

# WEEK 03

## Import CSV File

In [6]:

```
1 import pandas as pd
2 df = pd.read_csv(r"C:\Users\Anusha V\Desktop\diabetes.csv")
```

In [7]:

```
1 df.head(5)
```

Out[7]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunc
0	6	148	72	35	0	33.6	0.
1	1	85	66	29	0	26.6	0.
2	8	183	64	0	0	23.3	0.
3	1	89	66	23	94	28.1	0.
4	0	137	40	35	168	43.1	2.



In [5]:

```
1 df.tail(5)
```

Out[5]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFur
763	10	101	76	48	180	32.9	
764	2	122	70	27	0	36.8	
765	5	121	72	23	112	26.2	
766	1	126	60	0	0	30.1	
767	1	93	70	31	0	30.4	



In [10]:

```
1 df.describe()
```

Out[10]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	Diab
count	768.000000	768.000000	768.000000	768.000000	768.000000	768.000000	
mean	3.845052	120.894531	69.105469	20.536458	79.799479	31.992578	
std	3.369578	31.972618	19.355807	15.952218	115.244002	7.884160	
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	
25%	1.000000	99.000000	62.000000	0.000000	0.000000	27.300000	
50%	3.000000	117.000000	72.000000	23.000000	30.500000	32.000000	
75%	6.000000	140.250000	80.000000	32.000000	127.250000	36.600000	
max	17.000000	199.000000	122.000000	99.000000	846.000000	67.100000	



# Aggregation functions

```
In [14]: 1 sum = df['Age'].sum()  
2 sum
```

```
Out[14]: 25529
```

```
In [16]: 1 mean = df['Age'].mean()  
2 mean
```

```
Out[16]: 33.240885416666664
```

```
meadian = df[
```

```
In [17]: 1 median = df['Age'].median()  
2 median
```

```
Out[17]: 29.0
```

```
In [20]: 1 sum_result = df['Age'].max()  
2 sum_result
```

```
Out[20]: 81
```

```
In [22]: 1 count = df['Age'].count()  
2 count
```

```
Out[22]: 768
```

```
In [23]: 1 std = df['Age'].std()  
2 std
```

```
Out[23]: 11.76023154067868
```

```
In [24]: 1 var = df['Age'].var()  
2 var
```

```
Out[24]: 138.30304589037365
```

```
In [26]: 1 mode = df['Age'].mode()  
2 mode
```

```
Out[26]: 0    22  
Name: Age, dtype: int64
```

```
In [32]: 1 pivot_table = df.pivot_table(index = 'Pregnancies', aggfunc='mean')
2 pivot_table
```

Out[32]:

	Age	BMI	BloodPressure	DiabetesPedigreeFunction	Glucose
Pregnancies					
0	27.603604	34.290090	67.153153	0.520838	123.000000
1	27.370370	31.372593	67.792593	0.486496	112.748148
2	27.194175	30.583495	63.252427	0.491660	110.796117
3	29.026667	30.425333	66.586667	0.432147	123.586667
4	32.779412	32.141176	70.029412	0.446353	125.117647
5	39.035088	33.192982	76.210526	0.396421	118.859649
6	39.340000	30.290000	68.420000	0.429520	120.800000
7	41.111111	32.631111	70.777778	0.443622	136.444444
8	45.368421	31.568421	75.184211	0.504711	131.736842
9	44.178571	31.707143	77.892857	0.550679	131.392857
10	42.666667	30.641667	70.208333	0.454167	120.916667
11	44.545455	38.563636	74.181818	0.522545	126.454545
12	47.444444	32.344444	76.333333	0.444333	113.555556
13	44.500000	35.000000	73.800000	0.463300	125.500000
14	42.000000	35.100000	70.000000	0.312000	137.500000
15	43.000000	37.100000	70.000000	0.153000	136.000000
17	47.000000	40.900000	72.000000	0.817000	163.000000



```
In [31]: 1 melt_df = pd.melt(df,id_vars=['Pregnancies'], var_name='Glucose')
2 melt_df
```

Out[31]:

	Pregnancies	Glucose	value
0	6	Glucose	148.0
1	1	Glucose	85.0
2	8	Glucose	183.0
3	1	Glucose	89.0
4	0	Glucose	137.0
...	...	...	...
6139	10	Outcome	0.0
6140	2	Outcome	0.0
6141	5	Outcome	0.0
6142	1	Outcome	1.0
6143	1	Outcome	0.0

6144 rows × 3 columns

```
In [38]: 1 from functools import reduce
2 def multiply(x, y):
3     return x * y
4 column_to_reduce = 'Age'
5 result = reduce(multiply, df['Age'])
6 print(result)
7
```

In [39]:

```

1 def map_grades(grade):
2     grade_maping = {
3         'A':90,
4         'B':80,
5         'C':70,
6         'D':60,
7         'E':50,
8         'F':40
9     }
10 LeagueIndex = 'Grade'
11 df['Glucose'] = df['BloodPressure'].map(map_grades)
12 print(df)

```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	\
0	6	None	72	35	0	33.6	
1	1	None	66	29	0	26.6	
2	8	None	64	0	0	23.3	
3	1	None	66	23	94	28.1	
4	0	None	40	35	168	43.1	
..	...	...	...	...	...	...	
763	10	None	76	48	180	32.9	
764	2	None	70	27	0	36.8	
765	5	None	72	23	112	26.2	
766	1	None	60	0	0	30.1	
767	1	None	70	31	0	30.4	
	DiabetesPedigreeFunction	Age	Outcome				
0	0.627	50	1				
1	0.351	31	0				
2	0.672	32	1				
3	0.167	21	0				
4	2.288	33	1				
..	...	...	...				
763	0.171	63	0				
764	0.340	27	0				
765	0.245	30	0				
766	0.349	47	1				
767	0.315	23	0				

[768 rows x 9 columns]

In [41]:

```

1 def score_condition(score):
2     return score > 5
3 Likes = 'BloodPressure'
4 filtered_df = df[df['Glucose'].apply(score_condition)]
5 print(filtered_df)

```

```

-
TypeError                                     Traceback (most recent call last)
t)
~\AppData\Local\Temp\ipykernel_6124\1532437874.py in <module>
      2     return score > 5
      3 Likes = 'BloodPressure'
----> 4 filtered_df = df[df['Glucose'].apply(score_condition)]
      5 print(filtered_df)

~\anaconda3\lib\site-packages\pandas\core\series.py in apply(self, func, convert_dtype, args, **kwargs)
    4431         dtype: float64
    4432         """
-> 4433         return SeriesApply(self, func, convert_dtype, args, kwargs).apply()
    4434
    4435     def _reduce(

~\anaconda3\lib\site-packages\pandas\core\apply.py in apply(self)
    1086         return self.apply_str()
    1087
-> 1088         return self.apply_standard()
    1089
    1090     def agg(self):

~\anaconda3\lib\site-packages\pandas\core\apply.py in apply_standard(self)
    1141             # List[Union[Callable[..., Any], str]]]]"; expected
ed
    1142             # "Callable[[Any], Any]"
-> 1143             mapped = lib.map_infer(
    1144                 values,
    1145                 f, # type: ignore[arg-type]

~\anaconda3\lib\site-packages\pandas\_libs\lib.pyx in pandas._libs.lib.map_infer()

~\AppData\Local\Temp\ipykernel_6124\1532437874.py in score_condition(score)
      1 def score_condition(score):
----> 2     return score > 5
      3 Likes = 'BloodPressure'
      4 filtered_df = df[df['Glucose'].apply(score_condition)]
      5 print(filtered_df)

TypeError: '>' not supported between instances of 'NoneType' and 'int'
```

In [43]:

```
1 modify_score = lambda x: x * 2
2 LeagueIndex = 'BloodPressure'
3 df['Glucose'] = df['Age'].apply(modify_score)
4 print(df)
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	\
0	6	100	72	35	0	33.6	
1	1	62	66	29	0	26.6	
2	8	64	64	0	0	23.3	
3	1	42	66	23	94	28.1	
4	0	66	40	35	168	43.1	
..	...	...	...	...	...	...	...
763	10	126	76	48	180	32.9	
764	2	54	70	27	0	36.8	
765	5	60	72	23	112	26.2	
766	1	94	60	0	0	30.1	
767	1	46	70	31	0	30.4	

	DiabetesPedigreeFunction	Age	Outcome
0	0.627	50	1
1	0.351	31	0
2	0.672	32	1
3	0.167	21	0
4	2.288	33	1
..	...	...	...
763	0.171	63	0
764	0.340	27	0
765	0.245	30	0
766	0.349	47	1
767	0.315	23	0

[768 rows x 9 columns]

```
In [45]: 1 sorted_df = df.sort_values(by='Age', ascending=True)
2 print(sorted_df)
3
```

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	\
255	1	42	64	35	0	33.6	
60	2	42	0	0	0	0.0	
102	0	42	96	0	0	22.5	
182	1	42	74	20	23	27.7	
623	0	42	70	27	115	43.5	
..	...	...	...	...	...	...	...
123	5	138	80	0	0	26.8	
684	5	138	82	0	0	0.0	
666	4	140	82	18	0	32.5	
453	2	144	0	0	0	19.6	
459	9	162	74	33	60	25.9	

	DiabetesPedigreeFunction	Age	Outcome
255	0.543	21	1
60	0.304	21	0
102	0.262	21	0
182	0.299	21	0
623	0.347	21	0
..	...	...	...
123	0.186	69	0
684	0.640	69	0
666	0.235	70	1
453	0.832	72	0
459	0.460	81	0

[768 rows x 9 columns]

In [46]:

```
1 grouped_data = df.groupby('Age')
2 result = grouped_data['Age'].mean()
3 print(result)
```

```
Age
21    21.0
22    22.0
23    23.0
24    24.0
25    25.0
26    26.0
27    27.0
28    28.0
29    29.0
30    30.0
31    31.0
32    32.0
33    33.0
34    34.0
35    35.0
36    36.0
37    37.0
38    38.0
39    39.0
40    40.0
41    41.0
42    42.0
43    43.0
44    44.0
45    45.0
46    46.0
47    47.0
48    48.0
49    49.0
50    50.0
51    51.0
52    52.0
53    53.0
54    54.0
55    55.0
56    56.0
57    57.0
58    58.0
59    59.0
60    60.0
61    61.0
62    62.0
63    63.0
64    64.0
65    65.0
66    66.0
67    67.0
68    68.0
69    69.0
70    70.0
72    72.0
81    81.0
Name: Age, dtype: float64
```

In [47]: 1 df.head(2)

Out[47]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction
0	6	100	72	35	0	33.6	0.
1	1	62	66	29	0	26.6	0.

## VISUALIZATION

### Bar Chat

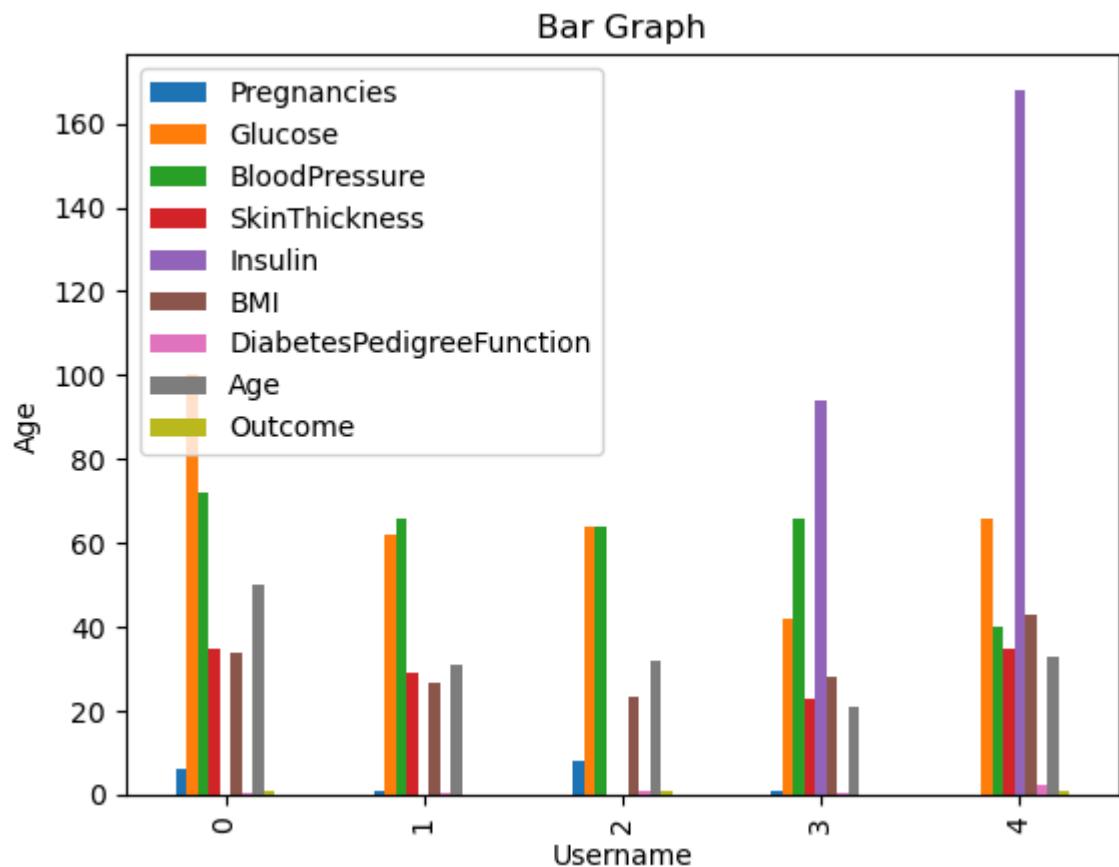
In [51]:

```

1 import matplotlib.pyplot as plt
2 bar_data = df.head(5)
3 plt.figure(figsize=(8, 6))
4 bar_data.plot(kind='bar')
5 plt.xlabel('Username')
6 plt.ylabel('Age')
7 plt.title('Bar Graph')
8 plt.show()

```

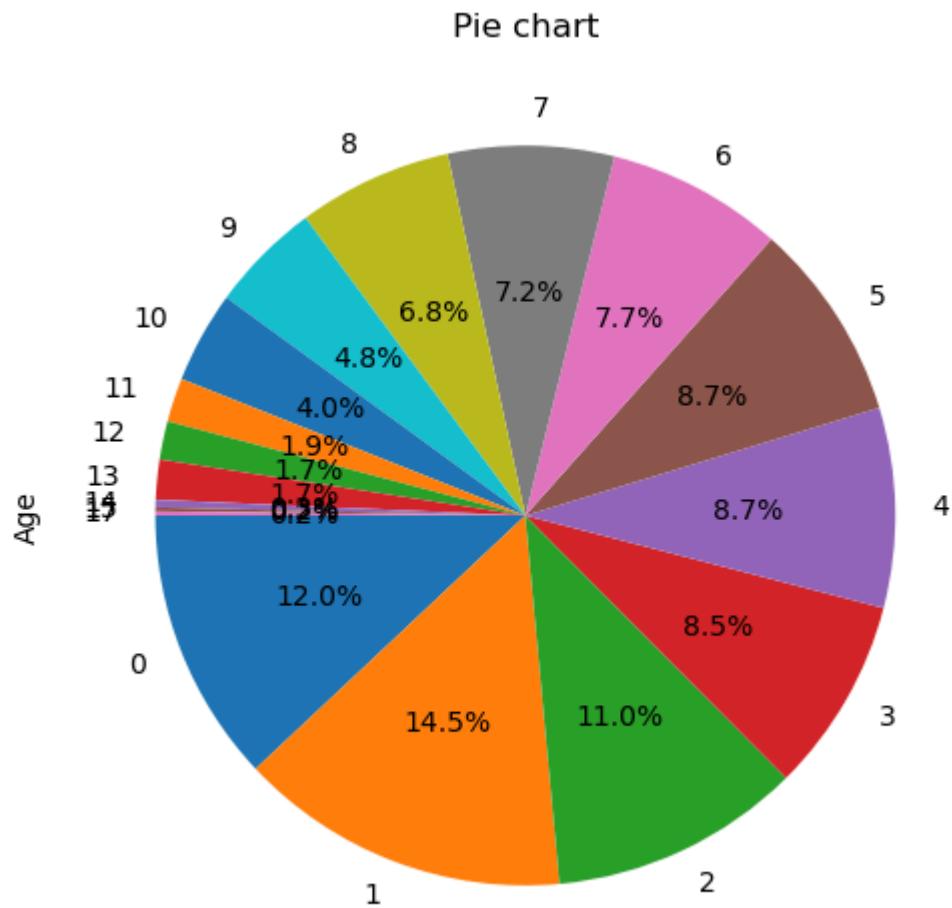
<Figure size 800x600 with 0 Axes>



### Pie Chart

In [53]:

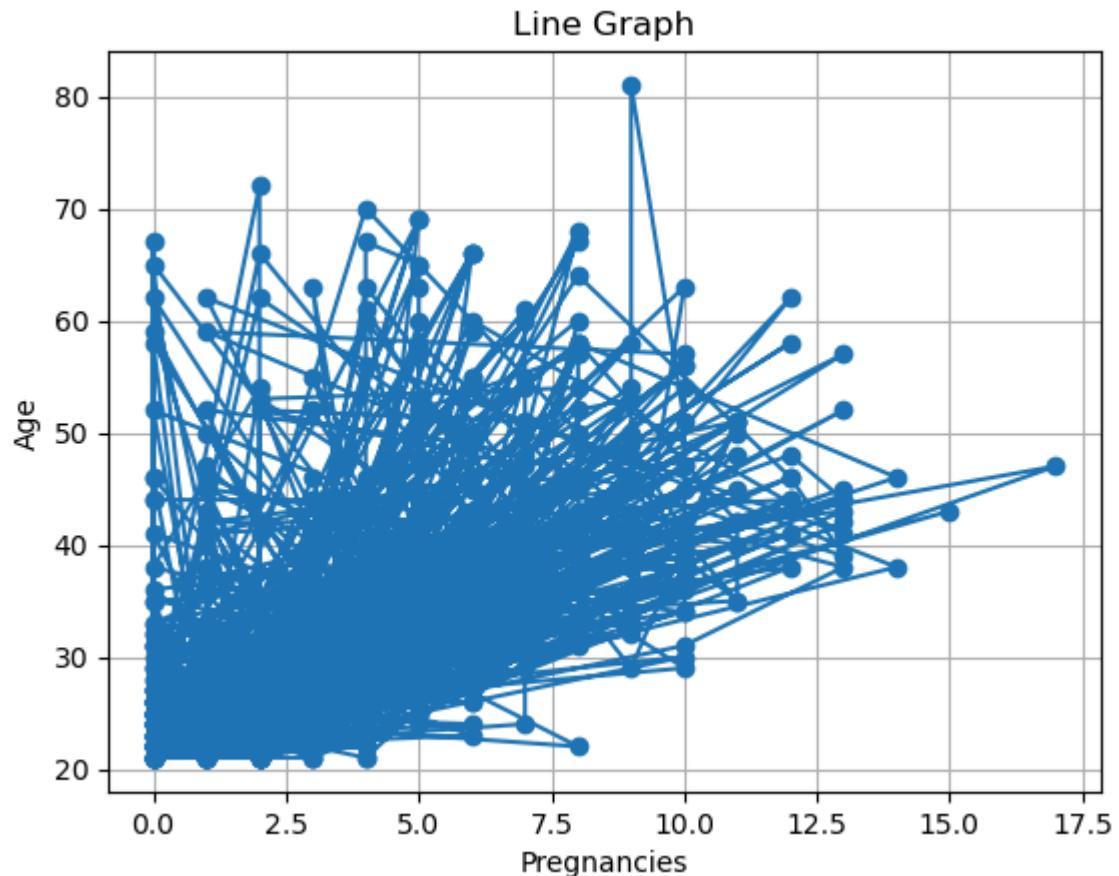
```
1 pie_data = df.groupby('Pregnancies')['Age'].sum()
2 plt.figure(figsize=(8, 6))
3 pie_data.plot(kind='pie', autopct='%1.1f%%', startangle=180)
4 plt.title('Pie chart')
5 plt.show()
6
```



## LINE GRAPH

In [56]:

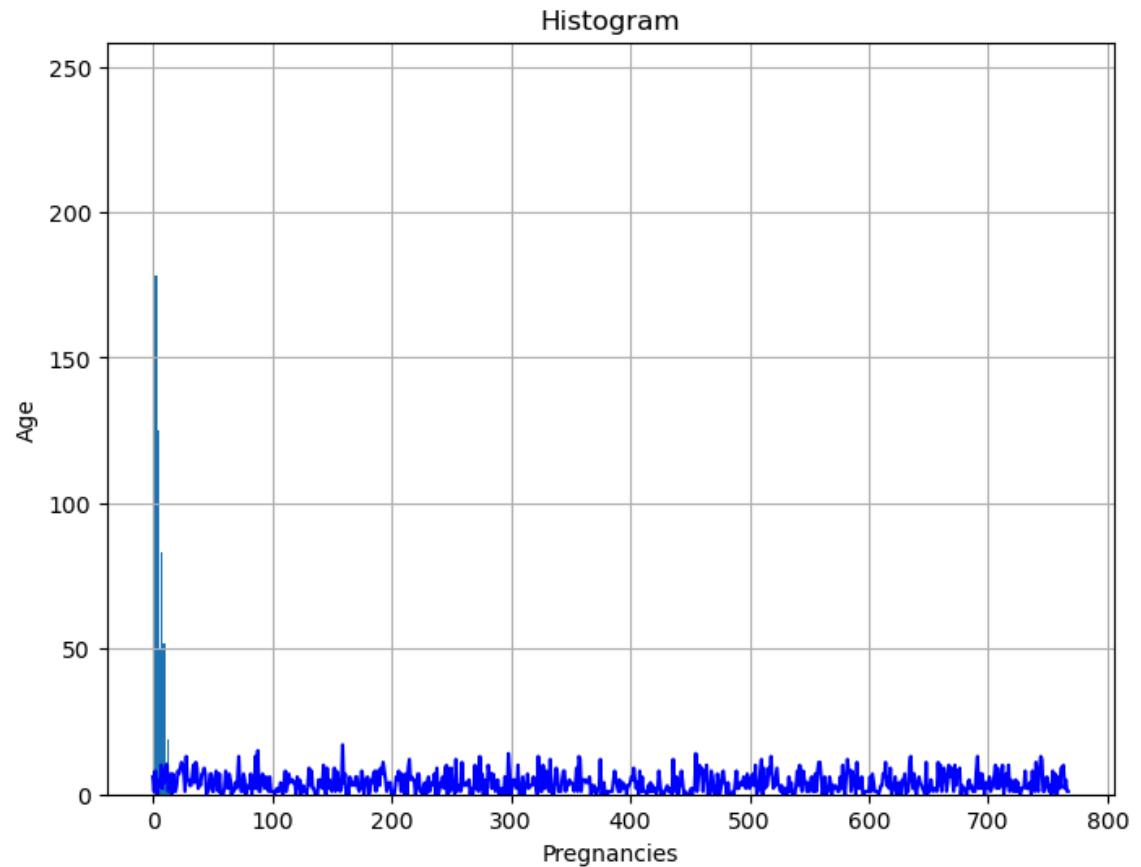
```
1 line_chart = df.head(2)
2 x_line = df['Pregnancies']
3 y_line = df['Age']
4 plt.plot(x_line, y_line, marker='o', linestyle='--')
5 plt.xlabel('Pregnancies')
6 plt.ylabel('Age')
7 plt.title('Line Graph')
8 plt.grid(True)
9 plt.show()
```



## HISTOGRAM

In [58]:

```
1 data_hist = df['Pregnancies']
2 plt.figure(figsize=(8, 6))
3 plt.hist(data_hist, bins=10)
4 plt.xlabel('Pregnancies')
5 plt.ylabel('Age')
6 plt.plot(df.index, df['Pregnancies'], label='Age', color='b')
7 # Display the plot
8 plt.grid(True) # Optional: Add grid line
9 plt.title('Histogram')
10 plt.show()
```

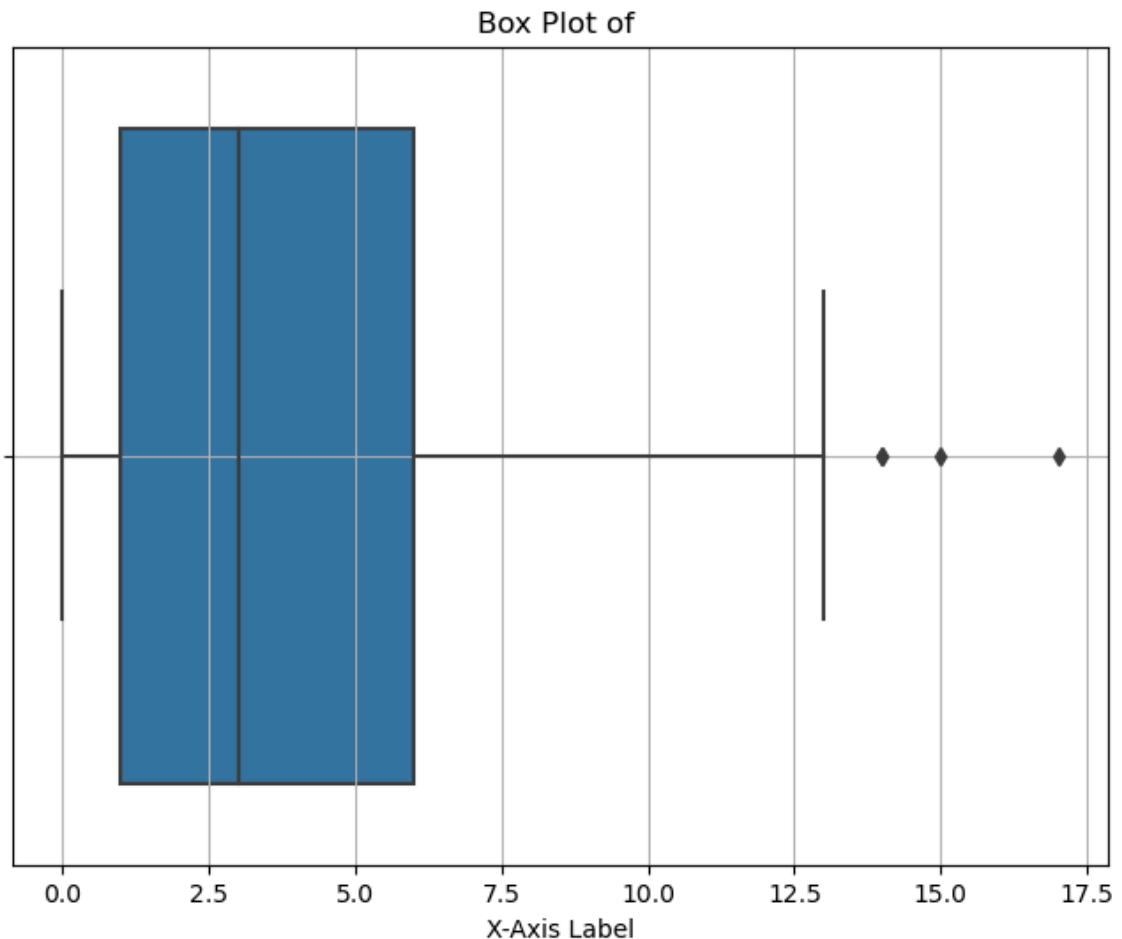


## BOX PLOT

In [59]:

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3 # Create a box plot
4 plt.figure(figsize=(8, 6))
5 sns.boxplot(x='Pregnancies', data=df)
6 plt.xlabel('X-Axis Label')
7 plt.title('Box Plot of ')
8 plt.grid(True)
9 plt.show()
```

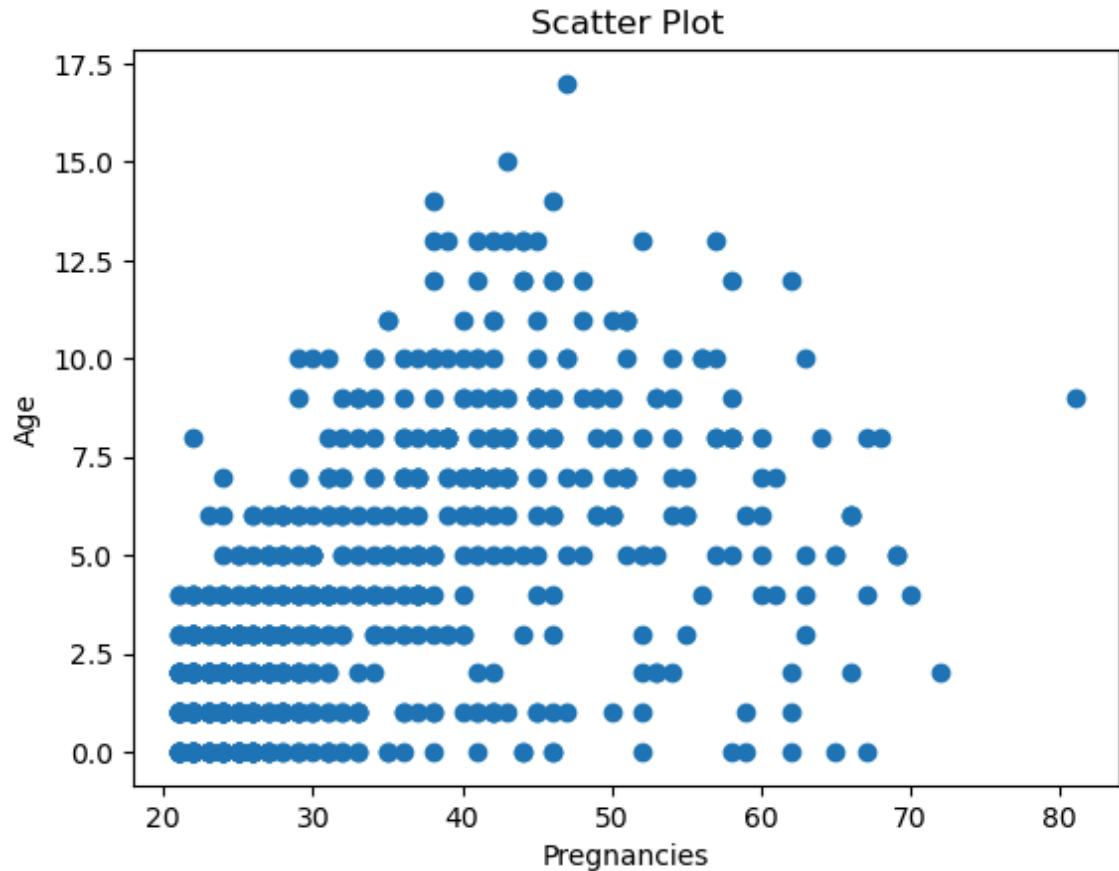
C:\Users\Anusha V\anaconda3\lib\site-packages\scipy\\_init\_.py:155: UserWarning: A NumPy version >=1.18.5 and <1.25.0 is required for this version of SciPy (detected version 1.26.2)  
warnings.warn(f"A NumPy version >={np\_minversion} and <{np\_maxversion}"



