

爬蟲基礎入門-Part 2-1

圖解功能介紹

處理型態



元素

字符串: '不論長怎樣' or "不論長怎樣"

整數: 10, 22, 34, ...

浮點數: 3.1415962, ...

籃子型態

List : [元素, 元素, 元素, 元素, 元素]

dict : { key : value, : , : }

元素級別 元素/籃子

item item item

set : { 元素, 元素, 元素, 元素 }

資料總攬, 不重複

位置

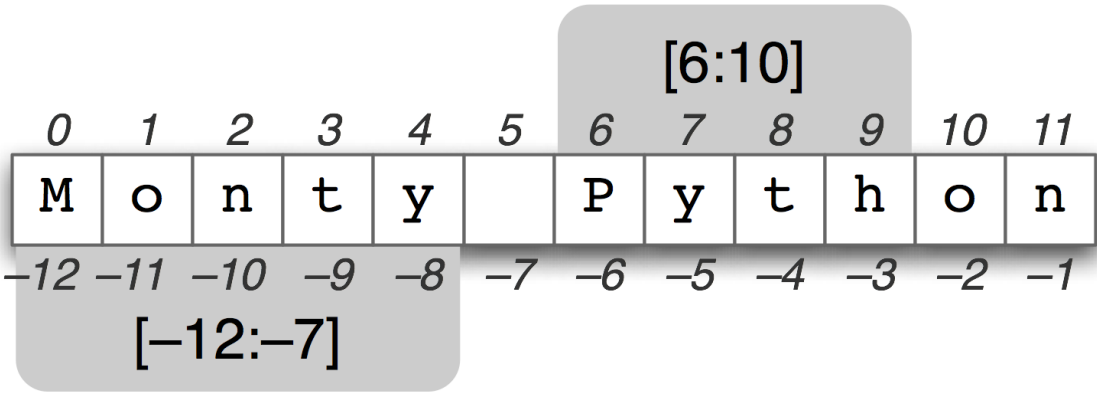
籃子/字符串

[某個位置]

籃子/字符串

[起點位置 : 終點位置]

注意: 終點數字本身不包含



舉例:

```
1 #位置
2 data_list = [10,20,30,40,50,60,70]
3 list_position = 2
4
5 print('第幾個位置:', list_position)
6 print('當個位置的值', data_list[list_position])
7 print('=====')
8 data_str1 = 'this is my data column'
9 data_str2 = "have other data column is '10'"
10 str_position = -10
11 print('str 選定的位置:', str_position)
12 print('字符串1當個位置的值:', data_str1[str_position])
13 print('字符串2當個位置的值:', data_str2[str_position])
14 print('=====')
15 #起點: 終點(不含)
16 str_p = 0
17 end_p = 3
18 my_filter_data = data_list[str_p:end_p]
19 print('我設定的起點/終點:', str_p, end_p)
20 print('我自己設定的區間:', my_filter_data)
21 print('=====')
22 #字符串裡面有變數
23 summary = f'我的數據{data_list}, \n它的型態是{type(data_list)}'
24 # summary = '我的數據{}, \n它的型態是{}'.format(data_list, type(data_list))
25 print(summary)
```

