可以匹配的子串
17 match = re.findall(pattern,string,re.I) #不区分大小写匹配字符串, 返回一个 match 对象
18 print(match) #输出匹配子串的列表
19 match = re.finditer(pattern, string, re.I) #返回一个由 match 对象组成的迭代器
20
21 #上面的代码可以通过编译模式字符串进行简化
22 pattern = re.compile(r'mr_\w+', re.I)
23 match = pattern.match(string)
24 match = pattern.match(string)
25 match = pattern.findall(string)
26 match = pattern.finditer(string)
27
28 #捕获组的使用
29 pattern = re.compile(r'(\w+)-(\d+)-(\w+)')
30 text = "abc-123-def"
31 match = pattern.match(text) #这里使用方法与上面相同
32 #match.group(0) 为 'abc-123-def', match.group(i) 为第 i 个捕获组捕获的内容
33 match = pattern.findall(text) #match = [('abc', '123', 'def')]
34

```

Table 4.3: matching pattern

标志	含义
A ASCII	只匹配 ASCII 范围内的字符
I IGNORECASE	不区分大小写
M MULTILINE	将 ^ 和 \$ 用于每一行的开头和结尾
S DOTALL	. 匹配所有字符, 包括换行符
X VERBOSE	忽略未转义的空格和注释

4.9.3 substitute string

sub 方法可以实现 vim 中批量搜索并替换字符串的操作。

模板为: `re.sub(pattern,repl,string,count,flags)`

- pattern: 模式字符串
- repl: 替换后的子字符串
- string: 原始字符串
- count: 可选, 最大替换次数, 默认为 0, 表示替换所有的匹配
- flags: 可选, 见 Table 4.3

Listing 4.21: string substitution

```

1 import re
2
3 pattern = r'1[34578]\d{9}' #模式字符串

```

```

4 string='电话号码是:13611111111'
5 result=re.sub(pattern,'1xxxxxxxxx',string)
6
7 pattern = re.compile(r'(1[34578]\d)\d{8}')
8 string = 'The phone numbers are:13611111111,15888888888,18333333333'
9 result = pattern.sub(r'\1*****', string) #这里\1 表示第一个捕获组
10 result = 'The phone numbers are:136*****,158*****,183*****'
11

```

4.9.4 split string

`re.split` 方法可以根据正则表达式切分字符串，匹配正则表达式的子串将被当作分隔符，并将分割结果以列表的形式返回。

模板为：`re.split(pattern,string,[maxsplit],[flags])`

- pattern: 模式字符串
- string: 待切分的字符串
- maxsplit: 可选，最大拆分次数
- flags: 可选，见 Table 4.3

Listing 4.22: string segmentation

```

1 import re
2
3 pattern=r'[?&]'
4 url='http://www.mingrisoft.com/login.jsp?username="mr"&pwd="mrsoft"'
5 result=re.split(pattern,url)
6 print(result)
7

```

4.10 date module

4.10.1 datetime module

datetime 用于时间和日期戳记录和转换

Listing 4.23: datetime & pytz

```

1 d = datetime.date(2020, 7, 24); print(d) #这里的 day 是一个 date 类型的变量
2 td = datetime.date.today() #输出今天的日期
3 print(d.year); print(d.month); print(d.day)
4 print(d.weekday()); print(tday.isoweekday())
5 #in weekday: Monday 0 Sunday 6, in isoweekday: Monday 1 Sunday 7
6 tdelta = datetime.timedelta(days = 7) #这里 tdelta 是 timedelta 类型的，可以参与加减法
7 date2 = date1 + timedelta; timedelta = date1 - date2
8 tdelta.total_seconds() #时间差转秒
9
10 t = datetime.time(9, 30, 45, 100000) #四个参数分别是 h m s μs
11 print(t.hour); print(t.minute); print(t.second); print(t.microsecond)
12 print(t.isoformat()) #以国际标准形式输出时间
13 print(t.strftime('%B %d, %Y')) #以自定义形式输出时间，这里可以参考 strftime 官方文档
14
15 dt = datetime.datetime.strptime('July 26, 2016', '%B %d, %Y') #根据模板由 string 转 datetime
16 dt = datetime.datetime(2020, 7, 8, 12, 10, 32, 100000, tz = pytz.UTC) #时区可选
17 t = dt.time(); d = dt.date(); print(dt.year)
18 sentence = f'Jack has a birthday on {birthday:%B %d, %Y}' #结合 fstring 使用
19

```

4.10.2 pytz module

pytz 用于管理时区 (time zone) 相关的信息。

Listing 4.24: pytz module

```

1 dt_today = datetime.datetime.today() #返回当前时区的时间
2 dt_now = datetime.datetime.now(tz = pytz.UTC) #返回指定时区的时间，默认没有时区信息
3 dt_utcnow = datetime.datetime.utcnow() #返回 utc 时间，不建议使用
4 dt_mtn = dt_now.astimezone(pytz.timezone('Asia/Shanghai'))
5 #需要带有时区参数的 datetime 变量才可以进行时区转换
6 for tz in pytz.all_timezones:
7     print(tz) #输出所有的时区名称
8 dt_mtn = pytz.timezone('America/New_York').localize(dt_today) #为时间添加时区
9 dt_time = datetime.datetime.fromtimestamp(mod_time) #将时间戳转换成 datetime 类型
10

```

4.10.3 calendar module

calendar 用于计算天数，星期数，星期几等数据。

Listing 4.25: calendar module

```

1 today = datetime.date.today()
2 days_in_current_month = calendar.monthrange(today.year, today.month)
3 #输出一个 tuple，第一个元素表示今天星期几 Monday0 Sunday6，第二个表示本月有几天
4

```

4.11 PIL module

PIL module 可以处理图片，module 中的具体内容见 [PIL 官方文档](#)。

Listing 4.26: PIL

```

1 from PIL import Image, ImageFilter
2
3 image1 = Image.open('pup1.jpg')
4 image1.save('pup1.png')
5 image1.show() #显示图片
6
7 image1.thumbnail((300, 300)) #压缩至 300*300 像素
8 image1.rotate(90) #逆时针旋转 90°
9 image1.convert(mode='L') #图片变黑白
10
11 image1.filter(ImageFilter.GaussianBlur(15)) #模糊图片，模糊程度可选
12

```

4.12 random module

Listing 4.27: random module

```

1 value = random.random() #生成 [0.0, 1.0) 的随机浮点数
2 value = random.uniform(1, 10) #生成 [1.0, 10.0] 的随机浮点数
3 value = random.randint(1, 10) #生成 [1, 10] 的随机整数
4 value = random.choice(['Red', 'Green', 'Blue']) #从 list 中随机取出一个
5 results = random.choices(['Red', 'Green', 'Blue'], k=10) #随机取 k 遍，输出一个 list
6 results = random.choices(['Red', 'Green', 'Blue'], weight=[18, 18, 2], k=10)
7 shuffled_list = list(range(2, 30)); random.shuffle(shuffled_list) #随机打乱 list
8 hand = random.sample(list(range(50)), k=5) #取 5 个样本，5 个样本不会重复
9 random.shuffle(my_list) #直接打乱列表，没有返回值
10

```

4.13 secrets module

secrets 可以生成随机长密码。

Listing 4.28: secrets module

```
1 import secrets
2 print(secrets.token_hex(16)) #16Bytes
3
```

4.14 cn2an module

cn2an module 可以实现中文数字到阿拉伯数字的转换。

Listing 4.29: cn2an

```
1 def chinese_to_int(chinese_str):
2     arabic_number = cn2an.cn2an(chinese_str, 'normal')
3     return arabic_number
4
```

4.15 Progress Bar

4.15.1 tqdm module

Listing 4.30: tqdm module

```
1 from tqdm import tqdm
2 import time
3
4 for i in tqdm(range(500), desc="Processing"):
5     time.sleep(0.01)
6
7 progress_bar = tqdm(total=2000, desc="Processing")
8 for i in range(200):
9     for j in range(10):
10         progress_bar.update(1)
11
```

5 Maths & Statistics

5.1 math

math module 提供了最基础的数学函数。

Listing 5.1: math module

```
1 math.ceil(val) #向上取整
2 math.floor(val) #向下取整
3 math.trunc(val) #向零取整
4 math.fabs(val) #取绝对值, 只用于浮点数, abs 通用于整数、浮点数、复数等
5
6 math.sqrt(val) #开平方根
7 math.pow(x, y) #计算 `x` 的 `y` 次方
8 math.log(x, base); math.log10(x); math.log2(x) #计算对数, 默认 base 为 e
9 math.e; math.exp(x) #得到 e 和 e 的乘方
10 math.sin(x); math.asin(x) #三角函数和反三角函数
11 math.sinh(x); math.asinh(x) #双曲函数和反双曲函数
12 math.gamma(x) #gamma 函数  $\Gamma(x)$ 
13
14 math.factorial(x) #返回 x 的阶乘
15 math.gcd(x, y) #返回 `x` `y` 的最大公约数
16 math.degrees(x); math.radians(x) #弧度和角度转换
17
```

5.2 numpy

5.2.1 定义矩阵

Listing 5.2: Define matrices using numpy

```

1 import numpy as np
2
3 a = np.array([0.1*i for i in range(100)])
4 a = np.array([[i+5*j for i in range(5)] for j in range(5)])
5 a = np.arange(25).reshape(5,5) #reshape 可以重新规定矩阵型号
6 a = np.linspace(0,10,100) #第三个参数是生成的列表长度
7 a = np.arange(0,10,0.1) #第三个参数是步长
8 a = np.logspace(0.9,10) #生成 10 的 0-9 次幂
9 a = np.eye(3) #生成三维单位阵
10 a = np.diag([1,2,3,4,5]) #生成对角阵
11 a = np.random.rand(2,3) #生成 2*3 随机矩阵,0-1 均匀分布
12 a = np.random.randn(2, 3)
13 a = np.random.random((2,3)) #生成 2*3 随机矩阵, 用元组表示大小
14 a = np.random.randint(low,high,size=(2,3)) #生成 2*3 随机整数矩阵
15

```

5.2.2 特殊函数

Listing 5.3: Advanced operations in numpy

```

1 X,Y=np.meshgrid(x,y) #将 x,y 扩展为一个矩阵
2 x,y=np.outer(x,y) #得到矩阵  $x^t y$ 
3 x=np.append(x1,x2) #拼接 array
4 a=X[1] #取出矩阵的一行
5 a=X[:,1] #取出矩阵的一列
6 X[:,0]=a #更改矩阵的一列
7 a=X[n,m] #取出满足条件的元素
8 (n,m)=np.where(X>1) #查找满足条件的元素坐标, n 和 m 都是 array
9 l=np.argwhere(X>1) #n 为坐标组成的二维 array
10

```

5.2.3 矩阵运算

Listing 5.4: calculation of matrix

