

```

1 ### Thong Ke Nhieu Chieu - BT ca nhan chuong 4
2 # Nguyen Duc Vu Duy - 18110004
3 #Import thu vien
4 import numpy as np
5
6 print("Exercise 4.26: ")
7 print("a) ")
8 # x=np.array([[1,18.95],[2,19],[3,17.95],[3,15.54],[4,14],[5,12.95],[6,8.94],[8,7.49],[9,6],[11,3.99]])
9 # x=pd.DataFrame([.48,40.53,2.19,.55,.74,.66,.93,.37,.22],[12.57,73.68,11.13,20.03,20.29,.78,4.64,.43,1.08]))
10 x=np.array([.48,12.57],[2.19,11.13],[.55,20.03],[.74,20.29],[.66,.78],[.93,4.64],[.37,.43],[.22,1.08]))
11 x_1,x_2=np.array(x[:,0]),np.array(x[:,1])
12
13 x_1_mean=np.sum(x_1)/8
14 x_2_mean=np.sum(x_2)/8
15
16 x_mean=np.array([x_1_mean],[x_2_mean]))
17
18 S=np.cov(x_1,x_2)
19 print('The Sample variance - covariance matrix S: ')
20 print(S)
21 print()
22
23 sq_distance=[]
24 for i in range(len(x)):
25     x_0=x[i].reshape(2,1)
26     x_left=np.dot((x_0-x_mean).T,np.linalg.pinv(S))
27     x_result=np.dot(x_left,(x_0-x_mean))
28     sq_distance.append(x_result)
29     print('The squared statistical distances of x_'+str(i+1), ': ',x_result)
30     print()

```

Exercise 4.26:

a)

The Sample variance - covariance matrix S:

```
[[ 0.37856429  1.03028214]
 [ 1.03028214 69.85975536]]
```

The squared statistical distances of x_1 : [[0.51814443]]

The squared statistical distances of x_2 : [[5.3838751]]

The squared statistical distances of x_3 : [[2.18499937]]

The squared statistical distances of x_4 : [[1.97289605]]

The squared statistical distances of x_5 : [[0.936943]]

The squared statistical distances of x_6 : [[0.39512842]]

The squared statistical distances of x_7 : [[1.22453865]]

The squared statistical distances of x_8 : `[[1.38347497]]`

```
1 print('b) ')
2 print('Squared Statistical distances of  $x_j$  in part a) are distributed as a chi square distribution')
3 print()
4 print('The observations falling within the estimated 50% probability contour of a bivariate normal distribution would have squared statisti
5 print()
6 print('the 50% percentile of chi square with 2 degree of freedom is 1.3863')
7 print()
8
9 for i in range(len(x)):
10     if sq_distance[i]<1.3863:
11         print('Observation x'+str(i+1)+' fall within the estimated 50%.')
12
13 print()
14 print('We have 5 among 10 observations. Thus, 50% of observations fall within the estimated 50%')
```

b)
Squared Statistical distances of x_j in part a) are distributed as a chi square distribution

The observations falling within the estimated 50% probability contour of a bivariate normal distribution would have squared statistical distance which is lower than the 50% perc
the 50% percentile of chi square with 2 degree of freedom is 1.3863

Observation x1 fall within the estimated 50%.
Observation x5 fall within the estimated 50%.
Observation x6 fall within the estimated 50%.
Observation x7 fall within the estimated 50%.
Observation x8 fall within the estimated 50%.

We have 5 among 10 observations. Thus, 50% of observations fall within the estimated 50%