## SCATTER PLOT

In [62]:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 x_column = 'Pregnancies'
4 y_column = 'Age'
5 # Create a scatter plot
6 plt.scatter(df['Age'], df['Pregnancies'])
7 # Add labels and title
8 plt.xlabel('Pregnancies')
9 plt.ylabel('Age')
10 plt.title('Scatter Plot')
11 # Show the plot
12 plt.show()
13
```



## BUBBLE CHART

In [67]:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 # Define the columns for x, y, and bubble size
5 x_column = 'Pregnancies' # Replace with the actual column name
6 y_column = 'Age' # Replace with the actual column name
7 size_column = 'Size_Column' # Replace with the actual column name
8 # Create a bubble chart
9 plt.figure(figsize=(10, 6))
10 sns.scatterplot(x=x_column, y=y_column, data=df, sizes=(20, 200), alpha=0.5)
11 # Customize the plot properties (title, labels, etc.)
12
13
14
```

```
File "C:\Users\Anusha V\AppData\Local\Temp\ipykernel_6124\1800711684.py", line 13
```

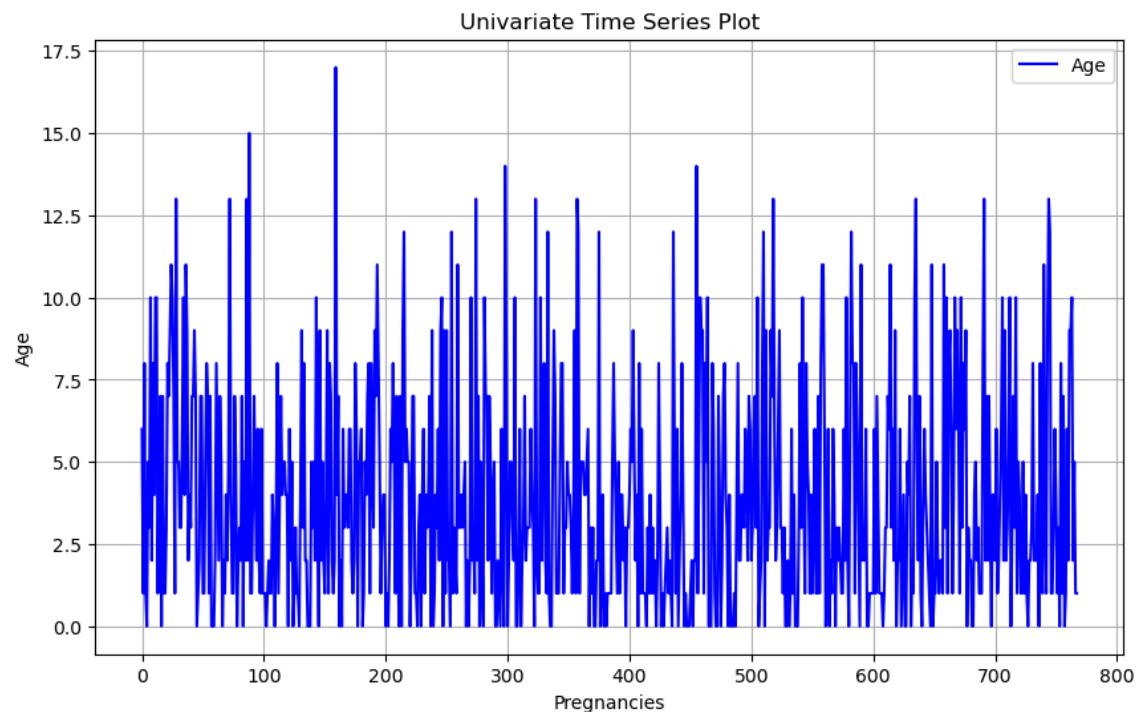
```
^
```

```
SyntaxError: positional argument follows keyword argument
```

## Univariate Time Series Plot

In [68]:

```
1 plt.figure(figsize=(10, 6)) # Optional: Set the figure size
2 # Plot the time series data
3 plt.plot(df.index, df['Pregnancies'], label='Age', color='b')
4 # Optional: Add Labels and a title
5 plt.xlabel('Pregnancies')
6 plt.ylabel('Age')
7 plt.title('Univariate Time Series Plot')
8 # Optional: Customize other plot properties
9 # Display the plot
10 plt.grid(True) # Optional: Add grid Lines
11 plt.legend() # Optional: Show the legend if you have multiple series
12 plt.show()
13 1
14
```



Out[68]: 1

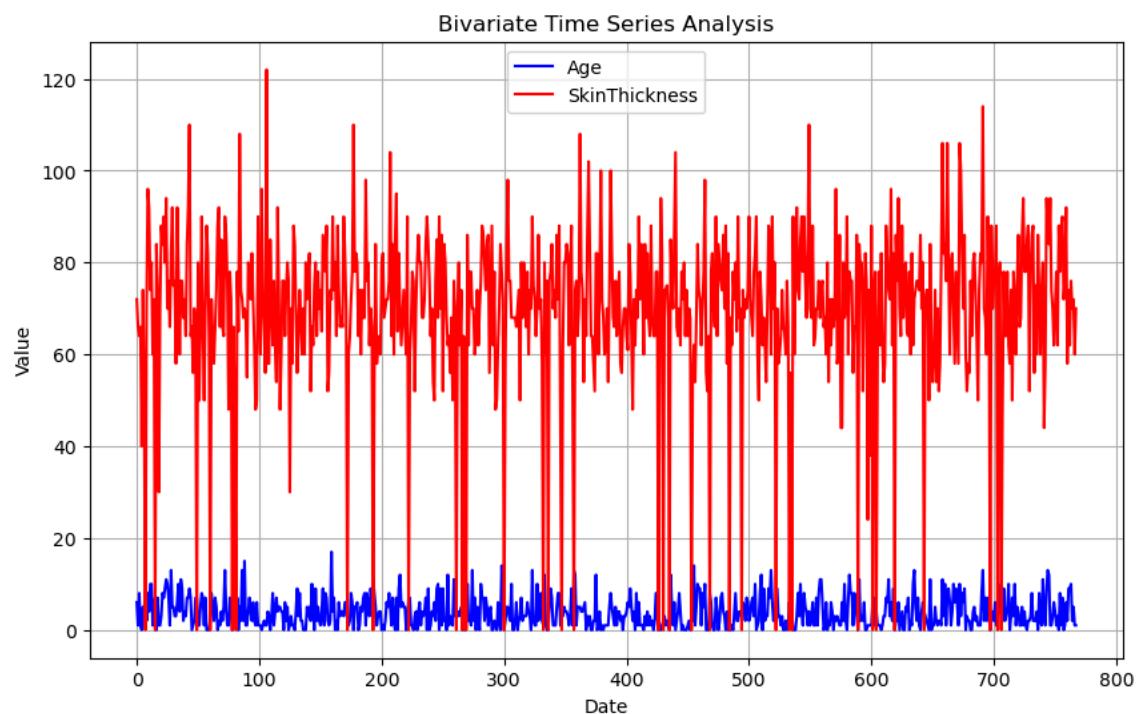
## Bivariate Time Series Analysis

In [69]:

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 # Plot the time series data
4 plt.figure(figsize=(10, 6))
5 plt.plot(df.index, df['Pregnancies'], label='Age', color='blue')
6 plt.plot(df.index, df['BloodPressure'], label='SkinThickness', color='red')
7 # Customize the plot
8 plt.xlabel('Date')
9 plt.ylabel('Value')
10 plt.title('Bivariate Time Series Analysis')
11 plt.legend()
12 plt.grid(True)
13 # Show the plot
14 plt.show()

```



In [71]:

```
1 df.head(1)
```

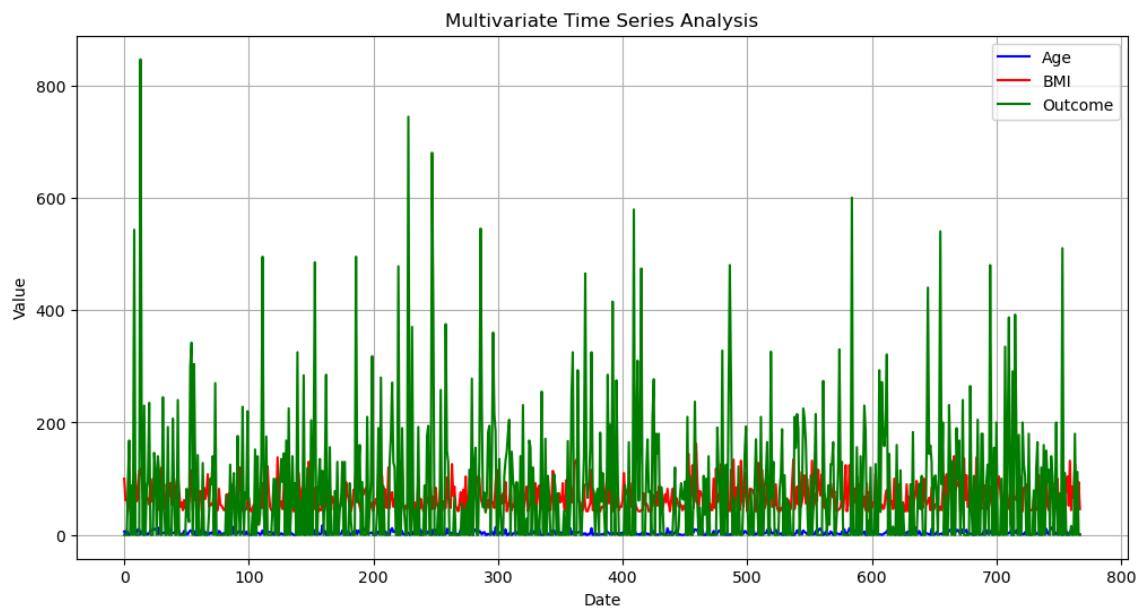
Out[71]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction
0	6	100	72	35	0	33.6	0.1478

## Multivariate Time Series Analysis

In [73]:

```
1 import matplotlib.pyplot as plt
2 # Plot the multivariate time series data
3 plt.figure(figsize=(12, 6))
4 plt.plot(df.index, df['Pregnancies'], label='Age', color='blue')
5 plt.plot(df.index, df['Glucose'], label='BMI', color='red')
6 plt.plot(df.index, df['Insulin'], label='Outcome', color='green')
7 # Customize the plot
8 plt.xlabel('Date')
9 plt.ylabel('Value')
10 plt.title('Multivariate Time Series Analysis')
11 plt.legend()
12 plt.grid(True)
13 # Show the plot
14 plt.show()
```



In [1]:

```
1 import pandas as pd
2 import numpy as np
3
4 # Create a Series from a list
5 my_list = [1, 2, 3, 4, 5]
6 series_from_list = pd.Series(my_list)
7 print(series_from_list)
8
9 # Create a Series from a Numpy array
10 my_array = np.array([10, 20, 30, 40, 50])
11 series_from_array = pd.Series(my_array)
12 print(series_from_array)
13
```

```
0    1
1    2
2    3
3    4
4    5
dtype: int64
0    10
1    20
2    30
3    40
4    50
dtype: int32
```

In [2]:

```
1 # Create a Series from a dictionary
2 my_dict = {'a': 1, 'b': 2, 'c': 3, 'd': 4, 'e': 5}
3 series_from_dict = pd.Series(my_dict)
4 print(series_from_dict)
5
6
```

```
a    1
b    2
c    3
d    4
e    5
dtype: int64
```

In [3]:

```
1 # You can also provide custom indices
2 custom_index_series = pd.Series([10, 20, 30], index=['a', 'b', 'c'])
3 print(custom_index_series)
4
```

```
a    10
b    20
c    30
dtype: int64
```

In [5]:

```
1 import pandas as pd
2 import numpy as np
3
4 # Create a DataFrame from a dictionary of lists
5 data_dict = {
6     'Name': ['Alice', 'Bob', 'Charlie'],
7     'Age': [25, 30, 35],
8     'City': ['New York', 'San Francisco', 'Los Angeles']
9 }
10
11 df_from_dict = pd.DataFrame(data_dict)
12 print(df_from_dict)
13
14
15
```

	Name	Age	City
0	Alice	25	New York
1	Bob	30	San Francisco
2	Charlie	35	Los Angeles

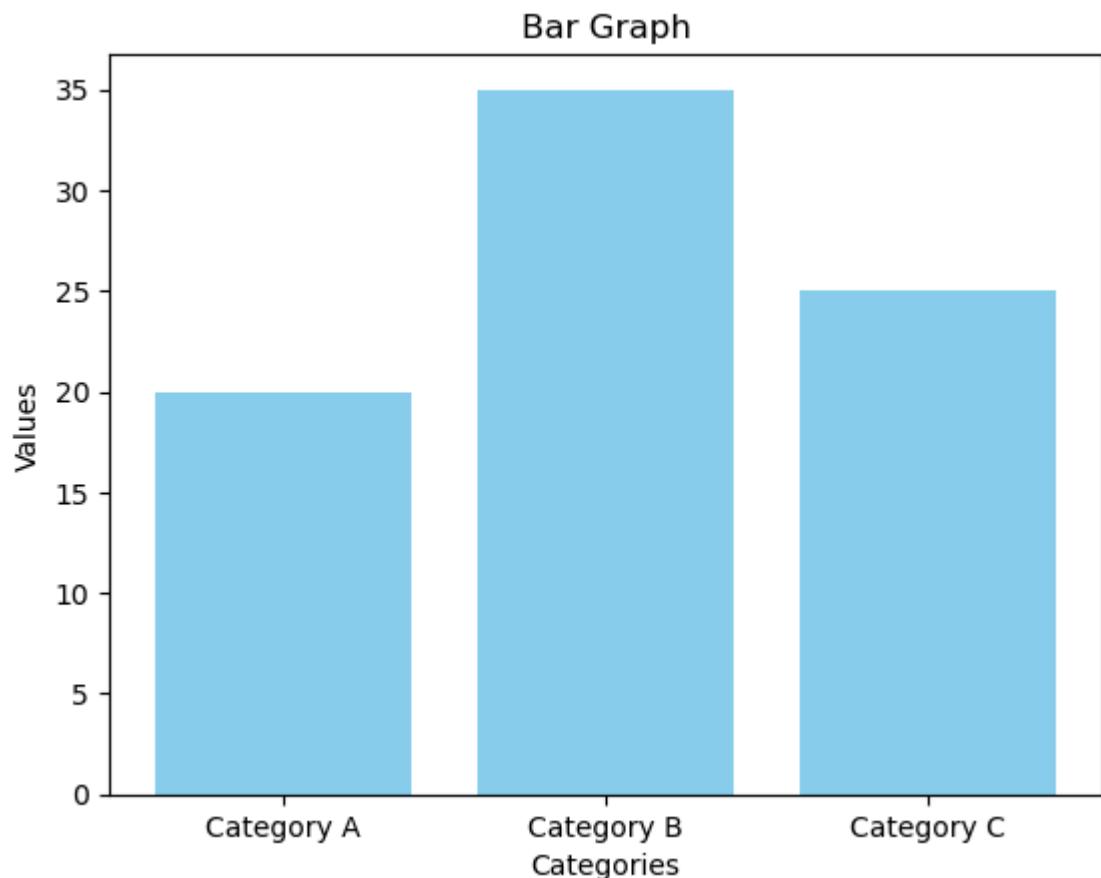
In [6]:

```
1 # Create a DataFrame from a dictionary of NumPy arrays
2 data_dict_np = {
3     'A': np.array([1, 2, 3]),
4     'B': np.array([4, 5, 6]),
5     'C': np.array([7, 8, 9])
6 }
7
8 df_from_dict_np = pd.DataFrame(data_dict_np)
9 print(df_from_dict_np)
```

	A	B	C
0	1	4	7
1	2	5	8
2	3	6	9

In [7]:

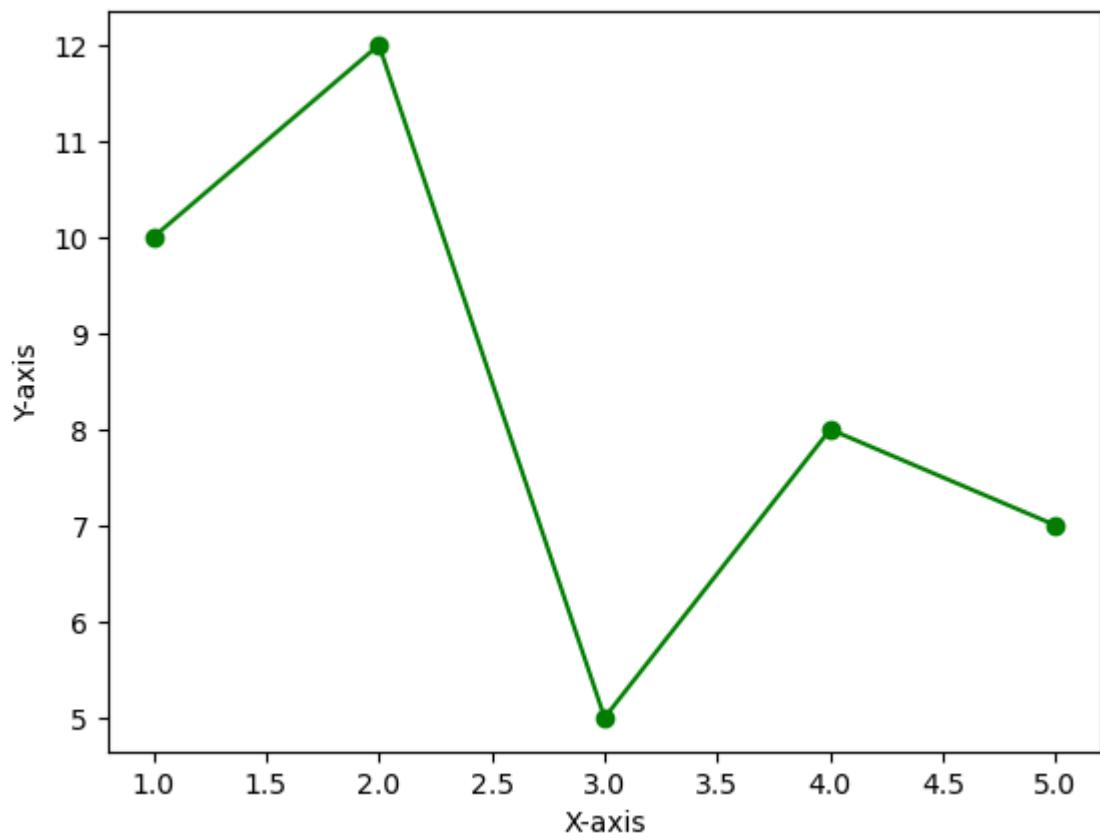
```
1 import matplotlib.pyplot as plt
2
3 categories = ['Category A', 'Category B', 'Category C']
4 values = [20, 35, 25]
5
6 plt.bar(categories, values, color='skyblue')
7 plt.title('Bar Graph')
8 plt.xlabel('Categories')
9 plt.ylabel('Values')
10 plt.show()
11
```



In [8]:

```
1 import matplotlib.pyplot as plt
2
3 x_values = [1, 2, 3, 4, 5]
4 y_values = [10, 12, 5, 8, 7]
5
6 plt.plot(x_values, y_values, marker='o', linestyle='-', color='green')
7 plt.title('Line Chart')
8 plt.xlabel('X-axis')
9 plt.ylabel('Y-axis')
10 plt.show()
11
```

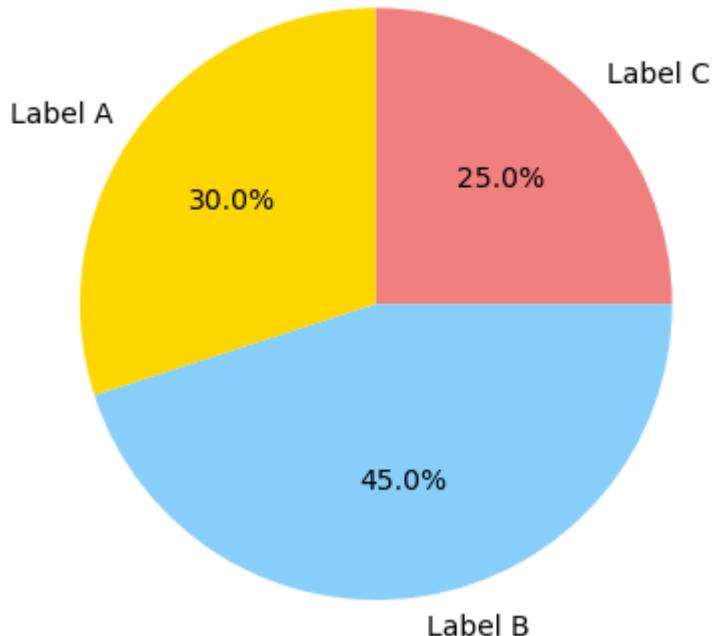
Line Chart



In [9]:

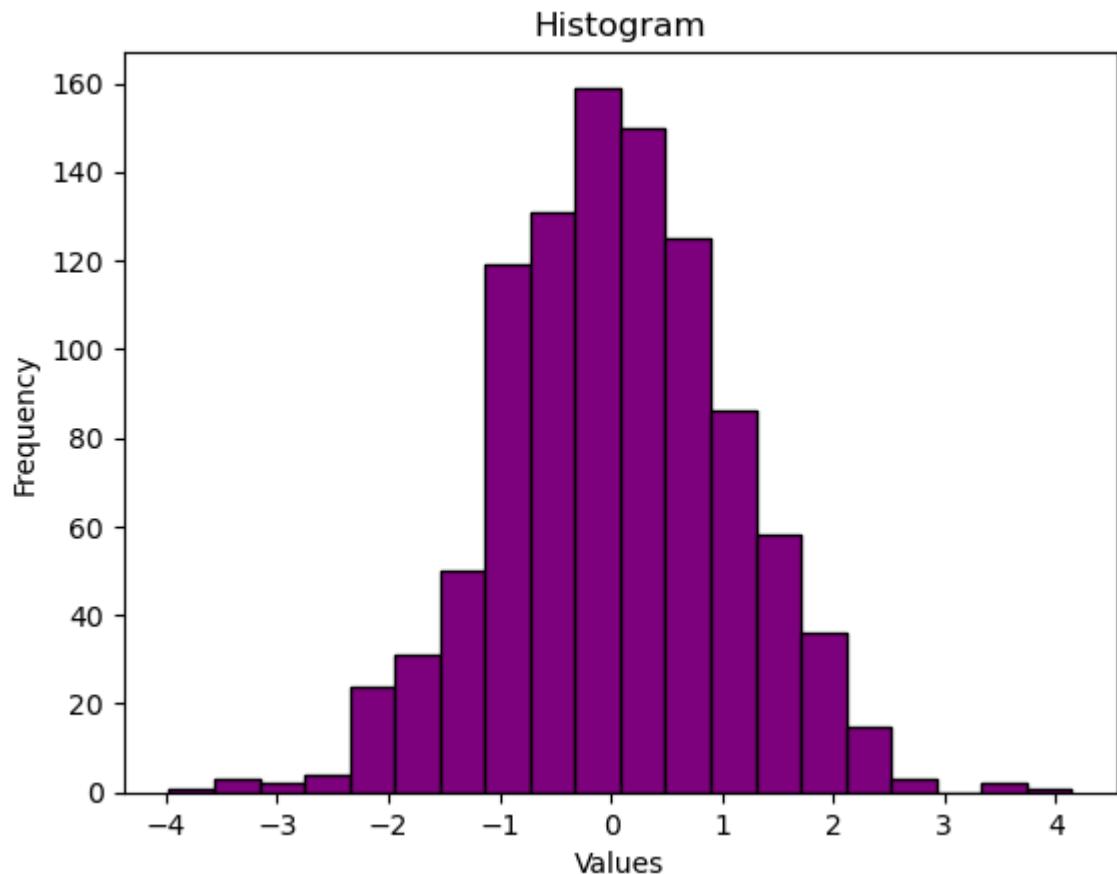
```
1 import matplotlib.pyplot as plt
2
3 labels = ['Label A', 'Label B', 'Label C']
4 sizes = [30, 45, 25]
5
6 plt.pie(sizes, labels=labels, autopct='%1.1f%%', startangle=90, colors=
7 plt.title('Pie Chart')
8 plt.show()
9
```

Pie Chart



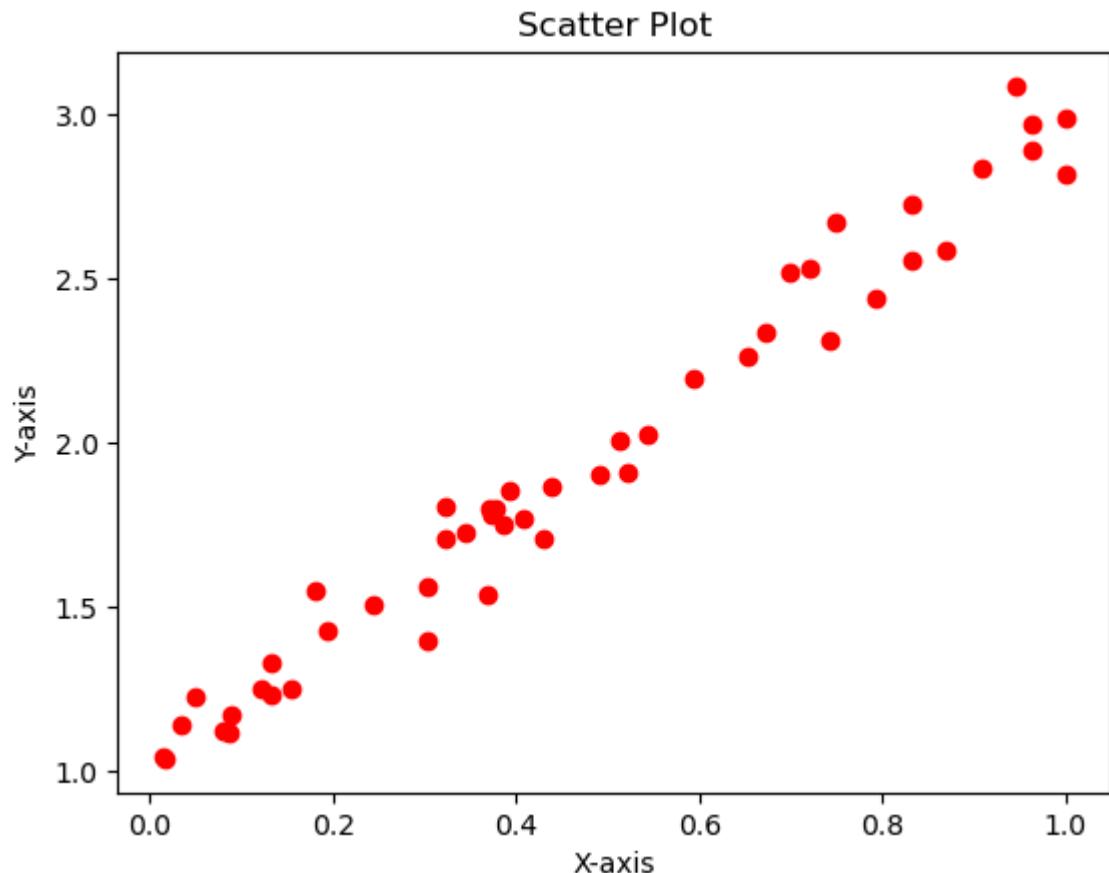
In [10]:

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 data = np.random.randn(1000)
5
6 plt.hist(data, bins=20, color='purple', edgecolor='black')
7 plt.title('Histogram')
8 plt.xlabel('Values')
9 plt.ylabel('Frequency')
10 plt.show()
11
```