```

1 c = np.dot(x,y) #矩阵乘法
2 c = x*y #对应元素相乘
3 c = np.transpose(a);c=a.T #矩阵转置
4 c = np.transpose(a, (2, 0, 1)) #调整矩阵轴的顺序
5 c = np.linalg.norm(x) #求 x 的  $L_2$  范数
6 result=np.linalg.inv(a) #求逆矩阵
7 result=np.linalg.det(a) #求行列式
8 result=np.linalg.matrix_rank(a) #求矩阵的秩
9 result = np.dot(a,np.linalg.inv(b)) #矩阵右除
10 result = np.dot(np.linalg.inv(a),(b)) #矩阵左除
11 matrix.sum(axis=0) #求和, axis 的轴会消失
12

```

5.3 Pandas

5.3.1 Series

Series 类型是一维数组，由 index 和 value 组成。

Listing 5.5: Series in Pandas

```

1 import pandas as pd
2
3 data=['A','B','C']; index=['a','b','c']
4 dict_data = {'a': 'A', 'b': 'B', 'c': 'C'}
5
6 series = pd.Series(data) #默认 index 从 0 开始编号
7 series_with_index = pd.Series(data,index=index) #规定 index
8 series_with_index = pd.Series(dict_data) #这一句与上一句相同
9 print(series.index,series.value) #会输出数组
10 print(series_with_index['a']) #调用 Series 的元素
11

```

5.3.2 DataFrame

DataFrame 的每列的名称为键，每个键对应一个数组，这个数组为值。

Listing 5.6: Create a dataframe

```

1 import pandas as pd
2
3 data = {'a':[1,2,3,4,5], 'b':[6,7,8,9,10], 'c':[11,12,13,14,15]}
4 index = ['A','B','C','D','E']
5 data_frame = pd.DataFrame(data) #创建 DataFrame 对象，默认 index 从 0 开始编号
6 data_frame = pd.DataFrame(data, index=index) #规定 index
7 data_frame = pd.DataFrame(data, columns=['a','b']) #指定列
8 print(data_frame)
9
10 df.set_index('a', inplace=True) #将 column `a` 设定成 index，仅用于创建时没有默认 index 的情况
11 df.reset_index(inplace=True) #将 index 重置为从 0 开始编号，仅用于创建时没有默认 index 的情况
12

```

5.3.3 读写数据

Pandas 模块可以将 csv 或 excel 文件转为 DataFrame 变量，也可以将 DataFrame 变量写入 csv 或 excel 文件。

Listing 5.7: Read and write files using Pandas

```

1 import pandas as pd
2
3 data = pd.read_csv('data.csv')
4 data = pd.read_excel('data.xlsx')
5 data = pd.read_csv('data.csv', index_col='Id') #设置指定 column 作为 index
6 data = pd.read_csv('data.csv', na_values=['NA', 'Missing']) #将给定值替换成 Na
7 d_parser = lambda x: pd.datetime.strptime(x, '%Y-%m-%d %I-%p')
8 data = pd.read_csv('data.csv', parse_date=['Date'], date_parser=d_parser)
9 #将 Date 列转为 datetime
10 data = pd.read_json('data.json', orient='records', lines=True) #读取列表形式的 json 文件
11
12 data.to_csv('data.csv', columns=['A','B'], index=False) #不写入行索引
13 data.to_csv('data.tsv', sep='\t') #写入 tsv 文件
14 data.to_excel('data.xlsx', columns=['A','B'], index=False) #写入 excel 文件
15 data.to_json('data.json', orient='records', lines=True) #以列表形式写入 json 文件,默认是字典形
16 式

```

5.3.4 数据库交互

pandas 中的数据也可以通过 SQL alchemy 进行快速的存储和访问。

Listing 5.8: pandas & database

```

1 from sqlalchemy import create_engine

```

```

1 import psycopg2
2
3 engine = create_engine('postgresql://dbuser:dbpass@localhost:5432/sample_db')
4
5

```

5.3.5 基本操作

Listing 5.9: Basic functions about DataFrame

```

1 A = data_frame.to_numpy(); data_frame = pd.DataFrame(A) #DataFrame 和 numpy 转换
2
3 pd.set_option('display.max_columns', 20) #设置最多显示列数
4 pd.set_option('display.max_rows', 20) #设置最多显示行数
5 print(data_frame.columns) #输出所有列
6 print(data_frame.shape) #输出行列数
7 print(data_frame.dtypes) #输出每一列的数据类型, 混合数据类型显示为 object
8 print(data_frame.info()) #输出 DataFrame 的基本信息
9 print(df.head(10)); print(df.tail(10)) #输出前/后 10 行
10
11 data_frame[['a', 'b']] #返回指定列组成的 DataFrame
12 data_frame.iloc[0] #输出第一行
13 data_frame.iloc[[0, 1]] #输出指定行
14 data_frame.iloc[[0, 1], [1, 2]] #输出第 1, 2 行第 2, 3 列
15 data_frame.loc[[0, 1], ['a', 'b']] #输出 index 为 0, 1, column 为 'a', 'b' 的 DataFrame
16 data_frame.loc[0:1, 'a':'b'] #与上面一句相同
17 #Loc 的变量也可以是 filter(True/False Series)
18 data_frame[data_frame > 0] #Filtering, 用 NaN 对齐
19 data_frame[data_frame['b'].isin([7, 8, 9])] #Filtering using list
20 data_frame[data_frame['b'].str.contains('Python', na=False)]
21 #提取 column 'b' 字符串中包含 'Python' 的行, NaN 不提取
22
23 #update column
24 data_frame.columns = ['Column_'+i for i in data_frame.columns]
25 data_frame.columns = df.columns.str.replace(' ', '_') #将 column 中的空格替换为下划线
26 data_frame.rename(columns={'a': 'Column_a', 'b': 'Column_b'}, inplace=True) #替换 columns 名
27
28 #update values
29 data_frame.loc['A':'B', 'a'] = [10, 11] #更新若干单元格
30 data_frame.loc['A':'B', 'a':'b'] = [[10, 11], [12, 13]] #更新若干单元格
31 data_frame.loc['B'] = [11, 12, 13] #更新一行
32 data_frame['d'] = ['AXIS', 'Root', 'TeX', 'Python', 'JavaScript'] #增添或修改一列数据
33 data_frame['d'] = data_frame['d'].str.lower() #全部改成小写
34 data_frame._append({'a': 17}, ignore_index=True) #插入一行, reset index 以避免冲突, 用 NaN 补齐
35 df_total = data_frame._append(df, ignore_index=True) #合并两个 DataFrame
36 df_total = data_frame.concat([df1, df2], axis='columns') #合并两个 DataFrame
37 df_total = data_frame.concat([df1, df2], axis='index') #合并两个 DataFrame
38
39 #drop values
40 data_frame.drop(index=['A', 'B'], inplace=True) #按 index 删除, inplace 表示直接替换原数据
41 data_frame.drop(columns=['a', 'b'], inplace=True) #按 column 删除
42

```

5.3.6 统计操作

Listing 5.10: Perform statistical operations using DataFrame

```

1 #数据预处理
2 data_frame['d_len'] = data_frame['d'].apply(len) #对每个元素进行函数操作
3 data_frame.apply(len, axis='rows') #对每一行进行 len 操作, 即获取 column 数
4 df_len = data_frame.applymap(len) #对 DataFrame 中每个元素进行 len 操作
5 sub_d = data_frame['d'].map({'Python': 'C++', 'JavaScript': 'Java'})
6 #替换 Series 的内容, 未匹配的替换成 NaN
7 sub_d = data_frame['d'].replace({'Python': 'C++', 'JavaScript': 'Java'})
8 #替换 Series 的内容, 未匹配的会保持原来的值
9
10 #处理空缺值

```

```

11 null_num = data_frame.isnull().sum() #统计空缺值数量, isnull 在空缺值返回 True, 否则返回 False
12 not_null_num = data_frame.nonnul().sum() #统计非空缺值数量
13 data_frame.dropna(axis='index', inplace=True) #删除包含空缺值的整行数据
14 data_frame.dropna(axis='index', how='all') #删除整行都空缺的数据
15 data_frame.dropna(axis='index', how='any', subset=['name', 'age']) #删除指定列有空缺的行
16 data_frame.fillna(0, inplace=True) #修改空缺值
17 data_frame.fillna(data_frame.mean()) #平均值填充
18 data_frame.fillna({'A':0, 'B':1, 'C':2}, inplace=True) #每列空缺值用指定的值代替
19
20 #处理不同数据类型
21 data_frame['age'] = data_frame['age'].astype(float) #全部转换成 float
22 data_frame['age'].unique() #返回所有出现的值
23
24 #基本运算 (逐项运算)
25 score = data_frame.a + data_frame.b - data_frame.c #可以直接做向量运算
26 filt = data_frame > 0
27
28 #字符串操作 (逐项操作)
29 df['school'] = df['school'].str.upper()
30 df['name'] = df['name'].str.split(' ')
31 df[['first_name', 'last_name']] = df['name'].str.split(' ', expand=True)
32 df['name'].str.contains('Jack') #匹配字符串, 返回一个 True-False Series
33
34 #常用的统计函数
35 data_frame.sort_index(inplace=True) #按 index 排序, NaN 会被当成最大值
36 data_frame.sort_values(by='a', ascending=False, inplace=True) #排序, ascending 为升序
37 data_frame.sort_values(by=['a', 'b'], ascending=[True, False]) #第一标准相同排第二标准, 以此类推
38 #以下函数在运行时会忽略 NaN
39 data_frame.nlargest(5, 'a') #取 'a' 最大的五行
40 data_frame.nlargest(5, columns=['a', 'b']) #取最大五项, 第一标准为 'a', 第二标准为 'b'
41 data_frame.nsmallest(5, columns=['a', 'b']) #取最小五项, 第一标准为 'a', 第二标准为 'b'
42 nums = df['a'].count() #统计非 NaN 的个数
43 nums = df['a'].value_counts() #统计 values 出现的频次, 返回一个 Series
44 nums = df['a'].value_counts(normalize=True) #统计 values 出现的频率
45 average = data_frame.mean() #求每列的平均值, 输出一个 Series
46 media = data_frame.median() #求每列的中位数
47 sum = data_frame.sum() #求每列的和
48 variance = data_frame.var() #样本方差
49 std = data_frame.std() #样本标准差
50 print(data_frame.describe()) #给出数据的所有常见统计参数
51
52 #分组
53 country_grp = df.groupby(['Country'])
54 print(country_grp.get_group('United States')) #查找一个 group, 返回一个 DataFrame
55 country_grp['SocialMedia'].value_counts().head() #统计每一个 group 的 SocialMedia, 返回 Series
56 country_grp['Salary'].agg(['median', 'mean']) #返回 DataFrame, columns=['median', 'mean']
57 country_grp['Language'].apply(lambda x: x.str.contains('French').sum()) #将函数作用于每个 group
58

```

5.3.7 处理 datetime

Listing 5.11: datetime in pandas

