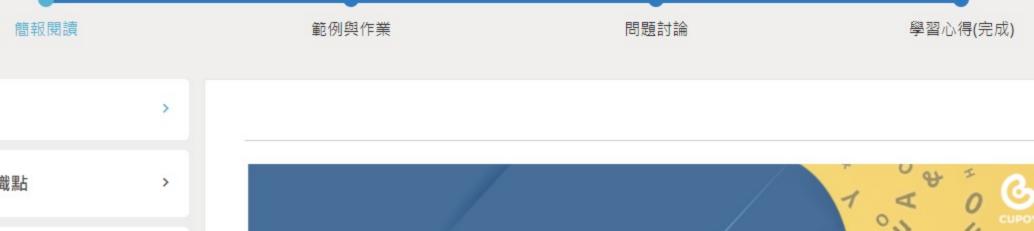
囯

E

anton

D15 pandas Split-Apply-Combine Strategy







(表1)

60

45

有一個函數 goupby 可以一行指令執行以上的邏輯(下圖 3)。

student_id

dtype: float64

sex

boy

girl

score_df.groupby('sex').mean()

69.500000

54.833333

60 45 50 boy 98 55 boy 70 89 boy 79 60 girl

68

70

math_score english_score chinese_score sex

70 boy

55 girl

77 girl

66.000000

61.833333

8	55	77	76	girl
9	25	57	60	girl
10	88	40	43	girl
運用索引	將資料分開再取	双平均(圖2)		
<pre>boy_score_df = girl_score_df = print(boy_score_ print(girl_score</pre>	score_df.loc df.mean())			
math_score english_score chinese_score dtype: float64	69.50 59.25 66.00			
math_score english_score chinese_score	54.833333 65.166667 61.833333			

運用 groupby 平均(圖3)

math_score english_score chinese_score

59.250000

65.166667

認識 Split-Apply-Combine 策略

以剛剛學生資料來分解一下 groupby 的邏輯過程 • Split: 將大的數據集拆成可獨立計算的小數據集(拆成男生、女生資料) • Apply:獨立計算各個小數據集(成績取平均) • Combine: 將小數據集運算結果合併 將 DataFrame 依照 A、B、C 拆成三個小數據集[split],各自計算總合[Apply],合併結果輸出 拆分成 A、B、C 小數據集的方法為 groupby SPLIT APPLY

2

2

Col 1 Col 2

SUM()

Col 1 Col 2

Col 1 Col 2

Col 1 Col 2

4

COMBINE

Col 1

Col 2

9

7



77

57

40

math_score english_score chinese_score

60

girl

2

43 girl

2

2

2

class

mean std

89

77

76

1.5 0.577350

60

55

77

76

60

43

girl

girl

girl

girl

girl

girl

chinese_score

mean

18.191115 66.000000 17.530925

std

	b	oy 1	7	4.000000	61.500000	62.5000	00	
		2	6	5.000000	57.000000	69.5000	00	
	g	irl 1	4:	2.000000	68.666667	65.6666	67	
		2	6	7.666667	61.666667	58.0000	00	
Group	oy 針對村	闌位做	多個	固分析				
								放成績平均以及標準差的計
		-	17.1.000		(中加入計算的) 維輯(M	iean,s	td),此時 groupby 自動會
生似多無過	支欄位(muli	upie coi	umn	S)				
		math s	cono	onalish scone	s chinaca scan		class	
	student id		core	eligitsii_score	e chinese_scor	e sex	CIASS	
	student_id	'			_	_		
	1		50	80	7	0 boy	1	
	2		60	45	5 5	0 boy	2	

69

79

68

70

77

57

40

std

english_score

mean

gi	rl 54.833	333 20.5	66153	65.166667	14.5	79666	61.833333	12.952477	1.5 0.54	17723
Groupby	同時針	計對多 個	固欄位	立做多個	国分;	析				
	以及最高	分的計算	算	I欄位做多	個分	析,例	如,學生	成績資料,	想針對性	別、班級做成
	score	_df.grou	ιрЬу(['	sex','cla	ass']).agg(['mean','	max'])		
			math	_score	е	nglish	_score	chinese_s	core	
			mean	ma	x m	ean	max	mean	max	
	sex	class								
	boy	1	74.00	00000	98 6	31.50000	00 80	62.500000	70	
		_	05.00			7 0000		00 500000		

• Split:將大的數據集拆成可獨立計算的小數據集 • Apply:獨立計算各個小數據集 • Combine: 將小數據集運算結果合併 • Groupby 可以同時針對多個欄位做多個分析

知識點回顧

groupby

groupby

4

• Groupby 可以拆成

6 df

0 -0.278565 1.267586

2 0.011435

1 -1.183920 -0.898350 a

98062 a	ne							
12								
trategy for Da	ta Mining							
tra	tegy for Dat	tegy for Data Mining						

a typical exploratory data analysis, we approach the problem by dividing the data set at some granular level and then aggregating the data at that granularity in order to understand the central tendency.

Similarly, a famous (must read) paper by, Hadley Wickham, outlines split-

analysis. Be it Marketing Segmentation, or any Behavioral Research, we use

apply-combine strategy as one of the most common strategies in data

Col 1 Col 2 1 3 **DataFrame** Col 1 Col 2 Col 1 Col 2

1

2

2 5 3

[Combine]

60 68 55 girl 2 45 70 77 girl

55

25

88

score_df.groupby(['sex','class']).mean()

10

sex class

6

8

10

math_score

mean

sex

70

56

60

45

55

25

88

score_df.groupby(['sex']).agg(['mean','std'])

std

boy 69.500000 20.680103 59.250000

均以	及最高	以 同時針 分的計算 及多分析	草	故多個	分析,例如	,學生	成績資料,
	score_	_df.grou	pby(['sex','	class	']).agg(['me	ean','	max'])
			math_score	2	english_so	ore	chinese_so
			mean	max	mean	max	mean
	sex	class					
	boy	1	74.000000	98	61.500000	80	62.500000
		2	65.000000	70	57.000000	69	69.500000
	girl	1	42.000000	56	68.666667	79	65.666667
		2	67.666667	88	61.666667	77	58.000000

參考資料

import pandas as pd

網站: python/pandas數據挖掘(十四)-groupby,聚合·分組級運算

df = pd.DataFrame({'key1':list('aabba'),

one

two

one

'key2': ['one','two','one','two','one'],

'data1': np.random.randn(5),

'data2': np.random.randn(5)})

-0.207110 b



this technique at some point during our analysis.

下一步:閱讀範例與完成作業