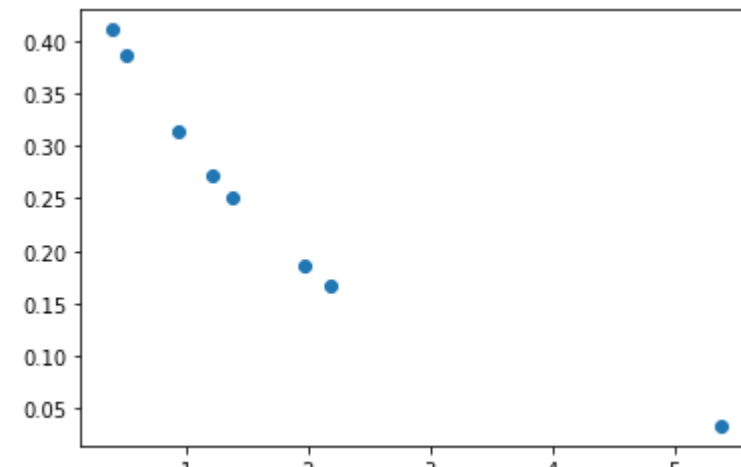
```
1 print('c) ')
2 print('The ordered distances in part a) ')
3 sq_distance_copy=sq_distance.copy()
4 sq_distance_order=sorted(sq_distance_copy)
5 print(sq_distance_order)
6 print()
7
8 #Import thu vien
9 from scipy.stats import chi2
10 import matplotlib.pyplot as plt
11
12 print('The chi square plot')
13 sq_distance_order_array=np.array(sq_distance_order).reshape(-1)
14 plt.plot(sq_distance_order_array, chi2.pdf(sq_distance_order_array, 2),'o')
15 plt.show()
16
```

c)

The ordered distances in part a)

```
[array([[0.39512842]]), array([[0.51814443]]), array([[0.936943]]), array([[1.22453865]]), array([[1.38347497]]), array([[1.97289605]]), array([[2.18499937]]), array([[5.3838751
```

The chi square plot



```
1 print('d) ')
2 print(' From part b, 50% of the observations fall within the estimated 50% probability contour of a bivariate normal distribution')
3 print()
4 print('Also, from part c, the line is nearly linear')
5 print()
6 print('Therefore, the data comes from a bivariate normal distribution')
```

d)

From part b, 50% of the observations fall within the estimated 50% probability contour of a bivariate normal distribution

Also, from part c, the line is nearly linear

Therefore, the data comes from a bivariate normal distribution

```
1 print('Exercise 4.29')
2 data=np.array([[8,98,7,2,12,8,2],[7,107,4,3,9,5,3],[7,103,4,3,5,6,3],[10,88,5,2,8,15,4],[6,91,4,2,8,10,3],[8,90,5,2,12,12,4],[9,84,7,4,12,1
3
4 x_5,x_6=np.array(data[:,4]),np.array(data[:,5])
5
6 x_5_mean=np.sum(x_5)/10
7 x_6_mean=np.sum(x_6)/10
8
9 x_mean_56=np.array([[x_5_mean],[x_6_mean]])
10
11 S_56=np.cov(x_5,x_6)
12 print('The Sample variance - covariance matrix S: ')
13 print(S_56)
14 print()
15
16 sq_distance_56=[]
17 for i in range(len(data)):
18     data_0=np.array([[data[i,4]],[data[i,5]]]).reshape(2,1)
19     data_left=np.dot((data_0-x_mean_56).T,np.linalg.pinv(S_56))
20     data_result=np.dot(data_left,(data_0-x_mean_56))
```

```

21 sq_distance_56.append(data_result)
22 print('The squared statistical distances of x_'+str(i+1), ': ',data_result)

```

Exercise 4.29

The Sample variance - covariance matrix S:

```

[[11.36353078  3.12659698]
 [ 3.12659698 30.97851336]]

