

Exp 1 or Linear regression

Apply least squares model for the regression of number of units sold on TV advertising budget for the Advertising data, for least squares coefficient estimates for simple linear regression

```
In [1]: 1 import pandas as pd
2 from sklearn.linear_model import LinearRegression
3 from sklearn.metrics import r2_score, mean_squared_error
4
5 df = pd.read_csv('Advertising.csv')
6
7 X = df['TV'].values.reshape(-1, 1)
8 y = df['Sales']
9
10 model = LinearRegression()
11 model.fit(X, y)
12 y_pred = model.predict(X)
13 r_squared = r2_score(y, y_pred)
14 r_squared
```

Out[1]: 0.8121757029987414

Multiple linear regression

```
In [2]: 1 import pandas as pd
2 from sklearn.linear_model import LinearRegression
3 from sklearn.metrics import r2_score, mean_squared_error
4 df = pd.read_csv('Advertising.csv')
5 X = df[['TV', 'Radio', 'Newspaper']]
6 y = df['Sales']
7 model = LinearRegression()
8 model.fit(X, y)
9 y_pred = model.predict(X)
10 r_squared = r2_score(y, y_pred)
11 r_squared
```

Out[2]: 0.9025912899684558

KNN

```
In [3]: 1 import pandas as pd
2 from sklearn.neighbors import KNeighborsRegressor
3 from sklearn.model_selection import train_test_split
4 df = pd.read_csv('Advertising.csv')
5
6 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
7 model = KNeighborsRegressor(n_neighbors=5)
8 model.fit(X_train, y_train)
9 y_pred = model.predict(X_test)
10 r_squared = r2_score(y_test, y_pred)
11 r_squared
```

Out[3]: 0.8991773755626823

Navi Bayesian

```
In [4]: 1 import pandas as pd
2 from sklearn.neighbors import KNeighborsRegressor
3 from sklearn.model_selection import train_test_split
4 from sklearn.metrics import r2_score
5
6 df = pd.read_csv('Advertising.csv')
7
8 X = df[['TV', 'Radio', 'Newspaper']]
9 y = df['Sales']
10
11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
12
13 model = KNeighborsRegressor(n_neighbors=5)
14 model.fit(X_train, y_train)
15
16 y_pred = model.predict(X_test)
17
18 r_squared = r2_score(y_test, y_pred)
19 print("R-squared score:", r_squared)
20
```

