

系級：四資工四乙 學號：B10615020 姓名：林哲旭

期末 Python 程式：FIFA19 選手能力評鑑

一、設計此程式的動機

身為一個電玩遊戲的愛好者，能夠精打細算的從遊戲資訊中挑取最高價值的角色是必要的技能，如今在大數據的課程學會使用 Python 製作簡易程式，不妨自行設計一個能夠跟自己興趣相符的程式，除了增加程式撰寫能力的基礎，還可以讓身在遊戲中的我獲得一份成就感，因此上 Kaggle 網站尋找到 FIFA19 足球遊戲的選手資料集，藉此來分析選手評分。

二、程式設計目標

撰寫 Python 程式，讀取 FIFA19 的選手相關資料，總共有 88 種角色屬性，包含人名、年紀、隊名、各種能力值等等，根據所有角色屬性把不重要的屬性篩選掉，例如：照片，國旗，背號等等。

ID	Name	Age	Photo	Nationality	Flag	Overall	Potential	Club	Club Logo	Value	Wage	Special	Preferred	International	Weak Foot
0	158023 L. Messi	31	https://cdn	Argentina	https://cdn	94	94	FC Barcel	https://cdn	€110.5M	€565K	2202	Left	5	4
1	20801 Cristiano R	33	https://cdn	Portugal	https://cdn	94	94	Juventus	https://cdn	€77M	€405K	2228	Right	5	4
2	190871 Neymar Jr	26	https://cdn	Brazil	https://cdn	92	93	Paris Saint	https://cdn	€118.5M	€290K	2143	Right	5	5
3	193080 De Gea	27	https://cdn	Spain	https://cdn	91	93	Manchest	https://cdn	€72M	€260K	1471	Right	4	3
4	192985 K. De Bru	27	https://cdn	Belgium	https://cdn	91	92	Manchest	https://cdn	€102M	€355K	2281	Right	4	5
5	183277 E. Hazard	27	https://cdn	Belgium	https://cdn	91	91	Chelsea	https://cdn	€93M	€340K	2142	Right	4	4
6	177003 L. Modrić	32	https://cdn	Croatia	https://cdn	91	91	Real Madi	https://cdn	€67M	€420K	2280	Right	4	4
7	176580 L. Suárez	31	https://cdn	Uruguay	https://cdn	91	91	FC Barcel	https://cdn	€80M	€455K	2346	Right	5	4
8	155862 Sergio Rai	32	https://cdn	Spain	https://cdn	91	91	Real Madi	https://cdn	€51M	€380K	2201	Right	4	3
9	200389 I. Oblak	25	https://cdn	Slovenia	https://cdn	90	93	Atlético M	https://cdn	€68M	€94K	1331	Right	3	3
10	188545 R. Lewand	29	https://cdn	Poland	https://cdn	90	90	FC Bayer	https://cdn	€77M	€205K	2152	Right	4	4
11	182521 T. Kroos	28	https://cdn	Germany	https://cdn	90	90	Real Madi	https://cdn	€76.5M	€355K	2190	Right	4	5
12	182493 D. Godin	32	https://cdn	Uruguay	https://cdn	90	90	Atlético M	https://cdn	€44M	€125K	1946	Right	3	3
13	168542 David Silv	32	https://cdn	Spain	https://cdn	90	90	Manchest	https://cdn	€60M	€265K	2115	Left	4	2
14	215914 N. Kanté	27	https://cdn	France	https://cdn	89	90	Chelsea	https://cdn	€63M	€225K	2189	Right	3	3
15	211110 P. Dybala	24	https://cdn	Argentina	https://cdn	89	94	Juventus	https://cdn	€89M	€205K	2092	Left	3	3
16	202126 H. Kane	24	https://cdn	England	https://cdn	89	91	Tottenham	https://cdn	€83.5M	€205K	2165	Right	3	4
17	194765 A. Griezma	27	https://cdn	France	https://cdn	89	90	Atlético M	https://cdn	€78M	€145K	2246	Left	4	3

