emborh AI共學社群 AI共學社群 > Python資料科學 > D15 pandas Spl ... mbine Strategy

D15 pandas Split-Apply-Combine Strategy



在數據分析中時常會分析不同族群的資料,例如,學生分數資料(如表 1),你想分析男生與女生的各科

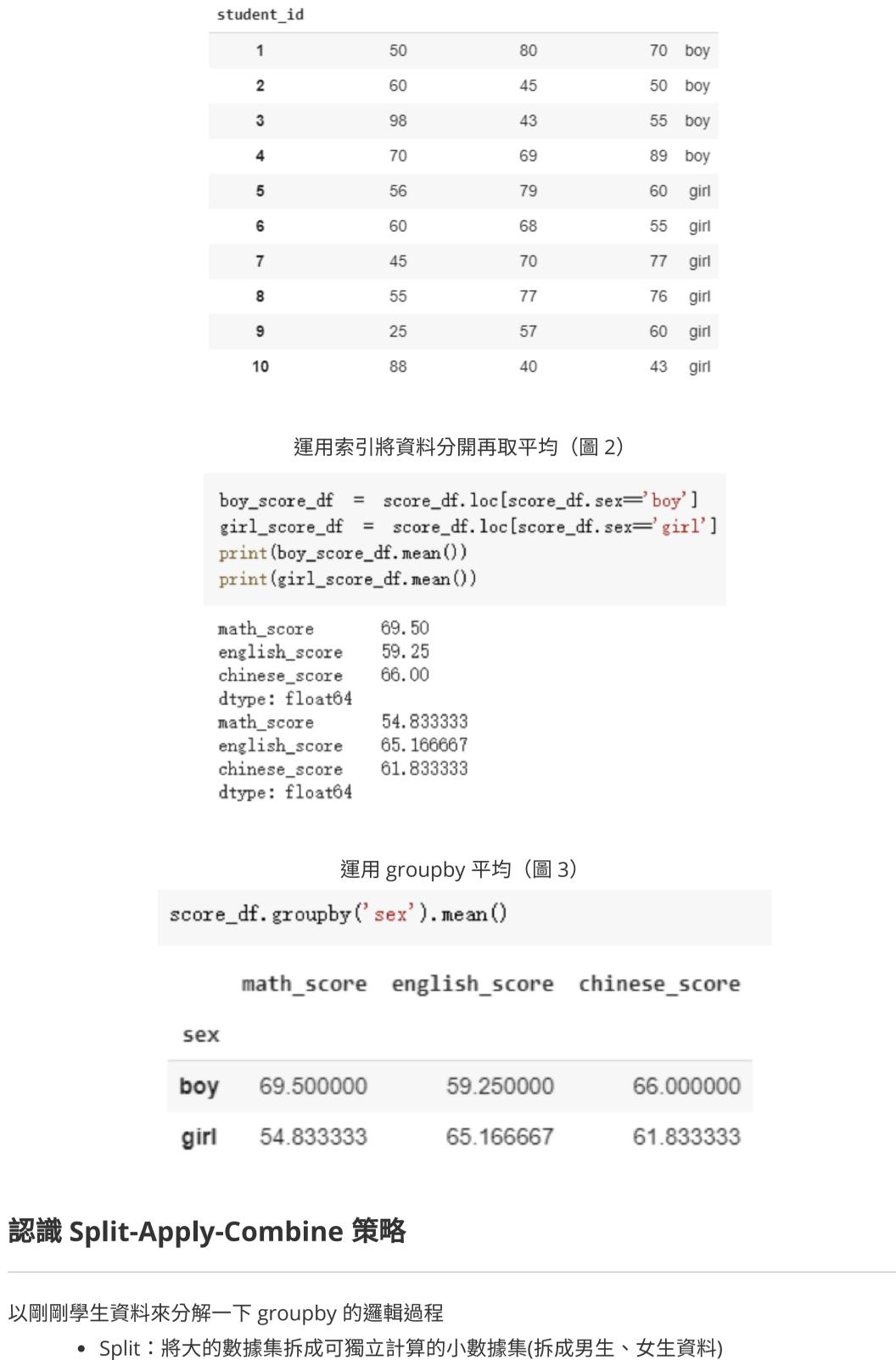
介紹如何透過 pandas groupby 函式實現資料科學的 Split-Apply-Combine 策略