R-squared score: 0.8991773755626823

Gradient Descent

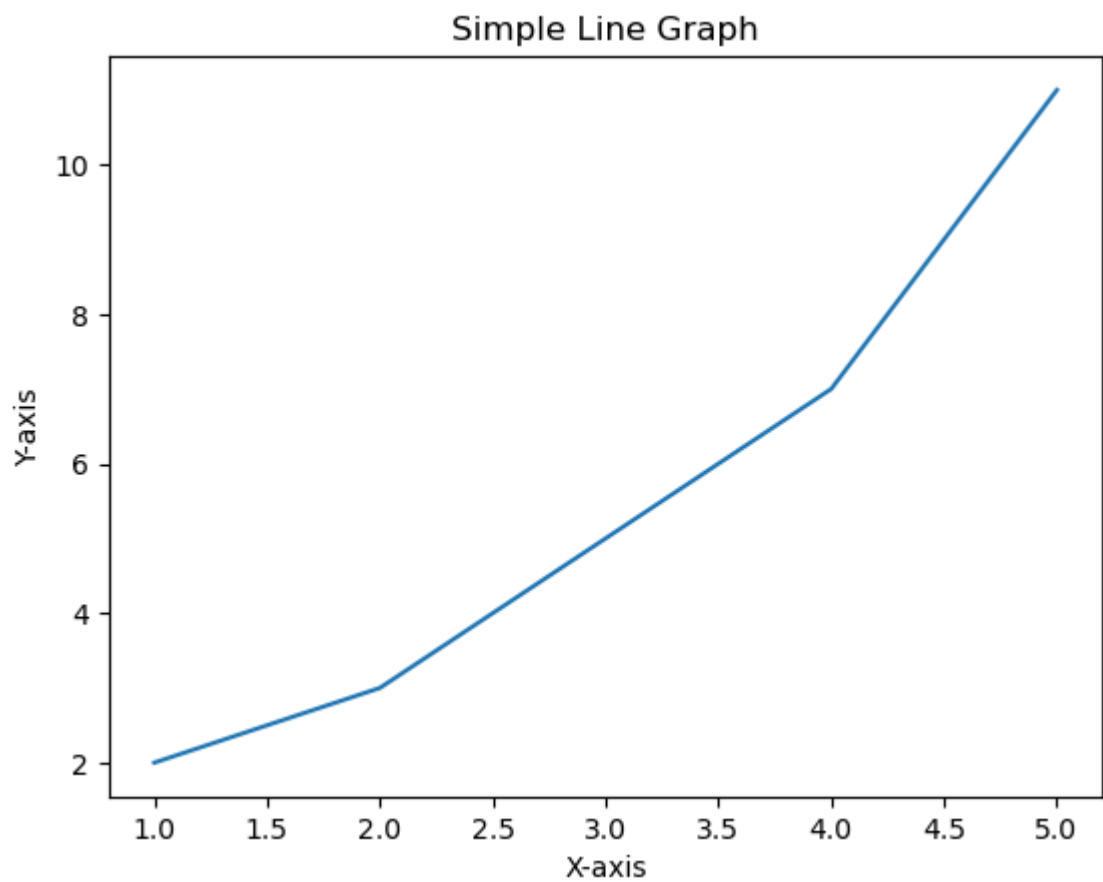
```
In [5]: 1 import pandas as pd
2 from sklearn.linear_model import SGDRegressor
3 from sklearn.preprocessing import StandardScaler
4 from sklearn.model_selection import train_test_split
5
6 df = pd.read_csv('Advertising.csv')
7
8 X = df[['TV', 'Radio', 'Newspaper']]
9 y = df['Sales']
10 scaler = StandardScaler()
11 X_train_scaled = scaler.fit_transform(X_train)
12 X_test_scaled = scaler.transform(X_test)
13 model = SGDRegressor(loss='squared_error')
14 model.fit(X_train_scaled, y_train)
15
16 y_test_pred = model.predict(X_test_scaled)
17 test_r_squared = r2_score(y_test, y_test_pred)
18
19 print("Testing R-squared:", test_r_squared)
```

