

Credit Card Task

By Gerard Utoware

Table of Contents

Overviews – Page 3

Q1 – Page 4

Q2 – Page 4

Q3 – Page 6

Overview of the tool used

The main advantage of python in comparison to other data analytic tools such as excel, would be Python's ability to handle large volumes of data without hindering productivity. Running scripts using the libraries on allows for more automation and ease of obtaining analysed data. With this task, the data was imported from python which ensures that data will not be lost or tampered with while performing data analysis. After defining the columns from the imported data, it was less tedious to run scripts than to perform individual analysis on excel.

Overview of the libraries used

Numpy, pandas and matplotlib were used in the script. All three python libraries allow different functions to be performed. Pandas allows data analysis scripts to be conducted on python. It utilises two other libraries, being numpy and matplotlib. Matplotlib allows for data visualisation through a variety of charts such as bar and pie charts. Numpy permits mathematical operations in python. The combination of these three libraries is what allowed the following data and charts to be formed and analysed. In addition to these, seaborn was installed in order for a heatmap to be printed and shown.

Overview of the dataset

Overall, the data contains fraudulent data that's been flagged up by the script using python. V8 and V28 are the two columns with the most outliers in the data with 79 and 75 fraudulent transactions respectively. V6 and V11 both only had 1 outlier in their columns. With that being said, the r coefficient is overall stays around 0 which means there is no strong positive or negative correlation in the dataset.

Q1&2

The screenshot displays a Jupyter Notebook environment with a dark theme. The top bar shows the file name 'helloworld' and the current directory 'SpeedyWork.py'. The left sidebar contains a 'Project' view showing the file structure and a 'Structure' view showing the code structure. The main area shows a Python script with the following code:

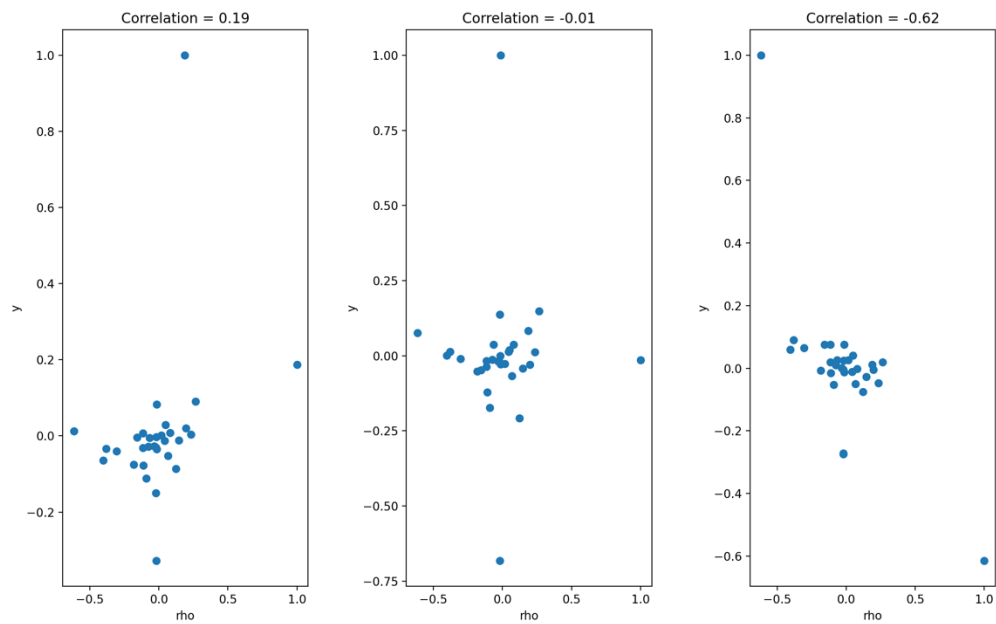
```
1 import matplotlib.pyplot as plt
2 import pandas as pd
3 import numpy as np
4 import pyod
5
6
7 #Written to import the csv file and to show the columns in the dataset
8 df = pd.read_csv('creditcard.csv')
9 print(df)
10 df.info()
11
12 #Detect outliers and valid transactions throughout the dataset by defining the IQR.
13 def find_outliers_IQR(df):
14     q1=df.quantile(0.25)
15     q3=df.quantile(0.75)
16     IQR=q3-q1
17     outliers = df[((df<(q1-1.5*IQR)) | (df>(q3+1.5*IQR)))]
18     valid=df[((df>(q1-1.5*IQR)) & (df<(q3+1.5*IQR)))]
19     print('outliers:')
20     print(outliers.count())
21     print('valid:')
22     print(valid.count())
23     return outliers
24
25 #This allows the number of outliers or valid transactions to be written down per column name
26 outliers = find_outliers_IQR(df[['Time','V1','V2','V3','V4','V5','V6','V7','V8','V9','V10','V11','V12','V13','V14','V15','V16','V17','V18','V19','V20','V21','V22','V23','V24','V25','V26','V27','V28','V29','V30','V31','V32','V33','V34','V35','V36','V37','V38','V39','V40','V41','V42','V43','V44','V45','V46','V47','V48','V49','V50','V51','V52','V53','V54','V55','V56','V57','V58','V59','V60','V61','V62','V63','V64','V65','V66','V67','V68','V69','V70','V71','V72','V73','V74','V75','V76','V77','V78','V79','V80','V81','V82','V83','V84','V85','V86','V87','V88','V89','V90','V91','V92','V93','V94','V95','V96','V97','V98','V99','V100','V101','V102','V103','V104','V105','V106','V107','V108','V109','V110','V111','V112','V113','V114','V115','V116','V117','V118','V119','V120','V121','V122','V123','V124','V125','V126','V127','V128','V129','V130','V131','V132','V133','V134','V135','V136','V137','V138','V139','V140','V141','V142','V143','V144','V145','V146','V147','V148','V149','V150','V151','V152','V153','V154','V155','V156','V157','V158','V159','V160','V161','V162','V163','V164','V165','V166','V167','V168','V169','V170','V171','V172','V173','V174','V175','V176','V177','V178','V179','V180','V181','V182','V183','V184','V185','V186','V187','V188','V189','V190','V191','V192','V193','V194','V195','V196','V197','V198','V199','V200','V201','V202','V203','V204','V205','V206','V207','V208','V209','V210','V211','V212','V213','V214','V215','V216','V217','V218','V219','V220','V221','V222','V223','V224','V225','V226','V227','V228','V229','V230','V231','V232','V233','V234','V235','V236','V237','V238','V239','V240','V241','V242','V243','V244','V245','V246','V247','V248','V249','V250','V251','V252','V253','V254','V255','V256','V257','V258','V259','V260','V261','V262','V263','V264','V265','V266','V267','V268','V269','V270','V271','V272','V273','V274','V275','V276','V277','V278','V279','V280','V281','V282','V283','V284','V285','V286','V287','V288','V289','V290','V291','V292','V293','V294','V295','V296','V297','V298','V299','V300','V301','V302','V303','V304','V305','V306','V307','V308','V309','V310','V311','V312','V313','V314','V315','V316','V317','V318','V319','V320','V321','V322','V323','V324','V325','V326','V327','V328','V329','V330','V331','V332','V333','V334','V335','V336','V337','V338','V339','V340','V341','V342','V343','V344','V345','V346','V347','V348','V349','V350','V351','V352','V353','V354','V355','V356','V357','V358','V359','V360','V361','V362','V363','V364','V365','V366','V367','V368','V369','V370','V371','V372','V373','V374','V375','V376','V377','V378','V379','V380','V381','V382','V383','V384','V385','V386','V387','V388','V389','V390','V391','V392','V393','V394','V395','V396','V397','V398','V399','V400','V401','V402','V403','V404','V405','V406','V407','V408','V409','V410','V411','V412','V413','V414','V415','V416','V417','V418','V419','V420','V421','V422','V423','V424','V425','V426','V427','V428','V429','V430','V431','V432','V433','V434','V435','V436','V437','V438','V439','V440','V441','V442','V443','V444','V445','V446','V447','V448','V449','V450','V451','V452','V453','V454','V455','V456','V457','V458','V459','V460','V461','V462','V463','V464','V465','V466','V467','V468','V469','V470','V471','V472','V473','V474','V475','V476','V477','V478','V479','V480','V481','V482','V483','V484','V485','V486','V487','V488','V489','V490','V491','V492','V493','V494','V495','V496','V497','V498','V499','V500','V501','V502','V503','V504','V505','V506','V507','V508','V509','V510','V511','V512','V513','V514','V515','V516','V517','V518','V519','V520','V521','V522','V523','V524','V525','V526','V527','V528','V529','V530','V531','V532','V533','V534','V535','V536','V537','V538','V539','V540','V541','V542','V543','V544','V545','V546','V547','V548','V549','V550','V551','V552','V553','V554','V555','V556','V557','V558','V559','V560','V561','V562','V563','V564','V565','V566','V567','V568','V569','V570','V571','V572','V573','V574','V575','V576','V577','V578','V579','V580','V581','V582','V583','V584','V585','V586','V587','V588','V589','V590','V591','V592','V593','V594','V595','V596','V597','V598','V599','V600','V601','V602','V603','V604','V605','V606','V607','V608','V609','V610','V611','V612','V613','V614','V615','V616','V617','V618','V619','V620','V621','V622','V623','V624','V625','V626','V627','V628','V629','V630','V631','V632','V633','V634','V635','V636','V637','V638','V639','V640','V641','V642','V643','V644','V645','V646','V647','V648','V649','V650','V651','V652','V653','V654','V655','V656','V657','V658','V659','V660','V661','V662','V663','V664','V665','V666','V667','V668','V669','V670','V671','V672','V673','V674','V675','V676','V677','V678','V679','V680','V681','V682','V683','V684','V685','V686','V687','V688','V689','V690','V691','V692','V693','V694','V695','V696','V697','V698','V699','V700','V701','V702','V703','V704','V705','V706','V707','V708','V709','V710','V711','V712','V713','V714','V715','V716','V717','V718','V719','V720','V721','V722','V723','V724','V725','V726','V727','V728','V729','V730','V731','V732','V733','V734','V735','V736','V737','V738','V739','V740','V741','V742','V743','V744','V745','V746','V747','V748','V749','V750','V751','V752','V753','V754','V755','V756','V757','V758','V759','V760
```