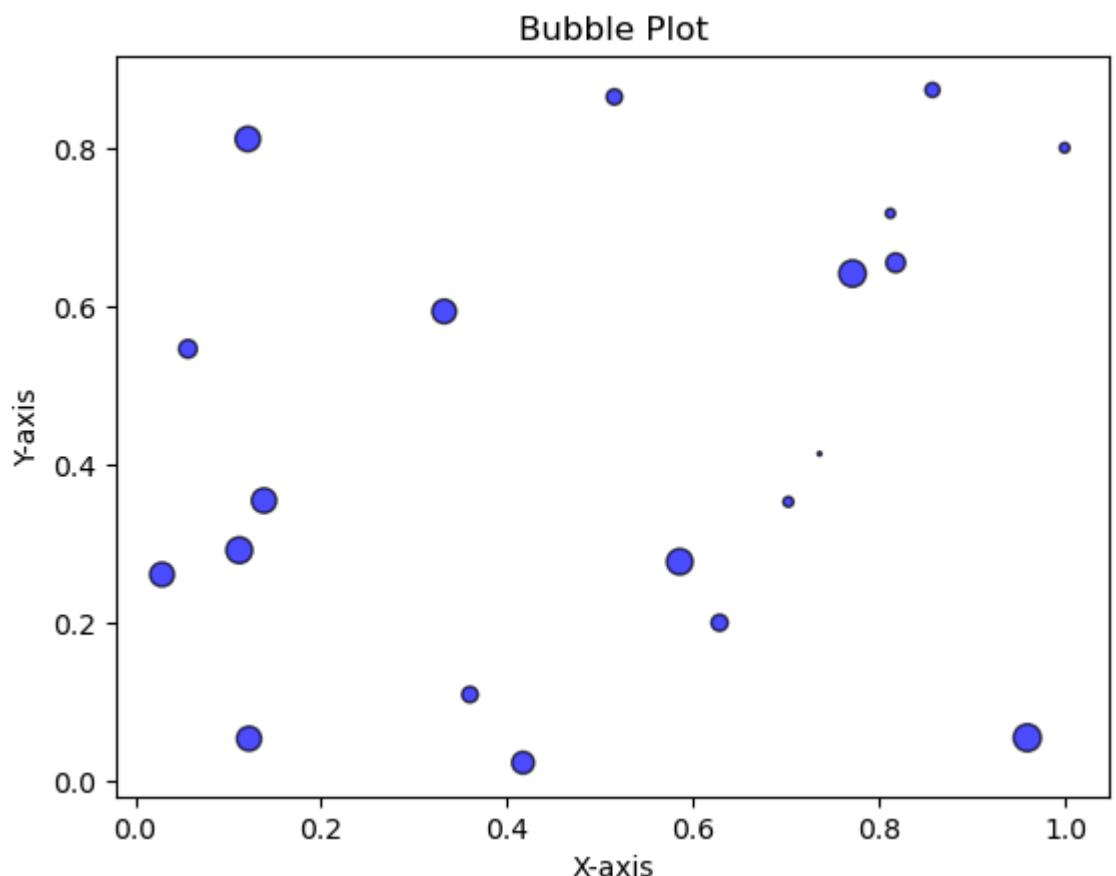
In [11]:

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 x_values = np.random.rand(50)
5 y_values = 2 * x_values + 1 + 0.1 * np.random.randn(50)
6
7 plt.scatter(x_values, y_values, color='red')
8 plt.title('Scatter Plot')
9 plt.xlabel('X-axis')
10 plt.ylabel('Y-axis')
11 plt.show()
12
```



In [12]:

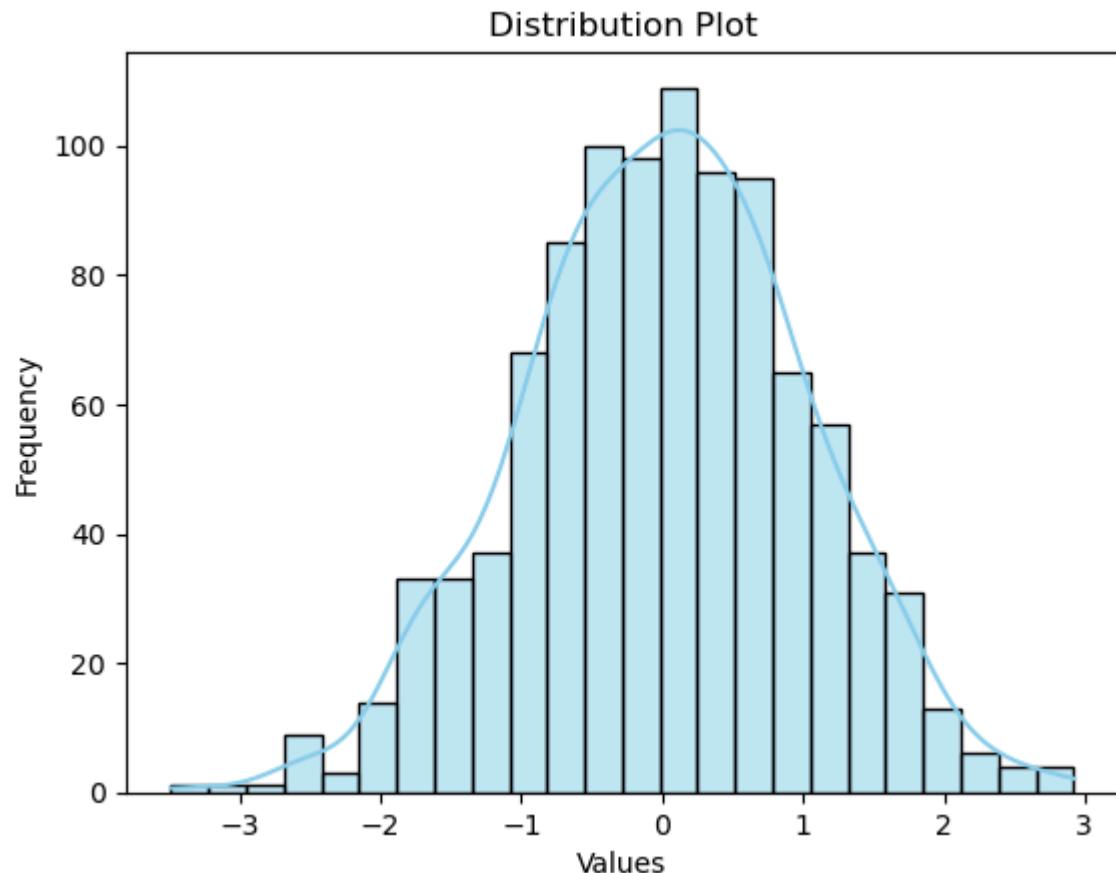
```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 # Generate random data
5 x_values = np.random.rand(20)
6 y_values = np.random.rand(20)
7 sizes = np.random.rand(20) * 100 # Random sizes for bubbles
8
9 # Scatter plot with bubbles
10 plt.scatter(x_values, y_values, s=sizes, alpha=0.7, c='blue', edgecolor='black')
11
12 plt.title('Bubble Plot')
13 plt.xlabel('X-axis')
14 plt.ylabel('Y-axis')
15 plt.show()
16
```



In [13]:

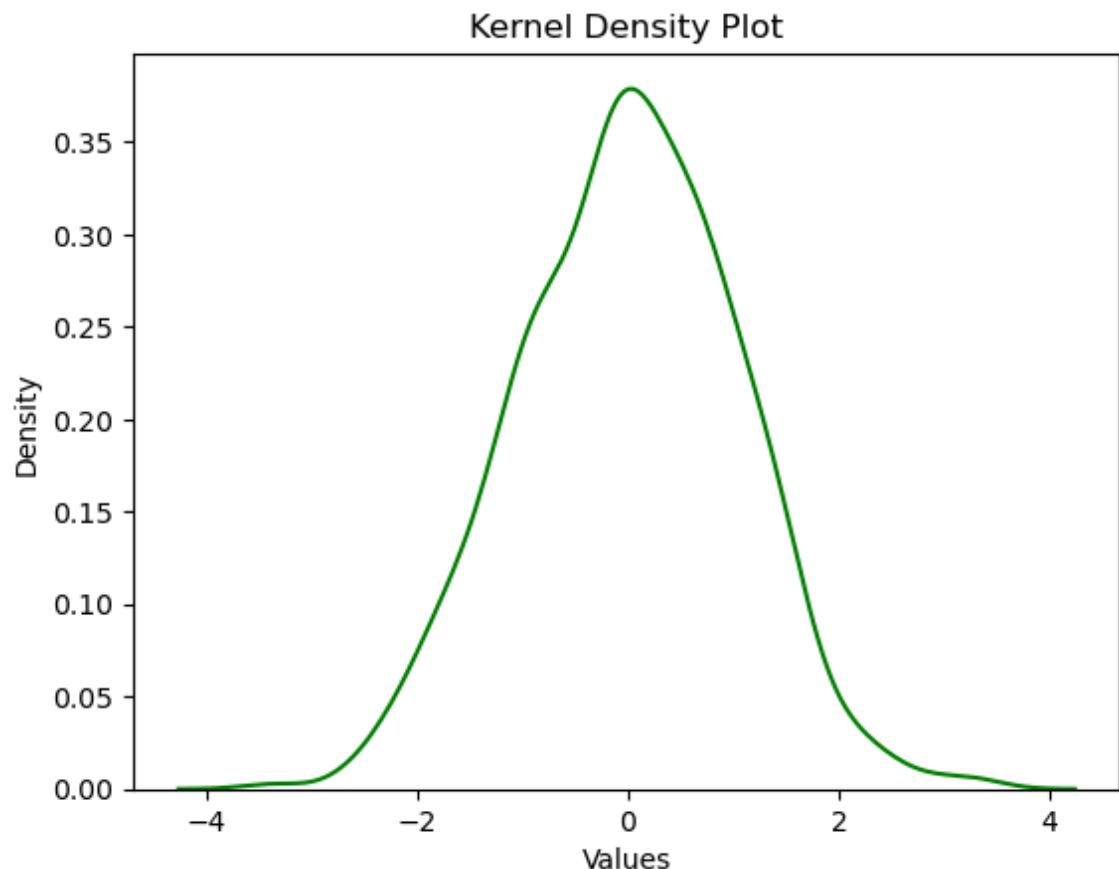
```
1 import seaborn as sns
2 import numpy as np
3
4 data = np.random.randn(1000)
5
6 sns.histplot(data, kde=True, color='skyblue')
7 plt.title('Distribution Plot')
8 plt.xlabel('Values')
9 plt.ylabel('Frequency')
10 plt.show()
11
```

C:\Users\Anusha V\anaconda3\lib\site-packages\scipy\\_init\_.py:155: UserWarning: A NumPy version >=1.18.5 and <1.25.0 is required for this version of SciPy (detected version 1.26.2)  
warnings.warn(f"A NumPy version >={np\_minversion} and <{np\_maxversion}"



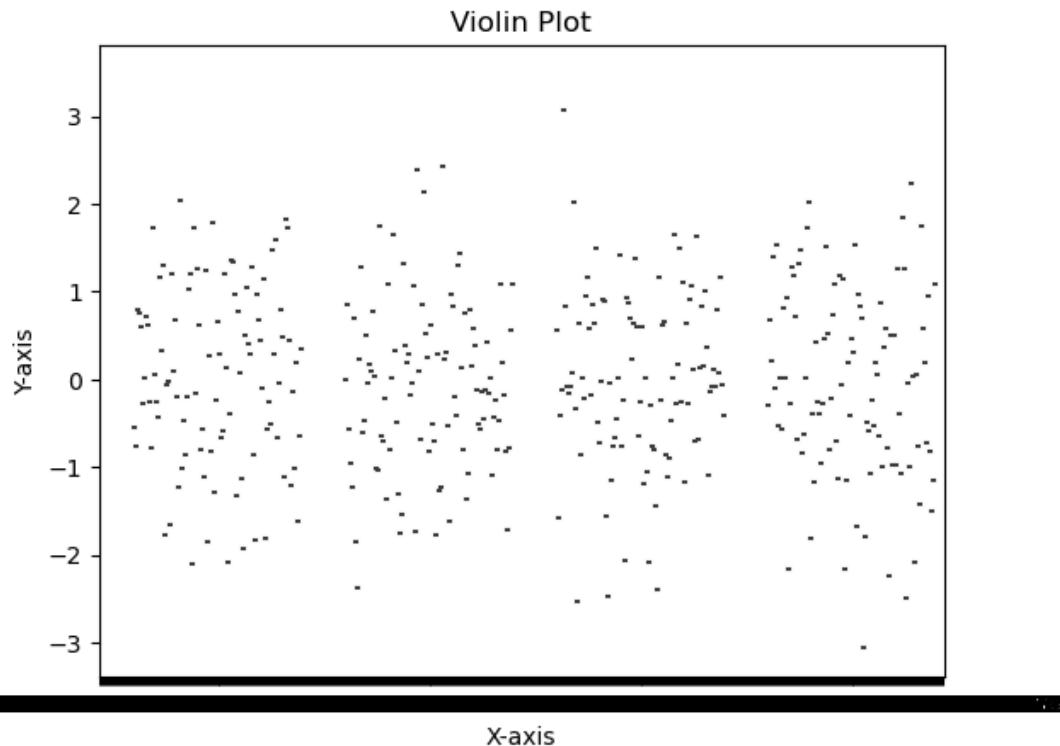
In [14]:

```
1 import seaborn as sns
2 import numpy as np
3
4 data = np.random.randn(1000)
5
6 sns.kdeplot(data, color='green')
7 plt.title('Kernel Density Plot')
8 plt.xlabel('Values')
9 plt.ylabel('Density')
10 plt.show()
11
```



In [15]:

```
1 import seaborn as sns
2 import numpy as np
3
4 data = np.random.randn(500, 2)
5
6 sns.violinplot(x=data[:, 0], y=data[:, 1], palette='muted')
7 plt.title('Violin Plot')
8 plt.xlabel('X-axis')
9 plt.ylabel('Y-axis')
10 plt.show()
11
```

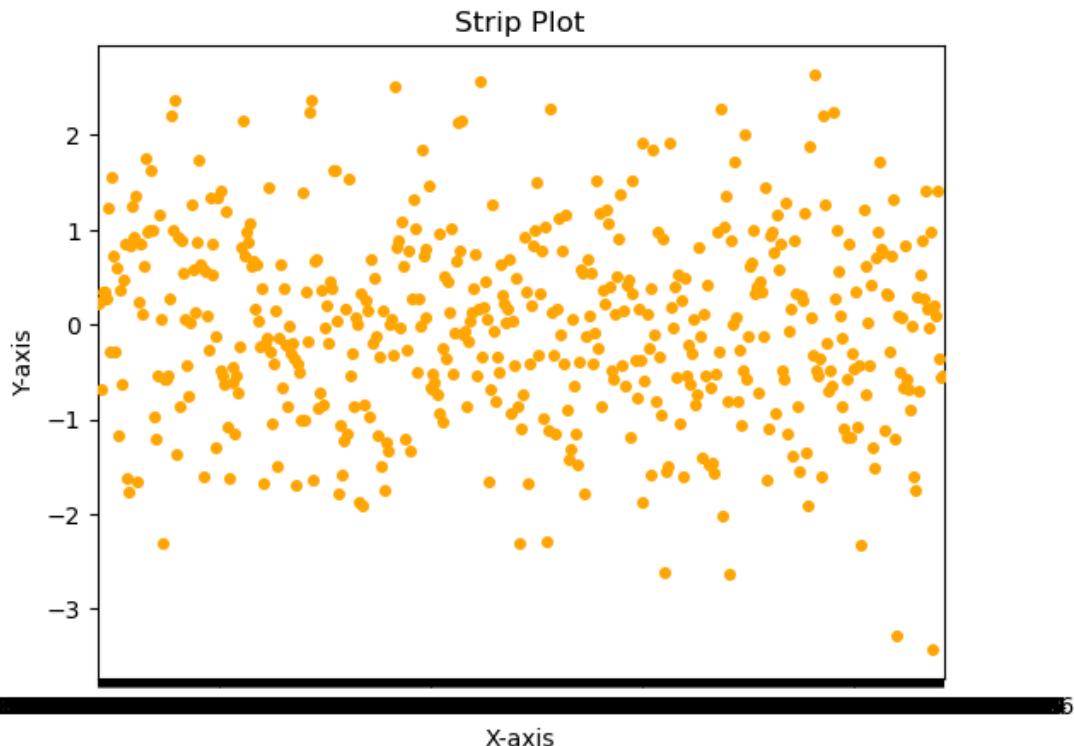


In [16]:

```

1 import seaborn as sns
2 import numpy as np
3
4 data = np.random.randn(500, 2)
5
6 sns.stripplot(x=data[:, 0], y=data[:, 1], jitter=True, color='orange')
7 plt.title('Strip Plot')
8 plt.xlabel('X-axis')
9 plt.ylabel('Y-axis')
10 plt.show()
11

```



## DATA SPLITTING

In [5]:

```

1 from sklearn.model_selection import train_test_split

```

```

C:\Users\Anusha V\anaconda3\lib\site-packages\scipy\__init__.py:155: UserWarning: A NumPy version >=1.18.5 and <1.25.0 is required for this version
of SciPy (detected version 1.26.2
    warnings.warn(f"A NumPy version >={np_minversion} and <{np_maxversion}"

```

In [15]:

```

1 import numpy as np
2 X = np.random.rand(100, 5)
3 y = np.random.randint(0, 2, size=100)

```

In [11]:

```

1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2
2 X_train, X_val, y_train, y_val = train_test_split(X_train, y_train, tes
3

```

```
In [12]: 1 print(f"Training set size: {len(X_train)}")  
2 print(f"Validation set size: {len(X_val)}")  
3 print(f"Test set size: {len(X_test)}")
```

Training set size: 491  
 Validation set size: 123  
 Test set size: 154

```
In [13]: 1 import pandas as pd  
2 df = pd.read_csv(r"C:\Users\Anusha V\Desktop\diabetes.csv")  
3 df.head(5)
```

Out[13]:

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction
0	6	148	72	35	0	33.6	0.
1	1	85	66	29	0	26.6	0.
2	8	183	64	0	0	23.3	0.
3	1	89	66	23	94	28.1	0.
4	0	137	40	35	168	43.1	2.



```
In [18]: 1 X = df.drop('Outcome', axis=1) # Features  
2 y = df['Outcome'] # Labels
```

```
In [19]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)  
2 X_train, X_val, y_train, y_val = train_test_split(X_train, y_train, test_size=0.25)  
3  
4 print(f"Training set size: {len(X_train)}")  
5 print(f"Validation set size: {len(X_val)}")  
6 print(f"Test set size: {len(X_test)}")
```

Training set size: 491  
 Validation set size: 123  
 Test set size: 154

## OUTLIERS

```
In [21]: 1 import pandas as pd  
2 from scipy import stats  
3 import matplotlib.pyplot as plt
```

In [20]:

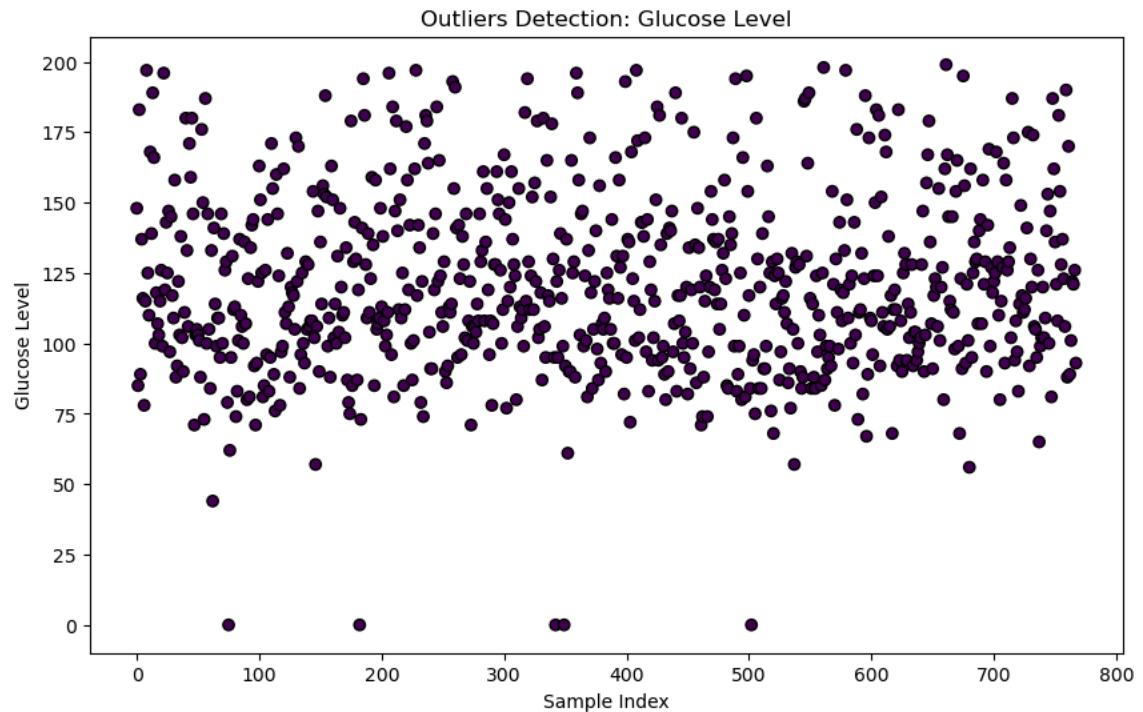
```
1 z_scores = np.abs(stats.zscore(df))
2
3 # Define a threshold for identifying outliers (commonly 3)
4 threshold = 3
5
6 # Find and print the indices of outliers
7 outlier_indices = np.where(z_scores > threshold)
8 outliers = set(zip(outlier_indices[0], outlier_indices[1]))
9
10 print("Indices of outliers:")
11 print(outliers)
```

Indices of outliers:

```
{(459, 7), (535, 2), (145, 5), (484, 2), (673, 5), (584, 4), (81, 5), (39
5, 6), (298, 0), (706, 5), (336, 2), (697, 2), (75, 1), (589, 2), (159,
0), (695, 4), (579, 3), (604, 2), (182, 1), (619, 2), (8, 4), (60, 5), (66
6, 7), (300, 2), (49, 2), (153, 4), (455, 0), (9, 5), (222, 2), (266, 2),
(332, 2), (15, 2), (81, 2), (453, 2), (370, 4), (753, 4), (706, 2), (468,
2), (621, 6), (370, 6), (371, 5), (426, 5), (494, 2), (60, 2), (655, 4),
(415, 4), (430, 2), (269, 2), (7, 2), (522, 5), (342, 1), (111, 4), (4,
6), (502, 1), (645, 4), (45, 6), (593, 6), (684, 5), (347, 2), (435, 2),
(286, 4), (13, 4), (445, 6), (426, 2), (703, 2), (409, 4), (684, 7), (349,
1), (88, 0), (228, 4), (330, 6), (522, 2), (533, 2), (643, 2), (601, 2),
(228, 6), (193, 2), (247, 4), (371, 6), (58, 6), (261, 2), (220, 4), (453,
7), (123, 7), (486, 4), (494, 5), (78, 2), (177, 5), (186, 4), (357, 2),
(49, 5), (445, 5), (172, 2)}
```

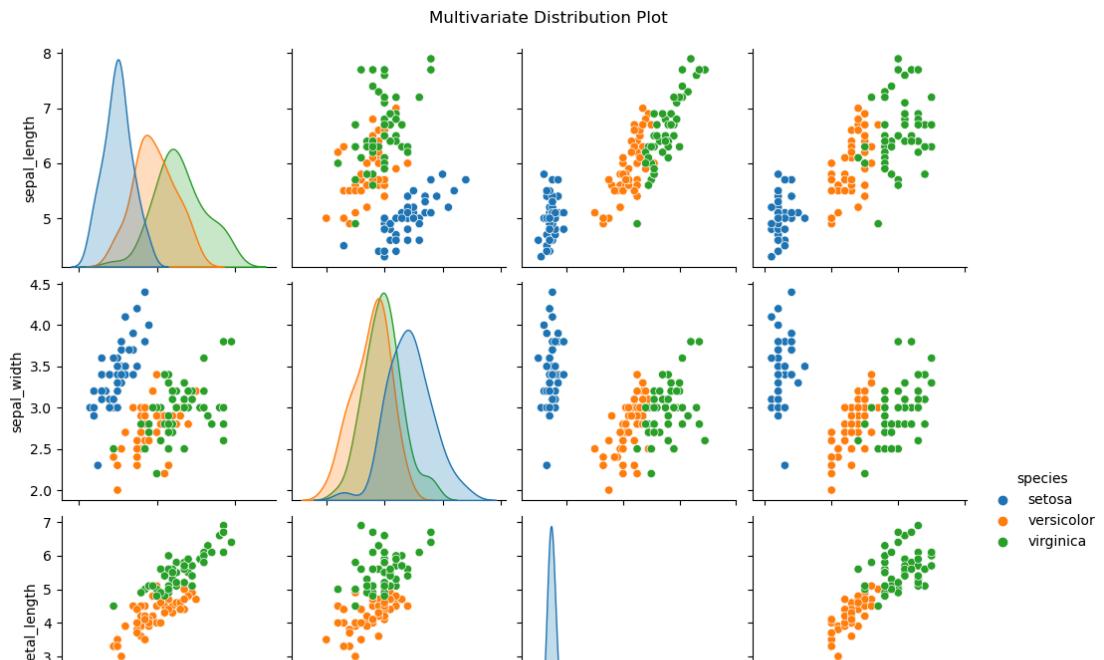
In [22]:

```
1 df = df.dropna()
2
3 # Function to detect outliers using Z-score
4 def detect_outliers_zscore(data_frame):
5     z_scores = np.abs(stats.zscore(data_frame))
6     outliers = (z_scores > 3).all(axis=1)
7     return outliers
8
9 # Identify outliers
10 outliers = detect_outliers_zscore(df[['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness', 'Insulin', 'BMI', 'DiabetesPedigreeFunction', 'Age']])
11
12 # Plot a graph of the outliers
13 plt.figure(figsize=(10, 6))
14 plt.scatter(df.index, df['Glucose'], c=outliers, cmap='viridis', marker='o')
15 plt.title('Outliers Detection: Glucose Level')
16 plt.xlabel('Sample Index')
17 plt.ylabel('Glucose Level')
18 plt.show()
```



In [29]:

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # Sample dataset
5 df = sns.load_dataset('iris') # You can replace this with your own data
6
7 # Multivariate Distribution Plot
8 sns.pairplot(df, hue='species', diag_kind='kde')
9 plt.suptitle('Multivariate Distribution Plot', y=1.02)
10 plt.show()
11
12 # Multivariate Comparison Plot
13 sns.violinplot(x='species', y='sepal_length', data=df)
14 plt.title('Multivariate Comparison Plot')
15 plt.show()
16
17 # Multivariate Relationship Plot
18 sns.scatterplot(x='sepal_length', y='petal_length', hue='species', data=df)
19 plt.title('Multivariate Relationship Plot')
20 plt.show()
21
22 # Multivariate Composition Plot
23 sns.jointplot(x='sepal_length', y='petal_length', kind='hex', data=df)
24 plt.suptitle('Multivariate Composition Plot', y=1.02)
25 plt.show()
26
27
```



In [30]:

```
1 # Initialize the model
2 model = LogisticRegression()
3
4 # Train the model
5 model.fit(X_train, y_train)
6
7 # Make predictions on the test set
8 y_pred = model.predict(X_test)
9
10 # Evaluate the model
11 accuracy = accuracy_score(y_test, y_pred)
12 precision = precision_score(y_test, y_pred)
13 recall = recall_score(y_test, y_pred)
14 f1 = f1_score(y_test, y_pred)
15 conf_matrix = confusion_matrix(y_test, y_pred)
16
17 # Print the evaluation metrics
18 print(f'Accuracy: {accuracy:.2f}')
19 print(f'Precision: {precision:.2f}')
20 print(f'Recall: {recall:.2f}')
21 print(f'F1 Score: {f1:.2f}')
22 print('Confusion Matrix:')
23 print(conf_matrix)
24
25
```

```
Accuracy: 0.75
Precision: 0.64
Recall: 0.67
F1 Score: 0.65
Confusion Matrix:
[[78 21]
 [18 37]]
```

```
C:\Users\Anusha V\anaconda3\lib\site-packages\sklearn\linear_model\_logistic.py:814: ConvergenceWarning: lbfsgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (`max_iter`) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html> (<https://scikit-learn.org/stable/modules/preprocessing.html>)

Please also refer to the documentation for alternative solver options:

[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression) ([https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression))

```
n_iter_i = _check_optimize_result()
```

In [32]:

```

1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.linear_model import LinearRegression
4 from sklearn.metrics import mean_absolute_error, mean_squared_error, r2
5
6 # Assuming 'Outcome' is the target variable
7 X = df.drop('Outcome', axis=1)
8 y = df['Outcome']
9
10 # Split the data into training and testing sets
11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
12
13 # Assuming Linear Regression as the model
14 model = LinearRegression()
15 model.fit(X_train, y_train)
16
17 # Make predictions on the test set
18 y_pred = model.predict(X_test)
19
20 # Calculate MAE (Mean Absolute Error)
21 mae = mean_absolute_error(y_test, y_pred)
22 print(f"MAE: {mae}")
23
24 # Calculate RMSE (Root Mean Squared Error)
25 rmse = np.sqrt(mean_squared_error(y_test, y_pred))
26 print(f"RMSE: {rmse}")
27
28 # Calculate R2 score
29 r2 = r2_score(y_test, y_pred)
30 print(f"R2 Score: {r2}")
31 0

```

Accuracy: 0.7467532467532467

Precision: 0.64

Recall: 0.67

F1 Score: 0.65

Confusion Matrix:

```

[[78 21]
 [18 37]]

```

C:\Users\Anusha V\anaconda3\lib\site-packages\sklearn\linear\_model\\_logistic.py:814: ConvergenceWarning: lbfgs failed to converge (status=1):  
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max\_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html> (<https://scikit-learn.org/stable/modules/preprocessing.html>)

Please also refer to the documentation for alternative solver options:

[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression) ([https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression))

n\_iter\_i = \_check\_optimize\_result(

In [38]:

```

1 from sklearn.metrics import confusion_matrix

```

```
In [39]: 1 import numpy as np
```

```
In [42]: 1 actual = np.array([1,0,1,1,0,1,0,1,1])
2 predicted = np.array([1,0,1,0,1,1,0,1,0])
3 cm = confusion_matrix(actual, predicted)
4 print("confusion matrix")
5 print(cm)
```

```
confusion matrix
[[2 1]
 [2 4]]
```

```
In [43]: 1 from sklearn.metrics import classification_report
2
3 # Example true labels and predicted labels
4 y_true = [1, 0, 1, 1, 0, 1, 0, 0, 1, 0]
5 y_pred = [1, 0, 1, 1, 0, 0, 1, 0, 1, 1]
6
7 # Generate a classification report
8 report = classification_report(y_true, y_pred)
9
10 # Print the classification report
11 print(report)
12
```

	precision	recall	f1-score	support
0	0.75	0.60	0.67	5
1	0.67	0.80	0.73	5
accuracy			0.70	10
macro avg	0.71	0.70	0.70	10
weighted avg	0.71	0.70	0.70	10

In [48]:

```

1 import numpy as np
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.ensemble import RandomForestRegressor
5 from sklearn.svm import SVR
6 from sklearn.linear_model import Lasso
7 from sklearn.metrics import mean_squared_error, r2_score
8 import matplotlib.pyplot as plt
9
10 # Generate sample data
11 np.random.seed(42)
12 X = np.random.rand(100, 1) * 10
13 y = 2 * X.squeeze() + np.random.randn(100) * 2 # True relationship: y
14
15 # Split the data into training and testing sets
16 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
17
18 # Random Forest Regression
19 rf_model = RandomForestRegressor(n_estimators=100, random_state=42)
20 rf_model.fit(X_train, y_train)
21 rf_pred = rf_model.predict(X_test)
22
23 # SVM Regression
24 svm_model = SVR(kernel='linear')
25 svm_model.fit(X_train, y_train)
26 svm_pred = svm_model.predict(X_test)
27
28 # Lasso Regression
29 lasso_model = Lasso(alpha=0.1)
30 lasso_model.fit(X_train, y_train)
31 lasso_pred = lasso_model.predict(X_test)
32
33 # Evaluate models
34 def evaluate_model(model_name, y_true, y_pred):
35     mse = mean_squared_error(y_true, y_pred)
36     r2 = r2_score(y_true, y_pred)
37     print(f"{model_name} - Mean Squared Error: {mse:.2f}, R2 Score: {r2:.2f}")
38
39 # Print evaluation metrics for each model
40 evaluate_model("Random Forest", y_test, rf_pred)
41 evaluate_model("SVM", y_test, svm_pred)
42 evaluate_model("Lasso", y_test, lasso_pred)
43
44

```

Random Forest - Mean Squared Error: 3.02, R2 Score: 0.92

SVM - Mean Squared Error: 2.52, R2 Score: 0.93

Lasso - Mean Squared Error: 2.60, R2 Score: 0.93

In [49]:

```

1 import numpy as np
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.ensemble import RandomForestRegressor
5 from sklearn.svm import SVR
6 from sklearn.linear_model import Lasso
7 from sklearn.metrics import mean_squared_error, r2_score
8 import matplotlib.pyplot as plt
9
10 # Generate sample data
11 np.random.seed(42)
12 X = np.random.rand(100, 1) * 10
13 y = 2 * X.squeeze() + np.random.randn(100) * 2 # True relationship: y
14
15 # Split the data into training and testing sets
16 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
17
18 # Random Forest Regression
19 rf_model = RandomForestRegressor(n_estimators=100, random_state=42)
20 rf_model.fit(X_train, y_train)
21 rf_pred = rf_model.predict(X_test)
22
23 # SVM Regression
24 svm_model = SVR(kernel='linear')
25 svm_model.fit(X_train, y_train)
26 svm_pred = svm_model.predict(X_test)
27
28 # Lasso Regression
29 lasso_model = Lasso(alpha=0.1)
30 lasso_model.fit(X_train, y_train)
31 lasso_pred = lasso_model.predict(X_test)
32
33 # Evaluate models
34 def evaluate_model(model_name, y_true, y_pred):
35     mse = mean_squared_error(y_true, y_pred)
36     r2 = r2_score(y_true, y_pred)
37     print(f"{model_name} - Mean Squared Error: {mse:.2f}, R2 Score: {r2}")
38
39 # Print evaluation metrics for each model
40 evaluate_model("Random Forest", y_test, rf_pred)
41 evaluate_model("SVM", y_test, svm_pred)
42 evaluate_model("Lasso", y_test, lasso_pred)
43
44 # Plot the true vs predicted values for each model
45 plt.scatter(X_test, y_test, color='black', label='True values')
46 plt.scatter(X_test, rf_pred, color='blue', label='Random Forest')
47 plt.scatter(X_test, svm_pred, color='red', label='SVM')
48 plt.scatter(X_test, lasso_pred, color='green', label='Lasso')
49 plt.legend()
50 plt.xlabel('X')
51 plt.ylabel('y')
52 plt.title('Regression Models Comparison')
53 plt.show()
54

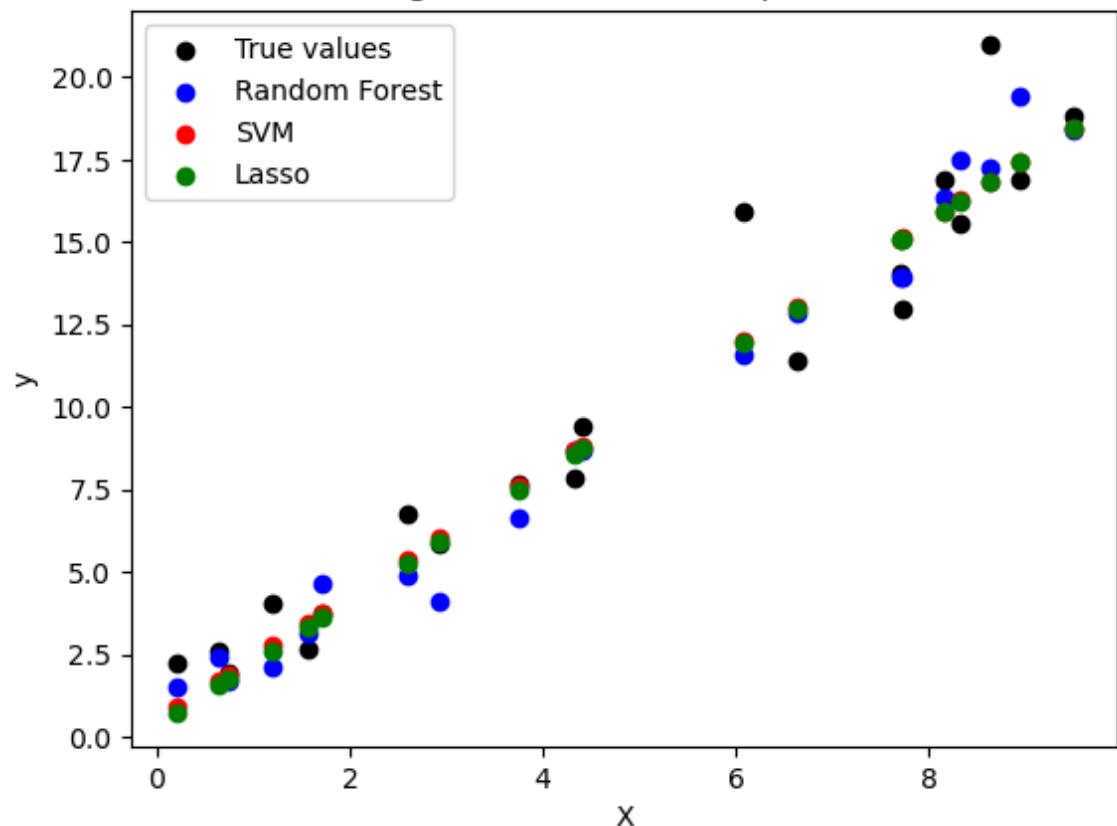
```

Random Forest - Mean Squared Error: 3.02, R2 Score: 0.92

SVM - Mean Squared Error: 2.52, R2 Score: 0.93

Lasso - Mean Squared Error: 2.60, R2 Score: 0.93

## Regression Models Comparison



In [50]:

```
1 import numpy as np
2 from sklearn.model_selection import train_test_split
3 from sklearn.ensemble import RandomForestRegressor
4 from sklearn.svm import SVR
5 from sklearn.linear_model import Lasso
6 from sklearn.metrics import mean_squared_error, r2_score
7
8 # Generate sample data
9 np.random.seed(42)
10 X = np.random.rand(100, 1) * 10
11 y = 2 * X.squeeze() + np.random.randn(100) * 2 # True relationship: y
12
13 # Split the data into training and testing sets
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
15
16 # Define and train models
17 models = {
18     "Random Forest": RandomForestRegressor(n_estimators=100, random_state=42),
19     "SVM": SVR(kernel='linear'),
20     "Lasso": Lasso(alpha=0.1)
21 }
22
23 # Evaluate models
24 for model_name, model in models.items():
25     model.fit(X_train, y_train)
26     y_pred = model.predict(X_test)
27     mse = mean_squared_error(y_test, y_pred)
28     r2 = r2_score(y_test, y_pred)
29     print(f"{model_name} - Mean Squared Error: {mse:.2f}, R2 Score: {r2:.2f}
```

Random Forest - Mean Squared Error: 3.02, R2 Score: 0.92

SVM - Mean Squared Error: 2.52, R2 Score: 0.93

Lasso - Mean Squared Error: 2.60, R2 Score: 0.93

In [51]:

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.model_selection import train_test_split
4 from sklearn.ensemble import RandomForestRegressor
5 from sklearn.svm import SVR
6 from sklearn.linear_model import Lasso
7 from sklearn.metrics import mean_squared_error, r2_score
8 import matplotlib.pyplot as plt
9
10 # Generate sample data
11 np.random.seed(42)
12 X = np.random.rand(100, 1) * 10
13 y = 2 * X.squeeze() + np.random.randn(100) * 2 # True relationship: y
14
15 # Split the data into training and testing sets
16 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
17
18 # Random Forest Regression
19 rf_model = RandomForestRegressor(n_estimators=100, random_state=42)
20 rf_model.fit(X_train, y_train)
21 rf_pred = rf_model.predict(X_test)
22
23 # SVM Regression
24 svm_model = SVR(kernel='linear')
25 svm_model.fit(X_train, y_train)
26 svm_pred = svm_model.predict(X_test)
27
28 # Lasso Regression
29 lasso_model = Lasso(alpha=0.1)
30 lasso_model.fit(X_train, y_train)
31 lasso_pred = lasso_model.predict(X_test)
32
33 # Evaluate models
34 def evaluate_model(model_name, y_true, y_pred):
35     mse = mean_squared_error(y_true, y_pred)
36     r2 = r2_score(y_true, y_pred)
37     print(f"{model_name} - Mean Squared Error: {mse:.2f}, R2 Score: {r2:.2f}")
38
39 # Print evaluation metrics for each model
40 evaluate_model("Random Forest", y_test, rf_pred)
41 evaluate_model("SVM", y_test, svm_pred)
42 evaluate_model("Lasso", y_test, lasso_pred)
43
44
```

Random Forest - Mean Squared Error: 3.02, R2 Score: 0.92

SVM - Mean Squared Error: 2.52, R2 Score: 0.93

Lasso - Mean Squared Error: 2.60, R2 Score: 0.93

In [52]:

```
1 from sklearn.datasets import load_iris
2 from sklearn.model_selection import train_test_split
3 from sklearn.svm import SVC
4 from sklearn.metrics import accuracy_score, precision_score, recall_score
5
6 # Load the Iris dataset (binary classification for simplicity)
7 iris = load_iris()
8 X, y = iris.data, iris.target
9 # Convert the problem into binary classification by considering only the first two classes
10 X, y = X[y != 2], y[y != 2]
11
12 # Split the data into training and testing sets
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
14
15 # Create and train a Support Vector Machine (SVM) classifier
16 classifier = SVC(kernel='linear', C=1.0)
17 classifier.fit(X_train, y_train)
18
19 # Make predictions on the test set
20 y_pred = classifier.predict(X_test)
21
22 # Evaluate the classifier
23 accuracy = accuracy_score(y_test, y_pred)
24 precision = precision_score(y_test, y_pred)
25 recall = recall_score(y_test, y_pred)
26 f1 = f1_score(y_test, y_pred)
27 conf_matrix = confusion_matrix(y_test, y_pred)
28
29 # Print the evaluation metrics
30 print(f"Accuracy: {accuracy:.2f}")
31 print(f"Precision: {precision:.2f}")
32 print(f"Recall: {recall:.2f}")
33 print(f"F1 Score: {f1:.2f}")
34 print("Confusion Matrix:")
35 print(conf_matrix)
36
```

```
Accuracy: 1.00
Precision: 1.00
Recall: 1.00
F1 Score: 1.00
Confusion Matrix:
[[12  0]
 [ 0  8]]
```

In [53]:

```
1 import numpy as np
2 from sklearn.datasets import make_classification
3 from sklearn.model_selection import train_test_split
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.multiclass import OneVsRestClassifier
6 from sklearn.metrics import accuracy_score, classification_report
7 from imblearn.over_sampling import SMOTE
8
9 # Generate sample data for binary classification
10 X_binary, y_binary = make_classification(n_samples=1000, n_features=20,
11
12 # Generate sample data for multi-class classification
13 X_multiclass, y_multiclass = make_classification(n_samples=1000, n_feat
14
15 # Generate sample data for multi-label classification
16 X_multilabel, y_multilabel = make_classification(n_samples=1000, n_feat
17
18 # Make a binary classification problem imbalanced
19 X_imbalanced, y_imbalanced = make_classification(n_samples=1000, n_feat
20
21 # Apply SMOTE for over-sampling in imbalanced classification
22 smote = SMOTE(random_state=42)
23 X_imbalanced_resampled, y_imbalanced_resampled = smote.fit_resample(X_i
24
25 # Define and train classifiers
26 classifiers = {
27     "Binary Classification": RandomForestClassifier(random_state=42),
28     "Multi-Class Classification": OneVsRestClassifier(RandomForestClass
29     "Multi-Label Classification": RandomForestClassifier(random_state=4
30     "Imbalanced Classification": RandomForestClassifier(random_state=42
31 }
32
33 # Train and evaluate classifiers
34 for problem_type, clf in classifiers.items():
35     if "Binary" in problem_type:
36         X_train, X_test, y_train, y_test = train_test_split(X_binary, y_b
37     elif "Multi-Class" in problem_type:
38         X_train, X_test, y_train, y_test = train_test_split(X_multiclass
39     elif "Multi-Label" in problem_type:
40         X_train, X_test, y_train, y_test = train_test_split(X_multilabel
41     elif "Imbalanced" in problem_type:
42         X_train, X_test, y_train, y_test = train_test_split(X_imbalance
43
44     clf.fit(X_train, y_train)
45     y_pred = clf.predict(X_test)
46
47     print(f"\n{problem_type} Evaluation:")
48     print(f"Accuracy: {accuracy_score(y_test, y_pred):.2f}")
49     print("Classification Report:")
50     print(classification_report(y_test, y_pred))
51
```



**Binary Classification Evaluation:****Accuracy:** 0.92**Classification Report:**

	precision	recall	f1-score	support
0	0.95	0.89	0.92	104
1	0.89	0.95	0.92	96
accuracy			0.92	200
macro avg	0.92	0.92	0.92	200
weighted avg	0.92	0.92	0.92	200

**Multi-Class Classification Evaluation:****Accuracy:** 0.84**Classification Report:**

	precision	recall	f1-score	support
0	0.75	0.85	0.80	61
1	0.89	0.83	0.86	81
2	0.88	0.84	0.86	58
accuracy			0.84	200
macro avg	0.84	0.84	0.84	200
weighted avg	0.85	0.84	0.84	200

**Multi-Label Classification Evaluation:****Accuracy:** 0.65**Classification Report:**

	precision	recall	f1-score	support
0	0.67	0.56	0.61	43
1	0.62	0.70	0.66	33
2	0.62	0.55	0.58	38
3	0.68	0.68	0.68	40
4	0.66	0.76	0.71	46
accuracy			0.65	200
macro avg	0.65	0.65	0.65	200
weighted avg	0.65	0.65	0.65	200

**Imbalanced Classification Evaluation:****Accuracy:** 0.98**Classification Report:**

	precision	recall	f1-score	support
0	0.98	0.98	0.98	180
1	0.98	0.98	0.98	179
accuracy			0.98	359
macro avg	0.98	0.98	0.98	359
weighted avg	0.98	0.98	0.98	359

In [1]:

```
1 from sklearn.datasets import make_classification
2 from sklearn.model_selection import train_test_split
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.multiclass import OneVsRestClassifier
5 from sklearn.ensemble import RandomForestClassifier
6 from sklearn.metrics import accuracy_score, classification_report
7
8 # Binary Classification
9 X_bin, y_bin = make_classification(n_samples=1000, n_features=20, n_informative=2, n_redundant=10, n_clusters_per_class=1, random_state=42)
10 X_bin_train, X_bin_test, y_bin_train, y_bin_test = train_test_split(X_bin, y_bin, test_size=0.2, random_state=42)
11
12 binary_classifier = LogisticRegression()
13 binary_classifier.fit(X_bin_train, y_bin_train)
14 y_bin_pred = binary_classifier.predict(X_bin_test)
15
16 print("Binary Classification Accuracy:", accuracy_score(y_bin_test, y_bin_pred))
17 print("Binary Classification Report:\n", classification_report(y_bin_test, y_bin_pred))
18
19 # Multi-class Classification
20 X_multi, y_multi = make_classification(n_samples=1000, n_features=20, n_informative=2, n_redundant=10, n_clusters_per_class=3, random_state=42)
21 X_multi_train, X_multi_test, y_multi_train, y_multi_test = train_test_split(X_multi, y_multi, test_size=0.2, random_state=42)
22
23 multi_classifier = OneVsRestClassifier(LogisticRegression())
24 multi_classifier.fit(X_multi_train, y_multi_train)
25 y_multi_pred = multi_classifier.predict(X_multi_test)
26
27 print("\nMulti-class Classification Accuracy:", accuracy_score(y_multi_test, y_multi_pred))
28 print("Multi-class Classification Report:\n", classification_report(y_multi_test, y_multi_pred))
29
30 # Multi-label Classification
31 X_multilabel, y_multilabel = make_classification(n_samples=1000, n_features=20, n_informative=2, n_redundant=10, n_clusters_per_class=2, random_state=42)
32 X_multilabel_train, X_multilabel_test, y_multilabel_train, y_multilabel_test = train_test_split(X_multilabel, y_multilabel, test_size=0.2, random_state=42)
33
34 multilabel_classifier = RandomForestClassifier()
35 multilabel_classifier.fit(X_multilabel_train, y_multilabel_train)
36 y_multilabel_pred = multilabel_classifier.predict(X_multilabel_test)
37
38 print("\nMulti-label Classification Accuracy:", accuracy_score(y_multilabel_test, y_multilabel_pred))
39 print("Multi-label Classification Report:\n", classification_report(y_multilabel_test, y_multilabel_pred))
40
41 # Imbalanced Classification
42 from imblearn.datasets import make_imbalance
43 from collections import Counter
44 from sklearn.svm import SVC
45
46 X_imbalanced, y_imbalanced = make_classification(n_samples=1000, n_features=20, n_informative=2, n_redundant=10, n_clusters_per_class=1, random_state=42)
47 X_imbalanced_train, X_imbalanced_test, y_imbalanced_train, y_imbalanced_test = train_test_split(X_imbalanced, y_imbalanced, test_size=0.2, random_state=42)
48
49 print("\nOriginal Class Distribution:", Counter(y_imbalanced_train))
50
51 imbalanced_classifier = SVC()
52 imbalanced_classifier.fit(X_imbalanced_train, y_imbalanced_train)
53 y_imbalanced_pred = imbalanced_classifier.predict(X_imbalanced_test)
54
55 print("Imbalanced Classification Accuracy:", accuracy_score(y_imbalanced_test, y_imbalanced_pred))
56 print("Imbalanced Classification Report:\n", classification_report(y_imbalanced_test, y_imbalanced_pred))
```

```
C:\Users\Anusha V\anaconda3\lib\site-packages\scipy\_init_.py:155: UserWarning: A NumPy version >=1.18.5 and <1.25.0 is required for this version of SciPy (detected version 1.26.2)
  warnings.warn(f"A NumPy version >={np_minversion} and <{np_maxversion}"
```

Binary Classification Accuracy: 0.82

Binary Classification Report:

	precision	recall	f1-score	support
0	0.86	0.78	0.82	104
1	0.78	0.86	0.82	96
accuracy			0.82	200
macro avg	0.82	0.82	0.82	200
weighted avg	0.82	0.82	0.82	200

Multi-class Classification Accuracy: 0.755

Multi-class Classification Report:

	precision	recall	f1-score	support
0	0.61	0.67	0.64	61
1	0.87	0.80	0.83	81
2	0.78	0.78	0.78	58
accuracy			0.76	200
macro avg	0.75	0.75	0.75	200
weighted avg	0.76	0.76	0.76	200

---

```
-
```

**TypeError** Traceback (most recent call last)  
t)  
~\AppData\Local\Temp\ipykernel\_6100\3238744082.py in <module>  
 29  
 30 # Multi-label Classification  
--> 31 X\_multilabel, y\_multilabel = make\_classification(n\_samples=1000, n  
 \_features=20, n\_informative=10, n\_classes=5, random\_state=42, n\_labels=3)  
 32 X\_multilabel\_train, X\_multilabel\_test, y\_multilabel\_train, y\_multi  
label\_test = train\_test\_split(X\_multilabel, y\_multilabel, test\_size=0.2, r  
andom\_state=42)  
 33

**TypeError:** make\_classification() got an unexpected keyword argument 'n\_labels'

In [6]:

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.metrics import accuracy_score, classification_report
4 from sklearn.datasets import make_classification
5
6 # Generate a hypothetical binary classification dataset
7 X, y = make_classification(n_samples=1000, n_features=20, n_informative=15,
8                            n_redundant=10, n_classes=2, random_state=42)
9
10 # Split the data into training and testing sets
11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
12
13 # Binary classification model (Logistic Regression)
14 model = LogisticRegression()
15 model.fit(X_train, y_train)
16
17 # Make predictions on the test set
18 predictions = model.predict(X_test)
19
20 # Evaluate the model
21 accuracy = accuracy_score(y_test, predictions)
22 report = classification_report(y_test, predictions)
23
24 print(f'Accuracy: {accuracy}')
25 print('Classification Report:\n', report)
26
```

Accuracy: 0.875

Classification Report:

	precision	recall	f1-score	support
0	0.88	0.89	0.89	111
1	0.86	0.85	0.86	89
accuracy			0.88	200
macro avg	0.87	0.87	0.87	200
weighted avg	0.87	0.88	0.87	200

In [3]:

```
1 from sklearn.datasets import load_iris
2 from sklearn.model_selection import train_test_split
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.metrics import accuracy_score, classification_report
5
6 # Sample data for multi-class classification (Iris dataset)
7 data = load_iris()
8 X, y = data.data, data.target
9
10 # Split the data into training and testing sets
11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
12
13 # Multi-class classification model (Logistic Regression)
14 model = LogisticRegression(multi_class='auto', max_iter=1000)
15 model.fit(X_train, y_train)
16
17 # Make predictions on the test set
18 predictions = model.predict(X_test)
19
20 # Evaluate the model
21 accuracy = accuracy_score(y_test, predictions)
22 report = classification_report(y_test, predictions)
23
24 print(f'Accuracy: {accuracy}')
25 print('Classification Report:\n', report)
26
```

Accuracy: 1.0

Classification Report:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	10
1	1.00	1.00	1.00	9
2	1.00	1.00	1.00	11
accuracy			1.00	30
macro avg	1.00	1.00	1.00	30
weighted avg	1.00	1.00	1.00	30

In [7]:

```

1 from sklearn.model_selection import train_test_split
2 from sklearn.multioutput import MultiOutputClassifier
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.metrics import accuracy_score, classification_report
5 from sklearn.datasets import make_multilabel_classification
6
7 # Generate a hypothetical multi-label classification dataset
8 X, y = make_multilabel_classification(n_samples=1000, n_features=20, n_labels=3)
9
10 # Split the data into training and testing sets
11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
12
13 # Multi-Label classification model (Random Forest)
14 model = MultiOutputClassifier(RandomForestClassifier(n_estimators=100))
15 model.fit(X_train, y_train)
16
17 # Make predictions on the test set
18 predictions = model.predict(X_test)
19
20 # Evaluate the model
21 accuracy = accuracy_score(y_test, predictions)
22 report = classification_report(y_test, predictions)
23
24 print(f'Accuracy: {accuracy}')
25 print('Classification Report:\n', report)
26
27

```

Accuracy: 0.69

Classification Report:

	precision	recall	f1-score	support
0	0.89	0.64	0.74	85
1	0.93	0.89	0.91	126
2	0.88	0.88	0.88	111
micro avg	0.90	0.82	0.86	322
macro avg	0.90	0.80	0.84	322
weighted avg	0.90	0.82	0.85	322
samples avg	0.79	0.72	0.74	322

C:\Users\Anusha V\anaconda3\lib\site-packages\sklearn\metrics\\_classification.py:1318: UndefinedMetricWarning: Precision and F-score are ill-defined and being set to 0.0 in samples with no predicted labels. Use `zero\_division` parameter to control this behavior.

\_warn\_prf(average, modifier, msg\_start, len(result))

C:\Users\Anusha V\anaconda3\lib\site-packages\sklearn\metrics\\_classification.py:1318: UndefinedMetricWarning: Recall and F-score are ill-defined and being set to 0.0 in samples with no true labels. Use `zero\_division` parameter to control this behavior.

\_warn\_prf(average, modifier, msg\_start, len(result))

In [5]:

```
1 from sklearn.datasets import make_classification
2 from sklearn.model_selection import train_test_split
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.metrics import accuracy_score, classification_report
5
6 # Generate imbalanced binary classification data
7 X, y = make_classification(n_samples=1000, n_features=20, n_informative=15,
8                            n_redundant=10, weights=[0.9, 0.1], random_state=42)
9
10 # Split the data into training and testing sets
11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
12
13 # Imbalanced classification model (Random Forest)
14 model = RandomForestClassifier(n_estimators=100)
15 model.fit(X_train, y_train)
16
17 # Make predictions on the test set
18 predictions = model.predict(X_test)
19
20 # Evaluate the model
21 accuracy = accuracy_score(y_test, predictions)
22 report = classification_report(y_test, predictions)
23
24 print(f'Accuracy: {accuracy}')
25 print('Classification Report:\n', report)
26
```

Accuracy: 0.98

Classification Report:

	precision	recall	f1-score	support
0	0.98	1.00	0.99	184
1	1.00	0.75	0.86	16
accuracy			0.98	200
macro avg	0.99	0.88	0.92	200
weighted avg	0.98	0.98	0.98	200

## vectroized operations

In [2]:

```
1 import numpy as np
2
3 # Creating two arrays
4 array1 = np.array([1, 2, 3, 4])
5 array2 = np.array([5, 6, 7, 8])
6
7 # Vectorized addition
8 result_addition = array1 + array2
9
10 # Vectorized multiplication
11 result_multiplication = array1 * array2
12
13 # Vectorized square root
14 result_sqrt = np.sqrt(array1)
15
16 # Print results
17 print("Vectorized Addition:", result_addition)
18 print("Vectorized Multiplication:", result_multiplication)
19 print("Vectorized Square Root:", result_sqrt)
20
```

```
Vectorized Addition: [ 6  8 10 12]
Vectorized Multiplication: [ 5 12 21 32]
Vectorized Square Root: [1.             1.41421356  1.73205081  2.           ]
```

In [3]:

```
1 import pandas as pd
2 from functools import reduce
3
4 # Creating a sample DataFrame
5 data = {'numbers': [1, 2, 3, 4, 5]}
6 df = pd.DataFrame(data)
7
8 # Using map with Lambda to square each element in the 'numbers' column
9 df['squared'] = df['numbers'].map(lambda x: x ** 2)
10
11 # Using filter with Lambda to keep only even numbers in the 'numbers' column
12 df_filtered = df[df['numbers'].apply(lambda x: x % 2 == 0)]
13
14 # Using reduce with Lambda to find the product of all elements in the 'numbers' column
15 product_of_numbers = reduce(lambda x, y: x * y, df['numbers'])
16
17 # Print results
18 print("Original DataFrame:")
19 print(df)
20
21 print("\nDataFrame with Squared Numbers:")
22 print(df)
23
24 print("\nDataFrame with Only Even Numbers:")
25 print(df_filtered)
26
27 print("\nProduct of Numbers using Reduce:", product_of_numbers)
28
```

Original DataFrame:

	numbers	squared
0	1	1
1	2	4
2	3	9
3	4	16
4	5	25

DataFrame with Squared Numbers:

	numbers	squared
0	1	1
1	2	4
2	3	9
3	4	16
4	5	25

DataFrame with Only Even Numbers:

	numbers	squared
1	2	4
3	4	16

Product of Numbers using Reduce: 120

```
In [5]: 1 pip install autograd  
2
```

```
Collecting autograd  
  Downloading autograd-1.6.2-py3-none-any.whl (49 kB)  
----- 49.3/49.3 kB 357.9 kB/s eta 0:  
00:00  
Requirement already satisfied: future>=0.15.2 in c:\users\anusha v\anaconda3  
\lib\site-packages (from autograd) (0.18.2)  
Requirement already satisfied: numpy>=1.12 in c:\users\anusha v\anaconda3  
\lib\site-packages (from autograd) (1.26.2)  
Installing collected packages: autograd  
Successfully installed autograd-1.6.2  
Note: you may need to restart the kernel to use updated packages.
```

In [6]:

```

1 import numpy as np
2 import autograd.numpy as ag_np
3 from autograd import grad
4
5 # Scalars
6 scalar = 5
7
8 # Vectors
9 vector = np.array([1, 2, 3])
10
11 # Matrices
12 matrix = np.array([[1, 2, 3], [4, 5, 6]])
13
14 # Tensors
15 tensor = np.array([[[1, 2, 3], [4, 5, 6]], [[7, 8, 9], [10, 11, 12]]])
16
17 # Gradient Computation using autograd
18 def scalar_function(x):
19     return 3 * x**2 + 2 * x + 1
20
21 # Compute the gradient of the scalar function
22 grad_scalar_function = grad(scalar_function)
23
24 # Test gradient computation
25 point = 2.0
26 gradient_at_point = grad_scalar_function(point)
27
28 # Print results
29 print("Scalar:", scalar)
30 print("Vector:", vector)
31 print("Matrix:")
32 print(matrix)
33 print("Tensor:")
34 print(tensor)
35
36 print("\nGradient of the scalar function f(x) = 3x^2 + 2x + 1 at x =",
```

```

Scalar: 5
Vector: [1 2 3]
Matrix:
[[1 2 3]
 [4 5 6]]
Tensor:
[[[ 1  2  3]
   [ 4  5  6]]

 [[ 7  8  9]
   [10 11 12]]]
```

Gradient of the scalar function  $f(x) = 3x^2 + 2x + 1$  at  $x = 2.0$  is 14.0

In [7]:

```

1 import pandas as pd
2 from sklearn.preprocessing import StandardScaler, LabelEncoder
3 from sklearn.model_selection import train_test_split
4
5 # Sample data
6 data = {'Feature1': [1, 2, 3, 4],
7          'Feature2': ['A', 'B', 'A', 'B'],
8          'Target': [0, 1, 0, 1]}
9
10 df = pd.DataFrame(data)
11
12 # Display original data
13 print("Original DataFrame:")
14 print(df)
15
16 # Handling missing values (if any)
17 df.fillna(value=0, inplace=True)
18
19 # Encoding categorical variables
20 label_encoder = LabelEncoder()
21 df['Feature2'] = label_encoder.fit_transform(df['Feature2'])
22
23 # Splitting the data into features (X) and target variable (y)
24 X = df.drop('Target', axis=1)
25 y = df['Target']
26
27 # Splitting the data into training and testing sets
28 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
29
30 # Standardizing features
31 scaler = StandardScaler()
32 X_train_scaled = scaler.fit_transform(X_train)
33 X_test_scaled = scaler.transform(X_test)
34
35 # Display preprocessed data
36 print("\nPreprocessed DataFrame:")
37 print("X_train_scaled:")
38 print(pd.DataFrame(X_train_scaled, columns=X.columns))
39 print("y_train:")
40 print(y_train)
41 print("X_test_scaled:")
42 print(pd.DataFrame(X_test_scaled, columns=X.columns))
43 print("y_test:")
44 print(y_test)
45

```

C:\Users\Anusha V\anaconda3\lib\site-packages\scipy\\_init\_.py:155: UserWarning: A NumPy version >=1.18.5 and <1.25.0 is required for this version of SciPy (detected version 1.26.2)

```
    warnings.warn(f"A NumPy version >={np_minversion} and <{np_maxversion}"
```

Original DataFrame:

	Feature1	Feature2	Target
0	1	A	0
1	2	B	1
2	3	A	0
3	4	B	1

Preprocessed DataFrame:

X\_train\_scaled:

	Feature1	Feature2
0	1.069045	1.414214
1	-1.336306	-0.707107
2	0.267261	-0.707107

y\_train:

3	1
0	0
2	0

Name: Target, dtype: int64

X\_test\_scaled:

	Feature1	Feature2
0	-0.534522	1.414214

y\_test:

1	1
---	---

Name: Target, dtype: int64

In [10]:

```
1 import pandas as pd
2
3 # Sample data with missing values
4 data = {'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Emily'],
5         'Age': [25, 30, None, 22, 35],
6         'Salary': [50000, None, 60000, 45000, 70000]}
7
8 df = pd.DataFrame(data)
9
10 # Display information about the DataFrame, including non-null counts
11 print("Original DataFrame Info:")
12 print(df.info())
13
14 # Check for missing values
15 print("\nMissing Values:")
16 print(df.isnull())
17
18 # Handling missing values
19 # Option 1: Drop rows with missing values
20 df_dropped = df.dropna()
21
22 # Option 2: Fill missing values with a specific value (e.g., mean)
23 df_filled = df.fillna(df.mean())
24
25 # Display processed DataFrames
26 print("\nDataFrame after Dropping Missing Values:")
27 print(df_dropped)
28
29 print("\nDataFrame after Filling Missing Values:")
30 print(df_filled)
31
```

Original DataFrame Info:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5 entries, 0 to 4
Data columns (total 3 columns):
 #   Column   Non-Null Count   Dtype  
--- 
 0   Name      5 non-null     object  
 1   Age       4 non-null     float64 
 2   Salary    4 non-null     float64 
dtypes: float64(2), object(1)
memory usage: 248.0+ bytes
None
```

Missing Values:

	Name	Age	Salary
0	False	False	False
1	False	False	True
2	False	True	False
3	False	False	False
4	False	False	False

DataFrame after Dropping Missing Values:

	Name	Age	Salary
0	Alice	25.0	50000.0
3	David	22.0	45000.0
4	Emily	35.0	70000.0

DataFrame after Filling Missing Values:

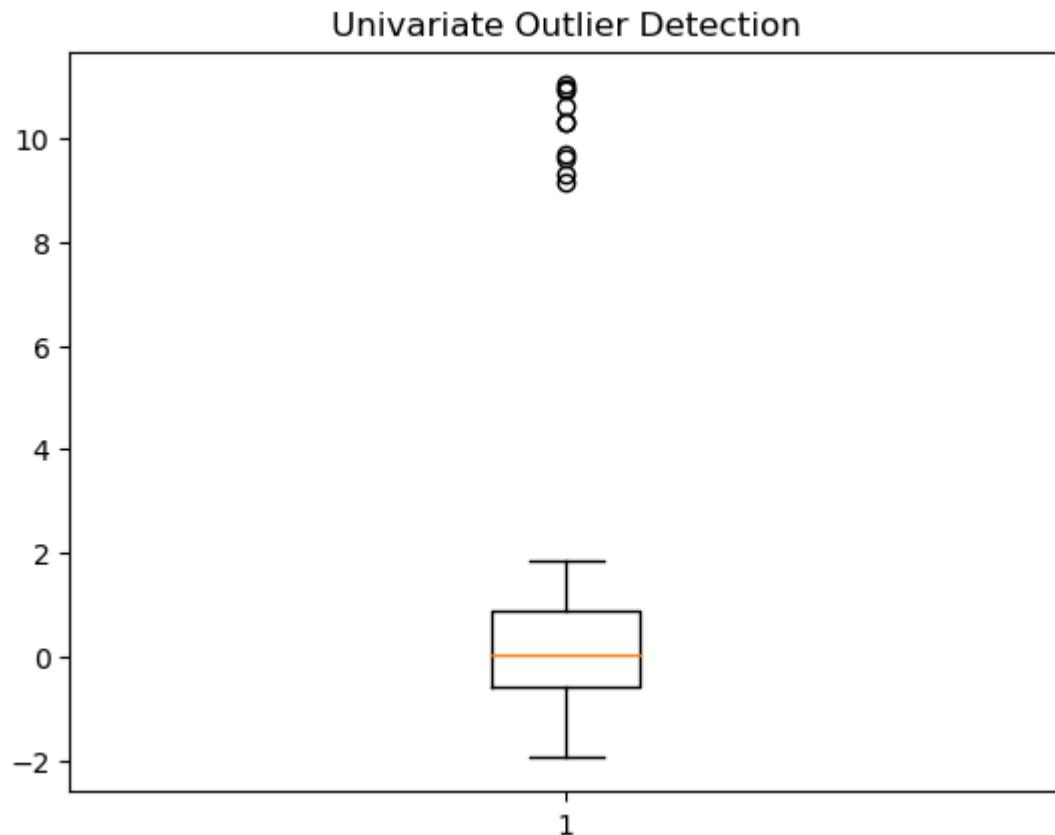
	Name	Age	Salary
0	Alice	25.0	50000.0
1	Bob	30.0	56250.0
2	Charlie	28.0	60000.0
3	David	22.0	45000.0
4	Emily	35.0	70000.0

```
C:\Users\Anusha V\AppData\Local\Temp\ipykernel_7028\26514610.py:23: Future
Warning: Dropping of nuisance columns in DataFrame reductions (with 'numer
ic_only=None') is deprecated; in a future version this will raise TypeErro
r. Select only valid columns before calling the reduction.
```

```
df_filled = df.fillna(df.mean())
```

In [11]:

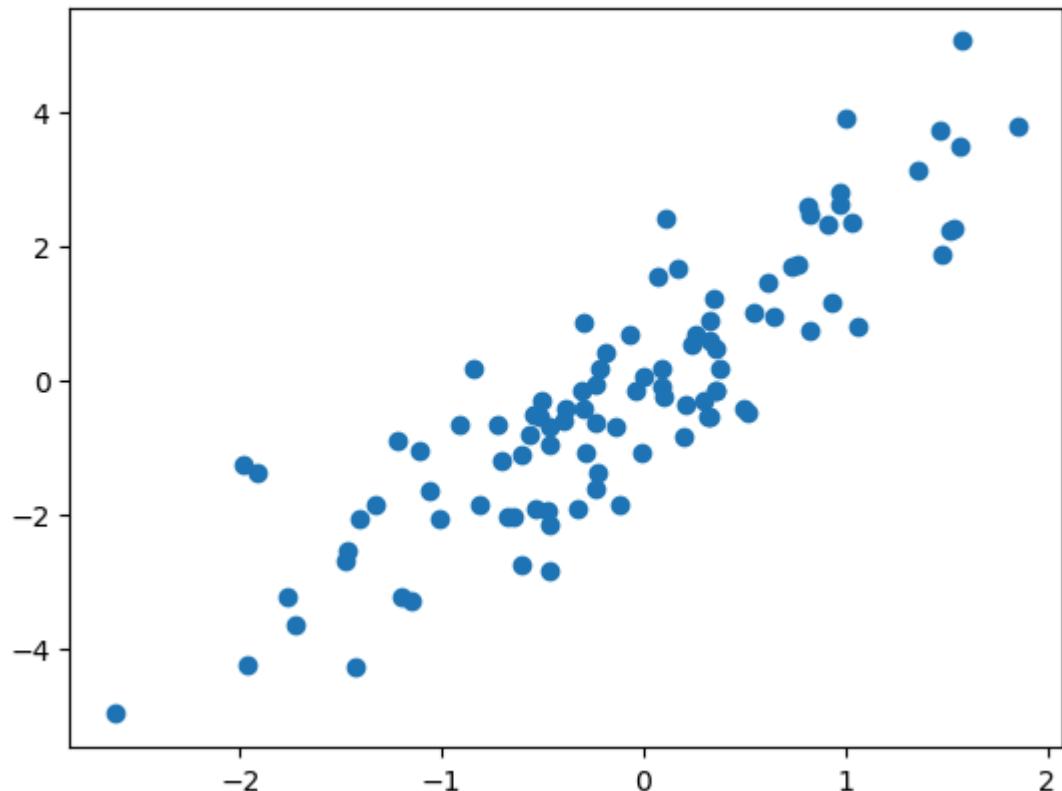
```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # Generating a sample univariate dataset
5 np.random.seed(42)
6 data_univariate = np.concatenate([np.random.normal(0, 1, 50), np.random.
7
8 # Plotting the univariate data
9 plt.boxplot(data_univariate)
10 plt.title("Univariate Outlier Detection")
11 plt.show()
12
```



In [12]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # Generating a sample bivariate dataset
5 np.random.seed(42)
6 data_bivariate_x = np.random.normal(0, 1, 100)
7 data_bivariate_y = data_bivariate_x * 2 + np.random.normal(0, 1, 100)
8
9 # Plotting the bivariate data
10 plt.scatter(data_bivariate_x, data_bivariate_y)
11 plt.title("Bivariate Outlier Detection")
12 plt.show()
13
```

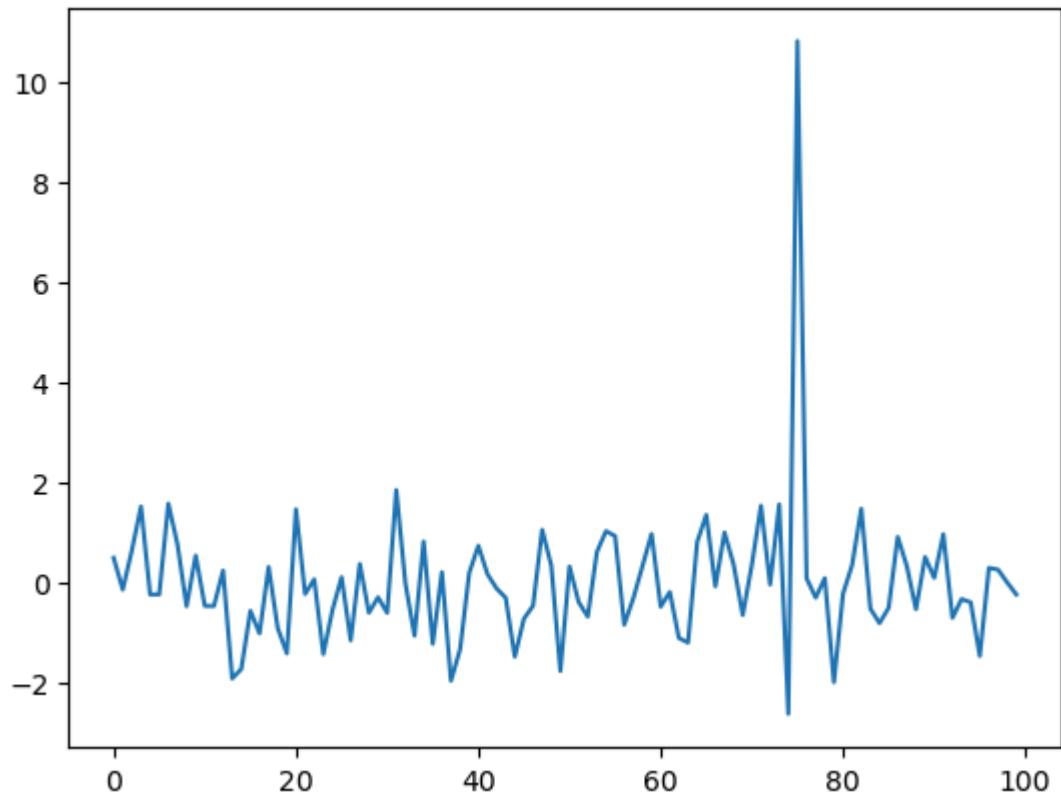
Bivariate Outlier Detection



In [13]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # Generating a sample time series dataset with an outlier
5 np.random.seed(42)
6 data_timeseries = np.random.normal(0, 1, 100)
7 data_timeseries[75] += 10 # Introducing an outlier
8
9 # Plotting the time series data
10 plt.plot(data_timeseries)
11 plt.title("Time Series Outlier Detection")
12 plt.show()
13
```

Time Series Outlier Detection



## Numerosity Data Reduction:

In [14]:

```

1 import numpy as np
2 from sklearn.datasets import make_classification
3 from sklearn.cluster import KMeans
4
5 # Generating a sample dataset for numerosity reduction
6 X, _ = make_classification(n_samples=1000, n_features=10, random_state=
7
8 # Applying K-Means clustering for numerosity reduction
9 kmeans = KMeans(n_clusters=5, random_state=42)
10 X_reduced = kmeans.fit_predict(X)
11
12 # Print the reduced dataset
13 print("Original Data Shape:", X.shape)
14 print("Reduced Data Shape:", X_reduced.shape)
15

```

Original Data Shape: (1000, 10)  
 Reduced Data Shape: (1000, )

## Dimensionality Data Reduction:

In [15]:

```

1 import numpy as np
2 from sklearn.datasets import make_classification
3 from sklearn.decomposition import PCA
4
5 # Generating a sample dataset for dimensionality reduction
6 X, _ = make_classification(n_samples=1000, n_features=10, random_state=
7
8 # Applying Principal Component Analysis (PCA) for dimensionality reduction
9 pca = PCA(n_components=3)
10 X_reduced = pca.fit_transform(X)
11
12 # Print the reduced dataset
13 print("Original Data Shape:", X.shape)
14 print("Reduced Data Shape:", X_reduced.shape)
15

```

Original Data Shape: (1000, 10)  
 Reduced Data Shape: (1000, 3)

## Data transformation in python code

In [16]:

```

1 import pandas as pd
2 from sklearn.preprocessing import StandardScaler, OneHotEncoder
3 from sklearn.compose import ColumnTransformer
4 from sklearn.pipeline import Pipeline
5
6 # Sample data
7 data = {'Name': ['Alice', 'Bob', 'Charlie', 'David'],
8         'Age': [25, 30, 22, 35],
9         'Gender': ['Female', 'Male', 'Male', 'Female'],
10        'Salary': [50000, 60000, 45000, 70000]}
11
12 df = pd.DataFrame(data)
13
14 # Display original data
15 print("Original DataFrame:")
16 print(df)
17
18 # Define columns for different transformations
19 numeric_features = ['Age', 'Salary']
20 categorical_features = ['Gender']
21
22 # Create transformers
23 numeric_transformer = StandardScaler()
24 categorical_transformer = OneHotEncoder()
25
26 # Combine transformers into a preprocessor using ColumnTransformer
27 preprocessor = ColumnTransformer(
28     transformers=[
29         ('num', numeric_transformer, numeric_features),
30         ('cat', categorical_transformer, categorical_features)
31     ])
32
33 # Define a pipeline with the preprocessor
34 pipeline = Pipeline(steps=[('preprocessor', preprocessor)])
35
36 # Fit and transform the data
37 transformed_data = pipeline.fit_transform(df)
38
39 # Display transformed data
40 transformed_df = pd.DataFrame(transformed_data, columns=['Age_scaled',
41 print("\nTransformed DataFrame:")
42 print(transformed_df)
43

```

Original DataFrame:

	Name	Age	Gender	Salary
0	Alice	25	Female	50000
1	Bob	30	Male	60000
2	Charlie	22	Male	45000
3	David	35	Female	70000

Transformed DataFrame:

	Age_scaled	Salary_scaled	Gender_Female	Gender_Male
0	-0.606092	-0.650945	1.0	0.0
1	0.404061	0.390567	0.0	1.0
2	-1.212183	-1.171700	0.0	1.0
3	1.414214	1.432078	1.0	0.0

# **Data transformation with binary coding, ranking transformation, discretization in python**

In [17]:

```
1 import pandas as pd
2 from sklearn.preprocessing import LabelBinarizer, MinMaxScaler, KBinsDiscretizer
3
4 # Sample data
5 data = {'Category': ['A', 'B', 'A', 'C', 'B'],
6         'Value': [10, 20, 15, 5, 25]}
7
8 df = pd.DataFrame(data)
9
10 # Display original data
11 print("Original DataFrame:")
12 print(df)
13
14 # Binary coding for categorical variable 'Category'
15 lb = LabelBinarizer()
16 binary_coding_result = lb.fit_transform(df['Category'])
17 binary_coding_df = pd.DataFrame(binary_coding_result, columns=[f'col_{i}' for i in range(len(lb.classes_))])
18
19 # Ranking transformation for numerical variable 'Value'
20 df['Value_rank'] = df['Value'].rank()
21
22 # Discretization for numerical variable 'Value'
23 discretizer = KBinsDiscretizer(n_bins=3, encode='ordinal', strategy='uniform')
24 df['Value_discretized'] = discretizer.fit_transform(df[['Value']])
25
26 # Display transformed data
27 print("\nDataFrame after Transformation:")
28 print(binary_coding_df)
29 print("\nDataFrame with Ranking Transformation:")
30 print(df[['Value', 'Value_rank']])
31 print("\nDataFrame with Discretization:")
32 print(df[['Value', 'Value_discretized']])
33
```

Original DataFrame:

	Category	Value
0	A	10
1	B	20
2	A	15
3	C	5
4	B	25

DataFrame after Transformation:

	A_binary	B_binary	C_binary
0	1	0	0
1	0	1	0
2	1	0	0
3	0	0	1
4	0	1	0

DataFrame with Ranking Transformation:

	Value	Value_rank
0	10	2.0
1	20	4.0
2	15	3.0
3	5	1.0
4	25	5.0

DataFrame with Discretization:

	Value	Value_discretized
0	10	0.0
1	20	2.0
2	15	1.0
3	5	0.0
4	25	2.0

## Data splitting using train\_test\_split

In [19]:

```
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3
4 # Sample data
5 data = {'Feature1': [1, 2, 3, 4, 5],
6         'Feature2': [10, 20, 30, 40, 50],
7         'Target': [0, 1, 0, 1, 1]}
8
9 df = pd.DataFrame(data)
10
11 # Splitting the data into features (X) and target variable (y)
12 X = df[['Feature1', 'Feature2']]
13 y = df['Target']
14
15 # Splitting the data into training and testing sets
16 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
17
18 # Display the split datasets
19 print("X_train:")
20 print(X_train)
21 print("\ny_train:")
22 print(y_train)
23 print("\nX_test:")
24 print(X_test)
25 print("\ny_test:")
26 print(y_test)
27
```

```
X_train:
  Feature1  Feature2
4          5        50
2          3        30
0          1        10
3          4        40

y_train:
4    1
2    0
0    0
3    1
Name: Target, dtype: int64

X_test:
  Feature1  Feature2
1          2        20

y_test:
1    1
Name: Target, dtype: int64
```

## Python code for underfitting and overfitting



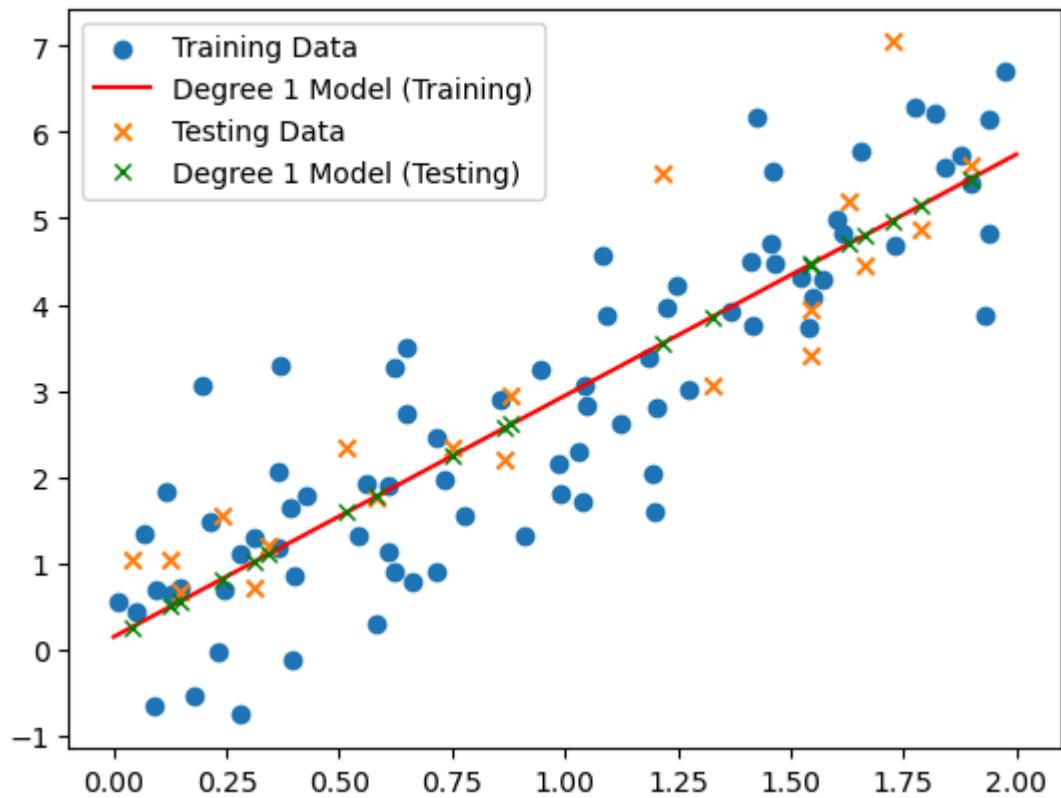
In [30]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.model_selection import train_test_split
4 from sklearn.linear_model import LinearRegression
5 from sklearn.preprocessing import PolynomialFeatures
6 from sklearn.metrics import mean_squared_error
7
8 # Generate synthetic data
9 np.random.seed(42)
10 X = 2 * np.random.rand(100, 1)
11 y = 3 * X + np.random.randn(100, 1)
12
13 # Split the data into training and testing sets
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
15
16 # Function to plot the data and model predictions
17 def plot_model(degree, X, y, X_test, y_test):
18     poly_features = PolynomialFeatures(degree=degree, include_bias=False)
19     X_poly = poly_features.fit_transform(X)
20
21     model = LinearRegression()
22     model.fit(X_poly, y)
23
24     # Training data
25     plt.scatter(X, y, label='Training Data')
26     x_range = np.linspace(0, 2, 100).reshape(-1, 1)
27     x_range_poly = poly_features.transform(x_range)
28     plt.plot(x_range, model.predict(x_range_poly), label=f'Degree {degree} {model}')
29
30     # Testing data
31     plt.scatter(X_test, y_test, label='Testing Data', marker='x')
32     x_test_poly = poly_features.transform(X_test)
33     plt.plot(X_test, model.predict(x_test_poly), 'x', label=f'Degree {degree} {model}')
34
35     plt.title(f'Polynomial Regression - Degree {degree}')
36     plt.legend()
37     plt.show()
38
39 # Try different polynomial degrees to observe underfitting and overfitting
40 degrees = [1, 3, 10]
41 for degree in degrees:
42     poly_features = PolynomialFeatures(degree=degree, include_bias=False)
43     X_poly = poly_features.fit_transform(X_train)
44
45     model = LinearRegression()
46     model.fit(X_poly, y_train)
47
48     # Training error
49     y_train_pred = model.predict(X_poly)
50     train_error = mean_squared_error(y_train, y_train_pred)
51
52     # Testing error
53     X_test_poly = poly_features.transform(X_test)
54     y_test_pred = model.predict(X_test_poly)
55     test_error = mean_squared_error(y_test, y_test_pred)
56
57     print(f'Degree {degree} - Training Error: {train_error:.3f}, Testing Error: {test_error:.3f}')
58
59     # Plot the model
60     plot_model(degree, X_train, y_train, X_test, y_test)
```

61

Degree 1 - Training Error: 0.848, Testing Error: 0.654

### Polynomial Regression - Degree 1



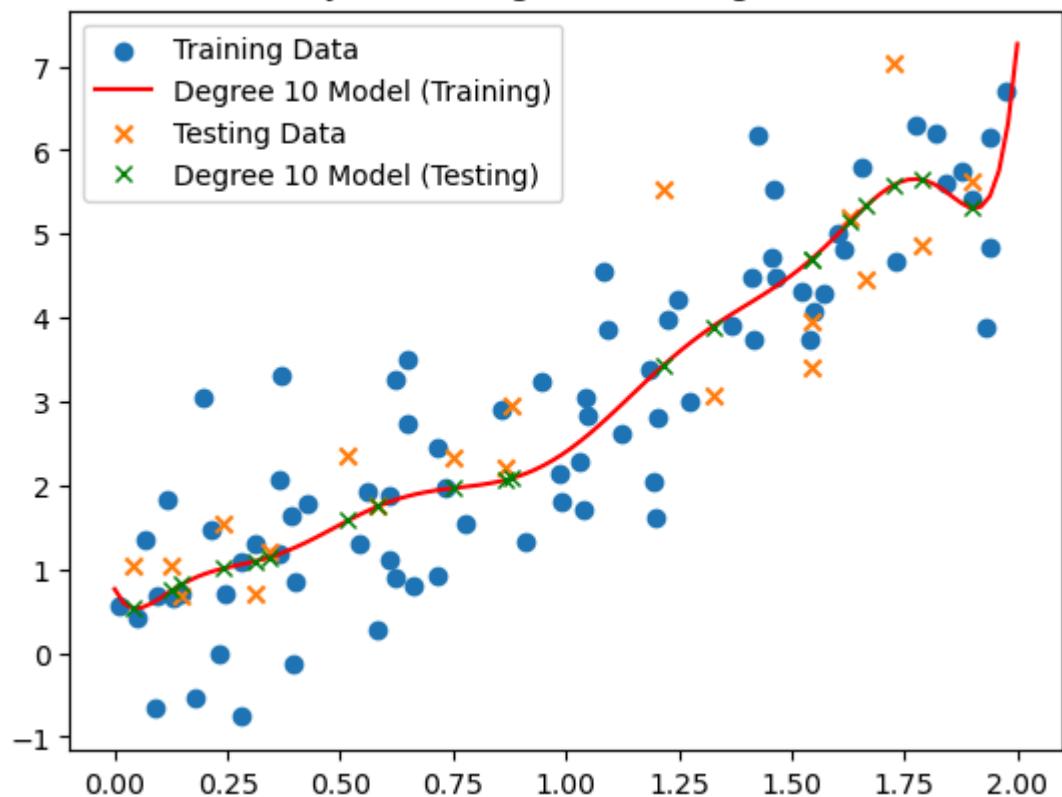
Degree 3 - Training Error: 0.808, Testing Error: 0.642

### Polynomial Regression - Degree 3



Degree 10 - Training Error: 0.770, Testing Error: 0.661

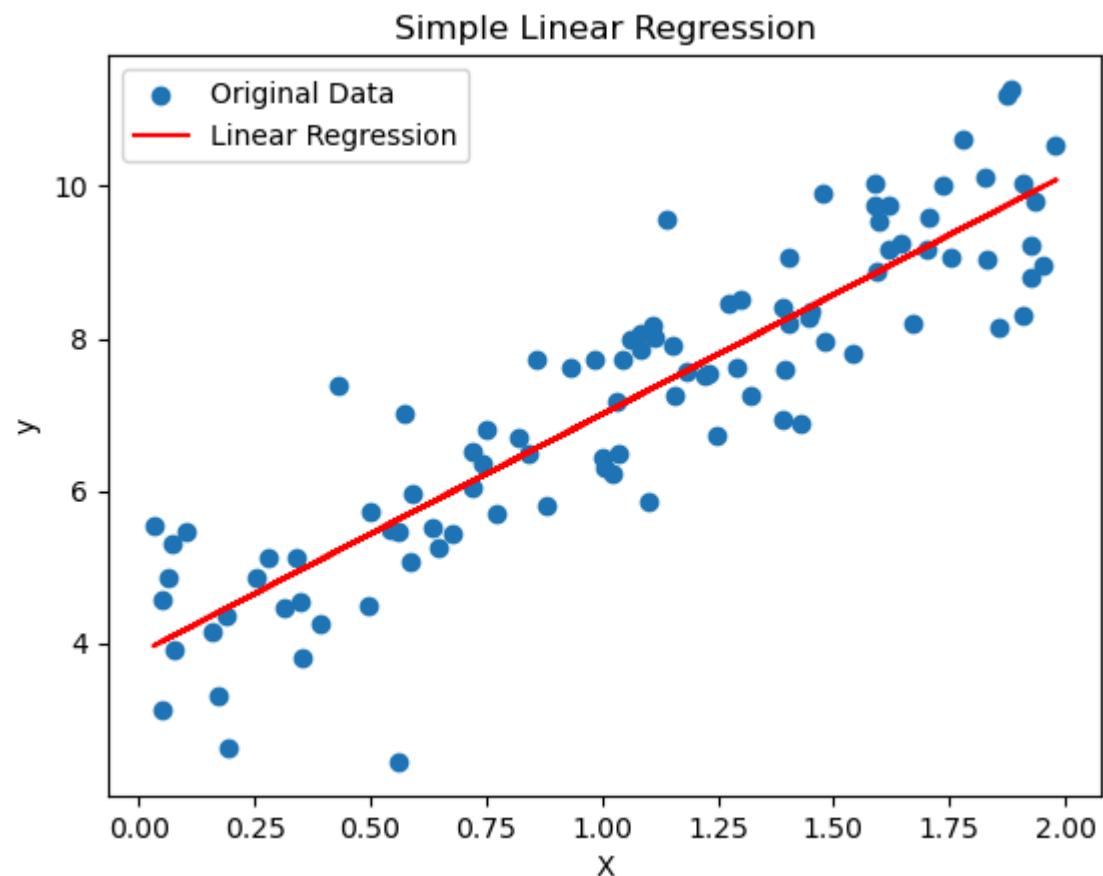
### Polynomial Regression - Degree 10



**simple linear regression model**

In [21]:

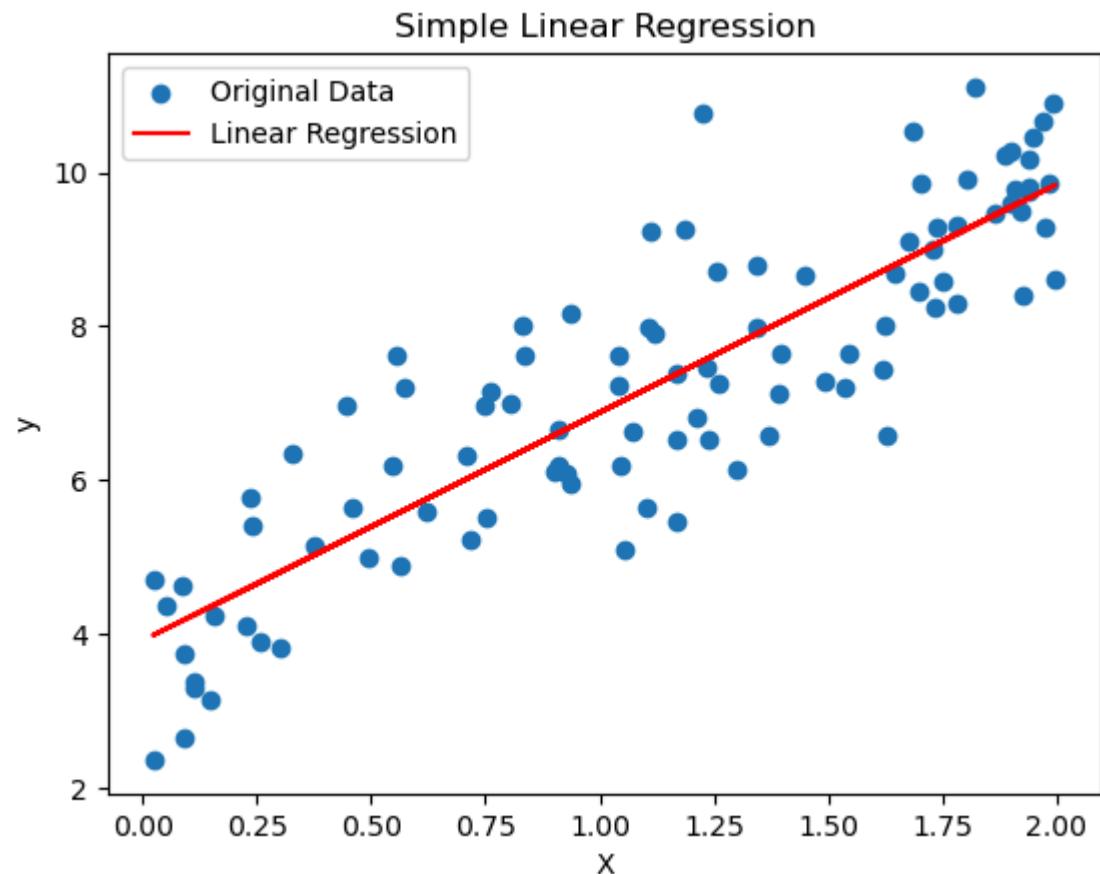
```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.linear_model import LinearRegression
4
5 # Generate sample data
6 X = 2 * np.random.rand(100, 1)
7 y = 4 + 3 * X + np.random.randn(100, 1)
8
9 # Train a simple linear regression model
10 lin_reg = LinearRegression()
11 lin_reg.fit(X, y)
12
13 # Plot the data and the Linear regression line
14 plt.scatter(X, y, label='Original Data')
15 plt.plot(X, lin_reg.predict(X), 'r', label='Linear Regression')
16 plt.xlabel('X')
17 plt.ylabel('y')
18 plt.title('Simple Linear Regression')
19 plt.legend()
20 plt.show()
21
```



## Simple Linear Regression:

In [22]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.linear_model import LinearRegression
4
5 # Generate sample data
6 X = 2 * np.random.rand(100, 1)
7 y = 4 + 3 * X + np.random.randn(100, 1)
8
9 # Train a simple linear regression model
10 lin_reg = LinearRegression()
11 lin_reg.fit(X, y)
12
13 # Plot the data and the Linear regression line
14 plt.scatter(X, y, label='Original Data')
15 plt.plot(X, lin_reg.predict(X), 'r', label='Linear Regression')
16 plt.xlabel('X')
17 plt.ylabel('y')
18 plt.title('Simple Linear Regression')
19 plt.legend()
20 plt.show()
21
```