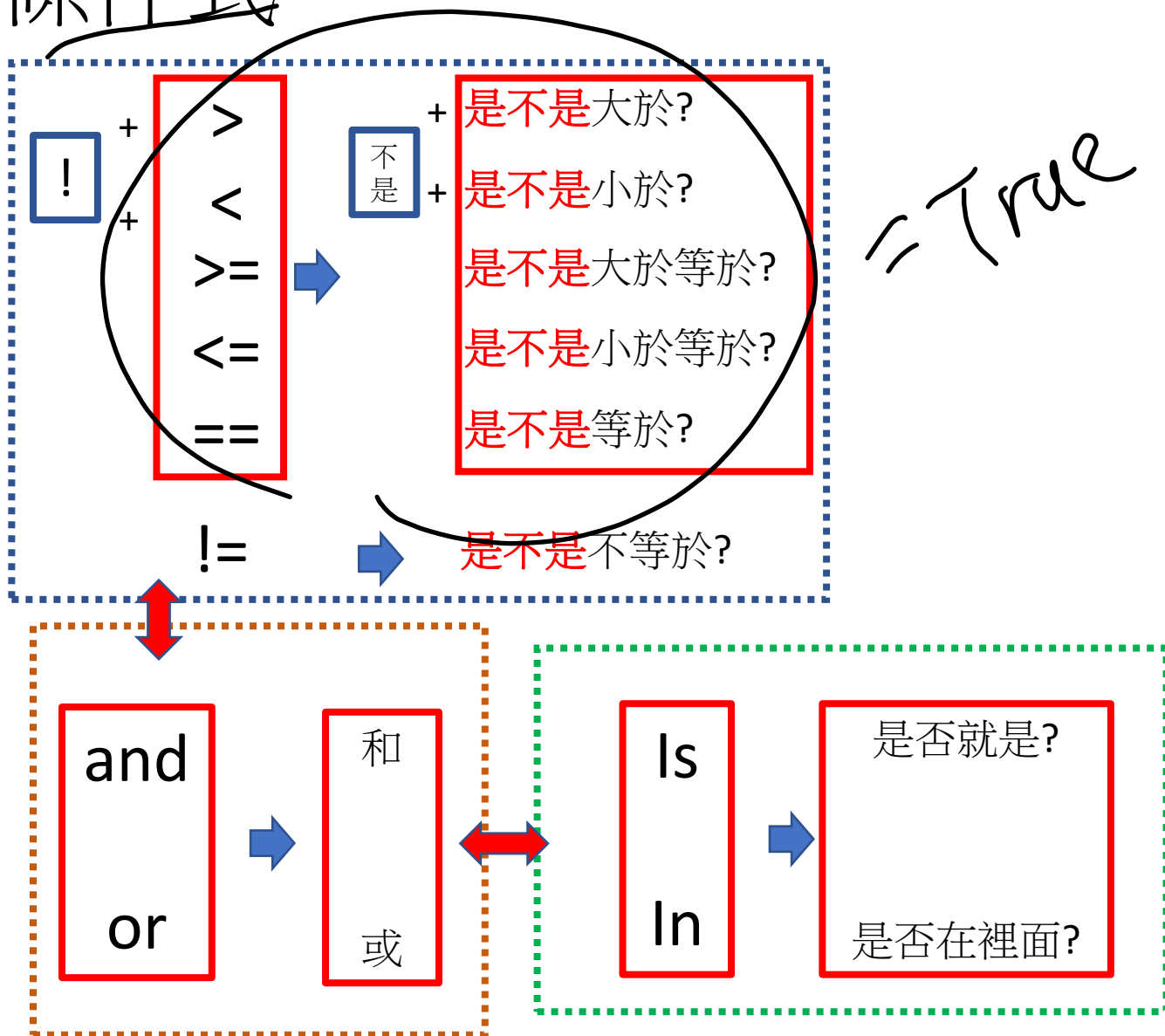
第幾個位置: 2
當個位置的值 30
=====

str 選定的位置: -10
字符串1當個位置的值: a
字符串2當個位置的值: m
=====

我設定的起點/終點: 0 3
我自己設定的區間: [10, 20, 30]
=====

我的數據[10, 20, 30, 40, 50, 60, 70],
它的型態是<class 'list'>

條件式



舉例:

```
1 #條件式-比大小
2 data_number1 = 100
3 data_number2 = 80
4 print('數據1是否大於數據2', data_number1 > data_number2)
5 print('數據1是否小於數據2', data_number1 < data_number2)
6 print('數據1是否等於數據2', data_number1 == data_number2)
7 print('數據1是否大於數據2', data_number1 != data_number2)
8 #條件式- in is
9 data_list1 = [1, 2, 3, 4, 5]
10 check_num = 1
11 print(f'確認是否{check_num}在{data_list1}裡面?', check_num in data_list1)
12 str1 = 'hello'
13 str2 = 'hello!!!'
14 print(f'確認字符串1{str1}是否就是字符串2{str2}?', str1 is str2)
15
16 #條件式 bool 取交集與聯集 TT=T TF=F FF=F
17 print(False and False)
18 print(True and False)
19 print(True and True)
20
21 print(False or False)
22 print(True or False)
23 print(True or True)
24
25 print(data_number1 == data_number2 and check_num < data_list1[-1])