	Fraudulent Transactions (Outliers)	Valid Transactions	Fractional Ratio
Time	0	284807	0
V1	7062	277745	3/121
V2	13526	271281	18/379
V3	3363	281444	3/254
V4	11148	273659	31/792
V5	12295	261842	27/602
V6	22965	275512	1/13
V7	8948	261842	19/575
V8	24134	275859	79/982
V9	8283	260673	17/552
V10	9496	276524	25/753
V11	780	284027	1/365
V12	15348	269459	9/167
V13	3368	281913	10/847
V14	14149	270658	31/624

V15	2894	281913	5/492
V16	8184	276623	5/174
V17	7420	277387	13/499
V18	7533	277274	26/983
V19	10205	274602	11/307
V20	27770	257037	43/441
V21	14497	270310	48/943
V22	1317	283490	4/865
V23	18541	266266	36/553
V24	4774	280033	3/179
V25	5367	279440	15/796
V26	5596	279211	19/967
V27	39163	245644	11/80
V28	30342	254465	75/704
Amount	31904	252903	41/366
Class	492	0	1

In order to obtain the outlier fractions per columns in the dataset, the number of fraudulent and valid transactions per column had to be ascertained first then tabulated. An outlier in statistics is defined as a point of data that falls more than 1.5 times the interquartile range above the third quartile or below the first quartile. These points on this dataset are more likely to be deemed as potential fraudulent and need to be flagged up. Using python, 'outliers' were defined and listed as shown in the screenshot above. Valid transactions were also printed so the total number of transactions could be identified to use to obtain the outlier fraction. Finally, the columns I wanted to show and find outliers for were listed individually.

Q3



```

#This allows the number of outliers or valid transactions to be written down per column name
outliers = find_outliers_IQR(df[['Time', 'V1', 'V2', 'V3', 'V4', 'V5', 'V6', 'V7', 'V8', 'V9', 'V10', 'V11', 'V12', 'V13', 'V14', 'V15', 'V16', 'V17', 'V18', 'V19', 'V20', 'V21', 'V22', 'V23']])

#Defining the R coefficient before printing
r_coe = df.corr(method='pearson')
print(df.corr(method='pearson'))
rho = np.corrcoef(r_coe)

#Plotting scatter graphs to show the change in the correlation coefficient throughout the data
fig, ax = plt.subplots(nrows=1, ncols=3, figsize=(12, 3))
for i in [0,1,2]:
    ax[i].scatter(rho[0,i], rho[0,i+1])
    ax[i].title.set_text('Correlation = ' + "{:.2f}".format(rho[0,i+1]))
    ax[i].set_xlabel='rho', ylabel='y'
fig.subplots_adjust(wspace=.4)
plt.show()