有一個函數 goupby 可以一行指令執行以上的邏輯(下圖 3)。

認識 groupby

(表 1) math_score english_score chinese_score sex student id

差異,前幾天有教到檢索可以將資料分成男生資料與女生資料,在將各資料算平均值(如圖 2),在這裡



[Combine]

DataFrame

Col 1 Col 2

• 拆分成 A、B、C 小數據集的方法為 groupby

SPLIT APPLY

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

SUM()

Col 1 Col 2

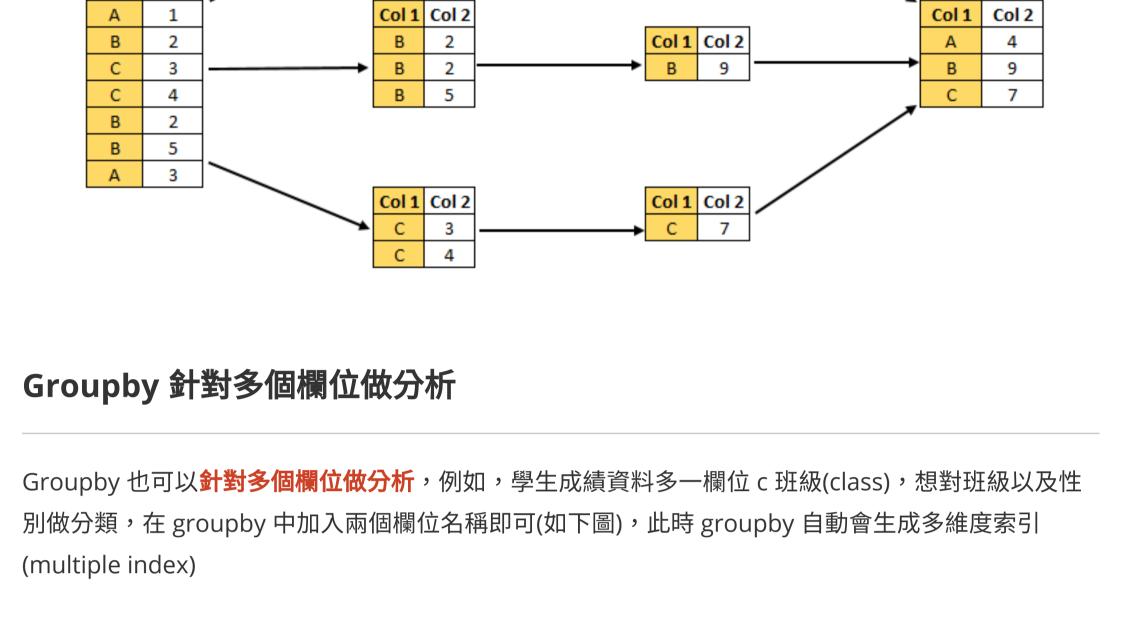
4

COMBINE

Col 1 Col 2

• 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



math_score english_score chinese_score sex class

80

45

43

69

79

70 boy

50 boy

55 boy

89 boy

60 girl

98 3 70

student_id

1

2

5

55 girl 60 6 68 77 girl 45 70 55 77 76 girl 8 25 9

50

60

56



	3	98	43	55	boy 1
	4	70	69	89	boy 2
	5	56	79	60	girl 1
	6	60	68	55	girl 2
	7	45	70	77	girl 1
	8	55	77	76	girl 2
	9	25	57	60	girl 1
	10	88	40	43	girl 2
score_df.groupby(['sex']).agg(['mean','std'])					

english_score

boy 69.500000 20.680103 59.250000 18.191115 66.000000 17.530925

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

max mean

math_score

67.666667

• Split:將大的數據集拆成可獨立計算的小數據集

mean

std

girl 54.833333 20.566153 65.166667 14.579666 61.833333 12.952477 1.5 0.547723

• Groupby 也可以**同時針對多個欄位做多個分析**,例如,學生成績資料,想針對性別、班級做成

45

50 boy

chinese_score

mean

std

english_score chinese_score

88 61.666667 77 58.000000 76

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

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

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

max mean

69 69.500000

class

mean std

max

1.5 0.577350

60

2

math_score

績平均以及最高分的計算

• 合併了多欄位以及多分析

sex

std

Groupby 同時針對多個欄位做多個分析

74.000000 98 61.500000 80 62.500000 70 70 57.000000 65.000000 girl 42.000000 56 68.666667 79 65.666667 77