```

1 d_range = pd.date_range('2024-01-01', period=6, freq='D') #区间生成日期 Series
2 df['Date'] = pd.to_datetime(df['Date'], format='%Y-%m-%d %I-%p') #string 转 datetime
3
4 df['Date'].dt.date() #对整行 datetime 进行操作
5 df['Date'].max(); df['Date'].min() #获取最早和最晚的时间
6 df.loc[df['Date'] >= pd.to_datetime('2020-5-20')]
7
8 df.set_index('Date', inplace=True) #下面将 Date 设置为 index
9 df['2019'] #得到 2019 年的数据
10 df['2019-02':'2020-01'] #得到指定时间区间内的数据
11 df['temperature'].resample('2D').max()
12 #将数据重采样为每两天的频率并求最大值, 重采样方式参考 date offset Series
13 df.resample('W').agg({'temperature': 'max', 'rain': 'sum', 'air_quality': 'mean'})
14 #重采样时, 对不同 column 执行不同的操作
15

```


6 visualization

6.1 matplotlib

6.1.1 画布预处理

Listing 6.1: basic setup of plt

```

1 from matplotlib import pyplot as plt
2
3 plt.figure(figsize=(10,20),facecolor,edgecolor)
4 plt.title("This is my title")
5 plt.xlabel("x"); plt.ylabel("y")
6 plt.xlim(2,22) #设置 x 的范围
7 plt.xscale('log') #设置 x 为对数轴
8 plt.xticks(ticks=[2*i+1 for i in range(10)],
9            labels=[2*i+1 for i in range(10)]) #设置坐标轴位置和标签
10 plt.legend(["line1", "line2"], loc="upper left") #显示图例, 以加入图表的顺序为序
11 plt.legend(loc=(0.03, 0.05)) #显示图例, 利用图表绘制时的 label, 直接以横纵坐标设定位置
12 plt.rcParams['font.sans-serif'] = ['SimHei'] #字体设置 (默认字体不支持中文)
13 plt.rcParams['axes.unicode_minus'] = False #正常显示负号
14
15 plt.style.use("seaborn-v0_8") #利用预设样式
16 plt.xkcd() #手绘样式
17 plt.tight_layout() #布局更紧凑, 可以避免 label 和 ticks 显示在画布范围外的问题
18
19 plt.grid(True) #开启网格线
20 plt.grid(axis=both, #axis=x or y or False
21         linestyle="dashed", #or dotted or dashdot
22         color="FFFFFF")
23 ) #添加网格线
24 plt.axhline(5, #水平参考线
25            xmin=0.1, xmax=0.9, #[0, 1] 范围内 i 昂对于 x 轴的归一化范围
26            color="red", linestyle="--", linewidth=0.5 #参考线样式
27            )
28 plt.axvline(10) #垂直参考线
29 plt.axhspan(5, 7, #水平参考区域
30            xmin=0.1, xmax=0.9, #[0, 1] 范围内 i 昂对于 x 轴的归一化范围
31            color="red", alpha=0.5, #填充样式, color 会同时作用于 facecolor 和 edgecolor
32            linestyle="--", linewidth=0.5 #区域边沿线样式
33            )
34 plt.axvspan(10, 12) #垂直参考区域
35
36 plt.annotate(text,xy=(5,10), #待注释点坐标
37             xytext=(7,12), #注释文本位置
38             color="FFFFFF",fontsize=16,
39             ha="center", #水平居中
40             va="bottom", #垂直对齐
41             arrowprops={"arrowstyle": "->", #or "-."
42                        "color": "FFFFFF"})
43 ) #显示注释点
44 plt.text(7,12,text) #显示无箭头注释
45
46 plt.savefig('plot.png') #保存图片
47 plt.show() #显示图片, 注意没有这句不会显示
48

```

6.1.2 子图表

figure 是一整张图片, plot 是其中的图。一个 figure 可以有多个 plot。

Listing 6.2: subplot

```

1 ax = plt.gca() #返回当前坐标轴
2 fig = plt.gcf() #返回当前 figure
3
4 fig, ax = plt.subplots( #创建 figure 和其中的 subplot, ax fig 是 list
5                        nrows=2, ncols=1,

```

```

6         sharex=True, #共用 x 轴 (xticks 只保留一个, xlabel 不会共用)
7     )
8
9     ax[0].plot(x, y)
10    ax[1].plot(x, y1)
11    ax[0].legend()
12
13    ax[0].set_title("Title of plot")
14    ax[0].set_xlabel("x")
15    ax[0].set_ylabel("y")
16
17    plt.tight_layout()
18    plt.show()
19    fig.savefig('test.jpg')
20

```

6.1.3 数据获取

Listing 6.3: fetch data

```

1  import numpy as np
2  import csv
3
4  x = np.arange(5)
5  y = np.array([3, 7, 8, 9, 13])
6  y2 = np.array([2, 4, 7, 11, 16])
7
8  with open('data.csv') as csv_file:
9      csv_reader = csv.DictReader(csv_file)
10     x = csv_reader.fieldnames
11     y = [float(i) for i in next(csv_reader).value()]
12     y1 = [float(i) for i in next(csv_reader).value()]
13

```

6.1.4 折线图

Listing 6.4: line plot

```

1  plt.plot(x, y, 'k--') #黑色虚线, 第三个变量是 format strings 样式, 见文档 Notes part
2  plt.plot(x, y,
3          color="#5a7d9a",
4          label="Line1",
5          linestyle="dashed", #or dotted or dashdot
6          #也可以设置为-- -. - : None ' ' ' '
7          linewidth=3,
8          marker=".", #or " " "o" "+" "x"
9          markersize=8,
10         markerfacecolor="blue",
11         markeredgecolor="cyan",
12         markeredgewidth=2
13     )
14     #面积图, y 轴是各组 y 数据的和
15     plt.stackplot(x, y1, y2, y3,
16                 color=["red", "yellow", "blue"],
17                 labels=["player1", "player2", "player3"])
18
19     plt.fill_between(x, y, alpha=0.25) #填充 x-y 下方的区域
20     plt.fill_between(x, y, 30, #填充 x-y 和 y=30 之间的区域
21                     where=(y > 30), #仅填充 y>30 的部分
22                     interpolate=True,
23     )
24     plt.fill_between(x, y, y1) #填充两条曲线之间的部分
25

```

6.1.5 柱状图

Listing 6.5: bar chart

```

1 plt.bar(x, y,
2         width=0.8, #柱形宽度
3         bottom=1, #柱形底部高度, 可以是和 x 一样长的列表, 用于 bar 的堆叠
4         hatch="/") #or "L" ""\ " //"
5 plt.bar(['Jack', 'Ron', 'Jane'], [93, 98, 96]) #这里 x 轴可以不是数
6 plt.barh(x,y) #横向柱状图, 适用于 bar 较多的情况
7
8 width = 0.25 #这是并列柱状图的画法
9 plt.bar(x-width/2, y, width=width, color="#008fd5", label="y")
10 plt.bar(x+width/2, y1, width=width, color="#e5ae38", label="y1")
11 #直方图
12 plt.hist(x, bins=5, #直方图, bins 可以是整数 (分组数量) 或列表
13         edgecolor="black"
14 )
15 plt.hist(x, bins=[10, 20, 30, 40]) #列表设定各 bins 的边界
16 plt.hist(x, bins=5,
17         log=True #y 轴以 10 的幂次表示, 用于数据差距比较大的情况
18 )
19

```

6.1.6 饼图

Listing 6.6: pie chart

```

1 #普通饼图
2 plt.pie(x, #不用换算成百分数, 原数据即可
3         labels=["一月", "二月", "三月"],
4         colors=["red", "blue", "yellow"],
5         wedgeprops={"edgecolor": "red"} #各区域之间的分割线
6         autopct="%0.1f%%", #一位小数百分比, %d%% 整数百分比, %0.1f 一位小数, %0.2f%% 两位小数百分比
7         explode=[0, 0.1, 0], #将第二块拉出 0.5
8         shadow=True, #加上阴影
9         startangle=90, #开始方向, 默认为右侧 (0), 逆时针为正方向
10 )
11 #圆环图
12 plt.pie(x, colors=["red", "blue", "yellow"],
13         shadow=True, labels=["一月", "二月", "三月"],
14         wedgeprops={"width": 0.6} #内圆半径, 外圆半径为 1
15 )
16

```

6.1.7 散点图

Listing 6.7: scatter plot

```

1 plt.scatter(x, y, s=100, #散点图, s 表示点的大小, 可以是和 x 一样长的 List
2            color="red", alpha=0.5, #填充样式
3            marker="o", #or " ", " + " " x " " ^ "
4            edgecolor="black", linewidth=2, linestyle="-.", #边缘线样式
5            )
6 plt.scatter(x, y, [i*50+100 for i in range(len(x))]) #气泡图
7
8 plt.scatter(x, y, c=[7, 5, 9, 2, 3], cmap='Greens') #通过 cmap 设置颜色
9 cbar = plt.colorbar(label='heat') #显示 colorbar 及其注释
10 #可以通过 print(plt.colormaps) 查看所有的可选 cmap
11 #误差棒图
12 plt.errorbar(x, y, #散点数组
13             yerr=[[0.1, 0.3, 0.2], [0.4, 0.1, 0.1]], #上置信度和置信度
14             xerr=[0.05, 0.04, 0.07], #上下置信度
15             fmt="o--", #散点样式和连线样式
16             color="red", #color of scatters

```

```

17         ecolor="blue", #color of the errorbars
18         elinewidth=3, #误差棒粗细
19         capsize=2, #误差棒横杠大小
20         capthick=2, #误差棒边界横杠厚度
21         uplims=False, #显示误差上界, False 显示, True 隐藏
22         lolims=[True, False, False], #只显示指定的下界
23     ) #其余选项和 plot, scatter 类似
24

```

6.1.8 箱型图

Listing 6.8: box plot

```

1  #箱型图
2  plt.boxplot(x, #x 是一个数组
3              showmeans=True, #显示均值
4              patch=True, #显示缺口箱图, 默认矩形箱图
5              vert=False, #False 为横向, True 为纵向
6              showfliers=True, #显示离群点, 默认显示
7              whis=1.5, #箱须延长的百分比范围
8              width=1, #箱体总宽度
9              flierprops={ #离群点样式
10                 "marker": "o", #or " ", " + " " x "
11                 "markerfacecolor": "red",
12                 "markeredgecolor": "black",
13                 "markersize": 8
14             },
15             patch_artist=True, #自定义箱型
16             boxprops={ #箱型样式
17                 "facecolor": "red",
18                 "edgecolor": "yellow",
19                 "linewidth": 4
20             },
21         )
22  #绘制多个箱型
23  plt.boxplot((x, y), #适用于 x 轴没有数值含义的时候
24              tick_labels=['x', 'y'], #均匀分布
25          )
26  plt.boxplot((x, y), #适用于 x 轴有数值
27              positions=(1, 2),
28              manage_ticks=True #刻度会根据 positions 自动调整
29          )
30

```

6.1.9 热力图

Listing 6.9: heat map

```

1  #绘制热力图
2  plt.imshow(x, #x 是个二维列表
3             cmap=plt.cm.cool #设置颜色
4         )
5  plt.colorbar(label="label") #显示图例
6

```

6.1.10 时间图像

Listing 6.10: Plotting Time Series Data

```

1  import pandas as pd
2  from datetime import datetime, timedelta
3  from matplotlib import pyplot as plt
4  from matplotlib import dates as mpl_dates

```