執行程式的時候，會偵測目的路徑中是否有指定的選手資料集(.csv)，如果系統沒有偵測到則會直接結束程式，如下圖所示，

```
directory is not found.  
  
An exception has occurred, use %tb to see the full traceback.  
  
SystemExit: 0
```

```

directory is found.
file is not found.

An exception has occurred, use %tb to see the full traceback.

SystemExit: 0

```

程式會自動在「D:\FIFA19」的工作目錄中，產生選手評分的預測結果(.csv)以及比較圖(.png)，最後可透過預測結果(.csv)來選取最高評分的角色。

H. Kane	ST	89	89
A. Griezmann	CAM	89	89
M. Neuer	GK	89	89
T. Courtois	GK	89	89
E. Cavani	LS	89	89
Coutinho	LW	88	88
P. Aubameyang	LM	88	88
S. Handanovic	GK	88	88
C. Eriksen	CAM	88	88
H. Lloris	GK	88	88

三、本程式所使用的套件功能介紹

本程式主要運用 3 個套件，其功能說明如下：

➤ **pandas** 套件

pandas 套件為資料處理套件，方便使用者在處理文書檔案(.csv, .json, .xls)非常的方便，無論是讀取、撰寫、修改，都能夠輕而易舉。

程式中使用以下幾種語法：

語法	功能
<code>read_csv("路徑")</code>	讀取本地端的檔案，或是網路上的來源檔，其中路徑的斜線與 windows 表示法會相反，需要多加注意。
<code>to_csv("路徑")</code>	用於輸出 Dataframe 為(.csv)檔，其中路徑的部分斜線與 windows 表示法不同只能有一個正斜線。

Dataframe	一種資料格式，用於檔案的存取，也可以用於操作新增、刪除、修改、更新，是非常方便管理資料的格式。
-----------	---

➤ **os 套件**

Python 提供 **os.path** 套件來進行檔案的操作，可以取得檔案路徑、檔案大小、建立目錄、刪除目錄、刪除檔案與執行命令...等操作。

程式中使用以下幾種語法：

語法	功能
<code>os.path.isdir("路徑")</code>	確認路徑中的資料夾是否存在。
<code>os.path.isfile ("路徑")</code>	確認路徑中的檔案是否存在。

➤ **sys 套件**

Python 提供 **sys** 套件用於系統上的執行指令，通常 **sys** 套件可以處理例外情況，或是強制性的執行某些關鍵指令，**sys** 中的部分指令具有不可回朔性，使用上必須多加小心。

程式中使用以下幾種語法：

語法	功能
<code>sys.exit(0)</code>	強制讓程式結束執行，並且視為正常結束。

➤ **sklearn 套件**

sklearn 提供許多預測模型，提供使用者只需要預處理資料集，就可以通過預測模型達到想要預測的效果，其中模型的實作相當複雜，因此開發者透過簡單的 **function** 窗口取代原先困難的內容，每個模型都有提供多樣的參數讓使用者根據不同的情形去調整模型，模型在訓練的過程中會不斷的反正修正，因此給予適當的資料集是相當重要的工作。

程式中使用以下幾種語法：

語法	功能
<code>train_test_split(資料集,答案集,訓練與預測資料分配比例)</code>	此功能相當的強大，能夠隨機的把整體資料集切割成訓練集以及預測集，並且使用者可以隨意決定切割的比例。
<code>tree.DecisionTreeClassifier()</code>	建立決策樹的分類器，括弧內有許多參數可以使用，此篇不再贅述。
<code>KNeighborsClassifier(n_neighbors=k)</code>	建立 KNN 的分類器，括弧內有許多參數可以使用，此篇不再贅述。

分類器變數.fit(訓練集,訓練集解答)	此功能用於訓練模型，給予訓練集以及解答，讓模型自行修正。
分類器變數.predict(預測集)	此功能用於已經訓練好的模型，接著使用想要預測的資料集讓模型預測，即可得知預測結果。

