emboh AI共學社群 我的

匰

D21 運用實際資料集進行資料視覺化練習





- 當你把數據轉換成了規範的格式,也已經採用了適當的統計和分析,接下來就是展示結果的時 候了,這時候數據可視化排上了用場 在可視化分析中,經常會遇到多個數據分布之間的比較,
- - 分布不同,用到的表達方式也不一樣。

自差異的細微差別。比如,在統計過程中,不同標準的數據集會有怎樣的差別,或者,如何通 過分析來改善評分功能。

可視化的好處

• 瞭解有關資料集屬性 # 我們可以使用 info()或是 descript() 方法瞭解有關資料集屬性的更多資訊。特別是行和列 的數量、列名稱、它們的數據類型和空值數

• 在對不同的分布數據進行比較時,通常有兩種形式,要麼突出異常值的差異,要麼突出它們各

- 導入數據集 Seaborn 在庫中附帶了幾個重要的數據集。安裝 Seaborn 後,數據集會自動下載。 您可以使用這些資料集中的任何一個進行學習。借助以下函數,您可以載入所需的數據集。
- load_dataset() • Seaborn 可以直接把 PANDAS 的 dataframe 當成資料匯入 • 本日範例,我們以Seaborn 內建的 IRIS 資料集做範例

導入必要的程式庫

import pandas as pd

import seaborn as sns

- from matplotlib import pyplot as plt #取得鳶尾花資料集
- df = sns.load_dataset('iris')
- <class 'pandas.core.frame.DataFrame'> RangeIndex: 150 entries, 0 to 149
- Data columns (total 5 columns): Column Non-Null Count Dtype

1 df.info()



sepal_length 150 non-null

float64

petal_length

sepal_length

sepal_width

petal_width

2 ·

petal_length

直方圖:

KDE:

setosa

setosa

#當一個或兩個正在研究的變數是分類的時,我們使用像條帶線()、swarmplot()等的圖。 # 查看到每個物種petal_length的差異。但是,散點圖的主要問題是散點圖上的點重疊。 sns.stripplot(x = "species", y = "petal_length", data = df) 6 petal_length w

分類式的變數 #上述散點圖的主要問題是散點圖上的點重疊。我們使用"抖動"參數來處理此類方案。 # 抖動會為數據添加一些隨機雜訊。此參數將沿分類軸調整位置。 sns.stripplot(x = "species", y = "petal_length", data = df, jitter=True)

versicolor

#另一個可以用作「抖動」的替代選項是函數群圖()。

此函數將散點圖的每個點都放在分類軸上,從而避免重疊點

sns.swarmplot(x = "species", y = "petal_length", data = df)

species

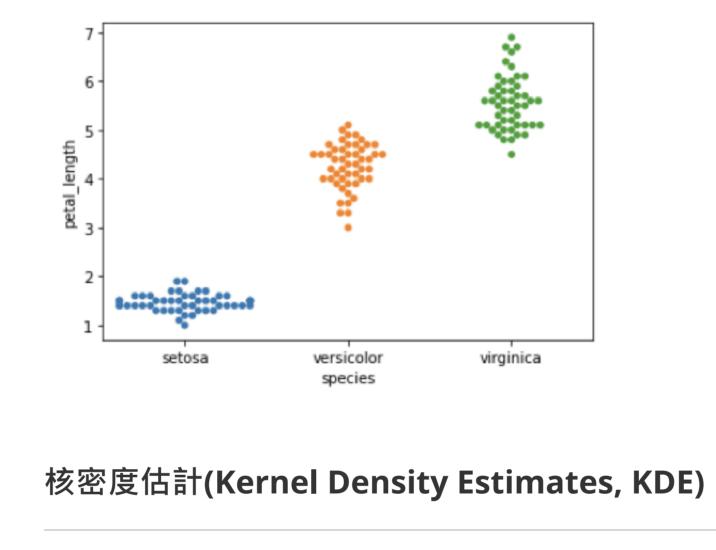
virginica

所謂核密度估計,就是採用平滑的峰值函數("核")來擬合觀察到的數據點,從而對真實的概率分佈曲線

versicolor

species

virginica



進行類比。以下面3個數據點(5,10,15)的一維數據集為例:

0.0

0.30

0.25

0.15

0.10

內核密度估計是估計變數分佈的非參數化方法。

sns.set_style("ticks")

7.5 -

th 6.5

<u>=</u>1 6.0

ਲ੍ਹੇ 5.5

5.0

4.5

4.5 -

4.0

sepal_width

2.5

2.0

2.5 -

2.0 hetal width 1.5

Plotly)

Seaborn

介紹一下這三種視覺化套件

網站:<u>blog.csdn</u>

分類專欄: python

import matplotlib.pyplot as plt import seaborn as sns

原创 u_7890 2019-04-1416:15:23 🧿 2860 🏚 收藏 2

df該dataframe的行名是中文,可以加上下面代碼中紅色這句,就可以顯示出中文

sns.set(font="simhei")#遇到標籤需要漢字的可以在繪圖前加上這句

3. Plotly

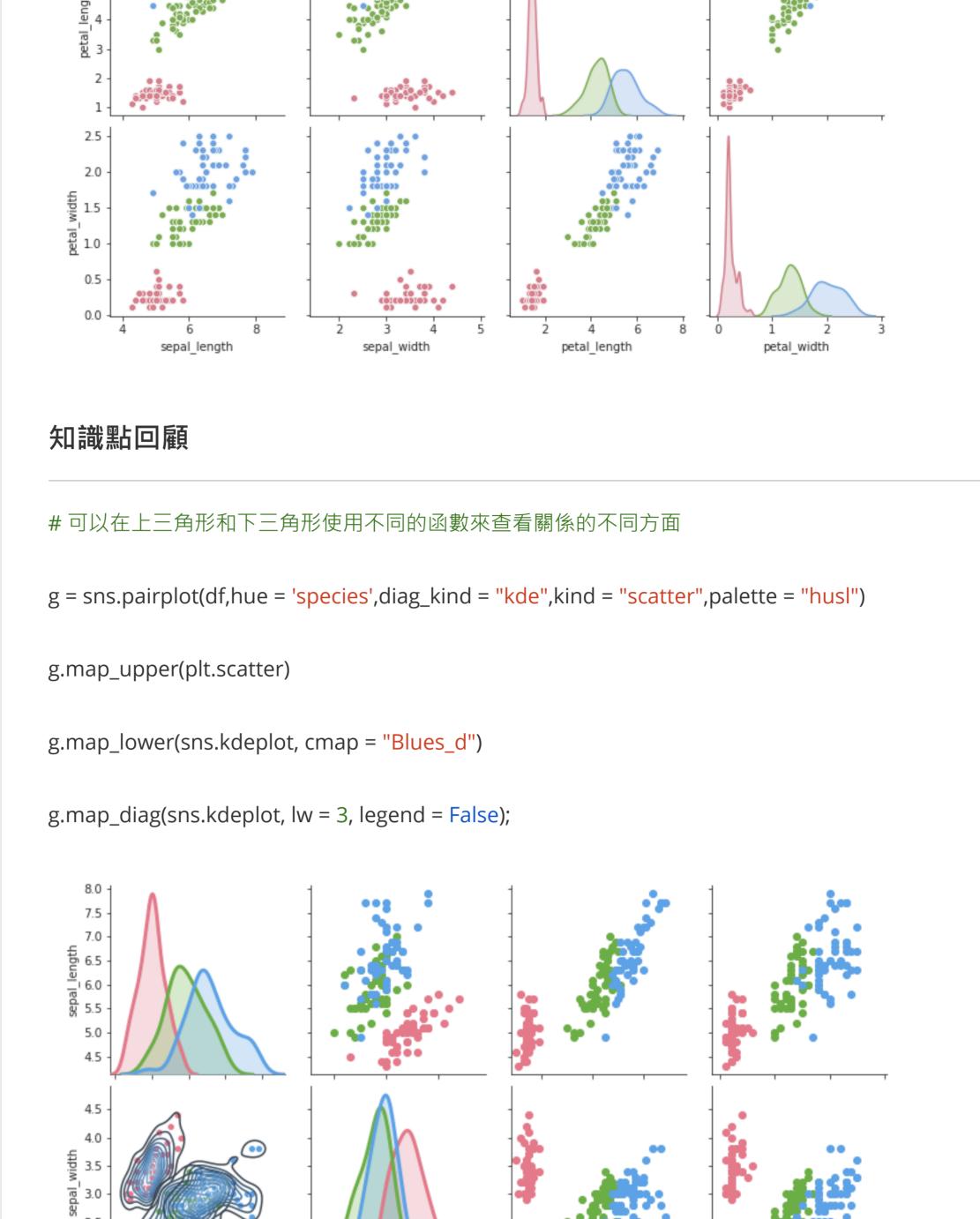
0.05 0.00 所有平滑的峰值函數均可作為KDE的核函數來使用,只要對歸一化后的KDE而言(描繪在圖上的是數據 點出現的概率值),該函數曲線下方的面積和等於1即可 — 只有一個數據點時,單個波峰下方的面積為 1,存在多個數據點時,所有波峰下方的面積之和為1。 概而言之,函數曲線需囊括所有可能出現的數據值的情況。 內核密度估計是用來繪製密度數據,這將更加準確地反映總體的基本變量。 1. 在數據點處為波峰 2. 曲線下方面積為1 #可以觀察每個情節的變化。繪圖採用矩陣格式,其中行名表示 x 軸,列名稱表示 y 軸。 # 對角線圖是內核密度圖,其中其他圖是散點圖