```

數據1是否大於數據2 True
數據1是否小於數據2 False
數據1是否等於數據2 False
數據1是否大於數據2 True
確認是否1在[1, 2, 3, 4, 5]裡面? True
確認字符串1hello是否就是字符串2hello!!!? False
False
False
True
False
True
True
False

判別式

If ----->如果

Elif----->再不然

Else----->沒有如果,一定要做

If  :

↔ 縮排

條件下才能處理的程序

elif  :

↔ 縮排

條件下才能處理的程序

else :

↔ 縮排

上述條件都沒有達成,
就要走這裡的處理程序

舉例:

```
1 #判別式 if elif else
2 data_number = 10
3 Threshold_value = 5
4 if data_number >= Threshold_value:
5     print(f'第一層判斷,{data_number} 比 {Threshold_value}大')
6 else:
7     print(f'第一層判斷,{data_number} 比 {Threshold_value}小')
8
9 print('=====')
10 data_number = 7
11 Threshold_value1 = 5
12 Threshold_value2 = 15
13 if data_number > Threshold_value1 and data_number < Threshold_value2:
14     print(f'你選的數字在於{Threshold_value1}到{Threshold_value2}之間')
15 elif data_number < Threshold_value1 or data_number > Threshold_value2:
16     print(f'你的數字要嘛比{Threshold_value1}小,要嘛比{Threshold_value2}大')
17 else:
18     print('你的數字剛好等於筏值')
19 |
```