```

```

The squared statistical distances of x_1 : [[98.11679939]]
The squared statistical distances of x_2 : [[118.36034603]]
The squared statistical distances of x_3 : [[139.74972261]]
The squared statistical distances of x_4 : [[110.48985158]]
The squared statistical distances of x_5 : [[116.33019929]]
The squared statistical distances of x_6 : [[92.48813839]]
The squared statistical distances of x_7 : [[88.96389426]]
The squared statistical distances of x_8 : [[52.39345208]]
The squared statistical distances of x_9 : [[98.83255908]]
The squared statistical distances of x_10 : [[91.79073439]]
The squared statistical distances of x_11 : [[116.60915696]]
The squared statistical distances of x_12 : [[99.68997693]]
The squared statistical distances of x_13 : [[68.85971073]]
The squared statistical distances of x_14 : [[104.65375632]]
The squared statistical distances of x_15 : [[116.33019929]]
The squared statistical distances of x_16 : [[110.76852113]]
The squared statistical distances of x_17 : [[126.31916125]]
The squared statistical distances of x_18 : [[86.98235906]]
The squared statistical distances of x_19 : [[103.86054199]]
The squared statistical distances of x_20 : [[118.36034603]]
The squared statistical distances of x_21 : [[95.53309159]]
The squared statistical distances of x_22 : [[127.86738944]]
The squared statistical distances of x_23 : [[88.93986219]]
The squared statistical distances of x_24 : [[146.46051769]]
The squared statistical distances of x_25 : [[93.18098112]]
The squared statistical distances of x_26 : [[127.86738944]]
The squared statistical distances of x_27 : [[98.83255908]]
The squared statistical distances of x_28 : [[122.07290619]]
The squared statistical distances of x_29 : [[113.60603634]]
The squared statistical distances of x_30 : [[134.72435136]]
The squared statistical distances of x_31 : [[109.79856446]]
The squared statistical distances of x_32 : [[98.11679939]]
The squared statistical distances of x_33 : [[125.5299788]]
The squared statistical distances of x_34 : [[92.70682395]]
The squared statistical distances of x_35 : [[129.32738441]]
The squared statistical distances of x_36 : [[110.76852113]]
The squared statistical distances of x_37 : [[125.5343714]]
The squared statistical distances of x_38 : [[81.05356433]]
The squared statistical distances of x_39 : [[97.95378394]]
The squared statistical distances of x_40 : [[122.19776596]]
The squared statistical distances of x_41 : [[95.16391276]]
The squared statistical distances of x_42 : [[135.31440391]]

```

```

1 print('b) ')
2 print('Squared Statistical distances of x_j in part a) are distributed as a chi square distribution')
3 print()
4 print('The observations falling within the estimated 50% probability contour of a bivariate normal distribution would have squared statisti
5 print()

```

```
6 print('the 50% percentile of chi square with 7 degree of freedom is 6.346')
7 print()
8
9 for i in range(len(data)):
10     if sq_distance_56[i]<6.346:
11         print('Observation x'+str(i+1)+' fall within the estimated 50%.')
12
13 print()
14 print('We have 0 among 42 observations. Thus, 0% of observations fall within the estimated 50%')
```

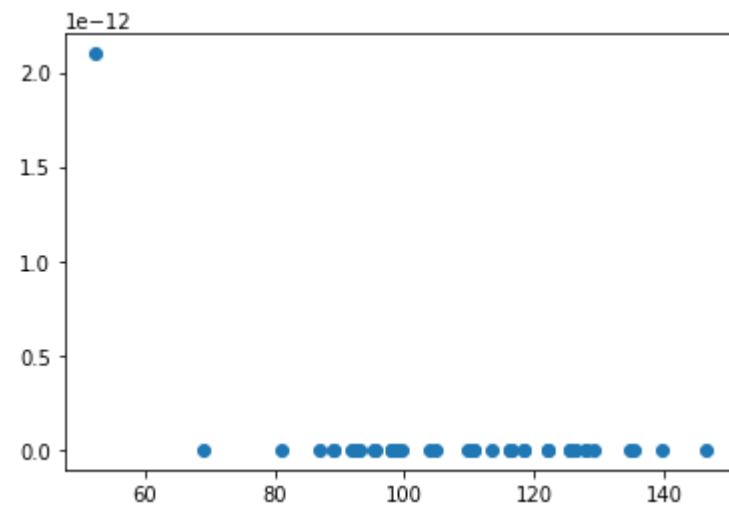
b)
Squared Statistical distances of x_j in part a) are distributed as a chi square distribution

The observations falling within the estimated 50% probability contour of a bivariate normal distribution would have squared statistical distance which is lower than the 50% percentile of chi square with 7 degree of freedom is 6.346

