

Q1.

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4
5 class datapreprocessing: 2 usages new *
6     def __init__(self,data): new *
7         self.data = data
8         self.normalized_data = None
9         self.standardized_data = None
10        self.iqr_data=None
11
12    def Normalized(self): 1 usage new *
13        self.normalized_data = (self.data - self.data.min())/(self.data.max()-self.data.min())
14
15    def Standardized(self): 1 usage new *
16        self.standardized_data = (self.data - self.data.mean())/(self.data.std())
17
18    def IQR(self): 1 usage new *
19        Q1 = self.data.quantile(0.25)
20        Q3 = self.data.quantile(0.75)
21        IQR = Q3-Q1
22        self.iqr_data = (self.data - Q1)/IQR
23
24    def Show_original(self): 1 usage new *
25        plt.figure(figsize=(6,4))
26        for column in self.data.columns:
27            plt.plot(*args: self.data[column],label=column)
28            plt.title("Original AAPL data set")
29            plt.xlabel(column)
30            plt.ylabel('Frequency')
31            plt.legend()
32            plt.show()
```

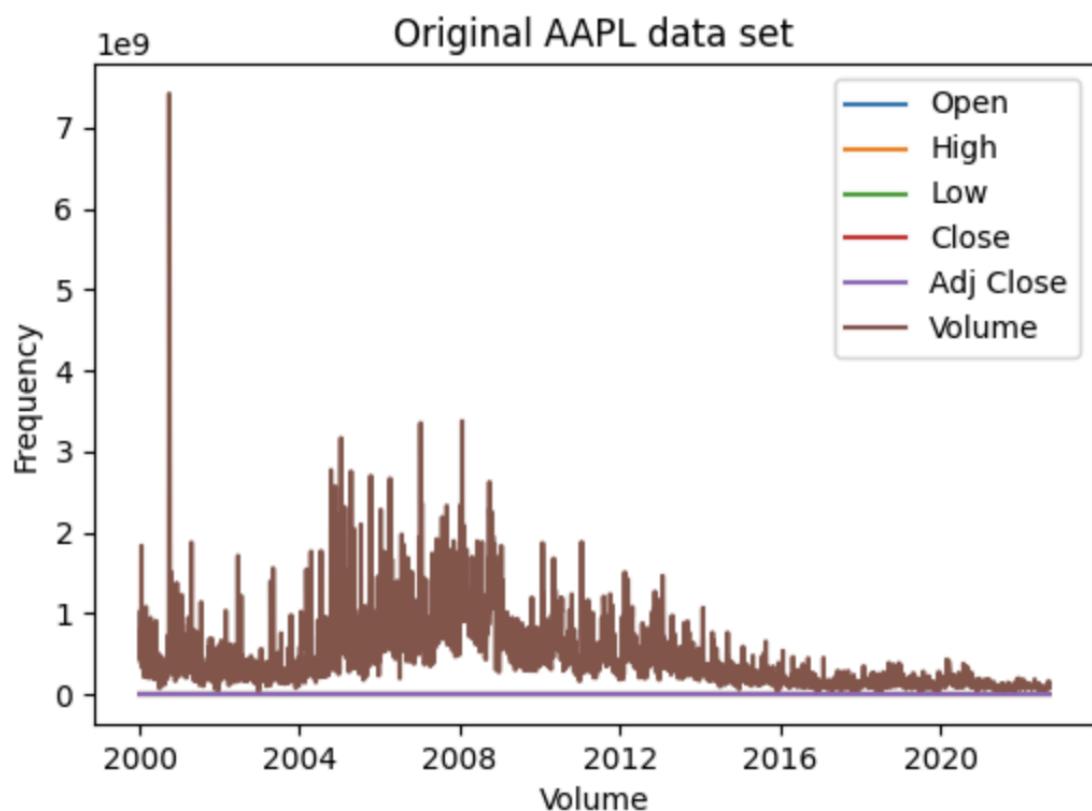
```
34     def Show_normalized(self): 1 usage  new *
35         if self.normalized_data is None:
36             self.Normalized()
37             plt.figure(figsize=(6, 4))