```

5
6     #获取时间
7     dates = pd.Series([datetime(2024, 6, 30), datetime(2024, 7, 5), datetime(2024, 7, 11)])
8     dates = pd.Series(["2024-6-30", "2024-7-5", "2024-7-11"])
9     dates = pd.to_datetime(dates).sort_values() #注意需要排序
10    y = [4, 6, 7]
11
12    plt.plot(dates, y)
13    plt.gcf().autofmt_xdate() #将日期斜方向写, gcf 表示当前图像
14
15    date_format = mpl_dates.DateFormatter('%b, %d %Y')
16    plt.gca().xaxis.set_major_formatter(date_format) #设置日期格式
17

```

6.1.11 极坐标图表

Listing 6.11: polar plots of plt

```

1     plt.polar(theta,r) #雷达图
2     plt.thetagrid(angles,labels) #角刻度标签
3     plt.rgrids(radii,rotation,labels) #r 方向刻度标签
4
5     ax=plt.axes(polar=True) #建立极坐标画布
6     ax.bar(x=theta,height=data,width=0.4,color="rainbow") #绘制南丁格玫瑰图
7     ax.bar(x=theta,height=100,width=0.4,color="white") #绘制中心空白
8     ax.text(angle,height,text) #添加注释
9     ax.grid(False)
10    plt.thetagrids(angles=[],labels=[]) #刻度标签
11    plt.rgrids(radii=[20],rotation,labels=['20'])
12

```

6.1.12 三维图表

Listing 6.12: 3D plots of plt

```

1     from mpl_toolkits.mplot3d import Axes3D
2     fig=plt.figure()
3     ax1=plt.axes(projection="3d")
4     ax1.scatter3D(x,y,z,cmap="blue")
5     ax1.plot3D(x,y,z,"gray")
6     ax1.plot_surface(X,Y,Z,rstride=0.1,cstride=0.1) #步长越短越清晰
7     ax1.contour(X,Y,Z,zdir='x',offset=-3,cmap="cold") #绘制等高线, 投影在 x=3 平面上
8     ax1.bar3d(X,Y,height,width,depth,Z,color="red",shade=True) #绘制柱状图,height 为柱底高度
9

```

6.1.13 动态实时图像

Listing 6.13: plotting live date

```

1     from matplotlib.animation import FuncAnimation
2     import itertools, random
3
4     x_vals, y_vals, index = [], [], itertools.count()
5     def animate(i):
6         x_vals.append(next(index))
7         y_vals.append(random.randint(0, 5))
8         plt.cla() #先清除上一帧图像
9         plt.plot(x_vals, y_vals)
10        plt.tight_layout()
11
12    ani = FuncAnimation(plt.gcf(), animate, interval=1000)
13
14    plt.tight_layout()

```

```

15 plt.show()
16

```

6.2 wordcloud

Listing 6.14: wordcloud

```

1 import matplotlib.pyplot as plt
2 import wordcloud as wc
3
4 text_data = """
5 Python is a popular programming language.
6 It is widely used for web development, data analysis, and artificial intelligence.
7 Word clouds are fun visualizations of text data.
8 Generate a word cloud using the wordcloud module.
9 """
10
11 # 生成词云对象
12 wordcloud = wc.WordCloud(width=800, height=400, background_color='white').generate(text_data)
13
14 # 显示词云图
15 plt.figure(figsize=(10, 5))
16 plt.imshow(wordcloud, interpolation='bilinear')
17 plt.axis('off')
18 plt.show()
19

```

7 website

7.1 请求连接

Python 自带的 `urllib` 和 `urllib3` 模块可以实现一些网络爬虫的常用操作。外接库中的 `requests` 库则更为常用。

7.1.1 urllib module

利用 `request` 模块可以实现 `get` 请求方式获取网页内容。

Listing 7.1: Fetch web content using GET method

```

1 import urllib.request #网络请求子模块
2
3 response = urllib.request.urlopen('http://www.baidu.com') #打开网页
4 html = response.read() #读取网页代码
5 response.close()
6
7 #也可以用 with 语句简化
8 with urllib.request.urlopen('http://www.baidu.com') as response:
9     html = response.read()
10 print(html)
11

```

也可以实现 `post` 请求方式获取网页内容。

Listing 7.2: Fetch web content using POST method

```

1 import urllib.parse #url 解析和引用模块
2 import urllib.request
3
4 #使用 urlencode 方法对数据进行处理, 并将处理后的数据设置为 utf-8 编码
5 data = bytes(urllib.parse.urlencode({'word':'hello'}), encoding='utf8')
6 with urllib.request.urlopen('http://httpbin.org/post', data=data) as response #打开网页
7     html = response.read()

```

```
8 print(html)
9
```

7.1.2 urllib3 module

urllib3 是一个更强大的 Python 库。

Listing 7.3: urllib3 module

```
1 import urllib3
2
3 #用 GET 的方式连接
4 http = urllib3.PoolManager() #创建对象, 用于处理连接和安全等细节
5 response = http.request('GET', 'https://www.baidu.com/') #连接网站
6 print(response.data) #输出读取内容
7
8 #用 POST 的方式连接
9 http = urllib3.PoolManager() #创建对象, 用于处理连接和安全等细节
10 response = http.request('POST', 'http://httpbin.org/post', fields={'word':'hello'})
11 print(response.data) #输出读取内容
12
```

7.1.3 requests module

另外有个更加人性化的第三方库 requests, 这个库更加常用。

Listing 7.4: requests module

```
1 import requests
2
3 data={'word':'hello'}
4 response = requests.get('http://www.baidu.com',params=data) #get 方法访问
5 response = requests.post('http://httpbin.org/post',data=data) #post 方法访问
6 if response.ok: #检验连接是否成功
7     print(response.content) #以字节流形式输出网页源码 (没有换行符和缩进, 可读性差)
8     print(response.text) #以文本形式输出网页源码 (带有换行符和缩进, 可读性好)
9
10 #下载文件
11 import zipfile
12 r = requests.get('<url>') #zip file 的 url
13 with open('data.zip', 'wb') as f: #建一个 data.zip 用于保存
14     f.write(r.content)
15
```

为了绕开反爬设计, 我们可以请求 header。

Listing 7.5: Handling request headers

```
1 import requests
2
3 url = 'http://www.bilibili.com/'
4 headers = {'User-Agent':'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36
5           (KHTML, like Gecko) Chrome/121.0.0.0 Safari/537.36'}
6 #这里 User-Agent 的内容要从网络监视器中复制过来
7 response = requests.get(url,headers=headers)
8 print(response.content)
9
```

7.2 处理 html 文件

BeautifulSoup module 用于处理 html 文件

Listing 7.6: Handling HTML

```

1 from bs4 import BeautifulSoup
2
3 with open('test.html') as html_file: #打开本地的 html 文件
4     soup = BeautifulSoup(html_file, 'lxml') #用 lxml 解析器解析 html 文件
5
6 source = requests.get('http://bilibili.com').text
7 soup = BeautifulSoup(source, 'lxml') #从网页上获取文件
8
9 print(soup.prettify()) #转换成易读的 html 文件
10 match = soup.title.text #得到文件中的某个部分
11 match_div = soup.find('div', class_='footer') #查找文件中第一个对应的内容
12 match_all = soup.find_all('div') #查找文件中所有对应的内容, 返回一个列表
13

```

7.3 Flask

Flask 是一个常见的网站开发框架。

7.3.1 create a website

Listing 7.7: create a web site

```

1 #my_app.py
2 from flask import Flask, abort
3 app = Flask(__name__)
4
5 @app.route("/") #root page
6 @app.route("/home") #两个网址共用一个网页
7 def hello():
8     return "<h1>Hello World!</h1>"
9
10 @app.route("/error") #两个网址共用一个网页
11 def hello():
12     abort(403) #返回 403 错误 (forbidden)
13
14 if __name__ == '__main__':
15     app.run(debug=True)
16

```

Listing 7.8: start a web site

```

1 export FLASK_APP=my_app.py
2 flask run #运行网站, 每次修改 python 文件需要重新启动
3 export FLASK_DEBUG=1
4 flask run #运行网站, 每次修改都会立即出现在网站上
5
6 python my_app.py #直接运行网站
7

```

7.3.2 template

template 就是每个 url 的 html 模板, 建立 template 可以避免将 html 语言写到 python 文件中造成混乱。

在下面的例子中, 我们将为 root page 创建一个 template/home.html 文件, 注意这个文件必须放在 template 文件夹下。

Listing 7.9: template/home.html

```

1 <!DOCTYPE html>
2 <html>
3 <head>
4     <meta charset="UTF-8"> <!-- 提供文档的元数据 -->

```



```

5      <title>网页标题</title>
6      <style> <!--定义内部CSS样式-->
7          .custom-text {
8              font-family: "Times New Roman", serif;
9              color: green;
10         }
11     </style> <!--更好的做法是将要用到的style都写到一个.css文件中
12 </head>
13 <body>
14     <header>
15         <h1>欢迎来到我的网站</h1>
16         <nav>
17             <ul>
18                 <li><a href="{ url_for('home') }}">主页</a></li>
19                 <li><a href="#about">关于我们</a></li>
20                 <li><a href="/contact">联系我们</a></li>
21             </ul>
22         </nav>
23     </header>
24     <main> <!--main part of body-->
25         <h1>这是一个<span style="color:red;">标题</span></h1> <!--h1至h6表示标题的等级-->
26         <p style="custom-text">这是一个段落。</p><br> <!--br为空格-->
27
28         <strong>这是重要的内容。</strong>
29         <small>这是一个较小的文本。</small>
30         <em>这是一个需要强调的文本。</em>
31         <mark>这是一个被标记的文本。</mark>
32         <p>H<sub>2</sub></sub>O 表示水。</p> <!--下标文本-->
33         <p>2<sup>2</sup> = 4</p> <!--上标文本-->
34         <p>这是一个 <code>print("Hello, World!")</code> 示例。</p>
35
36         <a href="https://www.example.com">访问示例网站</a>
37         <a href="#">点击这里回到顶部</a>
38         
39     </main>
40
41     <h2>无序列表</h2>
42     <ul> <!--ol有序无标记列表, li有序有标记列表-->
43         <li>苹果</li>
44         <li>香蕉</li>
45         <li>橙子</li>
46     </ul>
47
48     <div style="border: 1px solid black; padding: 10px;">
49         <h2>新闻标题</h2>
50         <p style="font-family: Arial, sans-serif; color: blue;">
51             这是新闻的内容部分。该部分使用`&lt;div&rt;`标签包裹, 以便进行样式处理和布局控制。
52         </p>
53     </div>
54     <div>
55         <button>Click me!</button> <!--按钮-->
56         <input type="text" name="username" placeholder="Enter name"> <!--文本输入框-->
57         <input type="password" name="password" placeholder="Password"> <!--密码输入框-->
58         <input type="submit" value="Submit"> <!--提交按钮-->
59         <input type="radio" name="gender" value="male"> <!--单选按钮-->
60         <input type="radio" name="gender" value="female"> <!--单选按钮-->
61         <input type="checkbox" name="agree" value="yes"> <!--复选框-->
62     </div>
63 </body>
64 </html>
65

```

html 文件中有一些特殊字符需要转义:

- 大于号小于号 < >: < >
- 书名号 《 》: ≪ &Rt;
- 与号 &: &
- 双引号": "

- 单引号': '
- 空格":
- 乘号 ×: ×
- 除号 ÷: ÷

Listing 7.10: using template in my_app.py