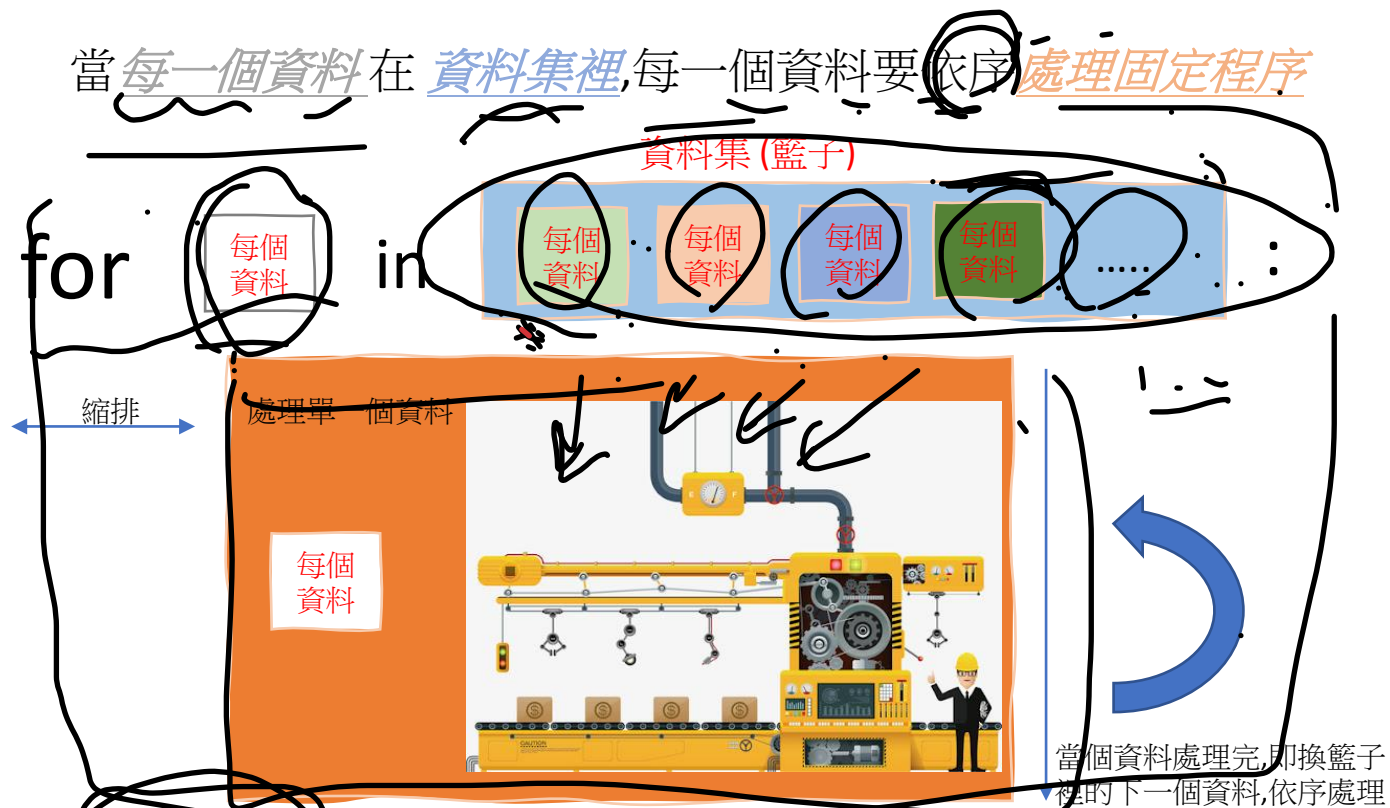
第一層判斷,10 比 5大

=====

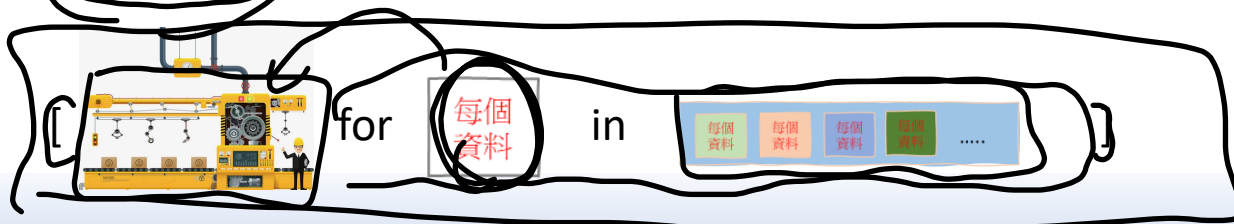
你選的數字在於5到15之間

有限迴圈 (再做一次)

當每一個資料在資料集裡, 每一個資料要依序處理固定程序



懶人寫法:



舉例:

```
1 #有限迴圈
2 data_list = [10,22,34,55,66,80,10]
3 data_treated = []
4 for i in data_list:
5     print('正在處理的數據:',i)
6     y = i+1
7     print('處理後的數據:',y)
8     data_treated.append(y)
9     print('-----每次循環結束-----')
10 print('正常循環處理後:',data_treated)
11 print('-----我是分隔線-----')
12 #懶人寫法, 上述也可以表示成:
13 data_treated2 = [i+1 for i in data_list]
14 print('用懶人方法寫的處理:',data_treated2)
15
16
17 print('-----我是分隔線-----')
18 data_sum = 0
19 processing = 1
20 for i in data_list:
21     print('正在處理的數據:',i)
22     data_sum = data_sum + i
23     print('當下循環的總和:',data_sum)
24     print(f'-----第{processing}次循環結束-----')
25     processing = processing+1
26
27
28 #學習重點:
29 #1. 學習到 for in 有限循環, 每次都是只處理籃子裡面的每一項元素
30 #2. 學習到 list 這個type 有一個功能叫做append, 可以在每次循環添加一些我處理完的資料
31 #3. 學習到 計數器的概念, 可以對數據做累加的寫法, 以及可以把處理完的值量代給自己本身, 供下一次循環使用
32 #4. 請懂懶人寫法
33 # 自我練習:
34 #嘗試 data_list 改為 range(起點, 終點, 間隔)
35
36
```