Testing R-squared: 0.9061334961284013

Matplotlib

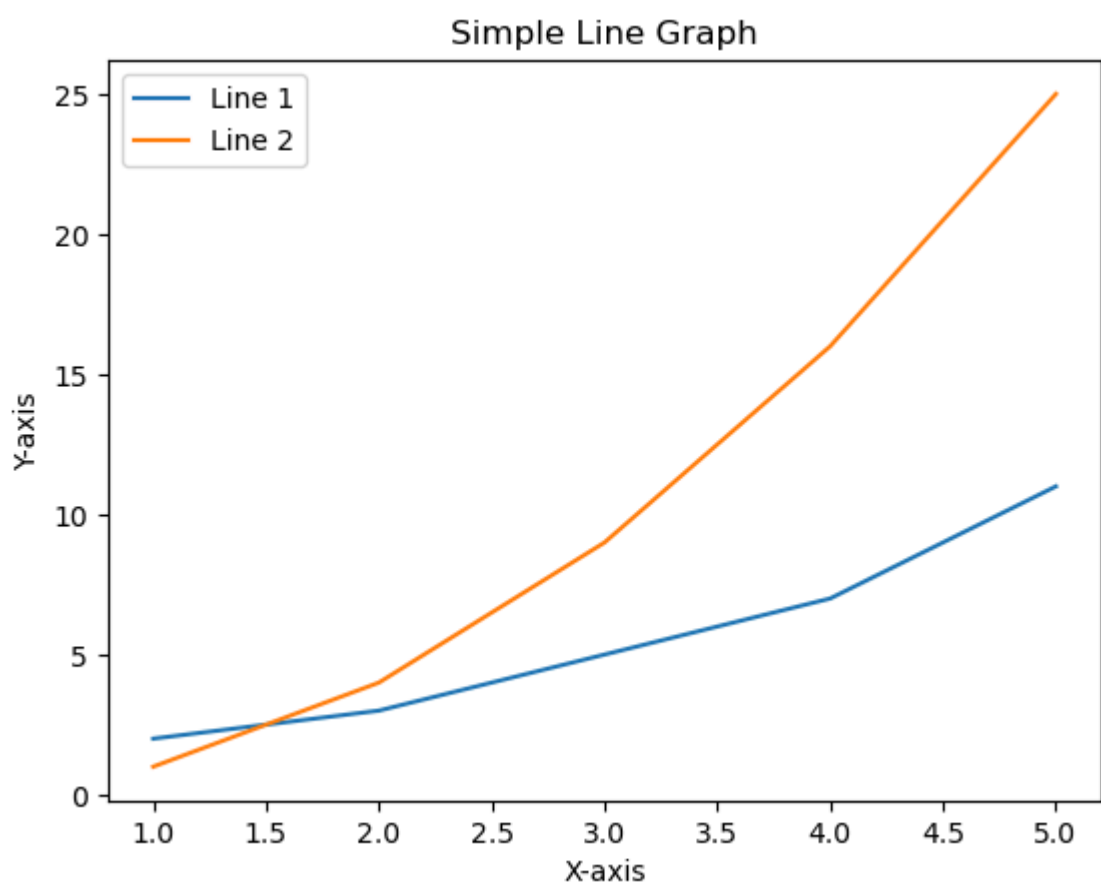
Simple line graph

```
In [6]: 1 import matplotlib.pyplot as plt
2
3
4 x = [1, 2, 3, 4, 5]
5 y = [2, 3, 5, 7, 11]
6
7 plt.plot(x, y)
8
9 plt.xlabel('X-axis')
10 plt.ylabel('Y-axis')
11 plt.title('Simple Line Graph')
12
13
14 plt.show()
15
```