```

1 @app.route("/")
2 def home():
3     return render_template('home.html')
4
5 @app.route("/post/<int:post_id>") #route 可以带有变量, 这里设置其为 int 类型
6 def post(post_id):
7     post = Post.query.get_or_404(post_id) #如果网页存在就返回 post_id, 否则返回 404
8     return render_template('post.html', post=post)
9

```

7.3.3 template inheritance

同一个网站中 template 往往有很多重合部分, template inheritance 可以实现统一设置模板。下面, 我们先创建一个 layout.html, 其中包含网页的公共部分。

Listing 7.11: layout.html

```

1 <head>
2     {% if title %}
3     <title>My app - {{ title }}</title>
4     {% else %}
5     <title>My app</title>
6     {% endif %}
7 </head>
8 <body>
9     {% block content %}{% endblock %} <!--需要替换的部分-->
10
11 </body>

```

在 home.html 中使用 layout.html。

Listing 7.12: Using template inheritance

```

1 {% extends "layout.html" %}
2 {% block content %}
3     {% for post in posts %} #读出posts的内容
4         <h1>{{ post.name }}</h1> #两个大括号表示代码结果要输出到网页上
5         <p>with number of {{ post.number }}</p>
6     {% endfor %} #结束读取
7 {% endblock content %}
8

```

7.3.4 bootstrap

bootstrap 定义了一些方法来优化网页, 参考bootstrap 官方文档。

Listing 7.13: bootstrap

```

1 <!doctype html>
2 <html lang="en">
3 <head>
4     <meta charset="utf-8">
5     <meta name="viewport" content="width=device-width, initial-scale=1">
6     <link
7         href="https://cdn.jsdelivr.net/npm/bootstrap@5.3.3/dist/css/bootstrap.min.css"

```

```

8         rel="stylesheet"
9         integrity="sha384-QWTKZyjpPEjISv5WaRU90FeRpok6YctnYmDr5pNlyT2bRjXh0JMhJy6hW+ALEwIH"
10        crossorigin="anonymous">
11
12    {% if title %}
13        <title>My app - {{ title }}</title>
14    {% else %}
15        <title>My app</title>
16    {% endif %}
17</head>
18<body>
19    <div class="container" mt-1> <!--利用bootstrap中的设定-->
20        <!--外间距:mt上间距, mb下间距, my上下间距-->
21        <!--外间距:ml左间距, mr右间距, mx左右间距-->
22        <!--内间距:pt左间距, pr右间距, px左右间距-->
23        <!--内间距:pl左间距, pr右间距, px左右间距-->
24        {% block content %}{% endblock %} <!--需要替换的部分-->
25    </div>
26
27    <script
28        src="https://cdn.jsdelivr.net/npm/bootstrap@5.3.3/dist/js/bootstrap.bundle.min.js"
29        integrity="sha384-YvpcrYf0tY3lHB60NNkmXc5s9fDVZLESaAA55NDz0xhy9GkcIdsIK1eN7N6jIeHz"
30        crossorigin="anonymous">
31    </script>
32</body>
33</html>
34

```

7.3.5 External CSS

External CSS(外部样式表) 将网页中文本、图片等元素所需样式都写在 static/main.css 文件里备用。main.css 中可以为每个板块类定义样式，这个样式会自动生效于网页中所有对应的类。

下面这个样例表示在 ‘article-metadata’ 类的元素内的链接 (<a> 标签) 当鼠标悬停 (hover) 时的样式。

Listing 7.14: static/main.css

```

1 .article-metadata a:hover {
2     color: #333;
3     text-decoration: none;
4 }

```

导入 External CSS 需要用到 url_for module

Listing 7.15: url_for

```

1 from flask import Flask, render_template, url_for
2
3 url_for('home') #返回 home 函数对应的 url
4 url_for('static', filename='picture1.jpg') #在 static 文件夹下 picture1.jpg 的地址
5 url_for('post', post_id=1) #返回 url: /post/1, 用于 decorator 里有变量的情况
6

```

在 layout.html 中调用 url_for 导入样式表。

Listing 7.16: layout.html

```

1 <head>
2     <link rel="stylesheet" type="text/css"
3         href="{{ url_for('static', filename='main.css') }}">
4     <!--static为默认目录, main.css只要在这个目录下面, url_for都会自动找到该文件-->
5 </head>
6

```

7.3.6 post data

template 中可以通过 post 传递变量，以实现信息交换。

Listing 7.17: post data in flask

```

1 posts = [ #posts 是一个 dict 组成的 list
2     {'name': 'Amy', 'number': 11},
3     {'name': 'Jack', 'number': 28}
4 ]
5
6 def home():
7     return render_template('home.html', posts=posts, title="Home page")
8

```

Listing 7.18: use posts in html

```

1 <head>
2     {% if title %}
3         <title>My app - {{ title }}</title>
4     {% else %}
5         <title>My app</title>
6     {% endif %}
7 </head>
8 <body>
9     {% for post i posts %} #读出posts的内容
10         <h1>{{ post.name }}</h1> #两个大括号表示代码结果要输出到网页上
11         <p>with number of {{ post.number }}</p>
12     {% endfor %} #结束读取
13 </body>
14

```

7.3.7 form & field

form 表示网页上的表单，field 表示提供用户输入的区域。一个 form 可以有多个输入。flask-wtf module 可用于处理表单。其中 wtforms 提供了各种表单区域，详见 [wtf field](#)。wtforms.validators 提供了检验输入法的方法，详见 [wtf validators](#)。

下面，我们建立一个 forms.py 文件来处理表单

Listing 7.19: form.py

```

1 from flask_wtf import FlaskForm
2 from wtforms import StringField, PasswordField, SubmitField, BooleanField
3     #这里常用的还有 TextareaField 文本编辑区域
4 from wtforms.validators import DataRequired, Length, Email, EqualTo, ValidationError
5 from models import User #models 将在下一目建立
6
7 class RegistrationForm(FlaskForm):
8     username = StringField('Username', validators=[DataRequired(), Length(min=2, max=20)])
9     email = StringField('Email', #这里的第一个变量就是下面 html 文件中的 label
10         validators=[DataRequired(), Email()])
11     password = PasswordField('Password', validators=[DataRequired()])
12     confirm_password = PasswordField('Confirm Password',
13         validators=[DataRequired(), EqualTo('password')])
14     remember = BooleanField('Remember Me')
15     submit = SubmitField('Sign Up')
16
17 def validate_username(self, username): #Validate 函数会在表单提交时自动调用
18     user = User.query.filter_by(username=username.data).first()
19     if user: #用户名已存在
20         raise ValidationError('That username is taken.')
21
22 def validate_email(self, email):
23     user = User.query.filter_by(email=email.data).first()
24     if user: #用户名已存在
25         raise ValidationError('That email is taken.')
26

```

创建 register.html 显示表单。<form> 用于建立表单，<fieldset> 用于将表单中的内容进行分组。

Listing 7.20: register.html

```

1  {% extends "layout.html" %}
2  {% block content %}
3      <div class="content-section">
4          <form method="POST" action=""> <!--action=""表示将表单数据提交到当前URL-->
5              {{ form.hidden_tag() }} <!--防止跨站请求伪造攻击-->
6              <fieldset class="form-group">
7                  <legend class="border-bottom mb-4">Join Today</legend>
8                  <!--Legend 一般用于fieldset内部，作为注释-->
9                  <div class="form-group">
10                     {{ form.username.label(class="form-control-label") }}
11                     {% if form.username.errors %}
12                         {{
13                             form.username(
14                                 class="form-control form-control-lg is invalid"
15                             )
16                         }}
17                     <div class="invalid-feedback">
18                         {% for errors in form.username.errors %}
19                             <span>{{ error }}</span>
20                         {% endfor %}
21                     </div>
22                     {% else %}
23                         {{ form.username(class="form-control form-control-lg") }}
24                     {% endif %}
25                 </div>
26                 <div class="form-group">
27                     {{ form.email.label(class="form-control-label") }}
28                     {% if form.email.errors %}
29                         {{
30                             form.email(
31                                 class="form-control form-control-lg is invalid"
32                             )
33                         }}
34                     <div class="invalid-feedback">
35                         {% for errors in form.email.errors %}
36                             <span>{{ error }}</span>
37                         {% endfor %}
38                     </div>
39                     {% else %}
40                         {{ form.email(class="form-control form-control-lg") }}
41                     {% endif %}
42                 </div>
43                 <div class="form-group">
44                     {{ form.password.label(class="form-control-label") }}
45                     {% if form.password.errors %}
46                         {{
47                             form.password(
48                                 class="form-control form-control-lg is invalid"
49                             )
50                         }}
51                     <div class="invalid-feedback">
52                         {% for errors in form.password.errors %}
53                             <span>{{ error }}</span>
54                         {% endfor %}
55                     </div>
56                     {% else %}
57                         {{ form.password(class="form-control form-control-lg") }}
58                     {% endif %}
59                 </div>
60                 <div class="form-group">
61                     {{ form.confirm_password.label(class="form-control-label") }}
62                     {% if form.confirm_password.errors %}
63                         {{
64                             form.confirm_password(
65                                 class="form-control form-control-lg is invalid"
66                             )
67                         }}

```

```

68         <div class="invalid-feedback">
69             {% for errors in form.confirm_password.errors %}
70                 <span>{{ error }}</span>
71             {% endfor %}
72         </div>
73     {% else %}
74         {{ form.confirm_password(class="form-control form-control-lg") }}
75     {% endif %}
76 </div>
77 </fieldset>
78 <div class="form-group">
79     {{ form.submit(class="btn btn-outline-info") }}
80 </div>
81 </form>
82 </div>
83 <div class="border-top pt-3">
84     <small class="text-muted">
85         Already Have An Account?
86         <a class="ml-2" href="{{ url_for('login') }}">Sign In</a>
87     </small>
88 </div>
89 {% endblock content %}
90

```

在 app.py 中设置密码 (防止用户信息泄露), 调用表单, 接收用户提交的表单, 并在处理表单后重定向。

Listing 7.21: Using forms in main.py

```

1  from forms import RegistrationForm
2  from flask import Flask, render_template, url_for, redirect
3
4  app.config['SECRET_KEY'] = '114514' #密码可以由 secrets module 生成
5
6  @app.route("/register", methods=['GET', 'POST']) #接收 GET 和 POST 表单
7  def register():
8      form = RegistrationForm()
9      if form.validate_on_submit(): #判断提交表单是否合理
10         flash(f'Account created for {form.username.data}!', 'success') #显示一次性消息
11         #消息类别有 "error" "danger" "warning" "success" 等
12         return redirect(url_for('home'))
13     return render_template('register.html', title='Register', form=form)
14

```

在 layout.html 中显示 flash 的内容并完成重定向。

Listing 7.22: layout.html