• Apply:獨立計算各個小數據集

1 import pandas as pd

data2

-0.207110 b

Split-Apply-Combine Strategy for Data Mining

1 -1.183920 -0.898350 a

• 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

sex class

- 知識點回顧 • Groupby 可以拆成
- 網站:<u>python/pandas數據挖掘(十四)-groupby,聚合,分組級運算</u>

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

key1 key2

one

two

one

data1 **0** -0.278565 1.267586 a

參考資料

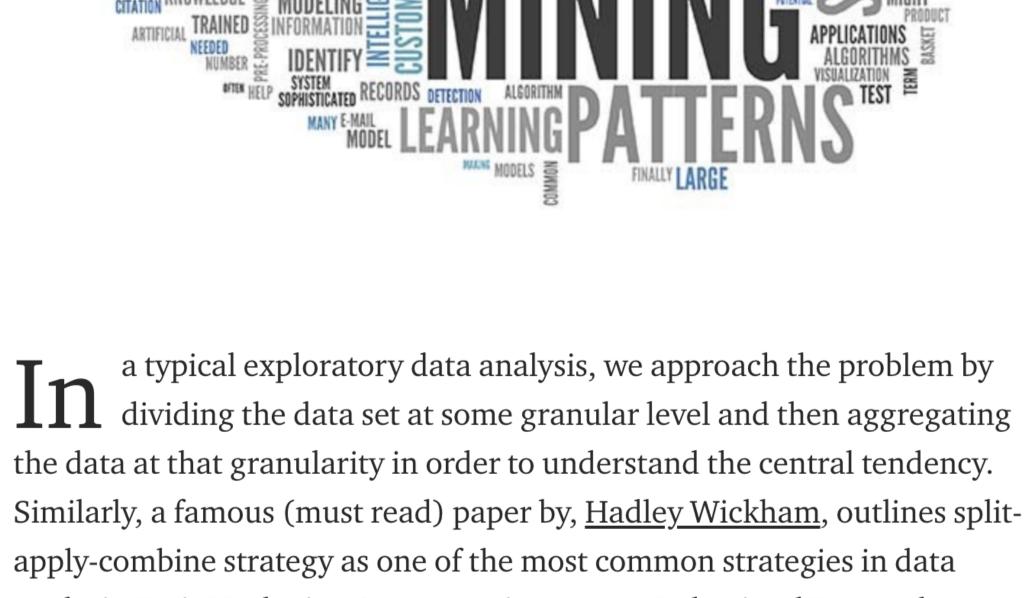
groupby

groupby

6 df

2 0.011435

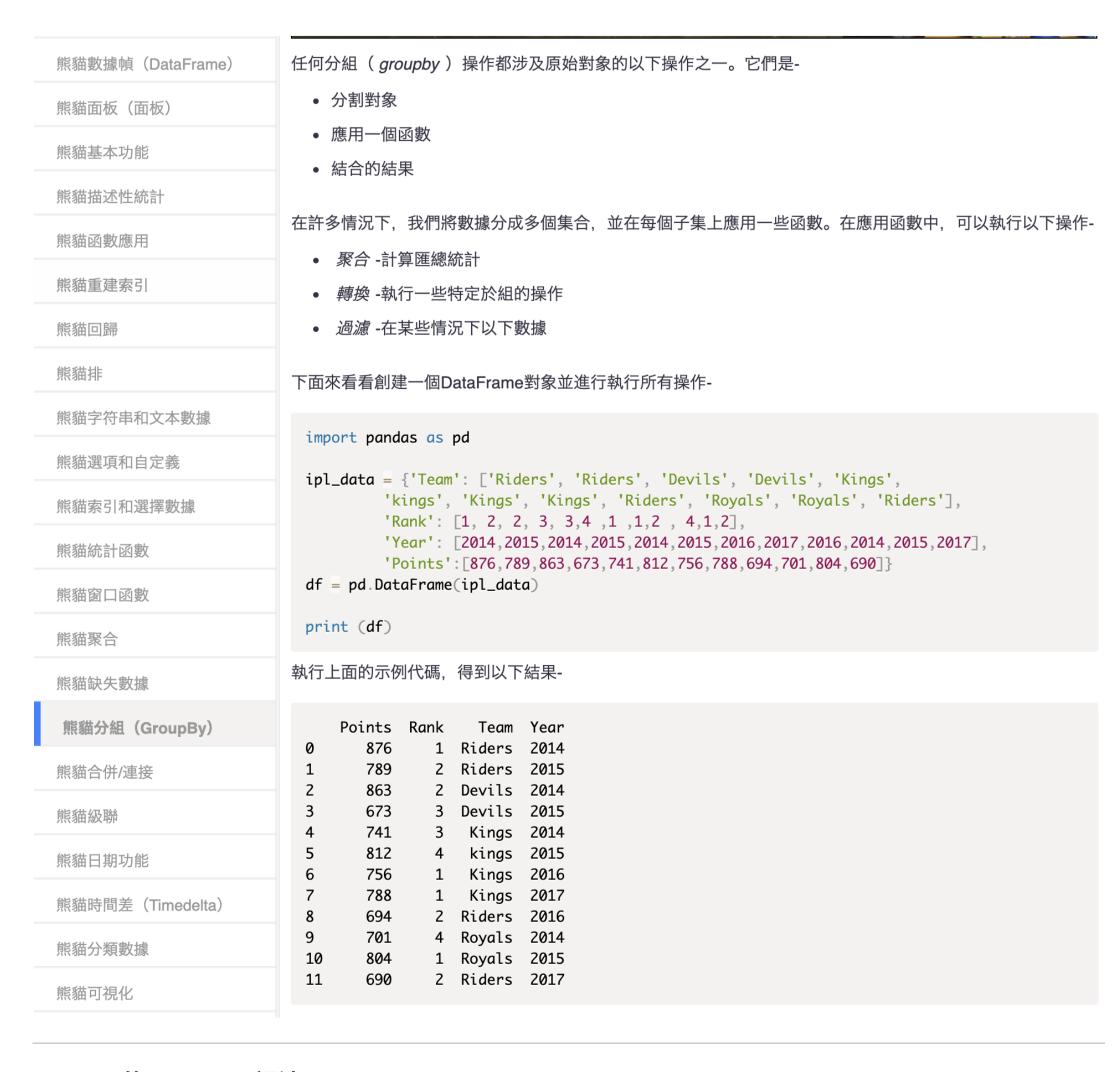
-1.706337 b 3 1.570595 two 4 1.149452 -1.098062 a one