sns.pairplot(df,hue = 'species',diag_kind = "kde",kind = "scatter",palette = "husl")

species

species

versicolor virginica



petal_length sepal_width petal_width 延伸閱讀 使用Seaborn 進行可視化 [資料分析&機器學習] 第2.5講:資料視覺化(Matplotlib, Seaborn, Plotly) 網站:<u>medium</u> • 針對 Matplotlib · Matplotlib & Pandas 帶入實例 • 針對 Seaborn · Seaborn & Pandas帶入實例

資料視覺化除了最後一步呈現你的成果之外,還可以在分析的過程中用資 料視覺化來看出一些insight,比方說用熱點圖來看你的Deep learning的 model是對圖片中哪一部分的看得較重要,或是可以降維之後將資料視覺 化去看資料在空間中的分佈,來決定下一步的分析要怎麼做。 Python資料視覺化主要有三大套件: 1. Matplotlib

其他還有像是Bokeh, ggplot...十幾種Python視覺化套件,以及更進階的BI Tool(Tableau, Spotfire, MicroStrategy)。如果你去大公司上班,ex: 台積電, Yahoo...都會購買這些要價不菲的BI Tool,對於資料處理&視覺化功能有更 豐富的運用。但對於一般使用者只要先掌握上述三個就夠了。今天就要來

[資料分析&機器學習]第2.5講:

資料視覺化(Matplotlib, Seaborn,

說明如何解決中文字在圖形得顯示問題 seaborn畫熱力圖時行名中包括中文顯示成方框

Seaborn 畫熱力圖時行名中包括中文顯示成方框

f, ax = plt.subplots(figsize=(10,10)) sns.heatmap(df, annot=False, ax=ax) 👍 點贊2 💬 評論 🙋 分享 🛊 收藏2 🗍 手機看 😩 打賞 … ■ 關注 🛑 一鍵三連 seaborn畫熱力圖坐標軸怎麼顯示中文和英文 熱力圖中橫坐標/擬坐標/包括中文,可用sns.set(font='LiSu')其中LiSu是字體。如何查看matplotlib包括哪些字體: from matplotlib.font_manager import fon... 【python】seaborn畫熱力圖時行名中包括中文顯示成方框 處理demo import matplotlib.pyplot as plt import <mark>seaborn</mark> as sns sns.set(font="simhei")#遇到標籤需要漢字的可以在繪圖前加上這句f, ax = plt.subplots(fi...

下一步:閱讀範例與完成作業