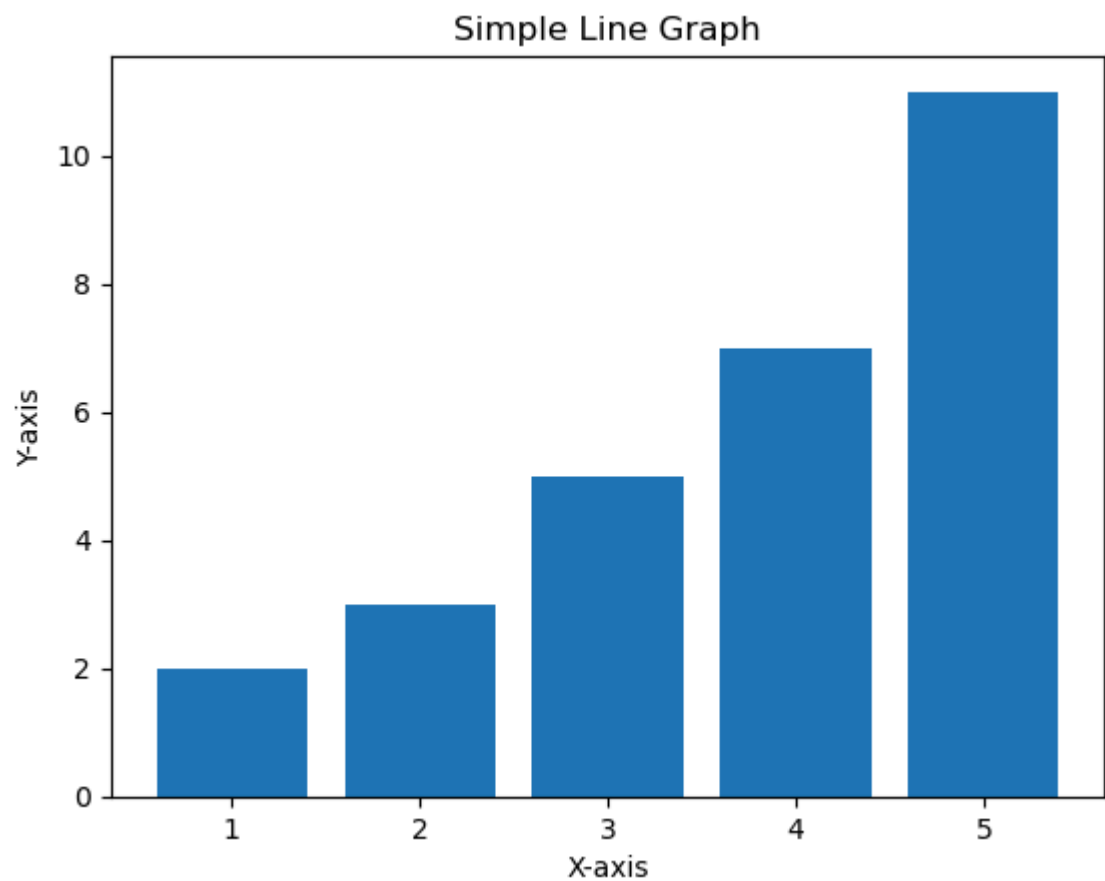
Multiple line graph

```
In [7]: 1 import matplotlib.pyplot as plt
2
3
4 x = [1, 2, 3, 4, 5]
5 y1 = [2, 3, 5, 7, 11]
6 y2 = [1, 4, 9, 16, 25]
7
8 plt.plot(x, y1, label='Line 1')
9 plt.plot(x, y2, label='Line 2')
10
11 plt.xlabel('X-axis')
12 plt.ylabel('Y-axis')
13 plt.title('Simple Line Graph')
14
15 plt.legend()
16
17 plt.show()
18
```



Bar graph

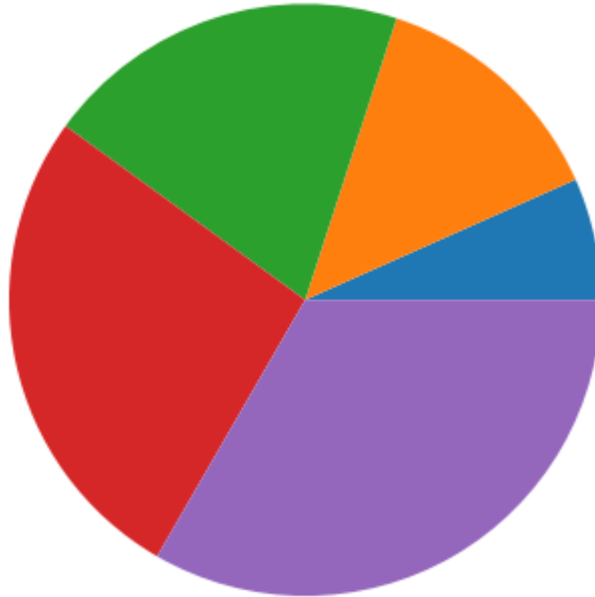
```
In [8]: 1 import matplotlib.pyplot as plt
2
3
4 x = [1, 2, 3, 4, 5]
5 y = [2, 3, 5, 7, 11]
6
7 plt.bar(x, y)
8
9 plt.xlabel('X-axis')
10 plt.ylabel('Y-axis')
11 plt.title('Simple Line Graph')
12
13
14 plt.show()
```