## Multiple Linear Regression:

In [23]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.linear_model import LinearRegression
4
5 # Generate sample data with multiple features
6 np.random.seed(42)
7 X_mult = 2 * np.random.rand(100, 2)
8 y_mult = 4 + 3 * X_mult[:, 0] + 5 * X_mult[:, 1] + np.random.randn(100)
9
10 # Train a multiple Linear regression model
11 lin_reg_mult = LinearRegression()
12 lin_reg_mult.fit(X_mult, y_mult)
13
14 # Display the coefficients and intercept
15 print("Coefficients:", lin_reg_mult.coef_)
16 print("Intercept:", lin_reg_mult.intercept_)
17
18 # Note: Visualization is not straightforward in multiple dimensions.
19
```

Coefficients: [3.1693339 5.17747302]

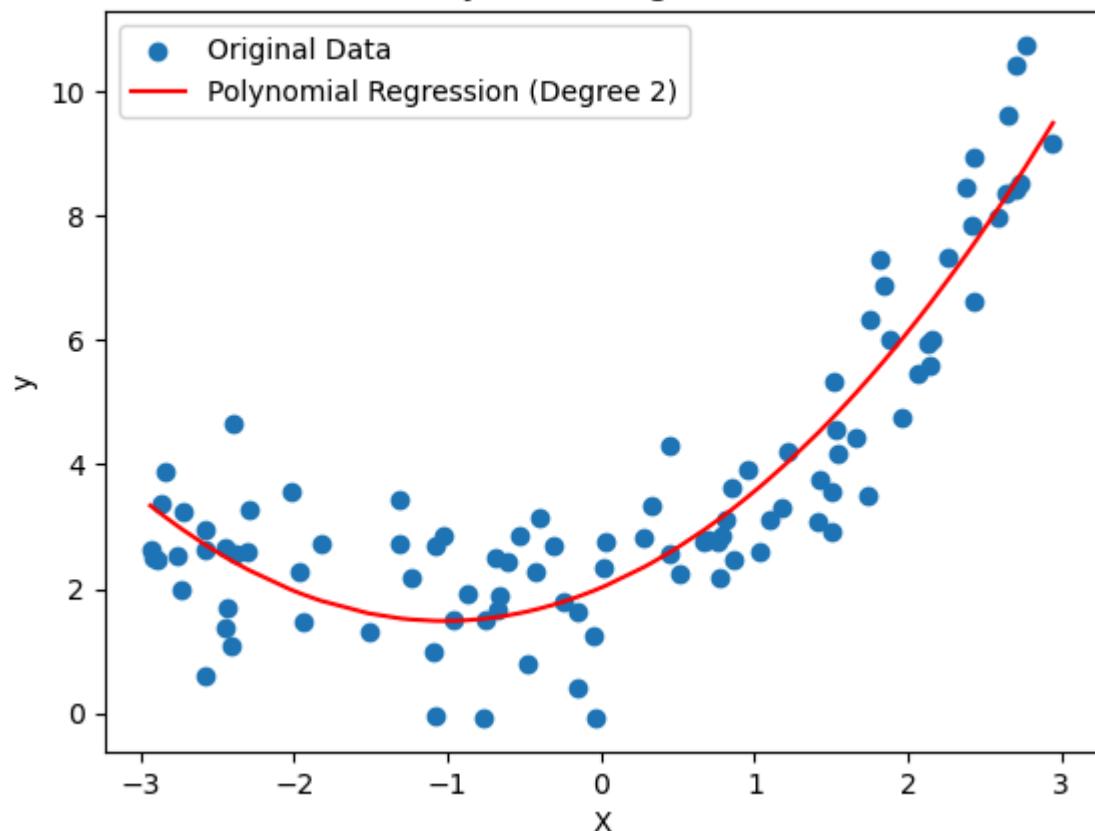
Intercept: 3.7722722641017086

## Polynomial Regression:

In [24]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.preprocessing import PolynomialFeatures
4 from sklearn.linear_model import LinearRegression
5
6 # Generate sample data
7 X_poly = 6 * np.random.rand(100, 1) - 3
8 y_poly = 0.5 * X_poly**2 + X_poly + 2 + np.random.randn(100, 1)
9
10 # Transform features to include polynomial terms up to degree 2
11 poly_features = PolynomialFeatures(degree=2, include_bias=False)
12 X_poly_transformed = poly_features.fit_transform(X_poly)
13
14 # Train a polynomial regression model
15 lin_reg_poly = LinearRegression()
16 lin_reg_poly.fit(X_poly_transformed, y_poly)
17
18 # Sort the data for better visualization
19 X_poly_sorted = np.sort(X_poly, axis=0)
20 y_poly_pred = lin_reg_poly.predict(poly_features.transform(X_poly_sorted))
21
22 # Plot the data and the polynomial regression line
23 plt.scatter(X_poly, y_poly, label='Original Data')
24 plt.plot(X_poly_sorted, y_poly_pred, 'r', label='Polynomial Regression')
25 plt.xlabel('X')
26 plt.ylabel('y')
27 plt.title('Polynomial Regression')
28 plt.legend()
29 plt.show()
30
```

Polynomial Regression



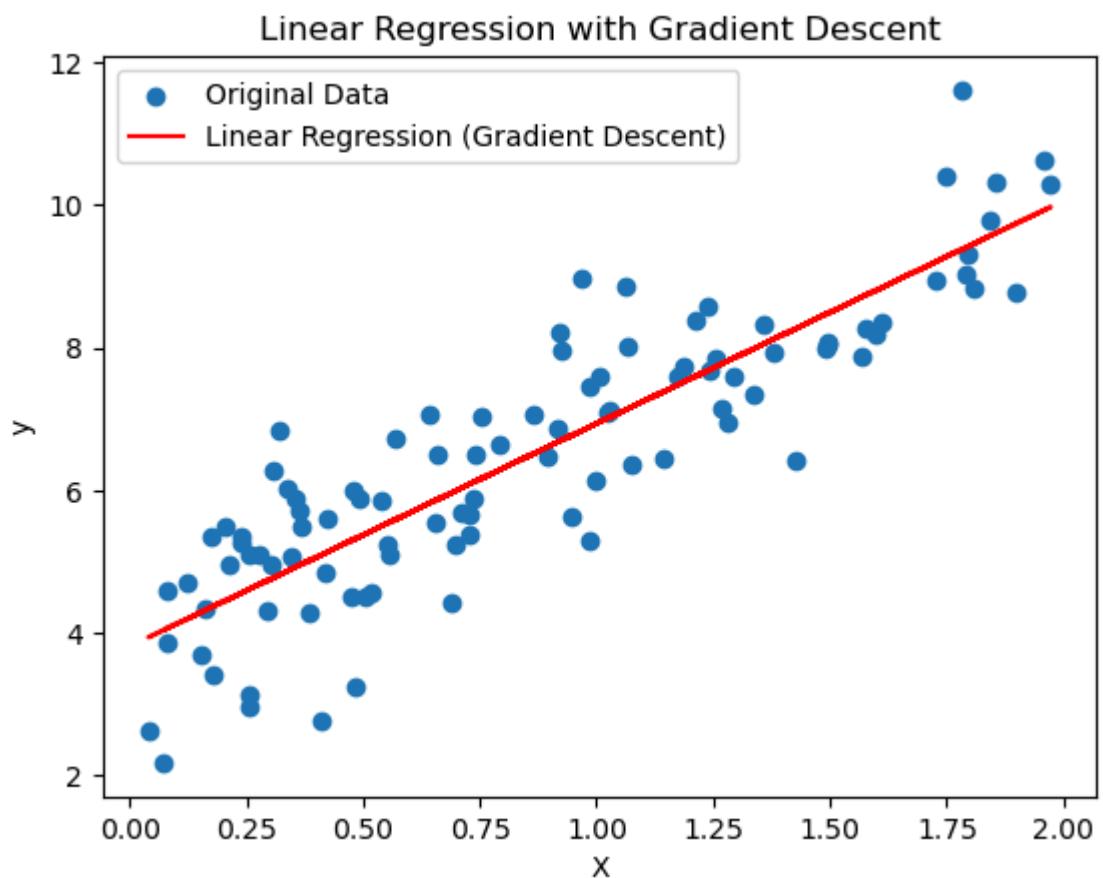
# Linear Regression with Gradient Descent:

In [25]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # Generate sample data
5 X = 2 * np.random.rand(100, 1)
6 y = 4 + 3 * X + np.random.randn(100, 1)
7
8 # Add bias term to features
9 X_b = np.c_[np.ones((100, 1)), X]
10
11 # Set hyperparameters
12 learning_rate = 0.01
13 n_iterations = 1000
14
15 # Initialize coefficients randomly
16 theta = np.random.randn(2, 1)
17
18 # Perform gradient descent
19 for iteration in range(n_iterations):
20     gradients = 2/100 * X_b.T.dot(X_b.dot(theta) - y)
21     theta = theta - learning_rate * gradients
22
23 # Display the final coefficients
24 print("Intercept:", theta[0][0])
25 print("Coefficient:", theta[1][0])
26
27 # Plot the data and the Linear regression Line
28 plt.scatter(X, y, label='Original Data')
29 plt.plot(X, X_b.dot(theta), 'r', label='Linear Regression (Gradient Descent)')
30 plt.xlabel('X')
31 plt.ylabel('y')
32 plt.title('Linear Regression with Gradient Descent')
33 plt.legend()
34 plt.show()
```

Intercept: 3.8131886040073186

Coefficient: 3.12272055734439



**Python code for R-square,RMSE, Gradient descent, optimize regression model**

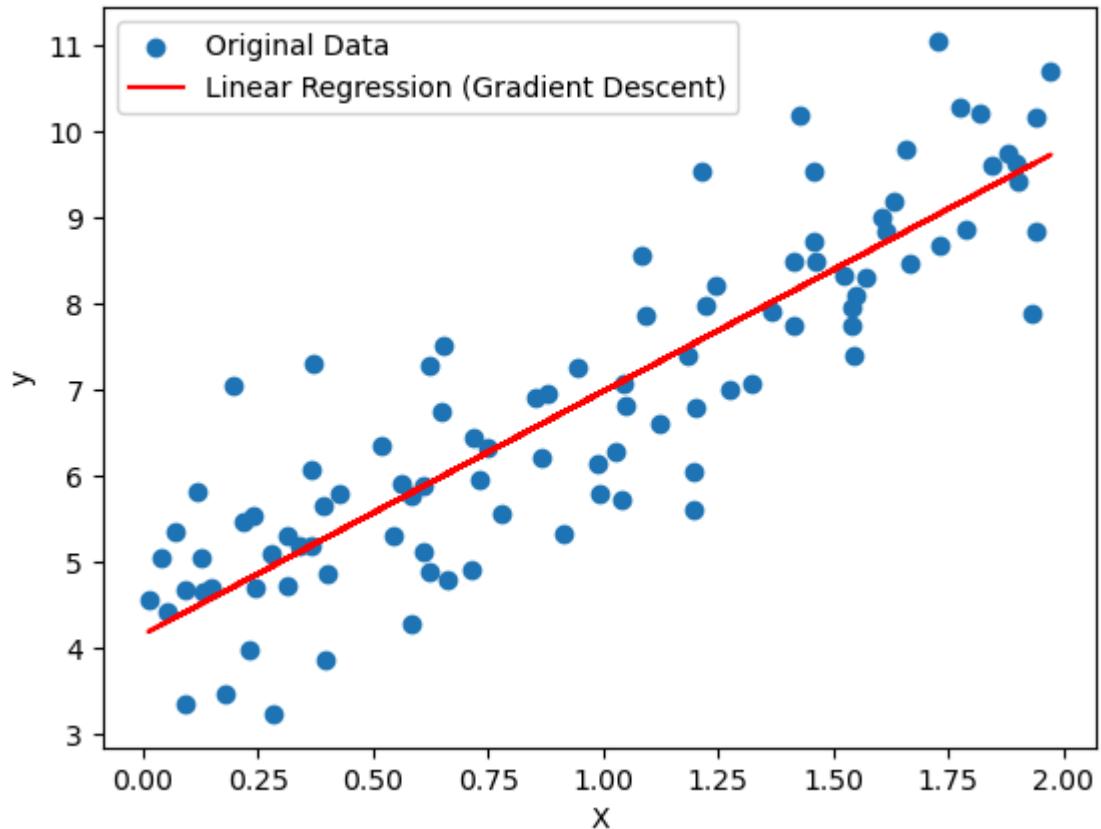
In [26]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.metrics import r2_score, mean_squared_error
4 from sklearn.linear_model import LinearRegression
5
6 # Generate sample data
7 np.random.seed(42)
8 X = 2 * np.random.rand(100, 1)
9 y = 4 + 3 * X + np.random.randn(100, 1)
10
11 # Add bias term to features
12 X_b = np.c_[np.ones((100, 1)), X]
13
14 # Initialize coefficients randomly
15 theta = np.random.randn(2, 1)
16
17 # Set hyperparameters
18 learning_rate = 0.01
19 n_iterations = 1000
20
21 # Perform gradient descent
22 for iteration in range(n_iterations):
23     gradients = 2/100 * X_b.T.dot(X_b.dot(theta) - y)
24     theta = theta - learning_rate * gradients
25
26 # Calculate predictions
27 y_pred = X_b.dot(theta)
28
29 # Calculate R-squared
30 r2 = r2_score(y, y_pred)
31
32 # Calculate Root Mean Squared Error (RMSE)
33 rmse = np.sqrt(mean_squared_error(y, y_pred))
34
35 # Display results
36 print("R-squared:", r2)
37 print("RMSE:", rmse)
38
39 # Plot the data and the optimized Linear regression line
40 plt.scatter(X, y, label='Original Data')
41 plt.plot(X, y_pred, 'r', label='Linear Regression (Gradient Descent)')
42 plt.xlabel('X')
43 plt.ylabel('y')
44 plt.title('Linear Regression with Gradient Descent')
45 plt.legend()
46 plt.show()
```

R-squared: 0.7689928256499341

RMSE: 0.8986467066353333

### Linear Regression with Gradient Descent



## Hold-out Method:

In [27]:

```

1 import numpy as np
2 from sklearn.model_selection import train_test_split
3 from sklearn.linear_model import LinearRegression
4 from sklearn.metrics import mean_squared_error
5
6 # Generate sample data
7 np.random.seed(42)
8 X = 2 * np.random.rand(100, 1)
9 y = 4 + 3 * X + np.random.randn(100, 1)
10
11 # Split the data into training and testing sets using hold-out method
12 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
13
14 # Train a Linear regression model
15 lin_reg = LinearRegression()
16 lin_reg.fit(X_train, y_train)
17
18 # Make predictions on the test set
19 y_pred = lin_reg.predict(X_test)
20
21 # Evaluate the model
22 mse_holdout = mean_squared_error(y_test, y_pred)
23 print("Hold-out Method MSE:", mse_holdout)
24

```

Hold-out Method MSE: 0.6536995137170021

## Leave-One-Out Cross-Validation (LOOCV):

In [28]:

```

1 from sklearn.model_selection import LeaveOneOut
2
3 # Perform Leave-one-out cross-validation
4 loo = LeaveOneOut()
5
6 mse_loocv = 0
7
8 for train_index, test_index in loo.split(X):
9     X_train_loo, X_test_loo = X[train_index], X[test_index]
10    y_train_loo, y_test_loo = y[train_index], y[test_index]
11
12    lin_reg_loo = LinearRegression()
13    lin_reg_loo.fit(X_train_loo, y_train_loo)
14
15    y_pred_loo = lin_reg_loo.predict(X_test_loo)
16    mse_loocv += mean_squared_error(y_test_loo, y_pred_loo)
17
18 mse_loocv /= len(X)
19 print("LOOCV MSE:", mse_loocv)
20

```

LOOCV MSE: 0.8389654996853929

## K-Fold Cross-Validation:

In [29]:

```

1 from sklearn.model_selection import KFold
2
3 # Perform k-fold cross-validation (k=5 in this example)
4 kf = KFold(n_splits=5, shuffle=True, random_state=42)
5
6 mse_kfold = 0
7
8 for train_index, test_index in kf.split(X):
9     X_train_kfold, X_test_kfold = X[train_index], X[test_index]
10    y_train_kfold, y_test_kfold = y[train_index], y[test_index]
11
12    lin_reg_kfold = LinearRegression()
13    lin_reg_kfold.fit(X_train_kfold, y_train_kfold)
14
15    y_pred_kfold = lin_reg_kfold.predict(X_test_kfold)
16    mse_kfold += mean_squared_error(y_test_kfold, y_pred_kfold)
17
18 mse_kfold /= kf.get_n_splits()
19 print("K-Fold Cross-Validation MSE:", mse_kfold)
20

```

K-Fold Cross-Validation MSE: 0.8364353497704167

# A python code on data exploration, preprocessing and splitting on Boston housing price from sci-kit learn

In [30]:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.datasets import load_boston
5 from sklearn.model_selection import train_test_split
6 from sklearn.preprocessing import StandardScaler
7
8 # Load the Boston Housing dataset
9 boston = load_boston()
10
11 # Create a DataFrame for easier exploration
12 boston_df = pd.DataFrame(data=boston.data, columns=boston.feature_names)
13 boston_df['TARGET'] = boston.target
14
15 # Explore the dataset
16 print("Head of the DataFrame:")
17 print(boston_df.head())
18
19 print("\nSummary Statistics:")
20 print(boston_df.describe())
21
22 # Visualize the correlation matrix
23 correlation_matrix = boston_df.corr()
24 plt.figure(figsize=(12, 8))
25 plt.imshow(correlation_matrix, cmap='coolwarm', interpolation='none', aspect='auto')
26 plt.colorbar()
27 plt.xticks(range(len(correlation_matrix)), correlation_matrix.columns)
28 plt.yticks(range(len(correlation_matrix)), correlation_matrix.columns)
29 plt.title("Correlation Matrix")
30 plt.show()
31
32 # Preprocessing: Standardize the features
33 scaler = StandardScaler()
34 X = scaler.fit_transform(boston.data)
35 y = boston.target
36
37 # Split the data into training and testing sets
38 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
39
40 # Display the shape of the training and testing sets
41 print("\nShape of Training Set (X_train, y_train):", X_train.shape, y_train.shape)
42 print("Shape of Testing Set (X_test, y_test):", X_test.shape, y_test.shape)
```

```
C:\Users\Anusha V\anaconda3\lib\site-packages\sklearn\utils\deprecation.py:87: FutureWarning: Function load_boston is deprecated; `load_boston` is deprecated in 1.0 and will be removed in 1.2.
```

The Boston housing prices dataset has an ethical problem. You can refer to the documentation of this function for further details.

The scikit-learn maintainers therefore strongly discourage the use of this dataset unless the purpose of the code is to study and educate about ethical issues in data science and machine learning.

In this special case, you can fetch the dataset from the original source::

```
import pandas as pd  
import numpy as np
```

## A python code on data exploration, preprocessing and splitting on cricket match result-past data from sci-kit learn

In [31]:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.model_selection import train_test_split
5 from sklearn.preprocessing import StandardScaler
6
7 # Sample cricket match dataset (hypothetical)
8 data = {
9     'Team1': ['India', 'Australia', 'England', 'India', 'Australia'],
10    'Team2': ['Australia', 'England', 'India', 'Australia', 'England'],
11    'Winner': ['India', 'Australia', 'India', 'India', 'England'],
12    'Margin': [50, 20, 30, 10, 15],
13    'Toss_Winner': ['India', 'Australia', 'England', 'Australia', 'England'],
14    'Toss_Decision': ['Bat', 'Field', 'Bat', 'Field', 'Field']
15 }
16
17 cricket_df = pd.DataFrame(data)
18
19 # Explore the dataset
20 print("Head of the DataFrame:")
21 print(cricket_df.head())
22
23 print("\nSummary Statistics:")
24 print(cricket_df.describe())
25
26 # Visualize the match margins
27 plt.figure(figsize=(8, 5))
28 plt.hist(cricket_df['Margin'], bins=10, edgecolor='black')
29 plt.title("Distribution of Match Margins")
30 plt.xlabel("Margin")
31 plt.ylabel("Frequency")
32 plt.show()
33
34 # Preprocessing: Convert categorical features to numerical using one-hot encoding
35 cricket_df_encoded = pd.get_dummies(cricket_df, columns=['Team1', 'Team2'])
36
37 # Split the data into features (X) and target variable (y)
38 X = cricket_df_encoded.drop('Winner', axis=1)
39 y = cricket_df_encoded['Winner']
40
41 # Standardize the features
42 scaler = StandardScaler()
43 X_standardized = scaler.fit_transform(X)
44
45 # Split the data into training and testing sets
46 X_train, X_test, y_train, y_test = train_test_split(X_standardized, y,
47
48 # Display the shape of the training and testing sets
49 print("\nShape of Training Set (X_train, y_train):", X_train.shape, y_train.shape)
50 print("Shape of Testing Set (X_test, y_test):", X_test.shape, y_test.shape)
```

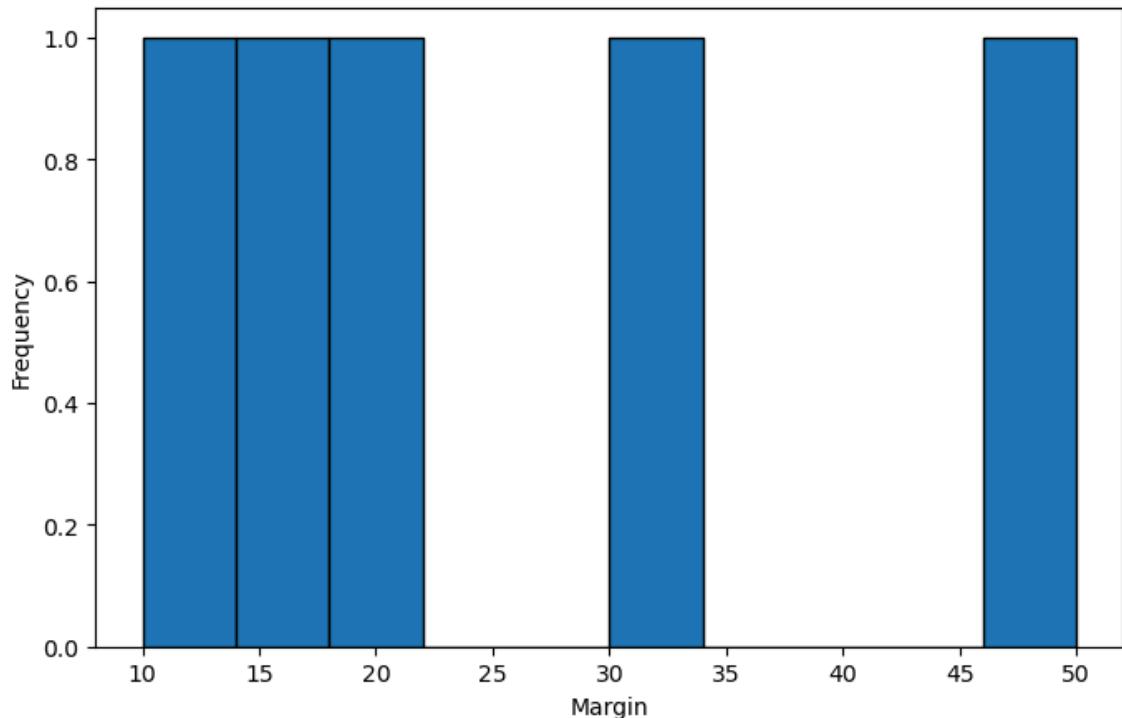
Head of the DataFrame:

	Team1	Team2	Winner	Margin	Toss_Winner	Toss_Decision
0	India	Australia	India	50	India	Bat
1	Australia	England	Australia	20	Australia	Field
2	England	India	India	30	England	Bat
3	India	Australia	India	10	Australia	Field
4	Australia	England	England	15	England	Field

Summary Statistics:

	Margin
count	5.000000
mean	25.000000
std	15.811388
min	10.000000
25%	15.000000
50%	20.000000
75%	30.000000
max	50.000000

Distribution of Match Margins



Shape of Training Set (X\_train, y\_train): (4, 8) (4,)  
 Shape of Testing Set (X\_test, y\_test): (1, 8) (1,)