We have 0 among 42 observations. Thus, 0% of observations fall within the estimated 50%

```
1 print('c) ')
2 print('The ordered distances in part a) ')
3 sq_distance_copy_56=sq_distance_56.copy()
4 sq_distance_order_56=sorted(sq_distance_copy_56)
5 print(sq_distance_order_56)
6 print('The chi square plot')
7 sq_distance_order_array_56=np.array(sq_distance_order_56).reshape(-1)
8 plt.plot(sq_distance_order_array_56, chi2.pdf(sq_distance_order_array_56, 2), 'o')
9 plt.show()
```

c)
The ordered distances in part a)
[array([[52.39345208]]), array([[68.85971073]]), array([[81.05356433]]), array([[86.98235906]]), array([[88.93986219]]), array([[88.96389426]]), array([[91.79073439]]), array([[91.80000000]])]
The chi square plot



```
1 print('Exercise 4.30')
```

```

2 print('a) ')
3
4 from sklearn.preprocessing import PowerTransformer
5 pt_1=PowerTransformer('box-cox')
6 pt_1.fit(x_1.reshape(-1,1))
7 print('The power transformation lambda_1 that makes the x_1 values approximately normal: ')
8 print(pt_1.lambdas_)
9 print()
10
11 transform_x_1=pt_1.transform(x_1.reshape(-1,1))
12
13 import statsmodels.api as sm
14 print('The Q-Q plot for the transformed data')
15 sm.qqplot(transform_x_1,line='r')
16 plt.title('the Q-Q plot')
17 plt.show()

```

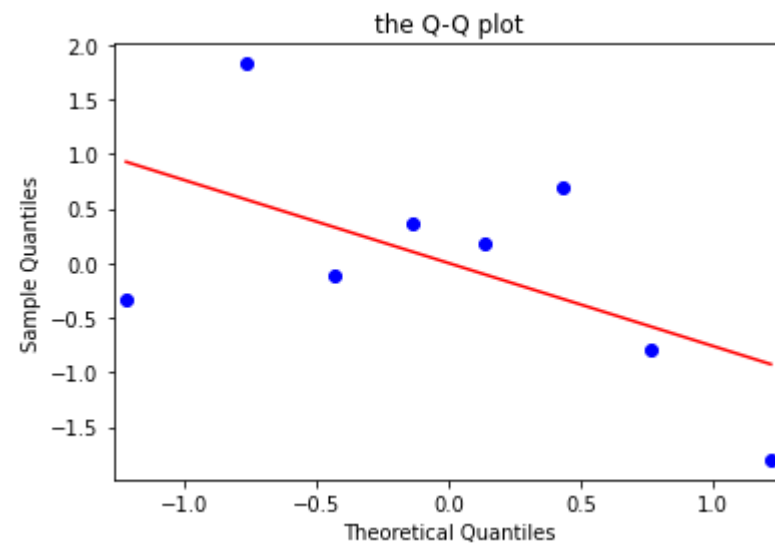
Exercise 4.30

a)

The power transformation lambda_1 that makes the x_1 values approximately normal:
[-0.23626729]

/usr/local/lib/python3.7/dist-packages/statsmodels/tools/_testing.py:19: FutureWarning: pandas.util.testing is deprecated. Use the functions in the public API at pandas.testing
import pandas.util.testing as tm

The Q-Q plot for the transformed data



```

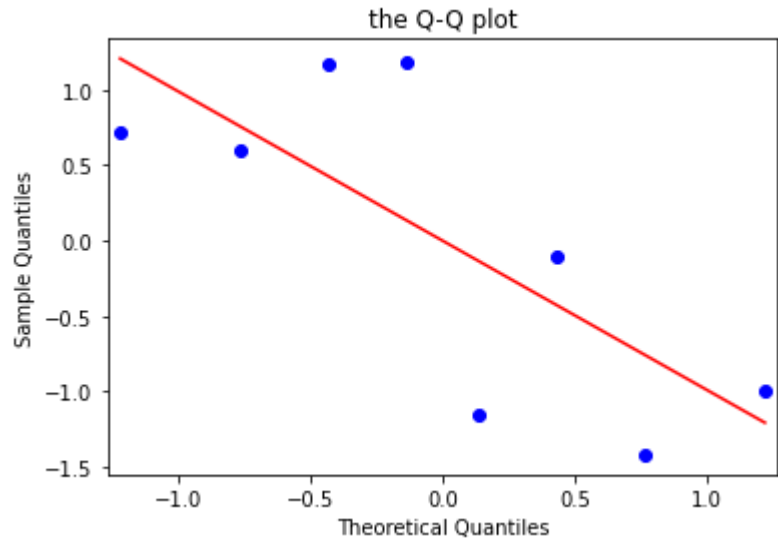
1 print('b) ')
2 pt_2=PowerTransformer('box-cox')
3 pt_2.fit(x_2.reshape(-1,1))
4 print('The power transformation lambda_2 that makes the x_2 values approximately normal: ')
5 print(pt_2.lambdas_)
6 print()
7
8 transform_x_2=pt_2.transform(x_2.reshape(-1,1))
9
10 print('The Q-Q plot for the transformed data')