38             for column in self.normalized_data.columns:
39                 plt.plot(*args: self.normalized_data[column], label=column)
40             plt.title("Normalized AAPL data set")
41             plt.xlabel(column)
42             plt.ylabel('Frequency')
43             plt.legend()
44             plt.show()
45
46     def Show_standardized(self): 1 usage  new *
47         if self.standardized_data is None:
48             self.Standardized()
49             plt.figure(figsize=(6, 4))
50             for column in self.standardized_data.columns:
51                 plt.plot(*args: self.standardized_data[column], label=column)
52             plt.title("Standardized AAPL data set")
53             plt.xlabel(column)
54             plt.ylabel('Frequency')
55             plt.legend()
56             plt.show()
57
58     def Show_IQR(self): 1 usage  new *
59         if self.iqr_data is None:
60             self.IQR()
61             plt.figure(figsize=(6, 4))
62             for column in self.iqr_data.columns:
63                 plt.plot(*args: self.iqr_data[column], label=column)
64             plt.title("IQR transformation AAPL data set")
65             plt.xlabel(column)
66             plt.ylabel('Frequency')
67             plt.legend()
68             plt.show()
```

Q2.

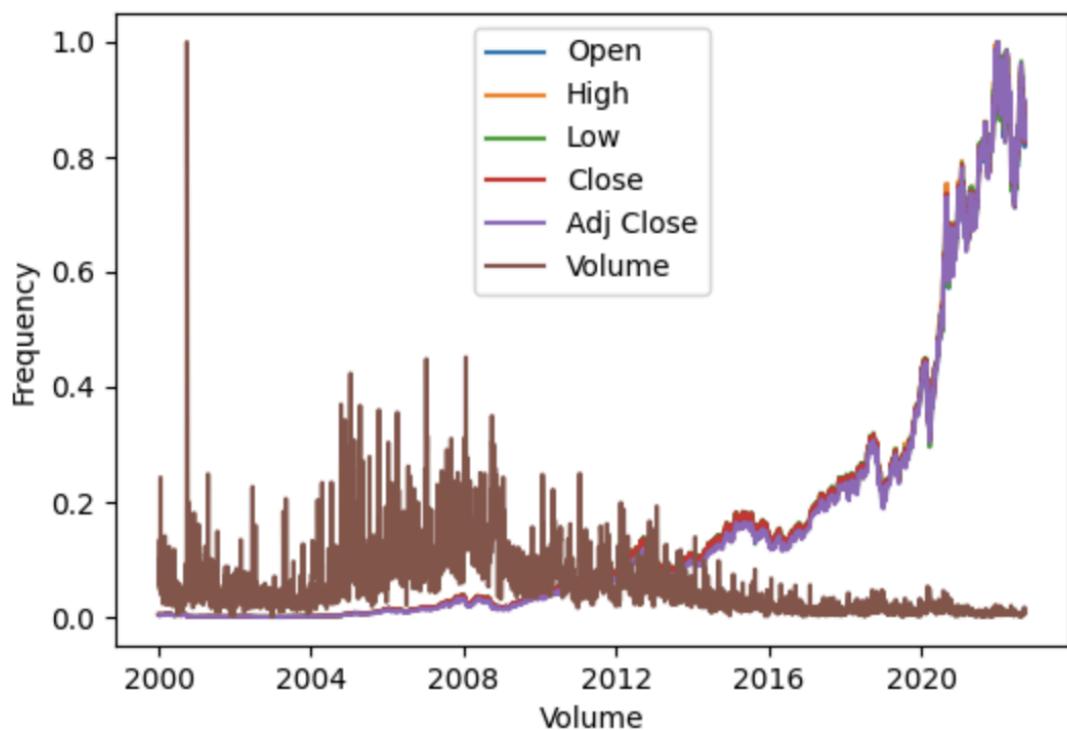
Code:

```
1 %%-----Q2-----  
2 import pandas as pd  
3 from fontTools.misc.bezierTools import epsilon  
4 from pandas_datareader import data  
5 import matplotlib.pyplot as plt  
6 import yfinance as yf  
7 from scipy.special import delta  
8  
9 from Home_Work_3_Q1 import datapreprocessing  
10 apple_stock = yf.download(tickers='AAPL', start="2000-01-01", end="2022-09-25")  
11 df = datapreprocessing(apple_stock)  
12 df.Show_original()  
13 df.Show_normalized()  
14 df.Show_standardized()  
15 df.Show_IQR()
```

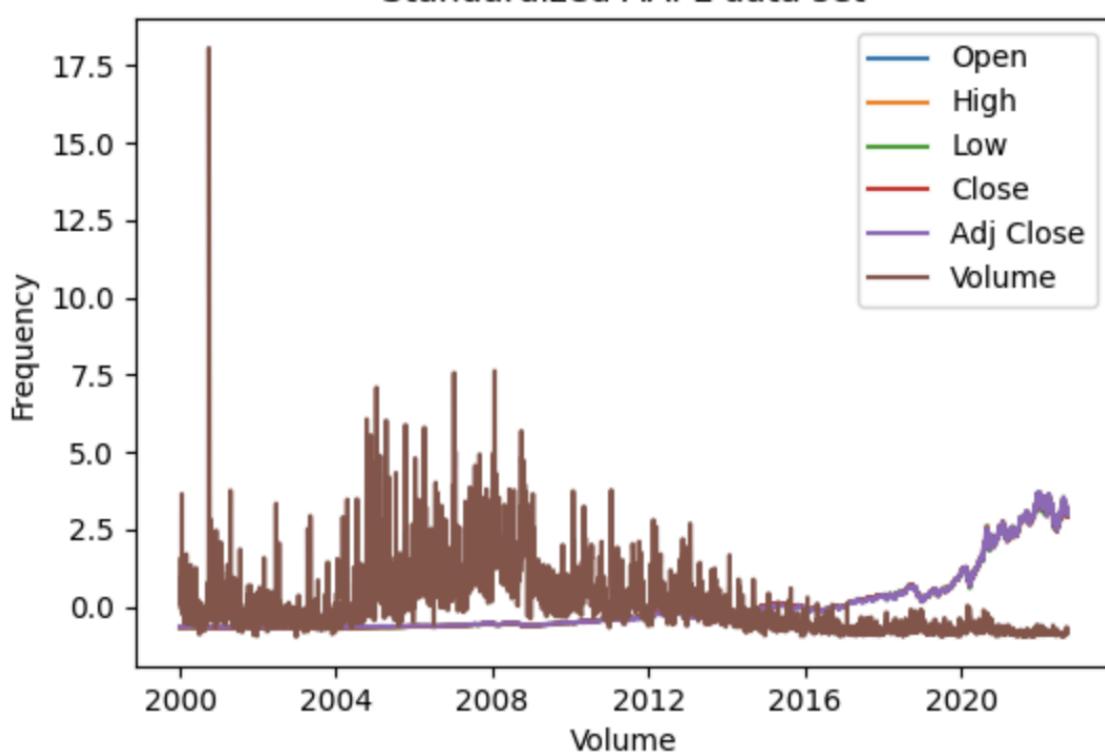
Output:

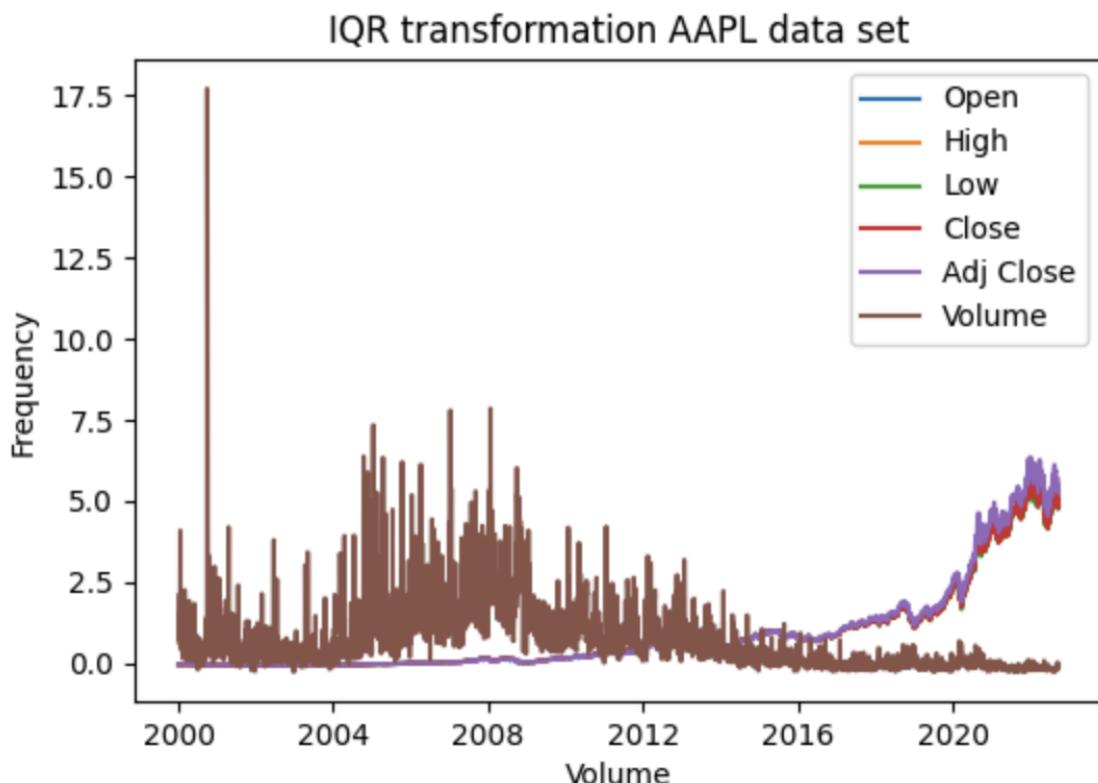


Normalized AAPL data set



Standardized AAPL data set





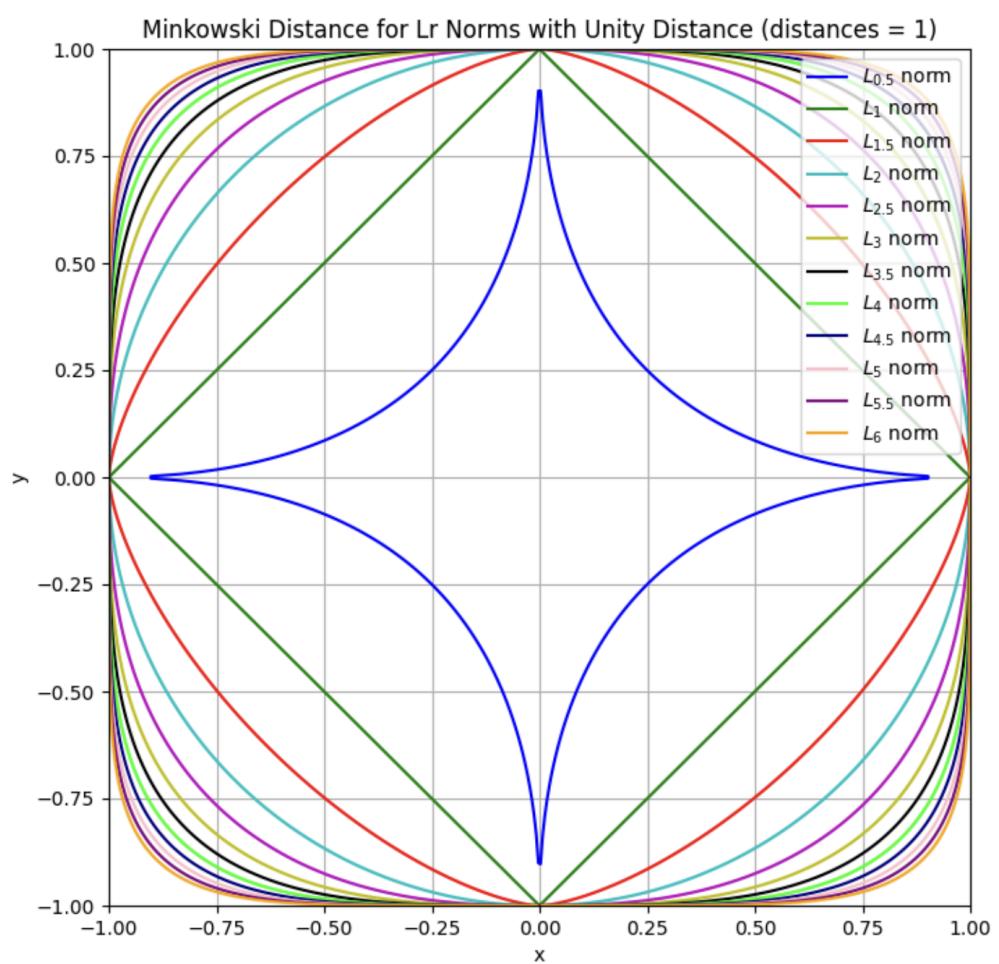
Q3.Code:

```

15 > %%-----Q3-----
16 import pandas as pd
17 import numpy as np
18 import matplotlib.pyplot as plt
19 r_values = [0.5,1,1.5,2,2.5,3,3.5,4,4.5,5,5.5,6]
20 colors = ['b', 'g', 'r', 'c', 'm', 'y', 'k', 'lime', 'navy', 'pink',
21 'purple', 'orange']
22 legend_handles = []
23 x = np.linspace(-1, stop: 1, num: 400)
24 y = np.linspace(-1, stop: 1, num: 400)
25 X,Y = np.meshgrid(*xi: x,y)
26 plt.figure(figsize=(8,8))
27 for r,color in zip(r_values, colors):
28     Z=(np.abs(X)**r+np.abs(Y)**r)**(1/r)
29     plt.contour(*args: X,Y,Z,levels=[1],colors=[color])
30     legend_handles.append(plt.Line2D(xdata: [0], ydata: [0], linestyle='--', color=color,
31 label=f'L_{r} norm'))
32 plt.title('Minkowski Distance for Lr Norms with Unity Distance (distances = 1)')
33 plt.xlabel('x')
34 plt.ylabel('y')
35 plt.legend(handles=legend_handles)
36 plt.grid(True)
37 plt.show()

```

Output:



Q4.

- 4) The Manhattan distance L_1 -norm, between two points (x_1, y_1) and (x_2, y_2)

$$d((x_1, y_1), (x_2, y_2)) = |x_1 - x_2| + |y_1 - y_2|$$

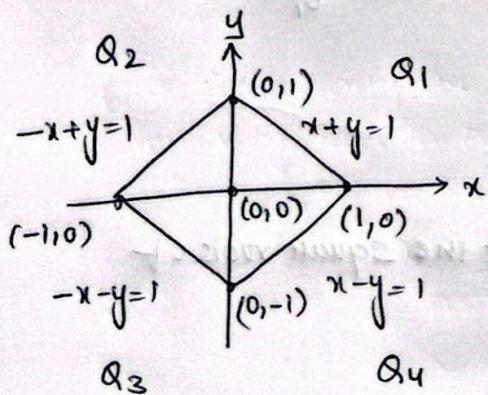
To find all points at a distance 1 unit from origin $(0,0)$ under L_1 -norm

$$|x - 0| + |y - 0| = 1$$

$$|x| + |y| = 1$$

considering all quadrants :-

First quadrant	Second quadrant	3rd quadrant	4th quadrant
$x \geq 0, y \geq 0$	$x \leq 0, y \geq 0$	$x \leq 0, y \leq 0$	$x \geq 0, y \leq 0$
$x+y=1$	$-x+y=1$	$-x-y=1$	$x-y=1$



plotting above 4 lines:-

$$x+y=1$$

$$-x+y=1$$

$$-x-y=1$$

$$x-y=1$$

we can see they form diamond shape centered $(0,0)$. The vertices of the diamond are $(1,0)$, $(0,1)$, $(-1,0)$ and $(0,-1)$.

: Therefore we have proven that L_1 Minkowski distance with unity distance forms a diamond shape.

Q5.

5) L_∞ Minkowski distance is defined as,

$$D(x, y) = \max(|x_1 - y_1|, |x_2 - y_2|, \dots, |x_n - y_n|)$$

In 2D, $D(x, y) = \max(|x_1 - y_1|, |x_2 - y_2|)$

→ Let us take a point $P(a, b)$

Now we need to find $q(x, y)$ such that L_∞ Minkowski distance between points $P, q = 1$

→ We look into all cases where $\max(|x-a|, |y-b|) = 1$.

$$\max(|x-a|, |y-b|) = 1$$

which refers as

$$|x-a| \text{ or } |y-b| \text{ or both } |x-a|, |y-b| = 1$$

① if $|x-a|=1$

$$x-a=1 \quad x-a=-1$$

$$x=a+1 \quad x=a-1$$

② if $|y-b|=1$

$$y-b=1 \quad b-y=1$$

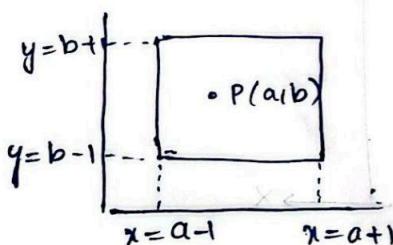
$$y=b+1 \quad y=b-1$$

We derive at four possible points

$$(a+1, b) \quad (a-1, b) \quad (a, b+1) \quad (a, b-1)$$

The above 4 points are at distance of 1 unit from $P(a, b)$

They form a square with centre $P(a, b)$.



Hence proved that L_∞ distance a square in 2 dimensions

Q6.**6.a)**

```
39 ▶ #%%-----Q6-----
40 import numpy as np
41 from prettytable import PrettyTable
42 import matplotlib.pyplot as plt
43 import pandas as pd
44 np.random.seed(5808)
45 x = np.random.normal(loc=1,np.sqrt(2), size=1000)
46 epsilon= np.random.normal(loc=2,np.sqrt(3), size=1000)
47 y = x+ epsilon
48 X = np.vstack((x,y)).T
49 n = X.shape[0]
50 X_centered = X- X.mean(axis=0)
51 cov_matrix = (1/(n-1))*(np.dot(X_centered.T,X_centered))
52 cov_matrix_table = PrettyTable()
53 cov_matrix_table.field_names =[ "Covariance Matrix Element","Value"]
54 for i in range(cov_matrix.shape[0]):
55     for j in range(cov_matrix.shape[1]):
56         element_name = f'Cov(x{i+1},x{j+1})'
57         value = "{:.2f}".format(cov_matrix[i,j])
58         cov_matrix_table.add_row([element_name,value])
59 print("Estimated Covariance Matrix:")
60 print(cov_matrix_table)
```

```
⌚ Estimated Covariance Matrix:
⌚ +-----+-----+
⌚ | Covariance Matrix Element | Value |
⌚ +-----+-----+
⌚ | Cov(x1,x1) | 2.09 |
⌚ | Cov(x1,x2) | 1.97 |
⌚ | Cov(x2,x1) | 1.97 |
⌚ | Cov(x2,x2) | 4.41 |
⌚ +-----+-----+
```

Justification of why the diagonal elements of calculated covariance matrix Σ is correct:

We know the diagonal elements of the covariance matrix represent the variance of the individual elements.

For x we know (given in question) it is 2 and for y when we calculate we get the variance as 4.41 (see the below attached code)and even when compare those values which we got in the covariance Cov(x1,x1) which is almost equal to 2 and Cov(x2,x2) which is almost equal to 4.41 below the python code to calculate the variance of y.