➤ matplotlib 套件

matplotlib 套件提供使用者用簡易的函數就可以畫出直線、圖形，圖表，用於實驗結果表示是非常方便的套件，畫出圖表後還可以儲存結果以便攜帶。

程式中使用以下幾種語法：

語法	功能
plt.figure(figsize=(x,y))	控制圖片大小。
plt.title("title",fontsize=x)	加入圖表標題。
plt.ylabel("y",fontsize=x) plt.xlabel("x",fontsize=x)	加入 X，Y 軸標題。
plt.xticks(fontsize=15) plt.yticks(fontsize=15)	控制 X，Y 軸刻度大小。
plt.plot(X 軸資料,Y 軸資料)	加入單一折線。
plt.legend()	加入圖例。
plt.savefig("路徑")	儲存圖表。
plt.show()	顯示圖表。

四、參考程式碼與程式碼解說

1	import pandas as pd
2	import matplotlib.pyplot as plt
3	from sklearn import tree
4	from sklearn.model_selection import train_test_split
5	from sklearn import metrics
6	from sklearn.neighbors import KNeighborsClassifier
7	from sklearn.svm import SVC
8	import os
9	import sys
此區塊主要是加入程式中需要的所有套件。	

1	#open file
2	#user can change path on his/her own
3	folderpath = 'D:/FIFA19'
4	if os.path.isdir(folderpath):
5	print("directory is found ◦ ")
6	else:
7	print("directory is not found ◦ ")
8	sys.exit(0)
9	
10	filepath = folderpath + '/data.csv'
11	if os.path.isfile(filepath):
12	print("file is found ◦ ")
13	else:
14	print("file is not found ◦ ")
15	sys.exit(0)
16	df = pd.read_csv("D:/FIFA19/data.csv")
此區塊主要是讀取選手資料集，當沒有路徑上的資料夾與檔案不存在時，將會跳出訊息告訴使用者需要檢查路徑。	
1	#drop non-useful attributes
2	drop_list = ["ID","Name","Photo","Nationality","Flag","Club","Club Logo","Special","Work
3	Rate", "Real Face","Jersey Number","Joined","Loaned From","Contract Valid Until",
4	"Crossing","Finishing","HeadingAccuracy","ShortPassing","Volleys","Dribbling",
5	"Curve","FKAccuracy","LongPassing","BallControl","Acceleration","SprintSpeed",
6	"Agility","Reactions","Balance","ShotPower","Jumping","Stamina","Strength",
7	"LongShots","Aggression","Interceptions","Positioning","Vision","Penalties",
8	"Composure","Marking","StandingTackle","SlidingTackle","GKDividing","GKHandling",
9	"GKKicking","GKPositioning","GKReflexes","LS","ST","RS","LW","LF","CF","RF","RW",
10	"LAM","CAM","RAM","LM","LCM","CM","RCM","RM","LWB","LDM","CDM","RDM","RWB",
11	LB","LCB","CB","RCB","RB","Overall","Position"]
12	fifa_data = df.drop(drop_list,axis = 1)
13	
14	#fill nan value to 0
15	fifa_data = fifa_data.fillna("0")
16	
17	#handle "€","M","K" in value,wage and release clause attributes
18	for i in range(len(fifa_data["Value"])):
19	temp_str = fifa_data["Value"][i]
20	temp_str = temp_str.replace("€","")

```

21     if "M" in temp_str:
22         temp_str = temp_str.replace("M", "")
23         temp_str = str(int(float(temp_str)*1000000))
24     elif "K" in temp_str:
25         temp_str = temp_str.replace("K", "")
26         temp_str = str(int(float(temp_str)*1000))
27     fifa_data.loc[i, "Value"] = temp_str
28
29 for i in range(len(fifa_data["Wage"])):
30     temp_str = fifa_data["Wage"][i]
31     temp_str = temp_str.replace("€", "")
32     if "M" in temp_str:
33         temp_str = temp_str.replace("M", "")
34         temp_str = str(int(float(temp_str)*1000000))
35     elif "K" in temp_str:
36         temp_str = temp_str.replace("K", "")
37         temp_str = str(int(float(temp_str)*1000))
38     fifa_data.loc[i, "Wage"] = temp_str
39
40 for i in range(len(fifa_data["Release Clause"])):
41     temp_str = fifa_data["Release Clause"][i]
42     temp_str = temp_str.replace("€", "")
43     if "M" in temp_str:
44         temp_str = temp_str.replace("M", "")
45         temp_str = str(int(float(temp_str)*1000000))
46     elif "K" in temp_str:
47         temp_str = temp_str.replace("K", "")
48         temp_str = str(int(float(temp_str)*1000))
49     fifa_data.loc[i, "Release Clause"] = temp_str
50
51 #handle Preferred Foot by translation in 0(Left) and 1(Right)
52 for i in range(len(fifa_data["Preferred Foot"])):
53     if "Left" in fifa_data["Preferred Foot"][i]:
54         fifa_data.loc[i, "Preferred Foot"] = "0"
55     else:
56         fifa_data.loc[i, "Preferred Foot"] = "1"
57
58 #handle body type by quantizing