```

```
11 sm.qqplot(transform_x_2,line='r')
12 plt.title('the Q-Q plot')
13 plt.show()

b)
The power transformation lambda_2 that makes the x_2 values approximately normal:
[0.23805971]
```

The Q-Q plot for the transformed data



```
1 print('c) ')
2 pt=PowerTransformer('box-cox')
3 pt.fit(x)
4 print('The power transformation lambda that makes the x_1, x_2 values jointly normal: ')
5 print(pt.lambdas_)
6 print()
7
8 print('The results are the same as those obtained in part a and b')
9
```

```
c)
The power transformation lambda that makes the x_1, x_2 values jointly normal:
[-0.23626729  0.23805971]
```

The results are the same as those obtained in part a and b

► Nhập dữ liệu

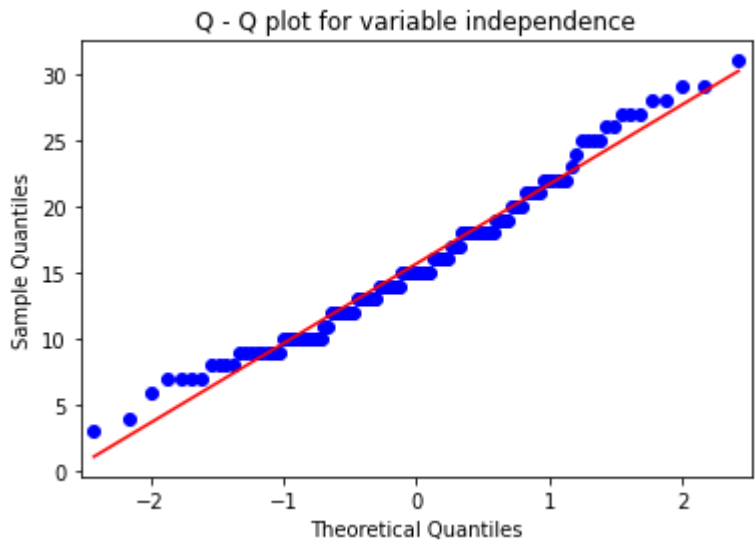
[] ↵ 1 cell hidden

▼ Bài 4.39

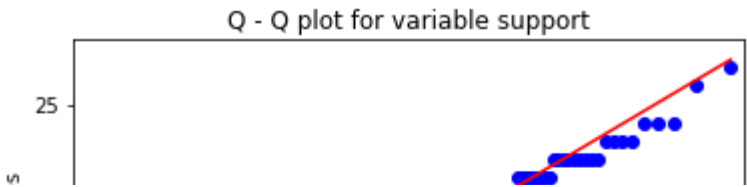
```
1 print('Exercise 4.39')
2 print('a) ')
```

```
3 dt_new=dt[:,0:5]
4 v_0,v_1,v_2,v_3,v_4=np.array([dt_new[:,0]],np.array([dt_new[:,1]]),np.array([dt_new[:,2]]),np.array([dt_new[:,3]]),np.array([dt_new[:,4]])
5
6 v_0_new=v_0.reshape(130,1)
7 v_1_new=v_1.reshape(130,1)
8 v_2_new=v_2.reshape(130,1)
9 v_3_new=v_3.reshape(130,1)
10 v_4_new=v_4.reshape(130,1)
11
12 v_0_sort=sorted(v_0_new.copy())
13 v_1_sort=sorted(v_1_new.copy())
14 v_2_sort=sorted(v_2_new.copy())
15 v_3_sort=sorted(v_3_new.copy())
16 v_4_sort=sorted(v_4_new.copy())
17 sm.qqplot(np.array([v_0_sort]).reshape(-1,1),line='r')
18 plt.title('Q - Q plot for variable independence')
19 plt.show()
20
```