```

1  {% with messages = get_flashed_messages(with_categories=true) %} <!-- success就是true-->
2      {% if messages %}
3          {% for category, message in messages %}
4              <div class="alert alert-{{ category }}">
5                  {{ messages }}
6              </div>
7          {% endfor %}
8      {% endif %}
9  {% endwith %}
10

```

7.3.8 flask-sqlalchemy

flask-sqlalchemy 可以用于存储网站运行产生的数据。
创建一个 models.py 文件, 用于处理数据库。

Listing 7.23: Using sqlalchemy in models.py

```

1  from flask_sqlalchemy import SQLAlchemy

```

```

2 from datetime import datetime
3 from app import db
4
5 app.config['SQLALCHEMY_DATABASE_URI'] = 'sqlite:///site.db' #///表示相对路径
6 db = SQLAlchemy(app)
7
8 class User(db.Model): #创建一个 User 表
9     id = db.Column(db.Integer, primary_key=True) #不重复的 key, 会自动生成, 从 1 开始
10    username = db.Column(db.String(20), unique=True, nullable=False) #用户名不重复且不为空
11    email = db.Column(db.String(120), unique=True, nullable=False)
12    image_file = db.Column(db.String(20), nullable=False, default='default.jpg')
13    #图片用 hash 值记录
14    password = db.Column(db.String(60), nullable=False) #密码用 hash 值记录
15    posts = db.relationship('Post', backref='author', lazy=True) #记录 Post 和 User 的关系
16    #这里 Post 指代下面的 Post 类型, lazy=True 会使实例直接带有 relationship 信息
17
18    def __repr__(self):
19        return f"User('{self.username}', '{self.email}', '{self.image_file}')"
20
21 class Post(db.Model): #创建一个 Post 表
22     id = db.Column(db.Integer, primary_key=True)
23     title = db.Column(db.String(100), nullable=False)
24     date_posted = db.Column(db.DateTime, nullable=False, default=datetime.utcnow)
25     #注意这里 utcnow 是一个函数而不是函数的返回值, 否则 default 会立即调用, 生成一个固定的时间
26     content = db.Column(db.Text, nullable=False) #长文本类型
27     user_id = db.Column(db.Integer, db.ForeignKey('user.id'), nullable=False)
28     #这里 user.id 指代 user 表中的 id 列 (表名是类名的小写)
29     #表示 user_id 是 user 表 id 列的一个外键
30
31    def __repr__(self):
32        return f"Post('{self.title}', '{self.date_posted}')"
33

```

在 app.py 文件中使用 models.py

Listing 7.24: Using models in main.py

```

1 from models import User, Post
2
3 app.config['SQLALCHEMY_DATABASE_URI'] = 'sqlite:///site.db' #///表示相对路径
4 db = SQLAlchemy(app) #database 和 app 绑定, 再次运行不会生成新的 database
5

```

在 shell 创建数据库文件, 并尝试加入数据。

Listing 7.25: create database file in terminal

```

1 >>> from app import db
2 >>> db.create_all()
3
4 >>> from app import User, Post
5 >>> user_1 = User(username='Jack', email='J@demo.com', password='password')
6 >>> db.session.add(user_1)
7 >>> db.session.commit()
8
9 >>> User.query.all() #输出所有 User
10 >>> User.query.first() #输出第一个 User
11 >>> User.query.filter_by(username='Jack').all()
12 >>> user = User.query.get(1) #按照 id 查找
13
14 >>> post_1 = Post(title='Blog 1', content='First Post Content!', user_id=user.id)
15 >>> db.session.add(post_1)
16 >>> db.session.commit()
17 >>> user.posts #此时可以输出用户的 posts
18 >>> db.query.first().author #利用 backref 查询作者
19
20 >>> db.drop_all() #删除数据库中的所有内容
21

```

将 website 产生的用户数据存入 database, 这里用到了 hash 生成, 见下一目。

```

1  @app.route("/register", methods=['GET', 'POST'])
2  def register():
3      form = RegistrationForm()
4      if form.validate_on_submit():
5          hash_password = bcrypt.generate_password_hash(form.password.data).decode('utf-8')
6          user = User(username=form.username.data,
7                      email=form.email.data,
8                      password=hashed_password)
9          db.session.add(user)
10         db.session.commit()
11         flash(f'Account created for {form.username.data}! You can log in now', 'success')
12         return redirect(url_for('login'))
13     return render_template('register.html', title='Register', form=form)
14

```

7.3.9 Bcrypt & LoginManager

flask-bcrypt.Bcrypt 可用于生成 hash 值。flask_login.LoginManager 用于控制用户登录

Listing 7.26: Using Bcrypt in app.py

```

1  from flask import request #这将帮助我们得到 GET 的内容，以便重定向
2  from flask_bcrypt import Bcrypt
3  from flask_login import LoginManager, UserMixin, login_user,
4                          current_user, logout_user, login_required
5  #这些变量可以在 html 中直接访问，无需通过 render_template 传递
6
7  bcrypt = Bcrypt(); bcrypt.init_app(app)
8  bcrypt = Bcrypt(app) #上面两句是一样的
9  login_manager = LoginManager(app)
10 login_manager.login_view = 'login' #指向 login 函数
11 #当没有 login 而访问需要 login 的 url 时，会被重定向到 login，同时通过 GET 传递重定向前的 route
12 login_manager.login_message_category = 'info' #info 是一个内置类
13 #这一句会改变重定向时提示的样式
14
15 pwd_hash = bcrypt.generate_password_hash('password').decode('utf-8')
16 #这个函数重复运行会生成不同的 hash 值
17 bcrypt.check_password_hash(pwd_hash, 'password') #返回 True
18 bcrypt.check_password_hash(pwd_hash, 'Testing') #返回 False
19
20 @login_manager.user_loader
21 def load_user(user_id):
22     return User.query.get(int())
23
24 class User(db.Model, UserMixin):
25     #definition of id etc.
26
27 @app.route("/login")
28 def login():
29     user = User() #add name email etc.
30     if current_user.is_authenticated:
31         return redirect(url_for('home'))
32     else:
33         login_user(user) #user 会载入 current_user，后面可访问，如 current_user.username
34         next_page = request.args.get('next') #得到重定向前的 route
35         return redirect(next_page) if next_page else redirect(url_for('home'))
36
37 @app.route("/account")
38 @login_required #需要 login 才能访问
39 def account():
40     form = UpdateAccountForm() #修改用户信息的表单
41     if form.validate_on_submit():
42         current_user.username = form.username.data #注意 form 要加 data 而 current_user 不用
43         current_user.email = form.email.data
44         db.session.commit() #直接修改 current_user 即可
45         flash('Your account has been updated!', 'success')
46         return redirect(url_for('account'))
47     elif request.method == 'GET':
48         form.username.data = current_user.username

```



```

49     form.email.data = current_user.email
50     image_file = url_for('static', filename='profile_pics/' + current_user.image_file)
51     return render_template('account.html', form=form, image_file=image_file)
52

```

7.3.10 FileField

FileField 允许网页传输文件 (图片、文本文件等)。FileAllowed 设置了允许传输的文件类型。

Listing 7.27: FileField

```

1  from flask_wtf.file import FileField, FileAllowed
2  import secrets
3
4  def save_picture(form_picture):
5      random_hex = secrets.token_hex(8) #生成随机文件名, 防止同名文件 bug
6      _, f_ext = os.path.splitext(form_picture.filename)
7      picture_fn = random_hex + f_ext
8      picture_path = os.path.join(app.root_path, 'static/profile_pics', picture_fn)
9      form_picture.save(picture_path)
10     return picture_fn
11
12     class UpdateAccountForm(FlaskForm):
13         picture = FileField('Update Profile Picture',
14                             validators=[FileAllowed(['jpg', 'png'])])
15
16     @app.route("/account", methods=['GET', 'POST'])
17         if form.validate_on_submit():
18             if form.picture.data:
19                 current_user.picture = save_picture(form.picture.data)
20

```

在 html 文件中调用 FileField。

Listing 7.28: Using FileField in html

```

1  <form method="POST" action="" enctype="multipart/form-data">
2  <!-- 这里必须给出编码方式enctype-->
3  <fieldset>
4      <div class="form-group">
5          {{ form.picture.label() }}
6          {{ form.picture(class="form-control-file") }}
7          {% if form.picture.errors %}
8              {% for error in form.picture.errors %}
9                  <span class="text-danger">{{ error }}</span><br>
10             {% endfor %}
11         {% endif %}
12     </div>
13 </fieldset>
14 </form>
15

```

7.3.11 Pagination

Pagination 可以减少 post 的内容, 提高网站运行速度, 同时避免页面过长, 增加可读性。为了在分页时网页的内容不混乱, 我们还需要对 post 的内容进行排序。

Listing 7.29: paginate

```

1  @app.route("/")
2  def home():
3      page = request.args.get('page', default=1, type=int)
4      posts = Post.query.order_by(Post.data).paginate(per_page=10, page=page)
5      #以 10 个为 1 页, 访问第 page 页
6      return render_template('home.html', posts=posts)
7

```

在 html 文件中，可以设置一个导航栏，导航到指定页码，方便随机访问。

Listing 7.30: Using pagination in home.html

```

1  {% for page_num in posts.items %}
2      <!--处理当前page的post内容-->
3  {% endfor %}
4
5  <!--导航栏-->
6  {% for page_num in posts.iter_pages(left_edfe=1, right_edge=1,
7                                     left_current=1, right_current=2) %}
8  <!--这是一个预设，仅开头结尾和当前page前后的page_num会显示，其他page_num会缩成一个None-->
9  <!--right_current包含current page，故要比left_current多1-->
10     {% if page_num %}
11         {% if posts.page == page_num %} <!--把当前page的button设计的不一样-->
12             <a
13                 class="btn btn-info mb-4"
14                 href="{{ url_for('home', page=page_num) }}"
15                 {{ page_num }}
16             </a>
17         {% else %}
18             <a
19                 class="btn btn-outline-info mb-4"
20                 href="{{ url_for('home', page=page_num) }}"
21                 {{ page_num }}
22             </a>
23         {% endif %}
24     {% else %}
25         ... <!--缩成None的page_num显示成...-->
26     {% endif %}
27 {% endfor %}
28

```

8 debug & tools

8.1 doctest

doctest 可以检查函数的输出，在代码注释样例中给出一组输入输出，若结果错误会报错。

Listing 8.1: doctest hello.py

```

1  from operator import floordiv, mod
2
3  def divide_exact(n, d):
4      """Return the quotient and remainder of dividing N by D
5      >>> q, r = divide_exact(2013, 10)
6      >>> q
7      201
8      >>> r
9      2
10     """
11     return floordiv(n, d), mod(n, d)
12

```

Listing 8.2: doctest bash