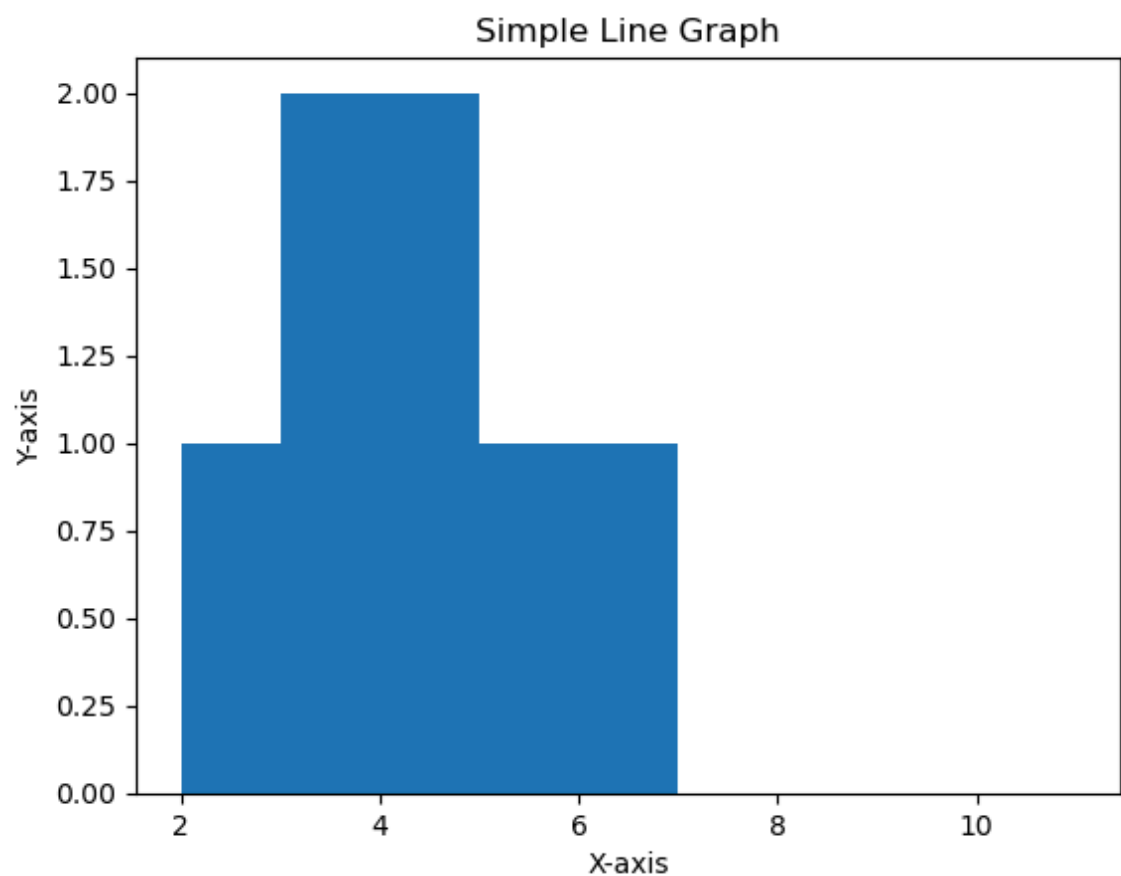
Pie chart

```
In [22]: 1 import matplotlib.pyplot as plt
          2 x = [1, 2, 3, 4, 5]
          3 plt.pie(x)
          4 plt.title('Simple Line Graph')
          5 plt.show()
```

Simple Line Graph



```
In [10]: 1 import matplotlib.pyplot as plt
2
3
4 x = [1, 2, 3, 4, 5]
5 y = [2, 3, 5, 7, 11]
6
7 plt.hist(x, y)
8
9 plt.xlabel('X-axis')
10 plt.ylabel('Y-axis')
11 plt.title('Simple Line Graph')
12
13
14 plt.show()
```