無限迴圈

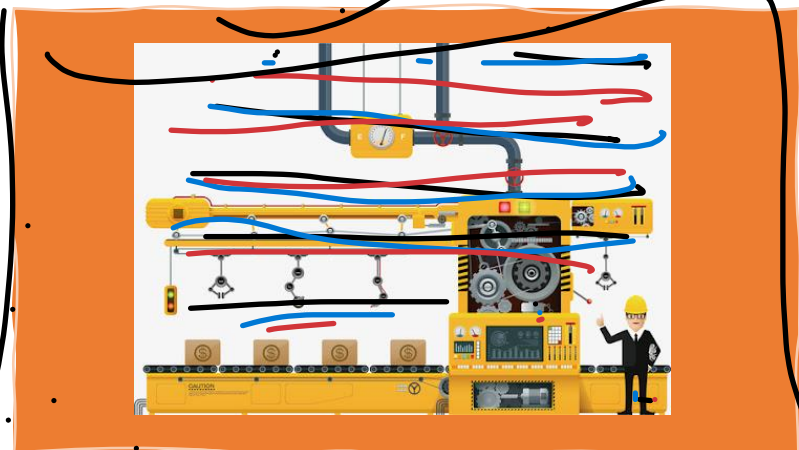
當條件成立時,直接處理固定程序,處理完後,重新繼續判斷,繼續迴圈,直到判斷式不成立:

while

條件式/判斷式

True
↓

縮排



處理完後再確定一次條件, 若依然成立, 即繼續下一個迴圈處理。

舉例:

```
1 #while 迴圈 無限迴圈 (條件迴圈)
2 my_data = 20
3 while True:
4     print('準備處理:', my_data)
5     my_data = my_data + 1
6     print('處理後:', my_data)
7
8 print('====我是分隔線====')
9 my_data = 20
10 processing = True
11 while processing:
12     print('準備處理:', my_data)
13     my_data = my_data + 1
14     print('處理後:', my_data)
15
16     if my_data == 5:
17         processing = False
18         print('====跳出循環====')
19     else:
20         print('還不用跳出循環')
21
22 #學習重點:
23 # 1. 學習到了while 跟for in 迴圈很像, 只是針對條件是否成立, 作為迴圈依據
24 # 2. 學習到了 可以在迴圈中埋入switch 中斷的開關, 並結合判斷式做數據處理, 防止無限迴圈
25 # 3. while 通常判斷式, 不寫即默認是True
```

```
====我是分隔線====
準備處理: 20
處理後: 21
還不用跳出循環
準備處理: 21
處理後: 22
還不用跳出循環
準備處理: 22
處理後: 23
還不用跳出循環
準備處理: 23
處理後: 24
還不用跳出循環
準備處理: 24
處理後: 25
還不用跳出循環
準備處理: 25
處理後: 26
====跳出循環====
```

自訂意函數1

input

$$y = x + 1$$

數據處理的世界

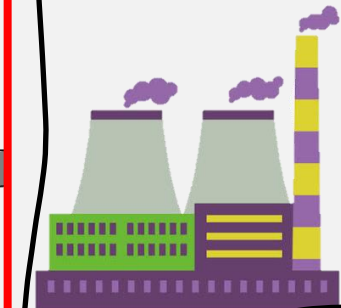
def

縮排

參數傳送員



處理程序 (變數不對外公開)