```

109 mean_x=1
110 variance_x=2
111 sample_size=1000
112 np.random.seed(5808)
113 x = np.random.normal(mean_x,np.sqrt(variance_x),sample_size)
114 mean_epsilon = 2
115 variance_epsilon = 3
116 epsilon = np.random.normal(mean_epsilon,np.sqrt(variance_epsilon),sample_size)
117 y = x + epsilon
118 variance_y = np.var(y)
119 print(f"Variance of y: {variance_y:.2f}")

```

Variance of y: 4.41

>>>

6.b)

```

62 eigen_values, eigen_vectors = np.linalg.eig(cov_matrix)
63 eigen_table = PrettyTable()
64 eigen_table.field_names = ["Eigen value","Eigen vector"]
65 for i in range(len(eigen_values)):
66     eigenvalue_formatted = "{:.2f}".format(eigen_values[i])
67     eigenvector_formatted = np.array2string(eigen_vectors[:,i],formatter={'float_kind' : lambda x: "{:.2f}".format(x)})
68     eigen_table.add_row([eigenvalue_formatted,eigenvector_formatted])
69 print("Eigen values & Eigen Vectors: ")
70 print(eigen_table)

+E-----+-----+
Eigen values & Eigen Vectors:
@+-----+-----+
@ | Eigen value | Eigen vector |
+-----+-----+
| 0.97 | [-0.87 0.50] |
| 5.54 | [-0.50 -0.87] |
+-----+-----+

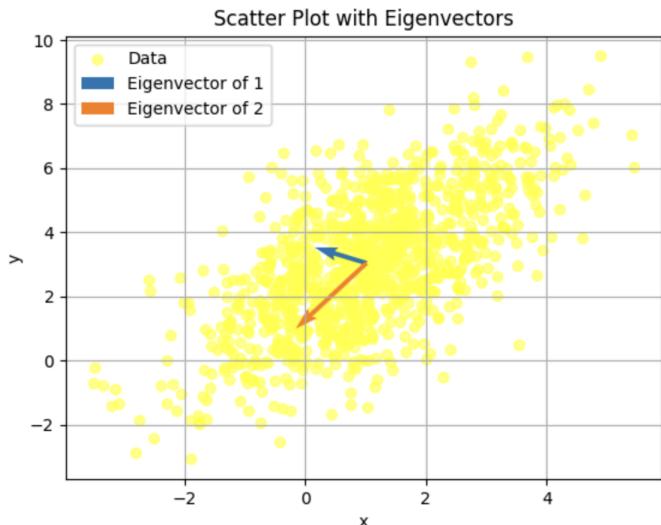
```

6.c)

```

71 plt.scatter(x, y, label='Data', alpha=0.5,color='yellow')
72 for i in range(len(eigen_values)):
73     plt.quiver(
74         *args: X.mean(axis=0)[0],
75         X.mean(axis=0)[1],
76         eigen_vectors[0, i] *np.sqrt(eigen_values[i]),
77         eigen_vectors[1, i] *np.sqrt(eigen_values[i]),
78         angles='xy',
79         scale_units='xy', scale=1,
80         label=f'Eigenvector of {i+1}',color=f'C{i}')
81 plt.xlabel('x')
82 plt.ylabel('y')
83 plt.title('Scatter Plot with Eigenvectors')
84 plt.legend()
85 plt.grid(True)
86 plt.show()

```



If I were to drop one feature (x or y), I would typically drop the feature corresponding to the smaller variance. Retaining the dimension with higher variance (larger eigenvalue) preserves more information. **So I will drop feature x in this case.**

- Eigenvector corresponding to the maximum eigenvalue represents the direction of maximum variance in the data. In other words, it points along the axis of the data's spread. This direction is often considered the "principal component" of the data.
- Eigenvector corresponding to the minimum eigenvalue represents the direction of minimum variance.

6.d)

```

88     eigen_values_scaled = eigen_values * (n - 1)
89     sqrt_eigenvalues_scaled = np.sqrt(eigen_values_scaled)
90     singular_values = np.linalg.svd(X_centered, compute_uv=False)
91     singular_values_table=PrettyTable()
92     singular_values_table.field_names = ["Singular Value"]
93     for i in range(len(singular_values)):
94         singular_value_formatted = "{:.2f}".format(singular_values[i])
95         singular_values_table.add_row([singular_value_formatted])
96     print("Singular Values Table:")
97     print(singular_values_table)
98     print("Scaled eigenvalues of  $X^T X$ :")
99     print(eigen_values_scaled)
100    print("Square roots of the scaled eigenvalues of  $X^T X$ :")
101    print(sqrt_eigenvalues_scaled)

```

