അവ AI共學社群 我的 AI共學社群 > Python資料科學 > D13 pandas 統計函式使用教學 D13 pandas 統計函式使用教學 囯 E 問題討論 簡報閱讀 範例與作業 學習心得(完成) > 重要知識點 Python資料科學程式馬拉松 統計函式 ▶ Pandas 統計函式使用教學 統計函式:平均值mean() > 統計函式:加總sum()·個、 數count() 陪跑專家: Hong 重要知識點 重要知識點 • 應用統計函式 • 自定義的行或列函式應用 統計函式 在生活中常聽到以下情況 1. 台灣平均薪資為 XXX 2. 今年指考最高分為 XXX 3. 今年台大最低入取分數為 XXX 4. 6 個標準差的良率 因為數據很多的情況下時常使用敘述統計量來描述數據的分佈與統計量,在資料分析中常拿來對資料 做初步的了解。接下來我們以 pandas 的 DataFrame 資料來做統計函式的介紹。 統計函式:平均值mean() 今天都以班上學生國文、英文、數學分數的資料(右表)為例子介紹各個統計函數。 首先是最常使用到的平均值 mean(), pandas 可針對指定欄位算平均值,如果沒指定會對全部欄位算 平均值。 [69] #指定欄位算平均 score\_df.math\_score.mean() 60.7 [70] #全欄位算平均 score\_df.mean() math\_score 60.7 english\_score 62.8 chinese\_score dtype: float64 score\_df math\_score english\_score chinese\_score student\_id 70 60 45 50 55 70 69 56 79 60 68 55 45 70 77 77 76 55 57 25 60 10 88 如果今天想要算每個學生的總平均分數怎麼辦? Pandas 統計函式中有參數 axis=0 為行運算, axis=1 為列運算,此參數適用在之後介紹的統計函式。 axis=1 col0 col1 col2 col3 col4 axis=1 row0 row1 axis=0 row2 axis=0 [71] #學生平均分數 score\_df.mean(axis=1) student\_id 66.666667 51.666667 65.333333 76.000000 65.000000 61.000000 64.000000 69.333333 47.333333 57.000000 10 dtype: float64 統計函式:加總sum(),個數count() 加總:計算總和,時常用在計算家庭開銷 個數:計算個數,時常用在出遊時的點名 以下利用加總算出學生3科總分,利用各數計算出應考人數 [74] #本次各科考試人數 [73] #學生3科總分數 score\_df.sum(axis=1) score\_df.count() student\_id math\_score 10 english\_score 10 200 1 chinese\_score 155 10 196 dtype: int64 4 228 195 6 183 192 208 8 9 142 10 171 dtype: int64 統計函式:中位數median() 中位數通常使用在有否贏過 50% 的數據,假如薪資中為數為 4 萬,超過 4 萬即為贏過 50% 的人,反 之亦然。 中位數:通過把所有觀察值高低排序後找出正中間的一個作為中位數。 如果觀察值有偶數個,則中位 數不唯一,通常取最中間的兩個數值的平均數作爲中位數。 以利用中位數算出各科中位數,如果今天數學考了60分超過了中位數的58分,我就可以說我數學贏 過了全班一半的同學。 [75] #各科中位數分佈 score\_df.median() 58.0 math\_score english\_score 68.5 chinese\_score 60.0 dtype: float64 統計函式:百分位數quantile() 百分位數使用在觀察數據百分比,最常運用到的是升學分數的百分位數。 百分位數:將一組數據從小到大排序,並計算相應的累計百分位,則某一百分位所對應數據的值就稱 為這一百分位的百分位數。如果百分位數設定在 50% 即為中位數。 以下計算 75% 的百分位數,如果我今天國文分數為 75分,我可以說我的國文贏過班上 75% 的同學 [77] #各科百分位數分佈(75%) score\_df.quantile(0.75) math\_score 67.50 english\_score 75.25 chinese\_score 74.50 Name: 0.75, dtype: float64 統計函式:最大值max()、最小值min() 最大最小值時常拿觀察極端值,也可以檢視資料的資料最小與最大分佈。 其中最小值常常拿來當通過門檻,例如:大學入取分數最低幾分。 以下計算全班各科最高與最低分: [79] #各科最小值 [78] #各科最大值 score\_df.min() score\_df.max() 25 math\_score 98 math\_score 40 english\_score english\_score 80 chinese\_score 43 chinese\_score 89 dtype: int64 dtype: int64 統計函式:標準差std(),變異數var() 標準差:在機率統計中最常使用作為測量一組數值的離散程度之用。一個較大的標準差,代表大部分 的數值和其平均值之間差異較大;一個較小的標準差,代表這些數值較接近平均值。 變異數:為標準差平方 以下計算出標準差,可以發現國文分數標準差比數學分數標準差來的小,所以國文的分散程度比較 小,也可以說國文分數較為集中。 [81] #各科變異數 [80] #各科標準差 score\_df.var() score\_df.std() math\_score 434.900000 math\_score 20.854256 english\_score 237.733333 english\_score 15.418603 chinese\_score 14.151953 200.277778 chinese\_score dtype: float64 dtype: float64 統計函式:相關係數corr() 相關係數:皮爾遜積矩相關係數 ( Pearson product-moment correlation coefficient ) 用於度量兩個 變數X和Y之間的相關程度(線性相依)。在自然科學領域中,該係數廣泛用於度量兩個變數之間的線 性相依程度。相關係數的值介於 -1 與 +1 之間,即 -1≤r≤+1。其性質如下: 1. 當 r>0 時,表示兩變數正相關,r<0 時,兩變數為負相關,r=0 時,表示兩變數間無線性相關關 2. 一般可按三級劃分: |r|<0.4 為低度線性相關; 0.4≤|r|<0.7 為顯著性相關; 0.7≤|r|<1為高度 線性相關。 1 8.0 0.4 -0.4-0.81 1 1 -1 -1 -1 可以發現說英文相對數學相關係數為 -0.53, 可以解釋說英文跟數學有負的高度線性相關, 可以說明此 班學生數學越高分英文越低分,另外國文相對英文相關係數為 0.68 為正向高度相關, 說明此班學生英 文越高分國文越高分。 [82] #各科之間的相關係數 score\_df.corr() math\_score english\_score chinese\_score -0.532708 math\_score 1.000000 -0.314552 english\_score -0.532708 1.000000 0.682340 1.000000 chinese\_score -0.314552 0.682340 自訂義的行或列函式應用 apply() 你有時候可能儲覺得說前面的統計函式不足以表達資料的特性,此時你可以使用 apply 做自定義的函 式。 像是學校最常使用的加分方式為開根號乘以十,例如:我考 49 分加分過後 √49 × 10 = 70,這種方程 式沒辦法在統計函式中算出來,需要藉由 apply 中 lambda 的函式達成。 其中 lambda x 相當於數學式中的 f(x) = √x × 10 [84] #各科開根號乘以十  $\rightarrow$  f(x) =  $\sqrt{x} \times 10$ score\_df.apply(lambda x : x\*\*(0.5)\*10) math\_score english\_score chinese\_score student\_id 89.442719 83.666003 70.710678 70.710678 77.459667 67.082039 98.994949 65.574385 74.161985 4 83.666003 83.066239 94.339811 88.881944 77.459667 74.833148 6 77.459667 82.462113 74.161985 7 67.082039 83.666003 87.749644 8 74.161985 87.749644 87.177979 77.459667 50.000000 75.498344 10 93.808315 63.245553 65.574385 apply 也適用先前統計函式,可以用下列程式碼看出兩個計算邏輯是等價的。 [87] #各科加總 [86] #各科加總apply score\_df.apply(sum, axis=1) score\_df.sum(axis=1) student\_id student\_id 200 1 200 1 2 155 2 155 3 3 196 196 228 228 5 195 195 183 6 6 183 7 192 192 8 8 208 208 9 9 142142 171 171 10 10 dtype: int64 dtype: int64 參考資料 Pandas 描述性統計 網站:程式教程網 有很多方法用來集體計算DataFrame的描述性統計信息和其他相關操作。 其中大多數是sum(), mean()等聚合函數,但其中一些,如sumsum(),產生一個相同大小的對象。 一般來說,這些方法 採用軸參數,就像ndarray.{sum,std,...},但軸可以通過名稱或整數來指定: • 數據幀(DataFrame) - 「index」(axis=0,默認), columns(axis=1) 下面創建一個數據幀(DataFrame),並使用此對象進行演示本章中所有操作。 示例 import pandas as pd import numpy as np #Create a Dictionary of series d = {'Name':pd.Series(['Tom', 'James', 'Ricky', 'Vin', 'Steve', 'Minsu', 'Jack', 'Lee', 'David', 'Gasper', 'Betina', 'Andres']), 'Age':pd.Series([25,26,25,23,30,29,23,34,40,30,51,46]), 'Rating':pd.Series([4.23,3.24,3.98,2.56,3.20,4.6,3.8,3.78,2.98,4.80,4.10,3 #Create a DataFrame df = pd.DataFrame(d) print df Pandas 函數應用 網站:程式教程網 表格函數應用 可以通過將函數和適當數量的參數作爲管道參數來執行自定義操作。 因此,對整個DataFrame執 行操作。 例如,爲DataFrame中的所有元素相加一個值2。 adder 函數 adder函數將兩個數值作爲參數相加並返回總和。 def adder(ele1,ele2): return ele1+ele2 現在將使用自定義函數對DataFrame進行操作。 df = pd.DataFrame(np.random.randn(5,3),columns=['col1','col2','col3']) df.pipe(adder,2) 下一步:閱讀範例與完成作業