```

1  python3 -m doctest hello.py
2  *****
3  File "/mnt/d/Desktop/python/program file/test/test_basic/Pythonproject
4  1/hello.py", line 8, in hello.divide_exact
5  Failed example:
6      r
7  Expected:
8      2
9  Got:
10     3
11  *****

```

Table 8.1: Types of Error

错误名	错误类型	例子
Syntax Error	语法错误	代码结构不是 Python
IndentationError	缩进错误	缩进不一致或缺少缩进
TypeError	对象类型错误	将字符串与数字相加
NameError	命名错误	使用了尚未定义的变量或函数
AttributeError	属性错误	使用了一个 class 没有的属性
IndexError	索引错误	访问超出 list 的索引位置
KeyError	键错误	访问字典中不存在的键
ValueError	值错误	将非数字字符串传给 int()
ImportError	导入模块错误	模块不存在或路径错误
ArithmeticError	数学错误	下面两个是他的子类
ZeroDivisionError	除以零	1/0
OverflowError	数值错误	算术结果超出数值类型范围
FileNotFoundError	文件不存在	以只读模式打开不存在的文件
IOError OSError	操作系统错误	文件操作出错
RuntimeError	标准错误以外的错误	'raise' 语句触发
AssertionError	断言错误	assert 语句触发
StopIteration	迭代器错误	迭代器没有更多项目供迭代

```

12 1 items had failures:
13   1 of 3 in hello.divide_exact
14 ***Test Failed*** 1 failures.
15

```

8.2 Try & Assert

try,except 用来处理可能出现报错的情况，并提供捕获错误的功能。
常见的错误类型如下，他们都有一个共同的父类 Exception

Listing 8.3: try except

```

1 try:
2     f = open('testfile.txt'); val = bad_val
3 except FileNotFoundError as e: #捕获文件打开错误，并执行以下部分
4     print(e) #输出 No such file or directory: 'testfile.txt'
5 except Exception: #捕获剩余的所有错误，exception 只会执行碰到的第一个
6     print('Sorry. Something went wrong.')
7 else: #若 try 没有出现错误则执行
8     print(f.read()); f.close()
9 finally: #无论是否出错都会执行
10    print('Executing Finally...')
11

```

raise 语句可以在特定情况下手动报出错误

Listing 8.4: Raise an error

```

1 try:
2     a = 2
3     if a == 2:
4         raise Exception
5 except Exception as e:
6     print(e)
7

```

Python 中的 assert 语句类似于 C++assert 断言，不需要导入库，可以与 try, except 结合使用。

Listing 8.5: assert in python

```

1 def area_square(r):
2     assert r > 0, 'A length must be positive'
3     return r * r
4
5 try:
6     area_square(-1)
7 except Exception as e:
8     print(e)
9

```

8.3 Unit Testing

unittest 是标准库中的一个 module，可以用于代码分块 debug。unittest 中有很多 assert 类型，参考 [unittest 文档](#)

Listing 8.6: Unit Testing

```

1 import unittest
2 import app #假设 app.py 是待测的文件
3
4 class TestApp(unittest.TestCase):
5
6     @classmethod
7     def setUpClass(cls): #setUpClass 会在所有 test 之前首先运行
8         print('setUpClass')
9
10    @classmethod
11    def tearDownClass(cls): #tearDownClass 会在所有 test 之后最后运行
12        print('tearDownClass')
13
14    def setUp(self): #setUp 会在每个 test 之前运行一次，注意此函数名不能自定义
15        self.eg1 = 10 #可以创建实例，初始化一些变量以简化 test 操作
16
17    def tearDown(self): #tearDown 会在每个 test 之后运行一次，注意此函数名不能自定义
18        pass #可以删除变量、实例以保证变量不影响下一个 test
19
20    def test_add(self): #检测 add 函数，这里函数名必须以 test 开头
21        result = app.add(self.eg1, 5) #检测 app.py 中的 add 函数
22        self.assertEqual(result, 15)
23
24    def test_div(self): #检测 div 函数
25        self.assertRaises(ValueError, app.div, 10, 0)
26        #检测 div(10, 0) 运行时是否会有 ValueError 报错
27        with self.assertRaises(ValueError):
28            app.div(10, 0) #这里利用文件管理器运行，与上一句效果相同
29
30    if __name__ == '__main__':
31        unittest.main() #运行所有的 test
32

```

用以下语句运行测试

Listing 8.7: run unittest

```

1 python -m unittest test_app.py
2

```

对于网站申请类型的脚本，我们不希望 test 的通过与否取决于网站是否能够连接。为此，可以使用 mock module。

Listing 8.8: mock module

```

1 ###main.py
2 import unittest
3 from unittest.mock import patch
4 from app import get_conn #假设需要检验 get_conn 函数

```

```

5
6 class TestApp(unittest.TestCase):
7
8     def test_get_conn(self):
9         with patch('app.requests.get') as mocked_get: #若连接成功,会继续运行
10             mocked_get.return_value.ok = True
11             mocked_get.return_value.text = "Success!" #捏造一个成功连接,返回 Success
12             result = get_conn() #试运行 get_conn
13             mocked_get.assert_called_with('http://company.com/Schafer/May') #检验 url 的正确性
14             self.assertEqual(result, 'Success') #检验 get_conn 的返回值
15
16 if __name__ == '__main__':
17     unittest.main()
18
19 ###app.py
20 import requests
21
22 def get_conn():
23     response = requests.get('http://company.com/Schafer/May')
24     if response.ok:
25         return response.text
26     else:
27         return "Fail to connect!"
28

```

8.4 jupyter notebook

jupyter notebook 可以用于代码笔记,代码汇报等,提供代码实时运行的功能。可以在 ubuntu 中打开。

Listing 8.9: run jupyter notebook

```

1 jupyter notebook
2

```

9 deep learning

9.1 pytorch 基本操作

9.1.1 数据操作

Listing 9.1: basic manipulation

```

1 import torch
2
3 x = torch.arange(12, dtype=torch.float32) #创建 1 至 12 的向量
4 x = torch.zeros((3, 4)) #零张量
5 x = torch.zeros_like(y) #创建同形状的零张量
6 x = torch.ones((3, 4)) #全为 1 的张量
7 x = torch.ones_like(y) #全为 1 的张量,形状同 y
8 x_repeated = torch.repeat_interleave(input, repeats, dim=None)
9     #input 为要进行重复的张量, repeats 为重复次数 (可以是整数或与 input 某维度长度相同的张量)
10     #dim 指定在哪个维度上重复元素,若未指定, input 会被展平
11 x_repeated = x.repeat(2, 1, 1) #在第 0 维重复 2 变,其它维不变
12 x = torch.randn(3, 5) #均值为 0,标准差为 1 的高斯分布
13 x = torch.tensor([[1, 3, 5, 6], [6, 9, 3, 5]]) #由 List 创建
14 x = torch(2.0) #创建一个标量
15
16 print(x.shape); print(x.numel()) #获取矩阵形状和总大小
17 print(len(x)) #获取向量的维度或矩阵列向量的维度
18 x = x.reshape(3, 4); x = x.reshape(3, 2, -1) #改变矩阵的形状,行优先
19 x = x.unsqueeze(0) #在第 0 维插入一个新的维度
20 x = x.permute((1, 0)) #将坐标轴顺序调换为 (1, 0)
21 A = torch.cat(x, y, dim=0) #将两个列数相同的张量纵向拼接

```

```

22 B = torch.cat(x, y, dim=1) #将两个行数相同的张量横向拼接
23 A = torch.stack(x, y, dim=0) #将两个形状相同的矩阵堆叠, 产生新维度 dim0
24
25 print(x[2, 3]); x[2, 3] = 9 #通过行列索引访问和修改
26 y = x[torch.tensor([1, 3, 5])] #通过一个 tensor 作为 index 取值
27 y = x[[0, 1], [1, 0]] #取出 (0, 1) 和 (1, 0) 两个元素
28 x[0: 2, :] = 12 #将前两行全部改成 12
29 x = x[None, :] #在 axis=0 处加入一个维度
30 x.fill_(10.0) #用 10.0 填充整个张量, 可用于修改 0-dim tensor
31
32 A = x.numpy(dtype=float); B = torch.tensor(A) #torch 和 numpy 转换
33 A = x.detach().numpy() #detach 可以丢掉当前计算图, 不参与梯度计算
34 x.detach_() #原地操作丢掉当前计算图
35 A = torch.tensor(data_frame.to_numpy()); B = pd.DataFrame(A.numpy()) #torch 和 DataFrame 转换
36 val = a.item(); val = float(a); val = int(a) #大小为 1 的张量转标量
37 x = x.type(y.dtype) #将 x 改为 y 的类型
38 x_long = x.long() #修改数据类型为 Long
39
40 print(x + y, x - y, x * y, x / y, x ** y) #同形状对应元素运算, 不同形状会广播
41 print(x == y, x > y) #构建二元张量
42 x[:] = exp(x) #求对应元素的 e 指数, 用切片表示法可以节省内存
43
44 x_sum = x.sum(); x_sum = x.sum(axis=[0, 1]) #求和, 生成单元素张量
45 x_sum_axis0 = A.sum(axis=0) #沿轴 0 降维, 逐列求和, 这里 axis 参数的轴会在降维后消失
46 sum_x = x.sum(axis=1, keepdims=True) #沿轴 1 求和降维, 但是轴数不变, 可以广播
47 x_mean_axis0 = A.mean(axis=0) #沿轴 0 降维, 逐列求均值
48 x_cumsum = torch.cumsum(x, dim=0) #计算张量在指定维度上的累积和, dim=0 即沿行累加
49 x_max = x.max() #返回一个 tensor, 最大值
50 x_argmax = x.argmax() #返回一个 tensor, 最大值对应的 index
51
52 sorted_tensor, indices = torch.sort(x, dim=1, descending=True)
53 #dim 默认为-1, indices 为排序后的张量中每个元素在原张量中 dim 维上的 index
54

```

注意, 直接赋值不会重新分配内存

Listing 9.2: equal and clone

```

1 a = torch.arange(10)
2 b = a; print(a is b, a == b) #返回 True, True
3 b = a.clone(); print(a is b, a == b) #返回 False, True
4

```

9.1.2 读写文件

Listing 9.3: File I/O

```

1 import torch
2 from torch import nn
3 from torch.nn import functional as F
4
5 x = torch.arange(4); y = torch.ones(3, 2)
6 torch.save([x, y], 'x-file') #写入文件
7 x2, y2 = torch.load('x-file') #读取文件
8
9 mydict = {'x': x, 'y': y} #以字典形式读写
10 torch.save(mydict, 'mydict')
11 mydict2 = torch.load('mydict')
12
13 net = Net() #nn.module 的子类
14 torch.save(net.state_dict(), 'net.params')
15 clone = Net(); clone.load_state_dict(torch.load('net.params'))
16 clone.eval()
17

```

9.1.3 线性代数

Listing 9.4: linear algebra

```

1  A = A.T #矩阵转置
2  result = torch.dot(x, y) #求向量点积
3  y = torch.mv(A, x) #求矩阵-向量积
4  C = torch.mm(A, B) #求矩阵-矩阵积
5  C = torch.matmul(A, B) #矩阵乘法, 可以处理多种输入形状, 与基本相同
6  C = A @ B #矩阵乘法, 在处理向量时会尝试将向量作为矩阵进行矩阵乘法, 而 matmul 会返回内积
7  l2 = torch.norm(x) #求向量的  $L_2$  范数
8  l1 = torch.abs(x).sum() #求向量的  $L_1$  范数
9

```

9.1.4 微分

Listing 9.5: differentiation