```

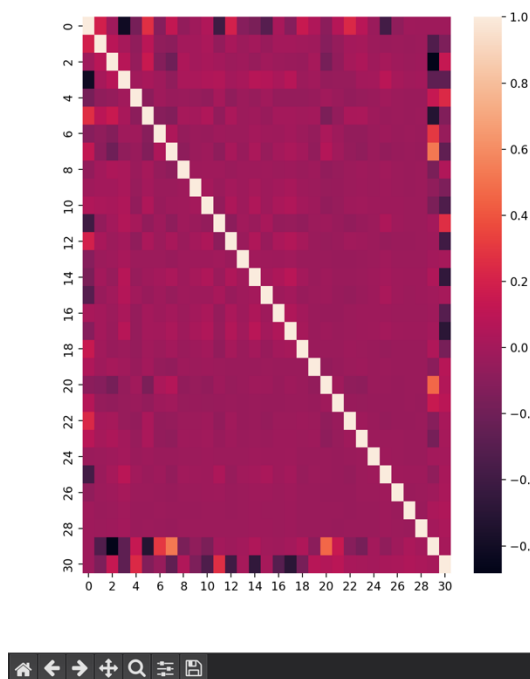
Run: SpeedyWork (1) x
dtype: Int64

	Time	V1	V2	...	V28	Amount	Class
Time	1.000000	1.173963e-01	-1.059333e-02	...	-9.412688e-03	-0.010596	-0.012323
V1	0.117396	1.000000e+00	3.777823e-12	...	-9.769215e-13	-0.227789	-0.101347
V2	-0.010593	3.777823e-12	1.000000e+00	...	2.525513e-12	-0.531409	0.091289
V3	-0.419618	-2.118614e-12	2.325661e-12	...	5.189189e-12	-0.210880	-0.192961
V4	-0.105260	-1.733159e-13	-2.314981e-12	...	-2.032372e-12	0.098732	0.133447
V5	0.173072	-3.473231e-12	-1.831952e-12	...	1.010196e-11	-0.386356	-0.094974
V6	-0.063016	-1.306165e-13	9.438444e-13	...	-6.069227e-13	0.215981	-0.043643
V7	0.084714	-1.116494e-13	5.403436e-12	...	2.958679e-13	0.397311	-0.187257
V8	-0.036949	2.114527e-12	2.133785e-14	...	1.866598e-12	-0.103079	0.019875
V9	-0.008660	3.016285e-14	3.238513e-13	...	-1.406856e-12	-0.044246	-0.097733
V10	0.030617	-2.615192e-12	1.463282e-12	...	5.116560e-12	-0.101502	-0.216883
V11	-0.247689	1.866551e-12	-8.314960e-13	...	-4.247931e-12	0.000104	0.154876
V12	0.124348	-1.238745e-12	6.139448e-13	...	-7.428113e-12	-0.009542	-0.260593

Version Control | Run | Python Packages | TODO | Python Console | Problems | Terminal | Services

PEP 8: E265 block comment should start with '#'

28:45 LF UTF-8 4 spaces Python 3.10 (helloworld)



```

helloworld SpeedyWork.py
SpeedyWork.py First.py
21 print(outliers.count())
22 print('valid:')
23 print(valid.count())
24 return outliers
25
26 #This allows the number of outliers or valid transactions to be written down
27 outliers = find_outliers_IQR(df[['Time','V1','V2','V3','V4','V5','V6','V7']
28
29 #Defining the R coefficient before printing
30 r_coe = df.corr(method='pearson')
31 print(df.corr(method='pearson'))
32 rho = np.corrcoef(r_coe)
33
34 #Plotting heatmap to show the change in the correlation coefficient through
35 sns.heatmap(rho)
36 plt.show()
37
Run: SpeedyWork (1)
0.019875
V9 -0.008660 3.016285e-14 3.238513e-13 ... -1.406856e-12 -0.044246
-0.097733
V10 0.030617 -2.615192e-12 1.463282e-12 ... 5.116568e-12 -0.101502
-0.216883
V11 -0.247689 1.866551e-12 -8.314960e-13 ... -4.247931e-12 0.000104
0.154876
V12 0.124348 -1.238745e-12 6.139448e-13 ... -7.428113e-12 -0.009542
-0.260593
V13 -0.065902 7.589589e-13 -1.181068e-12 ... -6.777880e-12 0.005293
-0.004570
V14 -0.098757 -1.871054e-13 -3.384604e-13 ... -1.700091e-12 0.033751
-0.302544
V15 -0.183453 -3.601390e-13 2.196083e-13 ... -4.214967e-12 -0.002986
-0.004223

```

These images depict the script and resulting plots or heatmap to show the correlation within the dataset. The R values change throughout the data, however, majority of the data show no correlation as the R is very close to 0. There are some points in the data where R changes are can reach as low as -0.6 or 0.6..