Exp 2

Compute t-statistic, Residual standard error, F-statistic and residual sum of squares (RSS) errors.


```
In [ ]: 1 import numpy as np
        2 import statsmodels.api as sm
        3 X = np.array([1, 2, 3, 4, 5])
        4 y = np.array([2, 3, 4, 5, 6])
        5 X = sm.add_constant(X)
        6 model = sm.OLS(y, X).fit()
        7 t_statistic = model.tvalues
        8 residual_standard_error = np.sqrt(model.mse_resid)
        9 f_statistic = model.fvalue
       10 rss = model.ssr
       11 print("T-Statistic:", t_statistic)
       12 print("Residual Standard Error:", residual_standard_error)
       13 print("F-Statistic:", f_statistic)
       14 print("Residual Sum of Squares (RSS):", rss)
```

Exp3

Do any of the predictors appear to be statistically significant? If so, which ones?

```
In [32]: 1 import numpy as np
2 import statsmodels.api as sm
3
4 # Sample data
5 X = np.array([1, 2, 3, 4, 5])
6 y = np.array([2, 3, 4, 5, 6])
7
8 # Add constant for intercept term
9 X = sm.add_constant(X)
10
11 # Fit linear regression model
12 model = sm.OLS(y, X).fit()
13
14 # Get model summary
15 print(model.summary())
16
17 # Extract p-values
18 p_values = model.pvalues[1:] # Exclude intercept term
19 print("P-Values:", p_values)
20
21 # Identify statistically significant predictors
22 significant_predictors = np.where(p_values < 0.05)[0]
23 print("Significant Predictor Indices:", significant_predictors)
24
```

OLS Regression Results

```

=====
====
Dep. Variable:          y      R-squared:
1.000
Model:                OLS      Adj. R-squared:
1.000
Method:              Least Squares      F-statistic:          6.085
e+32
Date:                Fri, 03 May 2024      Prob (F-statistic):          1.47
e-49
Time:                14:56:10      Log-Likelihood:          17
7.15
No. Observations:          5      AIC:          -3
50.3
Df Residuals:            3      BIC:          -3
51.1
Df Model:                1
Covariance Type:          nonrobust
=====
====
              coef      std err          t      P>|t|      [0.025      0.
975]
-----
----
const          1.0000      1.34e-16      7.44e+15      0.000          1.000
1.000
x1              1.0000      4.05e-17      2.47e+16      0.000          1.000
1.000
=====
====
Omnibus:          nan      Durbin-Watson:
1.000
Prob(Omnibus):          nan      Jarque-Bera (JB):
1.888
Skew:              1.500      Prob(JB):
0.389
Kurtosis:          3.250      Cond. No.
8.37
=====
====

```

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

P-Values: [1.46929683e-49]

Significant Predictor Indices: [0]

C:\Users\hites\anaconda3\Lib\site-packages\statsmodels\stats\stattools.py:

74: ValueWarning: omni_normtest is not valid with less than 8 observations; 5 samples were given.

warn("omni_normtest is not valid with less than 8 observations; %i "

Exp4

Compute the confusion matrix and overall fraction of correct predictions. Explain what the confusion matrix is telling you about the types of mistakes made by KNN

```
In [17]: 1 from sklearn.metrics import confusion_matrix, accuracy_score
2
3
4 y_true = [0, 1, 0, 1, 0, 1, 1, 0, 1, 0]
5 y_pred = [0, 1, 1, 1, 0, 0, 1, 0, 1, 1]
6
7 # Compute confusion matrix
8 cm = confusion_matrix(y_true, y_pred)
9 print("Confusion Matrix:")
10 print(cm)
11
12 # Compute overall fraction of correct predictions
13 accuracy = accuracy_score(y_true, y_pred)
14 print("Overall Accuracy:", accuracy)
15
16
```