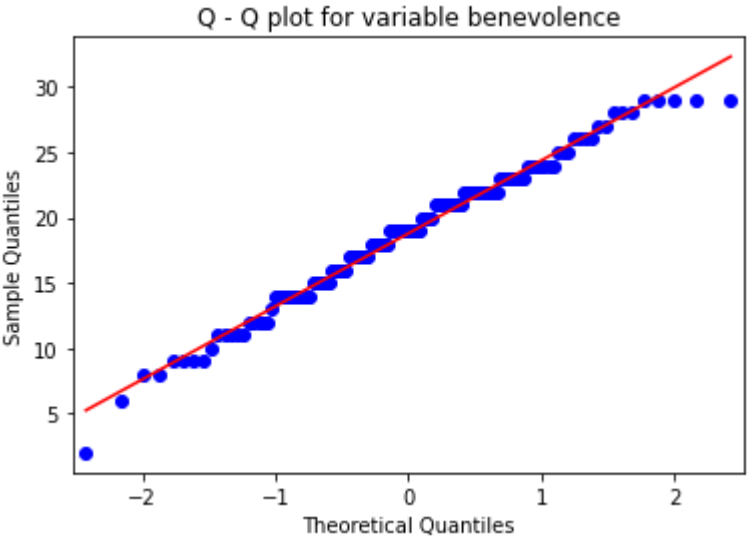
Exercise 4.39
a)



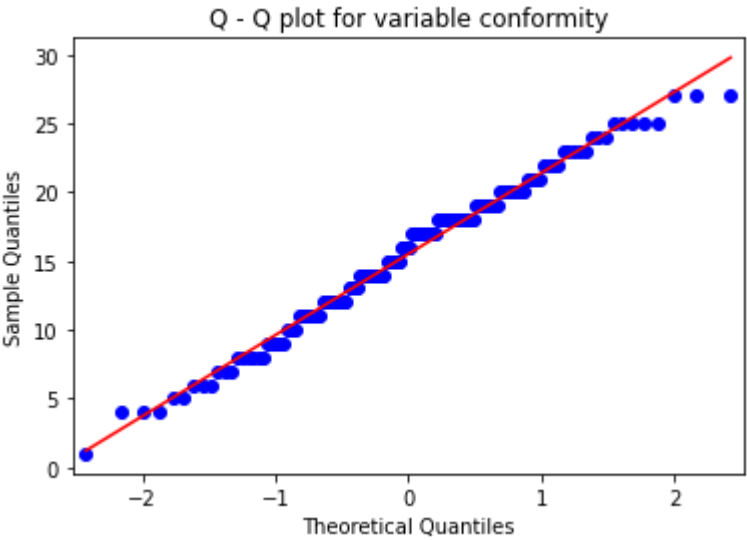
```
1 sm.qqplot(np.array([v_1_sort]).reshape(-1,1),line='r')
2 plt.title('Q - Q plot for variable support')
3 plt.show()
```

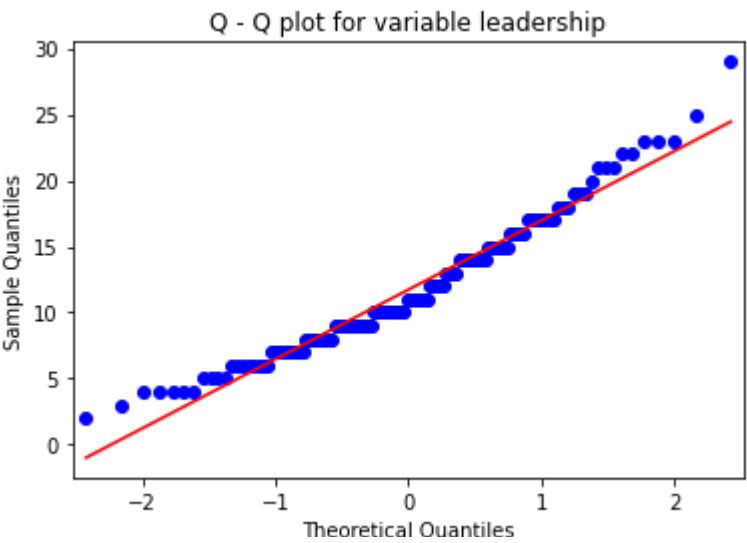
```
1 sm.qqplot(np.array([v_2_sort]).reshape(-1,1),line='r')
2 plt.title('Q - Q plot for variable benevolence')
3 plt.show()
```



```
1 sm.qqplot(np.array([v_3_sort]).reshape(-1,1),line='r')
2 plt.title('Q - Q plot for variable conformity')
3 plt.show()
```



```
1 sm.qqplot(np.array([v_4_sort]).reshape(-1,1),line='r')
2 plt.title('Q - Q plot for variable leadership')
3 plt.show()
```