```

```

59 monsters = ['C. Ronaldo','Messi','Neymar','Akinfenwa','Shaqiri','Courtois']
60 for i in range(len(fifa_data["Body Type"])):
61     if fifa_data["Body Type"][i] in monsters:
62         fifa_data.loc[i,"Body Type"] = "4"
63     elif fifa_data["Body Type"][i] == "Stocky":
64         fifa_data.loc[i,"Body Type"] = "3"
65     elif fifa_data["Body Type"][i] == "Normal":
66         fifa_data.loc[i,"Body Type"] = "2"
67     elif fifa_data["Body Type"][i] == "Lean":
68         fifa_data.loc[i,"Body Type"] = "1"
69     else:
70         fifa_data.loc[i,"Body Type"] = "0"
71
72 #handle Height by translating in CM
73 for i in range(len(fifa_data["Height"])):
74     temp_str = fifa_data["Height"][i]
75     temp_str = temp_str.replace("\'",".")
76     inch = float(temp_str) % 1 * 100
77     if inch > 12:
78         inch/=10
79     temp_str = str(float("{:.2f}".format((float(temp_str) // 1 * 12 + inch) * 2.54)))
80     fifa_data.loc[i,"Height"] = temp_str
81
82 #handle Weight by translating in KG
83 for i in range(len(fifa_data["Weight"])):
84     temp_str = fifa_data["Weight"][i]
85     temp_str = temp_str.replace("lbs","")
86     temp_str = str(float("{:.2f}".format(float(temp_str) * 0.453)))
87     fifa_data.loc[i,"Weight"] = temp_str

```

此區塊主要負責資料預處理，首先把不需要的 column drop，接著用 fillna 補上資料空值，因為模型需要數值化的資料因此需要把文字轉換成純數字，17-49 行用於處理資料中具有“€,K,M”文字的值；51-56 行用於處理轉換“Left,Right”為 0,1；58-70 行用於處理選手各種不同的體態字詞；72-80 行用於處理英呎轉成公分；82-87 用於處理磅轉為公斤。前述有提到預處理對於模型訓練非常重要，因此對於 drop_list 的挑選必須慎選。

```

1 #take overall as standard picking rule
2 fifa_ovr = df[["Overall","Name","Position"]]
3
4 #split fifa_data into training data and
5 train_data, test_data, train_ovr, test_ovr = train_test_split(fifa_data,fifa_ovr,test_size = 0.5)
6 test_name = pd.DataFrame({"Name":test_ovr["Name"]})
7 test_pos = pd.DataFrame({"Position":test_ovr["Position"]})
8 train_ovr = train_ovr.drop(["Name","Position"],axis = 1)
9 test_ovr = test_ovr.drop(["Name","Position"],axis = 1)
10
11 #copy data
12 train_Ddata = train_data
13 train_Dovr = train_ovr
14 test_Ddata = test_data
15 test_Dovr = test_ovr
16 test_Dname = test_name
17 test_Dpos = test_pos
18
19 #build "DecisionTree" classifier
20 clf = tree.DecisionTreeClassifier()
21 fifa_clf = clf.fit(train_Ddata, train_Dovr)
22
23 #predict
24 fifa_Dpredicted = fifa_clf.predict(test_Ddata)
25
26 #rearrange test standard overalls' index
27 test_Dovr.index = range(len(test_Dovr))
28 test_Dname.index = range(len(test_Dname))
29 test_Dpos.index = range(len(test_Dpos))
30
31 #accuracy
32 accuracy = metrics.accuracy_score(test_Dovr,fifa_Dpredicted)
33 print(accuracy)
34
35 #concat predicted answer and standard overall
36 df_Dpredicted = pd.DataFrame({'Score':fifa_Dpredicted})
37 Dresult = {
38     'Name': [],