```

> Singular Values Table:
<
+-----+
| Singular Value |
+-----+
|      74.39     |
|      31.11     |
+-----+
Scaled eigenvalues of X^T X:
[ 967.66474384 5534.26711837]
Square roots of the scaled eigenvalues of X^T X:
[31.10731014 74.392655  ]

```

- The **singular values** of a matrix (`X_centered`) are **the square roots of the eigenvalues** of the matrix $((X_{\text{centered}})^T) * (X_{\text{centered}})$. This is because, in Singular Value Decomposition (SVD), the matrix X is decomposed into $U\Sigma V^T$, and the matrix $X^T X = V\Sigma^2 V^T$. The eigenvalues of $X^T X$ are the squares of the singular values, which explains why the singular values are the

square roots of these eigenvalues.

5)

$$d. \quad \tilde{x} = U \Sigma V^T$$

U, V are orthogonal matrix ($U^T U = V^T V = I$)

$$U^T U = I$$

$$V^T V = I$$

$$\tilde{x}^T \tilde{x} = (U \Sigma V^T)^T (U \Sigma V^T)$$

$$= ((V^T)^T \Sigma^T U^T) \cdot (U \Sigma V^T)$$

$$= V \cdot \Sigma^T \cdot U^T \cdot U \cdot \Sigma \cdot V^T$$

$$= (V \cdot V^T) \cdot (U \cdot U^T) (\Sigma^T \cdot \Sigma)$$

$$= I \cdot I (\Sigma^T \cdot \Sigma) = \Sigma^2$$

$$\tilde{x}^T \tilde{x} = \Sigma^2$$

(Σ is diagonal matrix)

$$\downarrow \qquad \qquad \qquad \xrightarrow{\text{singular value } \sigma_i^2}$$

λ_i eigen values

of $\tilde{x}^T \tilde{x}$

\therefore singular values of \tilde{x} are the square roots of
the eigen values of $\tilde{x}^T \tilde{x}$.

6.e)

```

96 df = pd.DataFrame({'x': x, 'y': y})
97 correlation_matrix = df.corr()
98 correlation_matrix = correlation_matrix.round(2)
99 print("Correlation Matrix:")
100 print(correlation_matrix)
101 mean_x = np.mean(x)
102 mean_y = np.mean(y)
103 numerator_x_y = np.sum((x - mean_x) * (y - mean_y))
104 denominator_x = np.sqrt(np.sum((x - mean_x) ** 2))
105 denominator_y = np.sqrt(np.sum((y - mean_y) ** 2))
106 correlation_coefficient = numerator_x_y / (denominator_x * denominator_y)
107 print(f'Sample Pearson Correlation Coefficient between x and y (r):{ correlation_coefficient:.2f}')

```

Correlation Matrix:

	x	y
x	1.00	0.65
y	0.65	1.00

Sample Pearson Correlation Coefficient between x and y (r):0.65

Justification:(attached image below)

mathematical relation between covariance and correlation $\rho = \sigma_{xy} / (\sigma_x * \sigma_y)$

- ρ (rho) is the Pearson correlation coefficient.
- σ_{xy} is the covariance between X and Y.
- σ_x and σ_y are the standard deviations of X and Y, respectively.

6.e)

$$\text{co-relation coefficient} = \frac{\rho_{xy}}{\sigma_x \times \sigma_y}$$

$$\text{co-variance matrix} = \begin{bmatrix} 2.09 & 1.97 \\ 1.97 & 4.41 \end{bmatrix} \quad (\text{derived in 6.a})$$

$$\text{co-relation matrix} = \begin{bmatrix} 1.00 & 0.65 \\ 0.65 & 1.00 \end{bmatrix}$$

$$\rho(x, x) = \frac{\text{cov}(x, x)}{\sigma_x \cdot \sigma_x} = \frac{\sigma_x^2}{\sigma_x^2} = 1$$

$$\rho(y, y) = \frac{\text{cov}(y, y)}{\sigma_y \cdot \sigma_y} = \frac{\sigma_y^2}{\sigma_y^2} = 1$$

$$\rho(x, y) = \frac{\text{cov}(x, y)}{\sigma_x \cdot \sigma_y} = \frac{1.97}{\sqrt{2.09} \sqrt{4.41}} \approx 0.65$$

This matches the value in the correlation matrix, confirming that the correlation coefficient matrix is the normalized form of covariance matrix, using the standard deviations to scale the co-variance.

- The correlation matrix is derived by normalizing the co-variance matrix

Q.7)

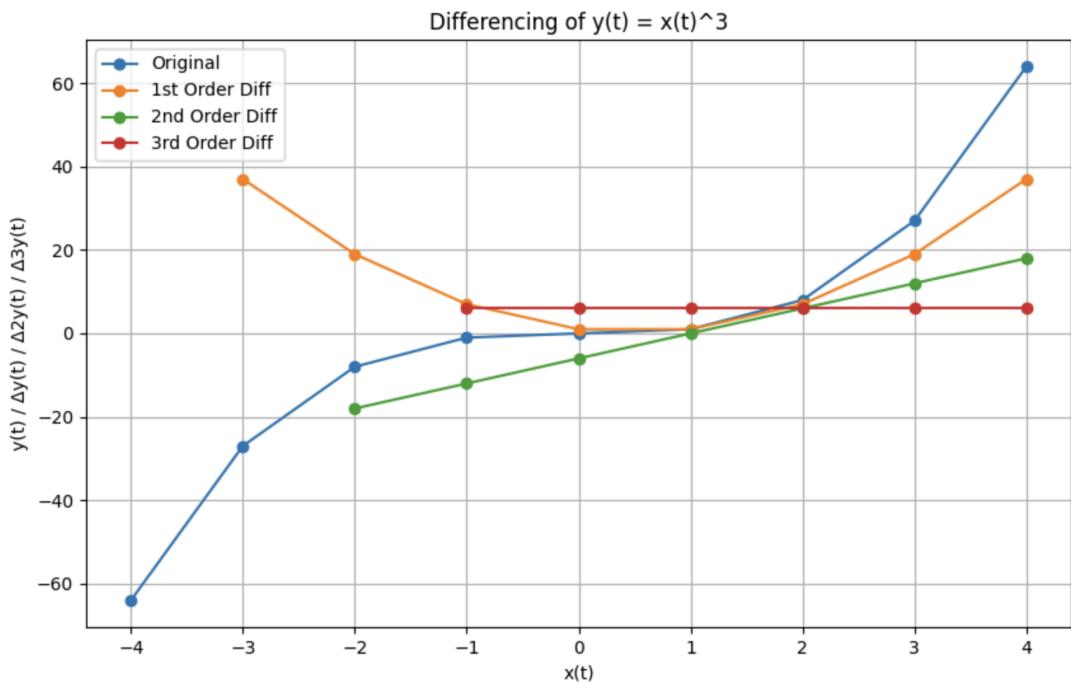
```