## A python code on data exploration, preprocessing and splitting on performance of a cricket player from sci-kit learn

In [32]:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.model_selection import train_test_split
5 from sklearn.preprocessing import StandardScaler
6
7 # Sample cricket player performance dataset (hypothetical)
8 data = {
9     'Matches_Played': [50, 48, 52, 45, 51],
10    'Runs_Scored': [1200, 1100, 1300, 1000, 1250],
11    'Batting_Average': [40.0, 38.5, 42.0, 35.0, 41.67],
12    'Strike_Rate': [85.0, 82.5, 87.5, 75.0, 80.0],
13    'Centuries': [5, 4, 6, 3, 5],
14    'Wickets_Taken': [20, 18, 22, 15, 21],
15    'Bowling_Average': [25.0, 27.0, 23.5, 30.0, 24.5],
16    'Player_Rating': [8.5, 8.2, 8.8, 7.9, 8.3]
17 }
18
19 cricket_player_df = pd.DataFrame(data)
20
21 # Explore the dataset
22 print("Head of the DataFrame:")
23 print(cricket_player_df.head())
24
25 print("\nSummary Statistics:")
26 print(cricket_player_df.describe())
27
28 # Visualize Runs Scored vs Batting Average
29 plt.scatter(cricket_player_df['Runs_Scored'], cricket_player_df['Batting_Average'])
30 plt.title("Runs Scored vs Batting Average")
31 plt.xlabel("Runs Scored")
32 plt.ylabel("Batting Average")
33 plt.show()
34
35 # Preprocessing: Standardize the features
36 scaler = StandardScaler()
37 X = scaler.fit_transform(cricket_player_df.drop('Player_Rating', axis=1))
38 y = cricket_player_df['Player_Rating']
39
40 # Split the data into training and testing sets
41 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
42
43 # Display the shape of the training and testing sets
44 print("\nShape of Training Set (X_train, y_train):", X_train.shape, y_train.shape)
45 print("Shape of Testing Set (X_test, y_test):", X_test.shape, y_test.shape)
```

Head of the DataFrame:

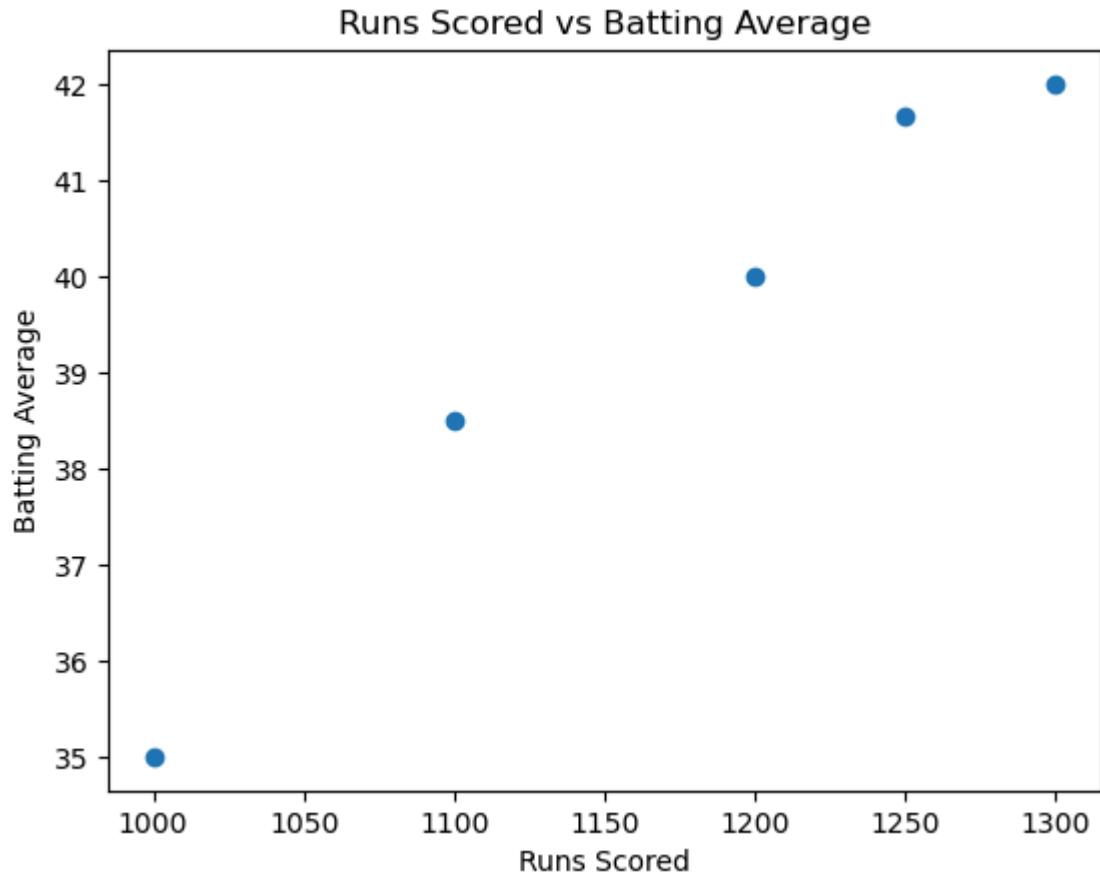
	Matches_Played	Runs_Scored	Batting_Average	Strike_Rate	Centuries	\
0	50	1200	40.00	85.0	5	
1	48	1100	38.50	82.5	4	
2	52	1300	42.00	87.5	6	
3	45	1000	35.00	75.0	3	
4	51	1250	41.67	80.0	5	

	Wickets_Taken	Bowling_Average	Player_Rating
0	20	25.0	8.5
1	18	27.0	8.2
2	22	23.5	8.8
3	15	30.0	7.9
4	21	24.5	8.3

Summary Statistics:

	Matches_Played	Runs_Scored	Batting_Average	Strike_Rate	Centurie
count	5.000000	5.000000	5.000000	5.000000	5.000000
0					
mean	49.200000	1170.000000	39.43400	82.000000	4.60000
0					
std	2.774887	120.415946	2.84768	4.808846	1.14017
5					
min	45.000000	1000.000000	35.00000	75.000000	3.00000
0					
25%	48.000000	1100.000000	38.50000	80.000000	4.00000
0					
50%	50.000000	1200.000000	40.00000	82.500000	5.00000
0					
75%	51.000000	1250.000000	41.67000	85.000000	5.00000
0					
max	52.000000	1300.000000	42.00000	87.500000	6.00000
0					

	Wickets_Taken	Bowling_Average	Player_Rating
count	5.000000	5.000000	5.000000
mean	19.200000	26.000000	8.340000
std	2.774887	2.573908	0.336155
min	15.000000	23.500000	7.900000
25%	18.000000	24.500000	8.200000
50%	20.000000	25.000000	8.300000
75%	21.000000	27.000000	8.500000
max	22.000000	30.000000	8.800000



Shape of Training Set (X\_train, y\_train): (4, 7) (4,)

Shape of Testing Set (X\_test, y\_test): (1, 7) (1,)

## A python code on data exploration, preprocessing and splitting on crop yield from sci-kit learn

In [33]:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.model_selection import train_test_split
5 from sklearn.preprocessing import StandardScaler
6
7 # Sample crop yield dataset (hypothetical)
8 data = {
9     'Temperature': [30, 32, 28, 35, 33],
10    'Rainfall': [100, 80, 120, 90, 110],
11    'Fertilizer_Use': [20, 25, 18, 30, 22],
12    'Crop_Yield': [3000, 3200, 2800, 3500, 3300]
13 }
14
15 crop_yield_df = pd.DataFrame(data)
16
17 # Explore the dataset
18 print("Head of the DataFrame:")
19 print(crop_yield_df.head())
20
21 print("\nSummary Statistics:")
22 print(crop_yield_df.describe())
23
24 # Visualize the relationship between Temperature and Crop Yield
25 plt.scatter(crop_yield_df['Temperature'], crop_yield_df['Crop_Yield'])
26 plt.title("Temperature vs Crop Yield")
27 plt.xlabel("Temperature (°C)")
28 plt.ylabel("Crop Yield (kg/ha)")
29 plt.show()
30
31 # Preprocessing: Standardize the features
32 scaler = StandardScaler()
33 X = scaler.fit_transform(crop_yield_df.drop('Crop_Yield', axis=1))
34 y = crop_yield_df['Crop_Yield']
35
36 # Split the data into training and testing sets
37 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
38
39 # Display the shape of the training and testing sets
40 print("\nShape of Training Set (X_train, y_train):", X_train.shape, y_train.shape)
41 print("Shape of Testing Set (X_test, y_test):", X_test.shape, y_test.shape)
```

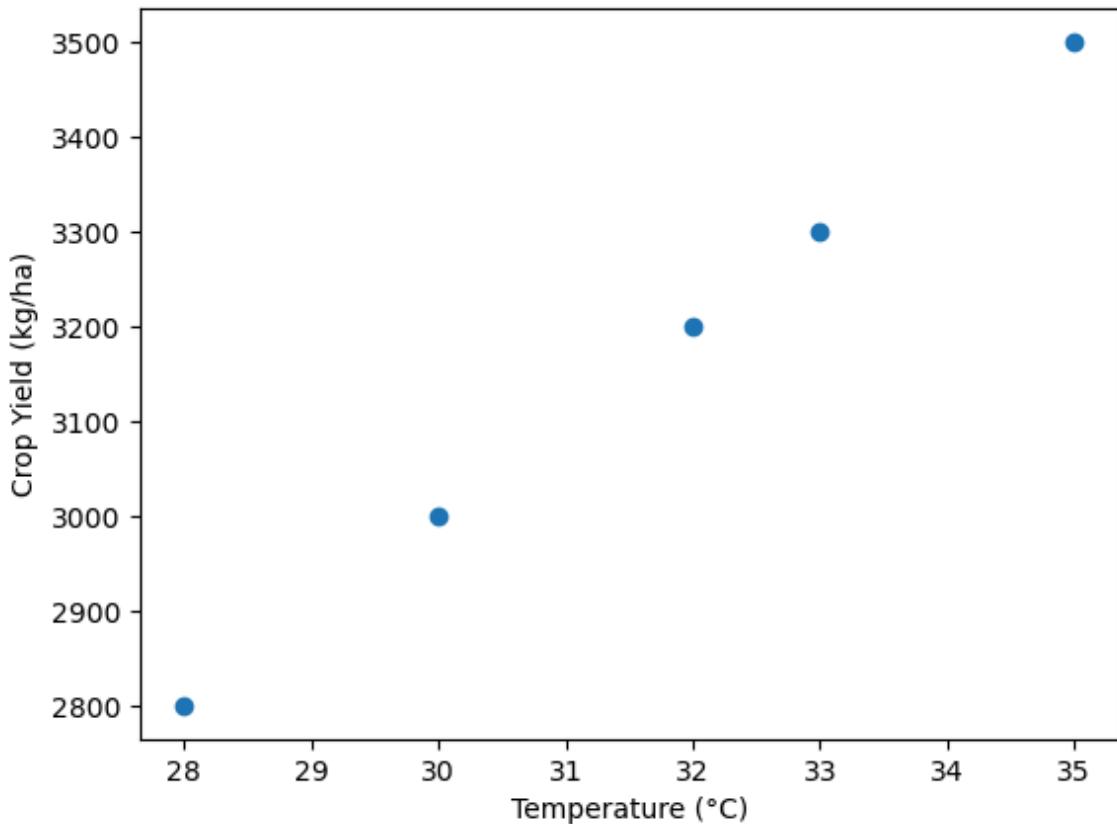
Head of the DataFrame:

	Temperature	Rainfall	Fertilizer_Use	Crop_Yield
0	30	100	20	3000
1	32	80	25	3200
2	28	120	18	2800
3	35	90	30	3500
4	33	110	22	3300

Summary Statistics:

	Temperature	Rainfall	Fertilizer_Use	Crop_Yield
count	5.000000	5.000000	5.000000	5.000000
mean	31.600000	100.000000	23.000000	3160.000000
std	2.701851	15.811388	4.690416	270.185122
min	28.000000	80.000000	18.000000	2800.000000
25%	30.000000	90.000000	20.000000	3000.000000
50%	32.000000	100.000000	22.000000	3200.000000
75%	33.000000	110.000000	25.000000	3300.000000
max	35.000000	120.000000	30.000000	3500.000000

Temperature vs Crop Yield



Shape of Training Set (X\_train, y\_train): (4, 3) (4, )

Shape of Testing Set (X\_test, y\_test): (1, 3) (1, )

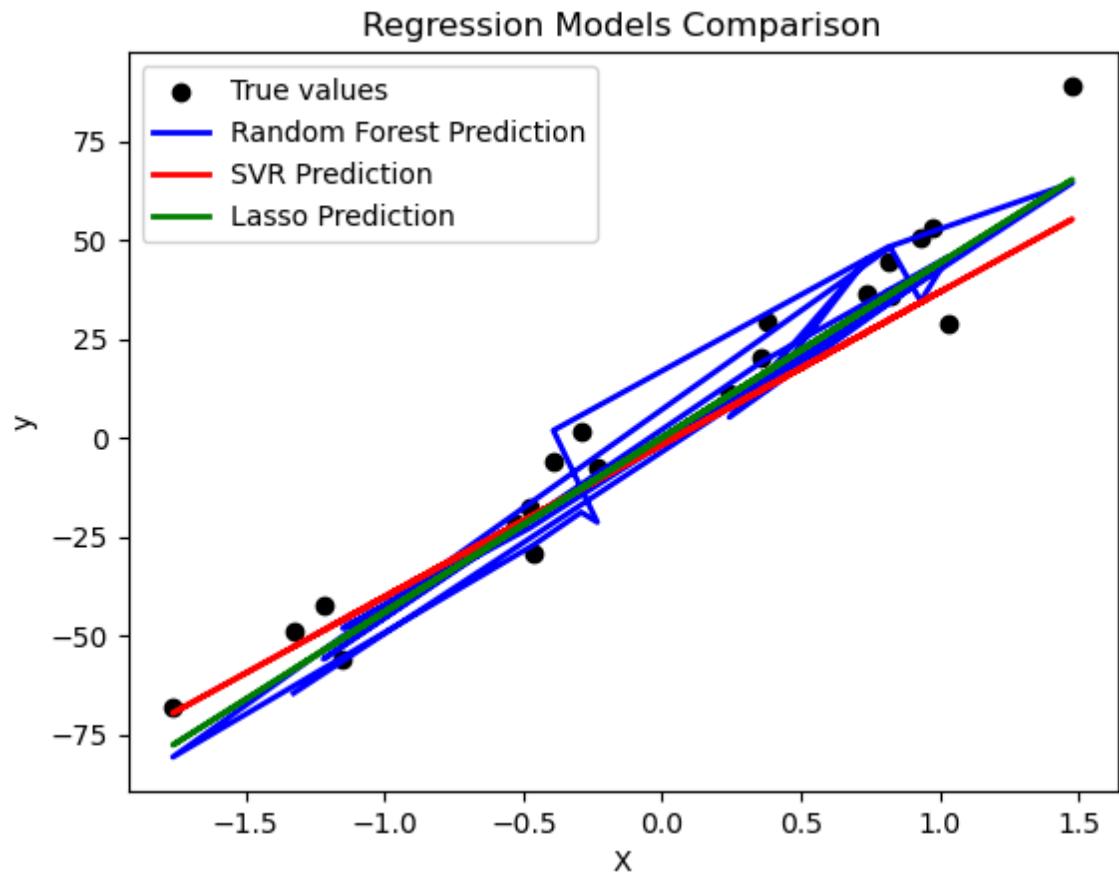


In [35]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.datasets import make_regression
4 from sklearn.model_selection import train_test_split
5 from sklearn.ensemble import RandomForestRegressor
6 from sklearn.svm import SVR
7 from sklearn.linear_model import Lasso
8 from sklearn.metrics import mean_squared_error
9 from sklearn.preprocessing import StandardScaler
10
11 # Generate a synthetic dataset
12 X, y = make_regression(n_samples=100, n_features=1, noise=10, random_state=42)
13
14 # Split the dataset into training and testing sets
15 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
16
17 # Standardize features
18 scaler = StandardScaler()
19 X_train_scaled = scaler.fit_transform(X_train)
20 X_test_scaled = scaler.transform(X_test)
21
22 # Random Forest Regression
23 rf_regressor = RandomForestRegressor(n_estimators=100, random_state=42)
24 rf_regressor.fit(X_train, y_train)
25
26 # Support Vector Regression (SVR)
27 svr_regressor = SVR(kernel='linear')
28 svr_regressor.fit(X_train_scaled, y_train)
29
30 # Lasso Regression
31 lasso_regressor = Lasso(alpha=0.1, random_state=42)
32 lasso_regressor.fit(X_train_scaled, y_train)
33
34 # Predictions
35 y_pred_rf = rf_regressor.predict(X_test)
36 y_pred_svr = svr_regressor.predict(X_test_scaled)
37 y_pred_lasso = lasso_regressor.predict(X_test_scaled)
38
39 # Evaluate the models
40 mse_rf = mean_squared_error(y_test, y_pred_rf)
41 mse_svr = mean_squared_error(y_test, y_pred_svr)
42 mse_lasso = mean_squared_error(y_test, y_pred_lasso)
43
44 print("Mean Squared Error (Random Forest):", mse_rf)
45 print("Mean Squared Error (SVR):", mse_svr)
46 print("Mean Squared Error (Lasso):", mse_lasso)
47
48 # Plot the results
49 plt.scatter(X_test, y_test, color='black', label='True values')
50 plt.plot(X_test, y_pred_rf, color='blue', linewidth=2, label='Random Forest')
51 plt.plot(X_test, y_pred_svr, color='red', linewidth=2, label='SVR Prediction')
52 plt.plot(X_test, y_pred_lasso, color='green', linewidth=2, label='Lasso')
53
54 plt.title('Regression Models Comparison')
55 plt.xlabel('X')
56 plt.ylabel('y')
57 plt.legend()
58 plt.show()
```

59

Mean Squared Error (Random Forest): 160.99855696169863  
Mean Squared Error (SVR): 154.7063747913098  
Mean Squared Error (Lasso): 104.41512794678081



**Python code for binary classification,multi-label classification,multi-class classification and imbalanced classification**



In [36]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.datasets import make_classification
4 from sklearn.model_selection import train_test_split
5 from sklearn.linear_model import LogisticRegression
6 from sklearn.multiclass import OneVsRestClassifier
7 from sklearn.svm import SVC
8 from sklearn.ensemble import RandomForestClassifier
9 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
10 from imblearn.over_sampling import SMOTE
11
12 # Generate synthetic datasets for binary and multiclass classification
13 X_bin, y_bin = make_classification(n_samples=1000, n_features=20, n_informative=15, n_redundant=5, n_clusters_per_class=1, random_state=42)
14 X_multi, y_multi = make_classification(n_samples=1000, n_features=20, n_informative=15, n_redundant=5, n_clusters_per_class=3, random_state=42)
15
16 # Binary Classification
17 X_bin_train, X_bin_test, y_bin_train, y_bin_test = train_test_split(X_bin, y_bin, test_size=0.2, random_state=42)
18
19 # Train a Logistic Regression model for binary classification
20 logreg_bin = LogisticRegression()
21 logreg_bin.fit(X_bin_train, y_bin_train)
22 y_bin_pred = logreg_bin.predict(X_bin_test)
23
24 # Evaluate the model for binary classification
25 print("Binary Classification:")
26 print("Accuracy:", accuracy_score(y_bin_test, y_bin_pred))
27 print("Classification Report:\n", classification_report(y_bin_test, y_bin_pred))
28 print("Confusion Matrix:\n", confusion_matrix(y_bin_test, y_bin_pred))
29
30 # Multiclass Classification
31 X_multi_train, X_multi_test, y_multi_train, y_multi_test = train_test_split(X_multi, y_multi, test_size=0.2, random_state=42)
32
33 # Train a Logistic Regression model for multiclass classification
34 logreg_multi = LogisticRegression(multi_class='ovr') # One-vs-Rest strategy
35 logreg_multi.fit(X_multi_train, y_multi_train)
36 y_multi_pred = logreg_multi.predict(X_multi_test)
37
38 # Evaluate the model for multiclass classification
39 print("\nMulticlass Classification:")
40 print("Accuracy:", accuracy_score(y_multi_test, y_multi_pred))
41 print("Classification Report:\n", classification_report(y_multi_test, y_multi_pred))
42 print("Confusion Matrix:\n", confusion_matrix(y_multi_test, y_multi_pred))
43
44 # Imbalanced Classification (using SMOTE for oversampling)
45 X_imbalanced, y_imbalanced = make_classification(n_samples=1000, n_features=20, n_informative=15, n_redundant=5, n_clusters_per_class=1, weights=[0.9, 0.1], random_state=42)
46
47 # Apply SMOTE to handle class imbalance
48 smote = SMOTE(sampling_strategy='auto', random_state=42)
49 X_resampled, y_resampled = smote.fit_resample(X_imbalanced, y_imbalanced)
50
51 # Split the resampled data into training and testing sets
52 X_resampled_train, X_resampled_test, y_resampled_train, y_resampled_test = train_test_split(X_resampled, y_resampled, test_size=0.2, random_state=42)
53
54 # Train a Random Forest model for imbalanced classification
55 rf_imbalanced = RandomForestClassifier(random_state=42)
56 rf_imbalanced.fit(X_resampled_train, y_resampled_train)
57 y_resampled_pred = rf_imbalanced.predict(X_resampled_test)
58
59 # Evaluate the model for imbalanced classification
60 print("Accuracy:", accuracy_score(y_resampled_test, y_resampled_pred))
61 print("Classification Report:\n", classification_report(y_resampled_test, y_resampled_pred))
62 print("Confusion Matrix:\n", confusion_matrix(y_resampled_test, y_resampled_pred))
```

```
62 print("\nImbalanced Classification:")
63 print("Accuracy:", accuracy_score(y_resampled_test, y_resampled_pred))
64 print("Classification Report:\n", classification_report(y_resampled_test, y_resampled_pred))
65 print("Confusion Matrix:\n", confusion_matrix(y_resampled_test, y_resampled_pred))
66
67 # Plot the original and resampled class distribution
68 plt.figure(figsize=(10, 4))
69
70 plt.subplot(1, 2, 1)
71 plt.title("Original Class Distribution (Imbalanced)")
72 plt.hist(y_imbalanced, bins=[0, 1, 2], align='left')
73 plt.xlabel("Class")
74 plt.ylabel("Count")
75
76 plt.subplot(1, 2, 2)
77 plt.title("Resampled Class Distribution (Balanced)")
78 plt.hist(y_resampled, bins=[0, 1, 2], align='left')
79 plt.xlabel("Class")
80 plt.ylabel("Count")
81
82 plt.tight_layout()
83 plt.show()
```



**Binary Classification:**

Accuracy: 0.81

**Classification Report:**

	precision	recall	f1-score	support
0	0.88	0.77	0.82	112
1	0.75	0.86	0.80	88
accuracy			0.81	200
macro avg	0.81	0.82	0.81	200
weighted avg	0.82	0.81	0.81	200

**Confusion Matrix:**

```
[[86 26]
 [12 76]]
```

**Multiclass Classification:**

Accuracy: 0.6

**Classification Report:**

	precision	recall	f1-score	support
0	0.64	0.55	0.59	66
1	0.55	0.65	0.60	63
2	0.61	0.61	0.61	71
accuracy			0.60	200
macro avg	0.60	0.60	0.60	200
weighted avg	0.60	0.60	0.60	200

**Confusion Matrix:**

```
[[36 15 15]
 [10 41 12]
 [10 18 43]]
```

**Imbalanced Classification:**

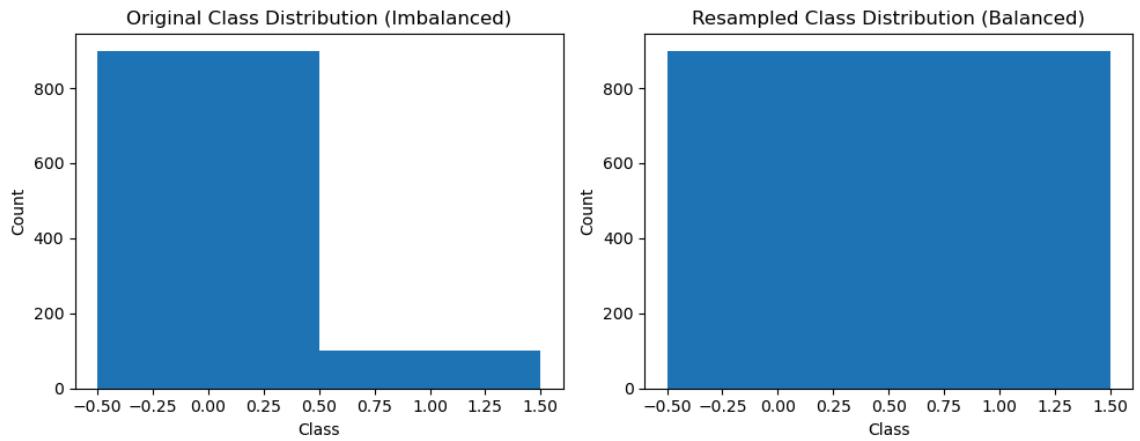
Accuracy: 0.9777777777777777

**Classification Report:**

	precision	recall	f1-score	support
0	0.99	0.97	0.98	179
1	0.97	0.99	0.98	181
accuracy			0.98	360
macro avg	0.98	0.98	0.98	360
weighted avg	0.98	0.98	0.98	360

**Confusion Matrix:**

```
[[173  6]
 [ 2 179]]
```



**Python code for evaluation matrices for classification, confusion matrix, accuracy, precession and recall, specificity,f1-score,AUC-ROC**

```
In [41]: 1 pip install scikit-learn matplotlib  
2
```

```
Requirement already satisfied: scikit-learn in c:\users\anusha v\anaconda3  
\lib\site-packages (1.0.2)  
Requirement already satisfied: matplotlib in c:\users\anusha v\anaconda3\l  
ib\site-packages (3.5.2)  
Requirement already satisfied: scipy>=1.1.0 in c:\users\anusha v\anaconda3  
\lib\site-packages (from scikit-learn) (1.9.1)  
Requirement already satisfied: threadpoolctl>=2.0.0 in c:\users\anusha v\ana  
conda3\lib\site-packages (from scikit-learn) (2.2.0)  
Requirement already satisfied: numpy>=1.14.6 in c:\users\anusha v\anaconda  
3\lib\site-packages (from scikit-learn) (1.24.4)  
Requirement already satisfied: joblib>=0.11 in c:\users\anusha v\anaconda3  
\lib\site-packages (from scikit-learn) (1.3.2)  
Requirement already satisfied: python-dateutil>=2.7 in c:\users\anusha v\ana  
conda3\lib\site-packages (from matplotlib) (2.8.2)  
Requirement already satisfied: pillow>=6.2.0 in c:\users\anusha v\anaconda  
3\lib\site-packages (from matplotlib) (9.2.0)  
Requirement already satisfied: kiwisolver>=1.0.1 in c:\users\anusha v\ana  
conda3\lib\site-packages (from matplotlib) (1.4.2)  
Requirement already satisfied: cycler>=0.10 in c:\users\anusha v\anaconda3  
\lib\site-packages (from matplotlib) (0.11.0)  
Requirement already satisfied: fonttools>=4.22.0 in c:\users\anusha v\ana  
conda3\lib\site-packages (from matplotlib) (4.25.0)  
Requirement already satisfied: pyparsing>=2.2.1 in c:\users\anusha v\ana  
conda3\lib\site-packages (from matplotlib) (3.0.9)  
Requirement already satisfied: packaging>=20.0 in c:\users\anusha v\anacon  
da3\lib\site-packages (from matplotlib) (21.3)  
Requirement already satisfied: six>=1.5 in c:\users\anusha v\anaconda3\lib  
\site-packages (from python-dateutil>=2.7->matplotlib) (1.16.0)  
Note: you may need to restart the kernel to use updated packages.
```

```
WARNING: Ignoring invalid distribution -umpy (c:\users\anusha v\anaconda3  
\lib\site-packages)  
WARNING: Ignoring invalid distribution -umpy (c:\users\anusha v\anaconda3  
\lib\site-packages)
```

In [27]:

```

1 from sklearn.model_selection import train_test_split
2 from sklearn.datasets import make_classification
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn.metrics import confusion_matrix, accuracy_score, precision_
5
6 # Generate synthetic dataset for demonstration
7 X, y = make_classification(n_samples=1000, n_features=20, n_classes=2,
8
9 # Split the dataset into training and testing sets
10 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2
11
12 # Initialize and train a classification model (Random Forest in this ex
13 model = RandomForestClassifier(n_estimators=100, random_state=42)
14 model.fit(X_train, y_train)
15
16 # Make predictions on the test set
17 y_pred = model.predict(X_test)
18
19 # Confusion matrix
20 conf_matrix = confusion_matrix(y_test, y_pred)
21 print("Confusion Matrix:")
22 print(conf_matrix)
23
24 # Accuracy
25 accuracy = accuracy_score(y_test, y_pred)
26 print(f"Accuracy: {accuracy:.4f}")
27
28 # Precision
29 precision = precision_score(y_test, y_pred)
30 print(f"Precision: {precision:.4f}")
31
32 # Recall
33 recall = recall_score(y_test, y_pred)
34 print(f"Recall: {recall:.4f}")
35
36 # Specificity (True Negative Rate)
37 specificity = conf_matrix[0, 0] / (conf_matrix[0, 0] + conf_matrix[0, 1]
38 print(f"Specificity: {specificity:.4f}")
39
40 # F1-score
41 f1 = f1_score(y_test, y_pred)
42 print(f"F1-score: {f1:.4f}")
43
44 # AUC-ROC
45 y_pred_prob = model.predict_proba(X_test)[:, 1]
46 roc_auc = roc_auc_score(y_test, y_pred_prob)
47 print(f"AUC-ROC: {roc_auc:.4f}")
48

```

Confusion Matrix:

```

[[88  5]
 [15 92]]