```



```

39     'Position': [],
40     'Score': [],
41     'Overall': []
42 }
43 df_Dresult = pd.DataFrame(Dresult)
44 for i in range(len(test_Dovr)):
45     df_Dresult = df_Dresult.append({
46         'Name': test_Dname['Name'][i],
47         'Position': test_Dpos['Position'][i],
48         'Score': df_Dpredicted['Score'][i],
49         'Overall': test_Dovr['Overall'][i]
50     },ignore_index=True)
51 df_Dresult.sort_values(by = 'Score',inplace=True,ascending=False,ignore_index=True)
52 print(df_Dresult)
53
54 #output result to csv
55 output_PATH = 'D:/FIFA19/fifa_D_result.csv'
56 df_Dresult.to_csv(output_PATH,index=0)

```

此區塊主要是使用決策數(DecisionTree)作為模型去預測資料，首先第 2 行先取出答案樣本，第 5 行切割資料為訓練集以及預測集，12-17 行為了保護資料而複製一份一樣的資料，接著第 20 行建立模型雛形，第 21 行訓練模型，再來交給第 24 行預測結果，因為切割訓練集以及預測集時為隨機切割，因此 index 不是從 1 開始所以 27-29 行用於重新定義 index，第 32 行計算預測結果跟答案樣本的相似度(準確性)，最後 35-56 行把分散的資料合併成一個 dataframe 且輸出至(.csv)檔。

```

1 #copy data
2 train_Kdata = train_data
3 train_Kovr = train_ovr
4 test_Kdata = test_Ddata
5 test_Kovr = test_ovr
6 test_Kname = test_name
7 test_Kpos = test_pos
8 train_Kovr = train_Kovr.values.ravel()
9
10 #Because we don't know how to choose K index in KNN algorithm,
11 #so we give a test to bulid many different K classifier.
12 accuracy = []
13 for k in range(1, 100):
14     knn = KNeighborsClassifier(n_neighbors=k)
15     knn.fit(train_Kdata, train_Kovr)
16     y_pred = knn.predict(test_Kdata)
17     accuracy.append(metrics.accuracy_score(test_Kovr, y_pred))
18 k_range = range(1,100)
19 plt.plot(k_range, accuracy)
20 plt.savefig("D:/FIFA19/K_selection.png")
21 plt.show()
22
23 #build "KNN" classifier
24 knn = KNeighborsClassifier(n_neighbors=1)
25 knn = knn.fit(train_Kdata, train_Kovr)
26
27 #predict
28 fifa_Kpredicted = knn.predict(test_Kdata)
29
30 #test standard overall rearrange index
31 test_Kovr.index = range(len(test_Kovr))
32 test_Kname.index = range(len(test_Kname))
33 test_Kpos.index = range(len(test_Kpos))
34
35 #accuracy
36 accuracy = metrics.accuracy_score(test_Kovr,fifa_Kpredicted)
37
38 #concate predicted answer and standard overall