return

輸出1

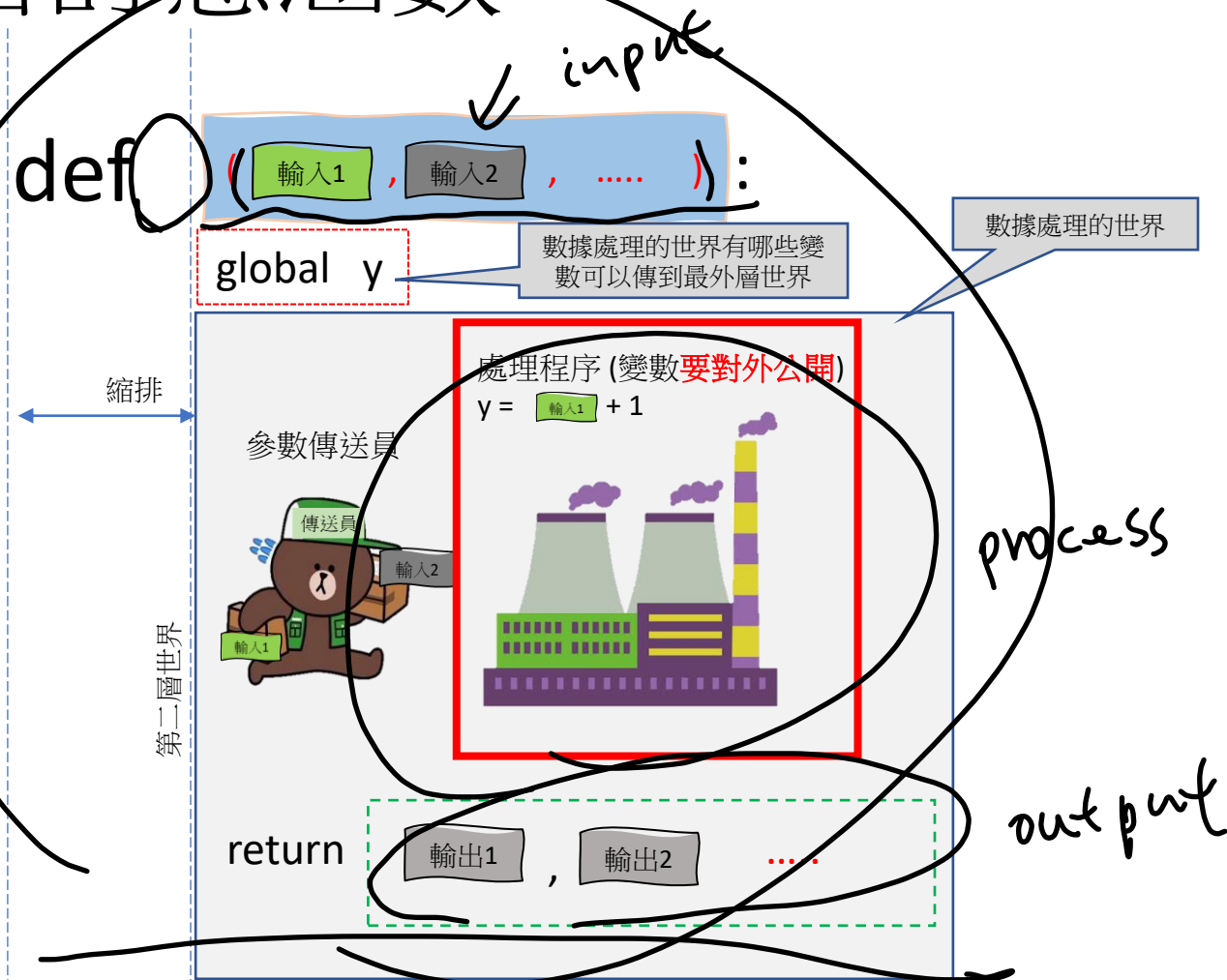
輸出2

.....