121 %%-----Q7-----
122 import numpy as np
123 import pandas as pd
124 import matplotlib.pyplot as plt
125 from prettytable import PrettyTable
126 x_t = np.arange(-4,5,1)
127 y_t = x_t ** 3
128
129 def first_order_diff(y): 1 usage  new *
130     return np.diff(y,n=1)
131 def second_order_diff(y): 1 usage  new *
132     return np.diff(y,n=2)
133 def third_order_diff(y): 1 usage  new *
134     return np.diff(y,n=3)
135
136 delta_y_t = first_order_diff(y_t)
137 delta2_y_t = second_order_diff(y_t)
138 delta3_y_t = third_order_diff(y_t)
139
140 first_order_diff = np.pad(delta_y_t.astype(float),(1,0),mode='constant',constant_values=(np.nan,))
141 second_order_diff = np.pad(delta2_y_t.astype(float),(2,0),mode='constant',constant_values=(np.nan,))
142 third_order_diff = np.pad(delta3_y_t.astype(float),(3,0),mode='constant',constant_values=(np.nan,))
143
144 data = {'x(t)': x_t, 'y(t)': y_t, 'Δy(t)': first_order_diff, 'Δ2y(t)':second_order_diff, 'Δ3y(t)': third_order_diff}
145 df = pd.DataFrame(data)
146 print(df)
147 plt.figure(figsize=(10, 6))
148 plt.plot(*args: x_t, y_t, label='Original', marker='o')
149 plt.plot(*args: x_t, first_order_diff, label='1st Order Diff', marker='o')
150 plt.plot(*args: x_t, second_order_diff, label='2nd Order Diff', marker='o')
151 plt.plot(*args: x_t, third_order_diff, label='3rd Order Diff', marker='o')
152 plt.xlabel('x(t)')
153 plt.ylabel('y(t) / Δy(t) / Δ2y(t) / Δ3y(t)')
154 plt.title('Differencing of y(t) = x(t)^3')
155 plt.legend()
156 plt.grid(True)
157 plt.show()

```

	x(t)	y(t)	Δy(t)	Δ2y(t)	Δ3y(t)
0	-4	-64	NaN	NaN	NaN
1	-3	-27	37.0	NaN	NaN
2	-2	-8	19.0	-18.0	NaN
3	-1	-1	7.0	-12.0	6.0
4	0	0	1.0	-6.0	6.0
5	1	1	1.0	0.0	6.0
6	2	8	7.0	6.0	6.0
7	3	27	19.0	12.0	6.0
8	4	64	37.0	18.0	6.0

>>>



$$\Delta y(t) = y(t+1) - y(t)$$

$$\Delta^2 y(t) = \Delta y(t+1) - \Delta y(t)$$

$$\Delta^3 y(t) = \Delta^2 y(t+1) - \Delta^2 y(t)$$

- This is the cubic function $y(t)=x^3$ (blue colored graph) which curves sharply on either side of $x=0$ with negative values for negative $x(t)$ and positive values for positive $x(t)$.
- 1st Order Diff (orange colored graph) shows the rate of change of the cubic function. It starts off with large values (because the cubic curve is steep) and decreases toward 0 at $x=0$, then increases again for positive $x(t)$.
- 2nd Order Diff (green colored graph) shows how the rate of change (slope) is itself changing. The negative values show that the slope is decreasing as x approaches 0, and the positive values show that the slope increases after 0 as it is a 2nd order diff of cubic function hence it's a linear graph.
- 3rd Order Diff (red coloured graph) is constant, indicating that the change in acceleration is uniform and 3rd order diff of x^3 is 6.