```

1  x = torch.arange(4.0); x.requires_grad_(True) #注意这里必须是 float 向量而不是 int 向量
2  x = torch.arange(4.0, requires_grad=True) #两句等效, 创建储存梯度的空间
3  y = 2 * torch.dot(x, x) #标量公式, 不能是向量或张量, 可以用 sum()
4  y.backward() #调用反向传播函数
5  print(x.grad) #输出梯度向量
6  x.grad.zero_() #在默认情况下, PyTorch 会累加梯度, 我们需要清除之前的值
7  with torch.no_grad(): #在这个上下文中所有的操作不会计算梯度
8
9  y = x * x; u = y.detach() #求微分时将 y 是为一个常数, 丢弃计算图中 y 的信息
10 z = u * x; z.sum().backward()
11 print(x.grad) #x.grad == u, 全部为 True
12 x.grad.zero_()
13

```

9.1.5 概率

Listing 9.6: Probability and Statistics

```

1  import torch
2  from torch.distributions import multinomial
3  from d2l import torch as d2l
4
5  fair_probs = torch.ones([6]) / 6 #概率向量
6  sample_vec = multinomial.Multinomial(10, fair_probs).sample((500,))
7  #进行 500 次测试, 每次取 10 个样本, 得到五行六列的矩阵
8
9  print(dir(torch.distributions)) #查看所有的函数和类
10
11 X = torch.normal(0, 1, (3, 5)) #生成正态分布数据, 期望为 0, 标准差为 1, 3 行 5 列
12 X = torch.normal(0, 1, A.shape)
13

```

9.2 线性神经网络

9.2.1 线性回归

```

1  import torch
2  from d2l import torch as d2l
3  from torch import nn #neural network
4
5  true_w = torch.tensor([2, -3.4])
6  true_b = 4.2
7  features, labels = d2l.synthetic_data(true_w, true_b, 1000) #初始数据集
8

```

```

9     batch_size = 10
10    data_iter = d2l.load_array((features, labels), batch_size) #将初始数据转换成迭代器
11
12    net = nn.Sequential(nn.Linear(2, 1)) #权重矩阵的形状
13    net[0].weight.data.normal_(0, 0.01) #初始化权重, 更多初始化方法参考torch.nn.init
14    nn.init.normal_(net[0].weight, mean=0, std=0.01) #这两句效果相同
15    net[0].bias.data.fill_(0) #初始化偏置
16    net = nn.Sequential(nn.LazyLinear(256), nn.ReLU(), nn.LazyLinear(10))
17    #延迟初始化, 只给定层的输出维数, 在给定输入后根据输入的维数确定权重的 shape 并初始化
18
19    loss = nn.MSELoss() #定义损失函数为  $L_2$  范数, 更多损失函数参考torch.nn Loss-functions
20    #nn 中的损失函数由 reduction 选项, 可选 'mean', 'none', 'sum' 表示对损失值进行预操作, 默认 'mean'
21    #注意 MSELoss 计算平方误差时不带系数 1/2
22    trainer = torch.optim.SGD(net.parameters(), lr=0.03) #设置优化算法为随机梯度下降
23
24    num_epochs = 3 #迭代次数
25    for epoch in range(num_epochs): #迭代训练
26        for X, y in data_iter:
27            l = loss(net(X), y)
28            trainer.zero_grad()
29            l.backward()
30            trainer.step()
31        l = loss(net(features), labels)
32
33    w = net[0].weight.data #获取权重和偏置数据数据
34    b = net[0].bias.data
35    print(net[0].weight.grad) #获取梯度数据
36    print(net[0].state_dict()) #获取一层的所有数据
37    iter_para = net.parameters() #获取一个迭代器, 包含所有的参数 tensor
38    iter_name_para = net.name_parameters() #获取迭代一, 所有的 (name, parameter) 对
39    print(net.eval()) #输出所有层的树, 可视化
40

```

9.2.2 softmax 回归

由于计算精度的影响 softmax 函数会出现上溢或者下溢的情况, 解决这个问题的技巧是:

- 在计算 softmax 之前, 先从所有 o_k 中减去 $\max(o_k)$ 以解决上溢
- 将计算 softmax 和交叉熵结合在一起, 保留 softmax 作为输出, 并传递取对数的结果

$$\log(\hat{y}_j) = \log\left(\frac{\exp(o_j - \max(o_k))}{\sum_k \exp(o_k - \max(o_k))}\right) = o_j - \max(o_k) - \log\left(\sum_k \exp(o_k - \max(o_k))\right)$$

Pytorch 的 CrossEntropyLoss 可以接受网络的未归一化输出 \mathbf{o} , 并在其内部进行上述操作, 输出 softmax 的结果和交叉熵损失。

Listing 9.7: Softmax

```

1     import torch
2     from torch import nn
3     from d2l import torch as d2l
4
5     train_iter, test_iter = d2l.load_data_fashion_mnist(256) #加载 Fashion-MNIST 图片集
6
7     net = nn.Sequential(nn.Flatten(), nn.Linear(784, 10)) #创建权重变量
8
9     def init_weights(m):
10         if type(m) == nn.Linear:
11             nn.init.normal_(m.weight, std=0.01)
12
13     net.apply(init_weights) #初始化权重
14
15     loss = nn.CrossEntropyLoss(reduction='none')
16     trainer = torch.optim.SGD(net.parameters(), lr=0.1)
17     num_epochs = 10

```



```

18 d2l.train_ch3(net, train_iter, test_iter, loss, num_epochs, trainer)
19 d2l.plt.show()
20

```

9.2.3 多层向量机

Listing 9.8: activation function

```

1 import torch
2 from d2l import torch as d2l
3
4 x = torch.arange(-8.0, 8.0, 0.1, requires_grad=True)
5 y = torch.relu(x) #ReLU 函数, 仅保留正值
6 y = torch.sigmoid(x) #sigmoid 压缩函数
7 y = torch.tanh(x) #双曲正切函数也能将实数域压缩到 (0, 1)
8

```

Listing 9.9: Multilayer

```

1 import torch
2 from torch import nn
3 from d2l import torch as d2l
4
5 net = nn.Sequential(nn.Flatten(), #将矩阵展平成向量
6                     nn.Linear(784, 256), #可选 bias=False, 不带偏置
7                     nn.ReLU(),
8                     nn.Dropout(0.2), #标准暂退法正则化
9                     nn.Linear(256, 10))
10
11 def init_weights(m):
12     if type(m) == nn.Linear:
13         nn.init.normal_(m.weight, std=0.01)
14
15 net.apply(init_weights);
16
17 batch_size, lr, num_epochs = 256, 0.1, 10
18 loss = nn.CrossEntropyLoss(reduction='none')
19 trainer = torch.optim.SGD(net.parameters(), lr=lr)
20 #parameters 表示模型中的所有计算梯度的参数
21 #这里可以设置 L2 范数惩罚
22 trainer = torch.optim.SGD(
23     [{"params":net[0].weight, 'weight_decay':wd}, #给 weight 张量加上 L2 范数惩罚
24     {"params":net[0].bias}], #net 是一个 OrderedDict, 每个元素代表一层
25     lr = lr
26 )
27
28 train_iter, test_iter = d2l.load_data_fashion_mnist(batch_size)
29 d2l.train_ch3(net, train_iter, test_iter, loss, num_epochs, trainer)
30 d2l.plt.show()
31

```

9.3 获取数据集

9.3.1 线性回归数据集

Listing 9.10: Get linear data

```

1 def synthetic_data(w, b, num_examples): #@save
2     """生成y=Xw+b+噪声"""
3     X = torch.normal(0, 1, (num_examples, len(w)))
4     y = torch.matmul(X, w) + b
5     y += torch.normal(0, 0.01, y.shape)
6     return X, y.reshape((-1, 1))
7

```

```

8     def data_iter(batch_size, features, labels): #加载迭代器
9         num_examples = len(features)
10        indices = list(range(num_examples))
11        # 这些样本是随机读取的, 没有特定的顺序
12        random.shuffle(indices)
13        for i in range(0, num_examples, batch_size):
14            batch_indices = torch.tensor(
15                indices[i: min(i + batch_size, num_examples)])
16            yield features[batch_indices], labels[batch_indices]
17

```

9.3.2 Fashion-MNIST 图像分类数据集

Fashion-MNIST 有 10 个类别, 每个类别 6000 张图片的训练数据集和 1000 张图片的测试数据集。

Listing 9.11: Download Fashion-MNIST

```

1     import torch
2     import torchvision
3     from torch.utils import data
4     from torchvision import transforms
5     from d2l import torch as d2l
6
7     d2l.use_svg_display()
8
9     def get_DataLoader_mnist():
10         return 4
11
12     def load_data_fashion_mnist(batch_size, resize=None):
13         trans = [transforms.ToTensor()] #将图像由 PIL 类型转换为 32 位浮点数, 并除以 255 归一
14         if resize: #调整图像大小以保证大小统一
15             trans.insert(0, transforms.Resize(resize))
16         trans = transforms.Compose(trans)
17         mnist_train = torchvision.datasets.FashionMNIST(
18             root='../data', train=True, transform=trans, download=True
19         )
20         mnist_test = torchvision.datasets.FashionMNIST(
21             root='../data', train=False, transform=trans, download=True
22         )
23         return (
24             data.DataLoader(mnist_train, batch_size, shuffle=True,
25                             num_workers=get_dataloader_worker()),
26             data.DataLoader(mnist_train, batch_size, shuffle=False,
27                             num_workers=get_dataloader_worker())
28         ) #采用 4 个进程来读取数据
29

```

9.3.3 数据加载切分

- torch.utils.data 可以用于切分数据。
- torch.utils.data.DataLoader 接收数据集 (需要有 `__len__` 和 `__getitem__`), 并返回一个 iterator
- torch.utils.data.Dataset 是一个基类, 数据集类可以是他的子类 (当然也可以是 list, torch.tensor 等)

Listing 9.12: DataLoader

```

1     import torch
2     import torch.utils.data
3
4     # 自定义数据集
5     class CustomDataset(torch.utils.data.Dataset):

```

```
6         def __init__(self, data, labels):
7             self.data = data
8             self.labels = labels
9
10        def __len__(self):
11            return len(self.data)
12
13        def __getitem__(self, index):
14            # 返回指定索引的数据和标签
15            return self.data[index], self.labels[index]
16
17        # 创建数据和标签
18        data = torch.randn(100, 3) # 100 个样本, 每个样本 3 个特征
19        labels = torch.randint(0, 2, (100,)) # 100 个标签
20
21        # 创建自定义数据集对象
22        dataset = CustomDataset(data, labels)
23
24        # 使用 DataLoader 加载自定义数据集
25        dataloader = torch.utils.data.DataLoader(dataset, batch_size=10, shuffle=True)
26
```