舉例:

```
1 #自訂義函數
2 def my_add_fun(num1,num2):
3     result = num1+num1
4     return result
5 #調用
6 my_result = my_add_fun(1,3)
7 print('調用我自訂義的加法函數結果:',my_result)
8 print('-----我是分隔線-----')
9 #(把比大小套的判斷寫成def)
10 def my_judge( data_1 , threshold ):    #偷偷註記: data_1 , threshold 都是可被計算,被比較的數字
11     if data_1>threshold:
12         print(f'您的數字{data_1}大於{threshold}')
13     elif data_1<threshold:
14         print(f'您的數字{data_1}小於{threshold}')
15     else:
16         print(f'您的數字等於{threshold}')
17
18 #調用
19 print(my_judge(100,20))
20
21 print('-----我是分隔線-----')
22 #(加總算平均)
23 def my_cal_mean(data):
24     sum_acc = 0
25     for i in data:
26         sum_acc = sum_acc + i
27
28     mean = sum_acc/len(data)
29
30     return mean
31 #調用
32 data_list = [1,2,3,4,5,6,7,8,9,10]
33 result = my_cal_mean(data_list)
34 print(result)
35 print('-----我是分隔線-----')
36 # 學習重點
37 #學會自訂一函數, 函數一開始的輸入type 要設定好, 方便以當下的type做該功能的調用.
38 #學習自訂一函數,最後return 是一個回傳的動作, 代表之後調用,重代後的值就是這個return的值
39 #學習到有些定義函數,著重在處理數據而已, 不一定要return 回傳值出去
40
41
```


自訂意函數2



舉例:

```
1 #自訂義函數 將函數處理過程中的變數, 也可以傳到最外面, 當作全域變數
2 def my_cal_mean(data):
3     global test
4     sum_acc = 0
5     for i in data:
6         sum_acc = sum_acc + i
7         test = 1000
8     mean = sum_acc / len(data)
9
10    return mean
11 result = my_cal_mean([1, 2, 3, 4, 5])
12 print(result)
13 #嘗試把第三行註解掉, 跑跑看是否能得到內部參數
14 print(test)
```

3.0
1000

自訂意函數3

def (輸入1) :

處理程序 (變數不對外公開)



return

輸出1

=

lambda 輸入1 : 輸出1

舉例:

```
1 #自訂義函數懶人寫法
2
3 #原本正規寫法:
4 def my_add_fun(num1,num2):
5     result = num1+num2
6     return result
7 add_result = my_add_fun(1,2)
8 print('正規寫法的結果:',add_result)
9
10
11 def my_multiply_fun(num1,num2):
12     result = num1*num2
13     return result
14 multiply_result = my_multiply_fun(3,2)
15 print('正規寫法的結果:',multiply_result)
16 print('=====')
17 #懶人寫法 ex1:
18 result = (lambda num1,num2 :num1+num2)(1,2)
19 print('懶人寫法的結果:',result)
20 #懶人寫法 ex2:
21 multiply = (lambda x1,x2:x1*x2)(3,2)
22 print('懶人寫法的結果:',multiply)
23
```

正規寫法的結果: 3

正規寫法的結果: 6

=====

懶人寫法的結果: 3

懶人寫法的結果: 6

什麼是套件包?

你用conda創的環境

My_new_env

pip install

別人寫好的功能包

== 1.3.0

(不加就默認最新版)

pip install

別人寫好的功能包

有些功能會基於別人寫的功能
做疊加, 因此版本匹配很重要

Ex: pip install jupyter
pip install pandas
...

My_new_env2

pip install

別人寫好的功能包

== 2.1.0

(不加就默認最新版)

...

你的電腦環境

重點1:

我們在安裝這些套件,通常會用pip install 或者conda install 做安裝

重點2: 開發時,通常會針對專案需求,創立一個虛擬環境,做版本管控

重點3:我們在開發功能時,不免俗會直接拿別人已經寫好的功能包做功能疊加的應用 (指令集參照官網)

Coding time