Independence, support, leadership are the variables that are nonnormal

```
1 print('b) ')
2
3 v_0_mean=np.sum(v_0)/130
4 v_1_mean=np.sum(v_1)/130
5 v_2_mean=np.sum(v_2)/130
6 v_3_mean=np.sum(v_3)/130
7 v_4_mean=np.sum(v_4)/130
8
9 v_mean=np.array([[v_0_mean],[v_1_mean],[v_2_mean],[v_3_mean],[v_4_mean]])
10 data_cov = np.vstack([v_0,v_1,v_2,v_3,v_4])
11 S_new = np.cov(data_cov)
12 print('The Sample variance - covariance matrix S: ')
13 print(S_new)
14
15 sq_distance_new=[]
16 for i in range(len(dt)):
17     x_0=dt[i,0:5].reshape(5,1)
18     x_left=np.dot((x_0-v_mean).T,np.linalg.pinv(S_new))
19     x_result=np.dot(x_left,(x_0-v_mean))
20     sq_distance_new.append(x_result)
21
22 print()
23 print('the 50% percentile of chi square with 5 degree of freedom is 4.3515')
24 print()
25
26 count=0
27 for i in range(len(dt)):
28     if sq_distance_new[i]<4.3515:
29         count+=1
30
31 print('There are '+str(count/130)+ ' percentage of observations whose squared distances less than 4.3515')
32 print('Therefore, all five variables do not have multivariate normality')
```

b)

The Sample variance - covariance matrix S:

```
[[ 34.75020871 -4.27668456 -18.07179487 -15.97286822  5.71645796]
 [ -4.27668456 17.51341682  0.41979726 -7.86821705 -8.72331544]
 [-18.07179487  0.41979726 29.84472272  9.34883721 -13.94215862]
 [-15.97286822 -7.86821705  9.34883721 33.04263566 -9.94186047]
 [  5.71645796 -8.72331544 -13.94215862 -9.94186047 26.95796064]]
```

the 50% percentile of chi square with 5 degree of freedom is 4.3515

There are 0.49230769230769234 percentage of observations whose squared distances less than 4.3515

Therefore, all five variables do not have multivariate normality

```
1 print('c) ')
2 pt_new_0=PowerTransformer('box-cox')
3 pt_new_0.fit(v_0.reshape(-1,1))
4 transform_v_0=pt_new_0.transform(v_0.reshape(-1,1))
5 print('The transformation to make independence variable normal: ',pt_new_0.lambdas_)
6 print()
7
8 pt_new_1=PowerTransformer('box-cox')
9 pt_new_1.fit(v_1.reshape(-1,1))
10 transform_v_1=pt_new_1.transform(v_1.reshape(-1,1))
11 print('The transformation to make support variable normal: ',pt_new_1.lambdas_)
12 print()
13
14 pt_new_4=PowerTransformer('box-cox')
15 pt_new_4.fit(v_4.reshape(-1,1))
16 transform_v_4=pt_new_4.transform(v_4.reshape(-1,1))
17 print('The transformation to make leadership variable normal: ',pt_new_4.lambdas_)
18 print()
```

c)

The transformation to make independence variable normal: [0.52377241]

The transformation to make support variable normal: [1.39626145]

The transformation to make leadership variable normal: [0.38154699]

✓ 0s completed at 11:44 PM

● ✕