```

39	df_Kpredicted = pd.DataFrame({'Score':fifa_Kpredicted})
40	Kresult = {
41	'Name': [],
42	'Position': [],
43	'Score': [],
44	'Overall': []
45	}
46	df_Kresult = pd.DataFrame(Kresult)
47	for i in range(len(test_Kovr)):
48	df_Kresult = df_Kresult.append({
49	'Name': test_Kname['Name'][i],
50	'Position': test_Kpos['Position'][i],
51	'Score': df_Kpredicted['Score'][i],
52	'Overall': test_Kovr['Overall'][i]
53	},ignore_index=True)
54	df_Kresult.sort_values(by = 'Score',inplace=True,ascending=False,ignore_index=True)
55	print(df_Kresult)
56	
57	#output result to csv
58	output_PATH = 'D:/FIFA19/fifa_K_result.csv'
59	df_Kresult.to_csv(output_PATH,index=0)
<p>此區塊主要是使用 K-近鄰演算法(KNN)作為模型去預測資料， 1-8 行為複製與前區塊一樣，因為 KNN 其中有一個參數 K 為鄰近數量，這個參數將會影響預測結果，必須謹慎選擇 K 值，因此 12-21 行為 K 值訓練，用 K=1 到 K=100 的方式逐次建立模型去訓練，藉此決定最後 K 取多少最為合理，此 K 值訓練的成果將會輸出成折線圖存至指定路徑，23-59 行與前區塊 19-56 行雷同，在此不多贅述。</p>	
1	#comparasion between DecisionTree and KNN
2	x_axis = []
3	D_line = []
4	K_line = []
5	O_line = []
6	for i in range(0,len(test_data),500):
7	x_axis.append(i)
8	D_line.append(fifa_Dpredicted[i])
9	K_line.append(fifa_Kpredicted[i])
10	O_line.append(test_ovr['Overall'][i])
11	plt.figure(figsize=(10,10))
12	plt.title("Comparasion of Original, DicisionTree, KNN",fontsize=15) # title

```

13 plt.ylabel("Accuracy",fontsize=15) # y label
14 plt.xlabel("Serial number of data",fontsize=15) # x label
15 plt.xticks(fontsize=15)
16 plt.yticks(fontsize=15)
17 plt.plot(x_axis,O_line,linewidth = 5,label = 'Original')
18 plt.plot(x_axis,D_line,'--',color = 'red',label = 'DecisionTree')
19 plt.plot(x_axis,K_line,'--',color = 'green',label = 'KNN')
20 plt.legend()
21 plt.savefig("D:/FIFA19/Comparasion.png")
22 plt.show()

```

此區塊為最後總整理，這次的程式使用兩種預測模型，決策樹與 K-近鄰演算法，為了比較這兩種演算法的好壞與標準答案的相似度，使用折線圖表示並存在指定路徑中，經結果顯示決策樹的預測較為貼近標準答案，因為決策樹善於處理大量的資料且容易被解釋，KNN 則是懶散學習法，速度快但受限於 K 值，因此結果較為不理想。

五、心得

經由這次作業不但了解資料處理的重要性，且對於一些訓練模型更加的了解，某業界人士說：「現在的時代人人都可以做 AI，只要有一台電腦跟足夠的資料」，這句話我非常的認同，因為整份程式用於 AI 的部分只有模型訓練，但模型訓練只有少少幾行，資料的預處理則是多上好幾倍，對於模型而言好的資料可以帶給他非常優秀的成果，而垃圾資料只會吐出垃圾資訊，就是人口中所戲稱的 **Garbage in Garbage out(GIGO)**。身為資工系的學生，大多時間都是在專研演算法，熟悉的語言頂多 C 與 C++，Python 則是很少碰觸到，而在現在的趨勢下 Python 逐漸變成人人必備的程式語言，為了與時代接軌，在大數據的課程中學到 Python 各種相關知識之外，也學會不少套件的應用，相信這門課程對於未來研究所或是找工作會非常有幫助。