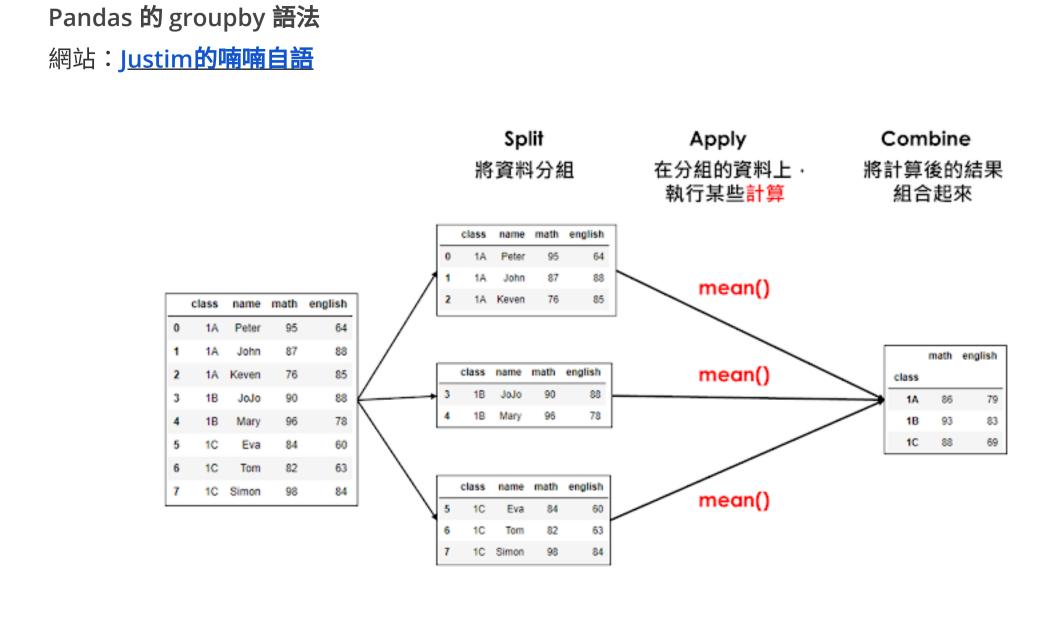


analysis. Be it Marketing Segmentation, or any Behavioral Research, we use this technique at some point during our analysis. 延伸閱讀

Pandas 分组(GroupBy)

網站:<u>易百教程</u>





下一步:閱讀範例與完成作業