```

Accuracy: 0.9000

Precision: 0.9485

Recall: 0.8598

Specificity: 0.9462

F1-score: 0.9020

AUC-ROC: 0.9379

# Python code for Hyper parameter tuning for decision tree classifier

In [39]:

```

1 import numpy as np
2 from sklearn.datasets import load_iris
3 from sklearn.model_selection import train_test_split, GridSearchCV
4 from sklearn.tree import DecisionTreeClassifier
5 from sklearn.metrics import accuracy_score
6
7 # Load Iris dataset
8 iris = load_iris()
9 X, y = iris.data, iris.target
10
11 # Split the dataset into training and testing sets
12 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
13
14 # Create a Decision Tree Classifier
15 dt_classifier = DecisionTreeClassifier()
16
17 # Define the hyperparameters to tune
18 param_grid = {
19     'criterion': ['gini', 'entropy'],
20     'max_depth': [None, 5, 10, 15],
21     'min_samples_split': [2, 5, 10],
22     'min_samples_leaf': [1, 2, 4]
23 }
24
25 # Use GridSearchCV to perform hyperparameter tuning
26 grid_search = GridSearchCV(estimator=dt_classifier, param_grid=param_gr
27 grid_search.fit(X_train, y_train)
28
29 # Display the best hyperparameters
30 best_params = grid_search.best_params_
31 print("Best Hyperparameters:")
32 print(best_params)
33
34 # Evaluate the model with the best hyperparameters on the test set
35 best_dt_model = grid_search.best_estimator_
36 y_pred = best_dt_model.predict(X_test)
37 accuracy = accuracy_score(y_test, y_pred)
38 print("\nAccuracy on Test Set:", accuracy)
39

```

Best Hyperparameters:

```
{'criterion': 'entropy', 'max_depth': None, 'min_samples_leaf': 4, 'min_samples_split': 2}
```

Accuracy on Test Set: 1.0

# Logistic regression model in python

In [42]:

```

1 import numpy as np
2 import pandas as pd
3 from sklearn.datasets import load_iris
4 from sklearn.model_selection import train_test_split
5 from sklearn.linear_model import LogisticRegression
6 from sklearn.metrics import accuracy_score, classification_report, conf
7
8 # Load the Iris dataset
9 iris = load_iris()
10 X, y = iris.data, iris.target
11
12 # Split the dataset into training and testing sets
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
14
15 # Create a Logistic Regression model
16 logreg_model = LogisticRegression()
17
18 # Train the model on the training set
19 logreg_model.fit(X_train, y_train)
20
21 # Make predictions on the test set
22 y_pred = logreg_model.predict(X_test)
23
24 # Evaluate the model
25 accuracy = accuracy_score(y_test, y_pred)
26 conf_matrix = confusion_matrix(y_test, y_pred)
27 class_report = classification_report(y_test, y_pred)
28
29 # Display the results
30 print("Logistic Regression Model Evaluation:")
31 print("Accuracy:", accuracy)
32 print("\nConfusion Matrix:")
33 print(conf_matrix)
34 print("\nClassification Report:")
35 print(class_report)
36

```

Logistic Regression Model Evaluation:

Accuracy: 1.0

Confusion Matrix:

```

[[10  0  0]
 [ 0  9  0]
 [ 0  0 11]]

```

Classification Report:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	10
1	1.00	1.00	1.00	9
2	1.00	1.00	1.00	11
accuracy			1.00	30
macro avg	1.00	1.00	1.00	30
weighted avg	1.00	1.00	1.00	30

```
C:\Users\Anusha V\anaconda3\lib\site-packages\sklearn\linear_model\_logistic.py:814: ConvergenceWarning: lbgfs failed to converge (status=1):  
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (`max_iter`) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html> (<https://scikit-learn.org/stable/modules/preprocessing.html>)

Please also refer to the documentation for alternative solver options:

[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression) ([https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression))

```
n_iter_i = _check_optimize_result(
```

## Python code for basic ensemble technique- Max voting, Averaging,Weighted average



In [43]:

```
1 import numpy as np
2 from sklearn.datasets import load_iris
3 from sklearn.model_selection import train_test_split
4 from sklearn.ensemble import RandomForestClassifier
5 from sklearn.linear_model import LogisticRegression
6 from sklearn.svm import SVC
7 from sklearn.metrics import accuracy_score
8
9 # Load the Iris dataset
10 iris = load_iris()
11 X, y = iris.data, iris.target
12
13 # Split the dataset into training and testing sets
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2
15
16 # Create three different models
17 model_rf = RandomForestClassifier(random_state=42)
18 model_lr = LogisticRegression(random_state=42)
19 model_svm = SVC(probability=True, random_state=42) # SVC requires prob
20
21 # Train the models
22 model_rf.fit(X_train, y_train)
23 model_lr.fit(X_train, y_train)
24 model_svm.fit(X_train, y_train)
25
26 # Make predictions on the test set
27 y_pred_rf = model_rf.predict(X_test)
28 y_pred_lr = model_lr.predict(X_test)
29 y_pred_svm = model_svm.predict(X_test)
30
31 # Ensembling techniques:
32
33 # Max Voting (Hard Voting)
34 y_pred_max_voting = np.array([np.argmax(np.bincount([pred_rf, pred_lr,
35
36 # Averaging
37 y_pred_averaging = (y_pred_rf + y_pred_lr + y_pred_svm) // 3
38
39 # Weighted Average
40 weights = {'rf': 2, 'lr': 1, 'svm': 1} # Adjust weights based on model
41 y_pred_weighted_average = (weights['rf'] * y_pred_rf + weights['lr'] *
42
43 # Evaluate individual models
44 acc_rf = accuracy_score(y_test, y_pred_rf)
45 acc_lr = accuracy_score(y_test, y_pred_lr)
46 acc_svm = accuracy_score(y_test, y_pred_svm)
47
48 # Evaluate ensemble methods
49 acc_max_voting = accuracy_score(y_test, y_pred_max_voting)
50 acc_averaging = accuracy_score(y_test, y_pred_averaging)
51 acc_weighted_average = accuracy_score(y_test, y_pred_weighted_average)
52
53 # Display results
54 print("Accuracy of Individual Models:")
55 print("Random Forest:", acc_rf)
56 print("Logistic Regression:", acc_lr)
57 print("Support Vector Machine:", acc_svm)
58
59 print("\nAccuracy of Ensemble Techniques:")
60 print("Max Voting (Hard Voting):", acc_max_voting)
61 print("Averaging:", acc_averaging)
```

```
62 | print("Weighted Average:", acc_weighted_average)
63 |
```

Accuracy of Individual Models:  
Random Forest: 1.0  
Logistic Regression: 1.0  
Support Vector Machine: 1.0

Accuracy of Ensemble Techniques:  
Max Voting (Hard Voting): 1.0  
Averaging: 1.0  
Weighted Average: 1.0

C:\Users\Anusha V\anaconda3\lib\site-packages\sklearn\linear\_model\\_logistic.py:814: ConvergenceWarning: lbfgs failed to converge (status=1):  
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max\_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html> (<https://scikit-learn.org/stable/modules/preprocessing.html>)

Please also refer to the documentation for alternative solver options:

[https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression) ([https://scikit-learn.org/stable/modules/linear\\_model.html#logistic-regression](https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression))

```
n_iter_i = _check_optimize_result(
```

## Python code for advanced ensemble technique-stacking, blending, bagging and boosting

```
In [45]: 1 pip install scikit-learn  
2
```

Requirement already satisfied: scikit-learn in c:\users\anusha v\anaconda3 \lib\site-packages (1.0.2)  
Note: you may need to restart the kernel to use updated packages.

Requirement already satisfied: joblib>=0.11 in c:\users\anusha v\anaconda3 \lib\site-packages (from scikit-learn) (1.3.2)

Requirement already satisfied: threadpoolctl>=2.0.0 in c:\users\anusha v\anaconda3\lib\site-packages (from scikit-learn) (2.2.0)

Requirement already satisfied: scipy>=1.1.0 in c:\users\anusha v\anaconda3 \lib\site-packages (from scikit-learn) (1.9.1)

Requirement already satisfied: numpy>=1.14.6 in c:\users\anusha v\anaconda3\lib\site-packages (from scikit-learn) (1.24.4)

WARNING: Ignoring invalid distribution -umpy (c:\users\anusha v\anaconda3 \lib\site-packages)



In [26]:

```
1 # Import necessary libraries
2 from sklearn.model_selection import train_test_split
3 from sklearn.ensemble import StackingClassifier, VotingClassifier, BaggingClassifier
4 from sklearn.tree import DecisionTreeClassifier
5 from sklearn.ensemble import RandomForestClassifier
6 from sklearn.linear_model import LogisticRegression
7 from sklearn.neighbors import KNeighborsClassifier
8 from sklearn.svm import SVC
9 from sklearn.metrics import accuracy_score
10
11 # Load or generate your dataset
12 # For simplicity, let's use the famous Iris dataset
13 from sklearn.datasets import load_iris
14 data = load_iris()
15 X, y = data.data, data.target
16
17 # Split the data into training and testing sets
18 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
19
20 # Base models
21 knn_model = KNeighborsClassifier(n_neighbors=3)
22 rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
23 svm_model = SVC(kernel='linear', probability=True)
24
25 # Stacking
26 stacking_model = StackingClassifier(
27     estimators=[('knn', knn_model), ('rf', rf_model), ('svm', svm_model)],
28     final_estimator=LogisticRegression()
29 )
30
31 # Blending (Voting Classifier)
32 blending_model = VotingClassifier(
33     estimators=[('knn', knn_model), ('rf', rf_model), ('svm', svm_model)],
34     voting='soft'
35 )
36
37 # Bagging
38 bagging_model = BaggingClassifier(
39     base_estimator=DecisionTreeClassifier(),
40     n_estimators=50,
41     random_state=42
42 )
43
44 # Boosting
45 boosting_model = AdaBoostClassifier(
46     base_estimator=DecisionTreeClassifier(max_depth=1),
47     n_estimators=50,
48     random_state=42
49 )
50
51 # Train models
52 stacking_model.fit(X_train, y_train)
53 blending_model.fit(X_train, y_train)
54 bagging_model.fit(X_train, y_train)
55 boosting_model.fit(X_train, y_train)
56
57 # Predictions
58 stacking_pred = stacking_model.predict(X_test)
59 blending_pred = blending_model.predict(X_test)
60 bagging_pred = bagging_model.predict(X_test)
61 boosting_pred = boosting_model.predict(X_test)
```

```
62  
63 # Evaluate models  
64 stacking_accuracy = accuracy_score(y_test, stacking_pred)  
65 blending_accuracy = accuracy_score(y_test, blending_pred)  
66 bagging_accuracy = accuracy_score(y_test, bagging_pred)  
67 boosting_accuracy = accuracy_score(y_test, boosting_pred)  
68  
69 # Print accuracies  
70 print(f"Stacking Accuracy: {stacking_accuracy}")  
71 print(f"Blending Accuracy: {blending_accuracy}")  
72 print(f"Bagging Accuracy: {bagging_accuracy}")  
73 print(f"Boosting Accuracy: {boosting_accuracy}")  
74
```

```
Stacking Accuracy: 1.0  
Blending Accuracy: 1.0  
Bagging Accuracy: 1.0  
Boosting Accuracy: 1.0
```

## Python code for random forest on breast cancer

In [47]:

```

1 import numpy as np
2 import pandas as pd
3 from sklearn.datasets import load_breast_cancer
4 from sklearn.model_selection import train_test_split
5 from sklearn.ensemble import RandomForestClassifier
6 from sklearn.metrics import accuracy_score, classification_report, conf
7
8 # Load the Breast Cancer dataset
9 breast_cancer = load_breast_cancer()
10 X, y = breast_cancer.data, breast_cancer.target
11
12 # Convert the data to a DataFrame for better visualization (optional)
13 feature_names = breast_cancer.feature_names
14 df = pd.DataFrame(data=np.c_[X, y], columns=np.append(feature_names, 't
15
16 # Split the dataset into training and testing sets
17 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2
18
19 # Create a Random Forest Classifier
20 rf_classifier = RandomForestClassifier(n_estimators=100, random_state=4
21
22 # Train the model on the training set
23 rf_classifier.fit(X_train, y_train)
24
25 # Make predictions on the test set
26 y_pred = rf_classifier.predict(X_test)
27
28 # Evaluate the model
29 accuracy = accuracy_score(y_test, y_pred)
30 conf_matrix = confusion_matrix(y_test, y_pred)
31 class_report = classification_report(y_test, y_pred)
32
33 # Display the results
34 print("Random Forest Classifier Evaluation:")
35 print("Accuracy:", accuracy)
36 print("\nConfusion Matrix:")
37 print(conf_matrix)
38 print("\nClassification Report:")
39 print(class_report)
40

```

Random Forest Classifier Evaluation:

Accuracy: 0.9649122807017544

Confusion Matrix:

```

[[40  3]
 [ 1 70]]

```

Classification Report:

	precision	recall	f1-score	support
0	0.98	0.93	0.95	43
1	0.96	0.99	0.97	71
accuracy			0.96	114
macro avg	0.97	0.96	0.96	114
weighted avg	0.97	0.96	0.96	114

# Evaluation matrices on inertia and dunn index

## Inertia (Within-Cluster Sum of Squares):

In [24]:

```
1 from sklearn.cluster import KMeans
2 from sklearn.datasets import make_blobs
3
4 # Generate sample data
5 X, y = make_blobs(n_samples=300, centers=4, random_state=42)
6
7 # Fit a KMeans model
8 kmeans = KMeans(n_clusters=4)
9 kmeans.fit(X)
10
11 # Get the inertia
12 inertia = kmeans.inertia_
13 print(f"Inertia: {inertia}")
```

Inertia: 564.9141808210252

## Dunn Index:

In [25]:

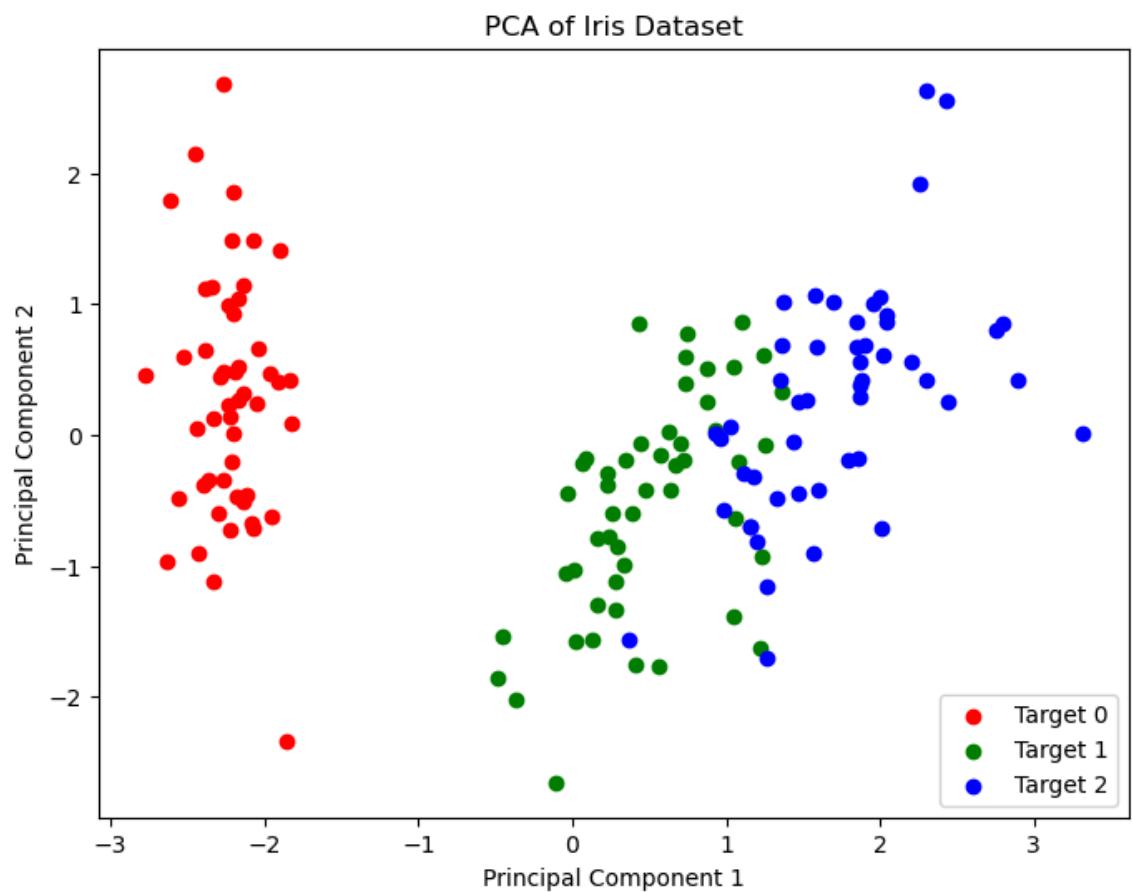
```
1 from sklearn.metrics import pairwise_distances
2 from sklearn.metrics import pairwise_distances_argmin_min
3 from sklearn.cluster import KMeans
4
5 def dunn_index(X, labels):
6     # Compute pairwise distances between points
7     dists = pairwise_distances(X)
8
9     # Compute minimum distance between points of different clusters
10    min_inter_cluster_distance = float('inf')
11    for i in range(len(labels)):
12        for j in range(i + 1, len(labels)):
13            if labels[i] != labels[j]:
14                min_inter_cluster_distance = min(min_inter_cluster_distance,
15
16    # Compute maximum intra-cluster distance
17    cluster_distances = []
18    for label in set(labels):
19        cluster_points = X[labels == label]
20        cluster_dists = pairwise_distances(cluster_points)
21        cluster_distances.append(cluster_dists.max())
22
23    # Compute Dunn Index
24    dunn_index = min_inter_cluster_distance / max(cluster_distances)
25    return dunn_index
26
27 # Fit a KMeans model
28 kmeans = KMeans(n_clusters=4)
29 kmeans.fit(X)
30
31 # Get the Dunn Index
32 labels = kmeans.labels_
33 dunn_index_value = dunn_index(X, labels)
34 print(f"Dunn Index: {dunn_index_value}")
```

Dunn Index: 0.22935905312943683

## Python code for dimensionality reduction using PCA

In [1]:

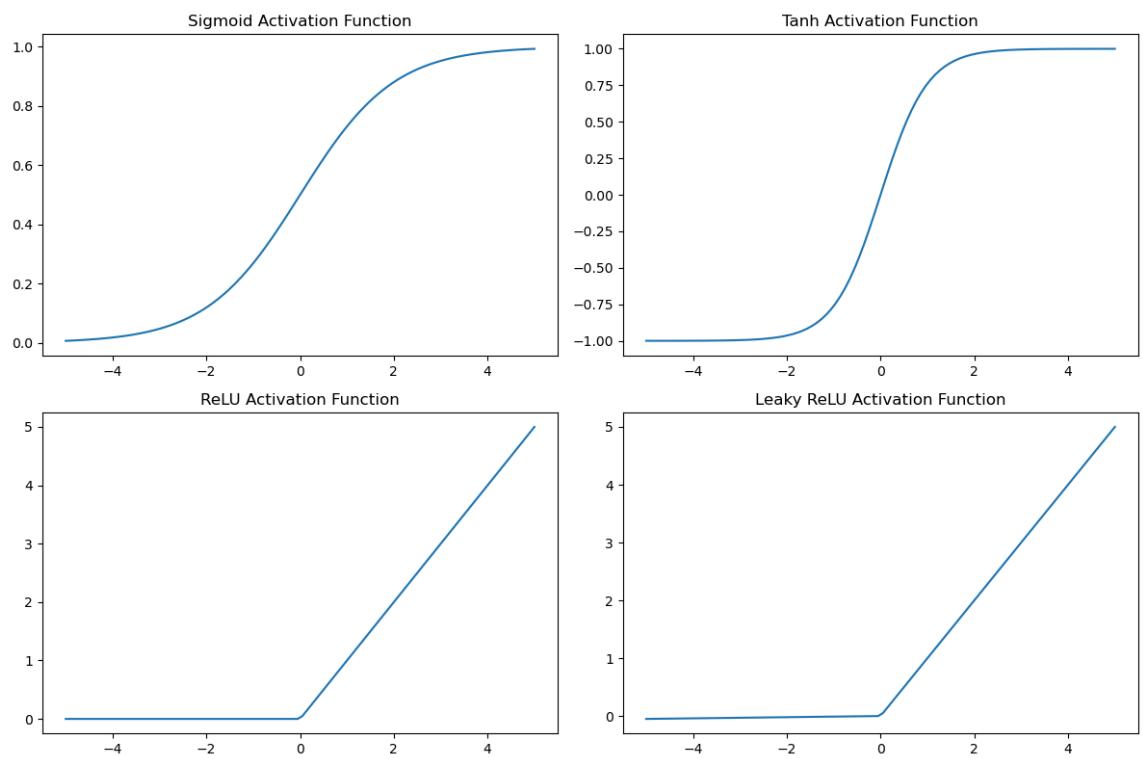
```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.datasets import load_iris
5 from sklearn.decomposition import PCA
6 from sklearn.preprocessing import StandardScaler
7
8 # Load the Iris dataset
9 iris = load_iris()
10 X = iris.data
11 y = iris.target
12
13 # Standardize the features (mean=0 and variance=1)
14 X_std = StandardScaler().fit_transform(X)
15
16 # Apply PCA with 2 components
17 pca = PCA(n_components=2)
18 X_pca = pca.fit_transform(X_std)
19
20 # Create a DataFrame for visualization
21 df_pca = pd.DataFrame(data=X_pca, columns=['Principal Component 1', 'Pr
22 df_pca['Target'] = y
23
24 # Visualize the PCA result
25 plt.figure(figsize=(8, 6))
26 targets = np.unique(y)
27 colors = ['r', 'g', 'b']
28
29 for target, color in zip(targets, colors):
30     indices_to_keep = df_pca['Target'] == target
31     plt.scatter(df_pca.loc[indices_to_keep, 'Principal Component 1'],
32                 df_pca.loc[indices_to_keep, 'Principal Component 2'],
33                 c=color, label=f'Target {target}')
34
35 plt.xlabel('Principal Component 1')
36 plt.ylabel('Principal Component 2')
37 plt.title('PCA of Iris Dataset')
38 plt.legend()
39 plt.show()
40
```



**Python code for activation function using sigmoid,Tanh,Relu,leaky relu**

In [2]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 def sigmoid(x):
5     return 1 / (1 + np.exp(-x))
6
7 def tanh(x):
8     return np.tanh(x)
9
10 def relu(x):
11     return np.maximum(0, x)
12
13 def leaky_relu(x, alpha=0.01):
14     return np.where(x > 0, x, alpha * x)
15
16 # Generate input values
17 x_values = np.linspace(-5, 5, 100)
18
19 # Apply activation functions
20 sigmoid_values = sigmoid(x_values)
21 tanh_values = tanh(x_values)
22 relu_values = relu(x_values)
23 leaky_relu_values = leaky_relu(x_values)
24
25 # Plot the activation functions
26 plt.figure(figsize=(12, 8))
27
28 plt.subplot(2, 2, 1)
29 plt.plot(x_values, sigmoid_values, label='Sigmoid')
30 plt.title('Sigmoid Activation Function')
31
32 plt.subplot(2, 2, 2)
33 plt.plot(x_values, tanh_values, label='Tanh')
34 plt.title('Tanh Activation Function')
35
36 plt.subplot(2, 2, 3)
37 plt.plot(x_values, relu_values, label='ReLU')
38 plt.title('ReLU Activation Function')
39
40 plt.subplot(2, 2, 4)
41 plt.plot(x_values, leaky_relu_values, label='Leaky ReLU')
42 plt.title('Leaky ReLU Activation Function')
43
44 plt.tight_layout()
45 plt.show()
46
```



**Python code for any data set dealing with missing values, category values, labelled encoding, one-hot coding**

In [3]:

```

1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.preprocessing import LabelEncoder, OneHotEncoder
4
5 # Load the Titanic dataset (a simplified version)
6 titanic_data = {
7     'Name': ['John', 'Jane', 'Bob', 'Alice', 'Charlie'],
8     'Age': [22, 30, None, 25, 32],
9     'Sex': ['male', 'female', 'male', 'female', 'male'],
10    'Survived': [1, 1, 0, 1, 0]
11 }
12
13 df = pd.DataFrame(titanic_data)
14
15 # Display the original dataset
16 print("Original Dataset:")
17 print(df)
18 print("\n")
19
20 # Dealing with missing values (replace with mean or median in a real scenario)
21 df['Age'].fillna(df['Age'].median(), inplace=True)
22
23 # Dealing with categorical values using Label Encoding
24 label_encoder = LabelEncoder()
25 df['Sex_LabelEncoded'] = label_encoder.fit_transform(df['Sex'])
26
27 # Dealing with categorical values using One-Hot Encoding
28 one_hot_encoder = OneHotEncoder(drop='first', sparse=False)
29 sex_one_hot_encoded = one_hot_encoder.fit_transform(df[['Sex']])
30 df_encoded = pd.concat([df, pd.DataFrame(sex_one_hot_encoded, columns=['Sex_Female'])], axis=1)
31
32 # Display the preprocessed dataset
33 print("Preprocessed Dataset:")
34 print(df_encoded)
35

```

Original Dataset:

	Name	Age	Sex	Survived
0	John	22.0	male	1
1	Jane	30.0	female	1
2	Bob	NaN	male	0
3	Alice	25.0	female	1
4	Charlie	32.0	male	0

Preprocessed Dataset:

	Name	Age	Sex	Survived	Sex_LabelEncoded	Sex_Female
0	John	22.0	male	1	1	1.0
1	Jane	30.0	female	1	0	0.0
2	Bob	27.5	male	0	1	1.0
3	Alice	25.0	female	1	0	0.0
4	Charlie	32.0	male	0	1	1.0

## Python code for n-grams

In [7]:

```

1 import nltk
2 from nltk import ngrams
3
4 # Sample text
5 text = "This is a sample sentence for generating n-grams."
6
7 # Tokenize the text into words
8 words = nltk.word_tokenize(text)
9
10 # Function to generate n-grams
11 def generate_ngrams(input_list, n):
12     return list(ngrams(input_list, n))
13
14 # Generate n-grams
15 ngram_size = 2# Change this value to generate different n-grams (e.g.,
16 result_ngrams = generate_ngrams(words, ngram_size)
17
18 # Display the generated n-grams
19 print(f"{ngram_size}-grams:")
20 for ngram in result_ngrams:
21     print(ngram)
22

```

2-grams:

('This', 'is')  
 ('is', 'a')  
 ('a', 'sample')  
 ('sample', 'sentence')  
 ('sentence', 'for')  
 ('for', 'generating')  
 ('generating', 'n-grams')  
 ('n-grams', '.')

## sentiment analysis on Spacey,spacy and keras

In [10]:

```

1 import spacy
2
3 # Load spaCy model
4 nlp = spacy.load("en_core_web_sm")
5
6 # Sample text
7 text = "I love using spaCy for NLP tasks."
8
9 # Analyze sentiment
10 doc = nlp(text)
11 sentiment = doc.sentiment
12
13 # Display sentiment score
14 print("Sentiment Score:", sentiment)
15

```

Sentiment Score: 0.0

```
In [13]: 1 pip install spacy tensorflow keras  
2
```

Requirement already satisfied: spacy in c:\users\anusha v\anaconda3\lib\\site-packages (3.5.3)  
Note: you may need to restart the kernel to use updated packages.

Requirement already satisfied: tensorflow in c:\users\anusha v\anaconda3\lib\\site-packages (2.14.0)  
Requirement already satisfied: keras in c:\users\anusha v\anaconda3\lib\\site-packages (2.14.0)  
Requirement already satisfied: catalogue<2.1.0,>=2.0.6 in c:\users\anusha v\anaconda3\lib\\site-packages (from spacy) (2.0.7)  
Requirement already satisfied: wasabi<1.2.0,>=0.9.1 in c:\users\anusha v\anaconda3\lib\\site-packages (from spacy) (0.9.1)  
Requirement already satisfied: preshed<3.1.0,>=3.0.2 in c:\users\anusha v\anaconda3\lib\\site-packages (from spacy) (3.0.6)  
Requirement already satisfied: numpy>=1.15.0 in c:\users\anusha v\anaconda3\lib\\site-packages (from spacy) (1.24.4)  
Requirement already satisfied: thinc<8.2.0,>=8.1.8 in c:\users\anusha v\anaconda3\lib\\site-packages (from spacy) (8.1.10)  
Requirement already satisfied: spacy-legacy<3.1.0,>=3.0.11 in c:\users\anusha v\anaconda3\lib\\site-packages (from spacy) (3.0.11)

```
In [ ]: 1
```

In [22]:

```
1 import spacy
2 from tensorflow.keras.models import Sequential
3 from tensorflow.keras.layers import Embedding, LSTM, Dense
4 from tensorflow.keras.preprocessing.text import Tokenizer
5 from tensorflow.keras.preprocessing.sequence import pad_sequences
6
7 # Load Spacy model
8 nlp = spacy.load("en_core_web_sm")
9
10 # Sample text for sentiment analysis
11 text = "ChatGPT is an amazing language model. I love using it!"
12
13 # Spacy for sentiment analysis
14 doc = nlp(text)
15 sentiment_score = doc.sentiment
16
17 print(f"Spacy Sentiment Score: {sentiment_score}")
18
19 # Define a simple LSTM model for sentiment analysis
20 model = Sequential()
21 model.add(Embedding(input_dim=10000, output_dim=16, input_length=100))
22 model.add(LSTM(64))
23 model.add(Dense(1, activation='sigmoid'))
24 model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
25
26 # Tokenize and pad the text
27 tokenizer = Tokenizer(num_words=10000, oov_token='<OOV>')
28 tokenizer.fit_on_texts([text])
29 sequences = tokenizer.texts_to_sequences([text])
30 padded_sequences = pad_sequences(sequences, maxlen=100, truncating='post')
31
32 # Make a prediction with the model
33 prediction = model.predict(padded_sequences)
34
35 print(f"Keras/TensorFlow Sentiment Prediction: {prediction[0][0]}")
36
37
```

```
Spacy Sentiment Score: 0.0
1/1 [=====] - 1s 617ms/step
Keras/TensorFlow Sentiment Prediction: 0.500891923904419
```

In [ ]:

1