Confusion Matrix:

[[3 2]

[1 4]]

Overall Accuracy: 0.7

Exp 5

Compute Mallow's Cp, Akaike information criterion (AIC), adjusted R Squared and Bayesian information criterion (BIC)

```
In [18]: 1 import numpy as np
2 import statsmodels.api as sm
3
4
5 X = np.array([[1, 2, 3, 4, 5], [2, 3, 4, 5, 6]]).T
6 y = np.array([2, 3, 4, 5, 6])
7
8 model = sm.OLS(y, sm.add_constant(X)).fit()
9
10 residuals = model.resid
11 p = len(model.params)
12 n = len(y)
13 rss = np.sum(residuals ** 2)
14
15 Cp = (1/n) * (rss + 2 * p * (rss/n))
16 AIC = (2 * p) + (n * np.log(rss/n))
17 adj_r_squared = 1 - ((rss/(n - p - 1)) / (np.var(y)))
18 BIC = n * np.log(rss/n) + p * np.log(n)
19
20
21 print("Mallow's Cp:", Cp)
22 print("AIC:", AIC)
23 print("Adjusted R-squared:", adj_r_squared)
24 print("BIC:", BIC)
25
```

Mallow's Cp: 1.8591479383796198e-29

AIC: -328.71653386035314

Adjusted R-squared: 1.0

BIC: -329.8882201230508

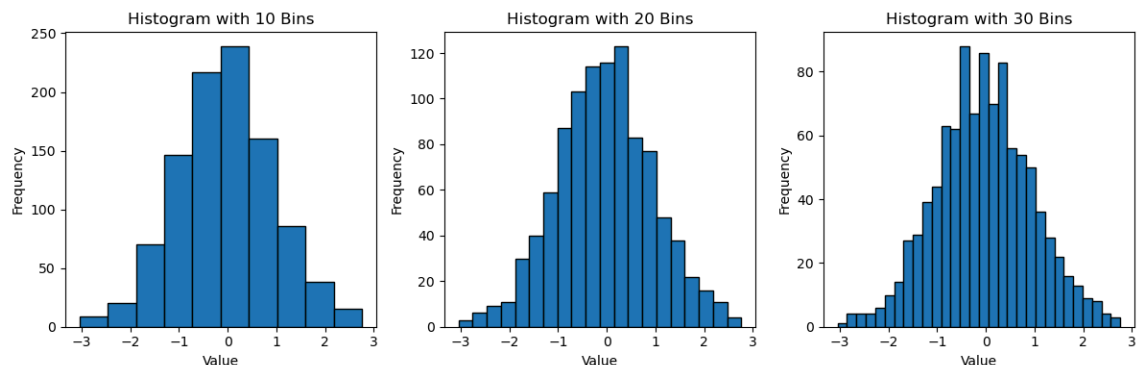
Exp6

Produce some histograms with differing numbers of bins for a few of the quantitative variables

```

In [23]: 1 import numpy as np
          2 import matplotlib.pyplot as plt
          3
          4 # Generate some example data
          5 np.random.seed(0)
          6 data = np.random.normal(loc=0, scale=1, size=1000) # Example data with
          7
          8 # Create histograms with different numbers of bins
          9 bins_list = [10, 20, 30] # Different numbers of bins
         10 titles = ['Histogram with 10 Bins', 'Histogram with 20 Bins', 'Histogram with 30 Bins']
         11
         12 plt.figure(figsize=(12, 4))
         13 for i, bins in enumerate(bins_list, start=1):
         14     plt.subplot(1, len(bins_list), i)
         15     plt.hist(data, bins=bins, edgecolor='black')
         16     plt.title(titles[i-1])
         17     plt.xlabel('Value')
         18     plt.ylabel('Frequency')
         19
         20 plt.tight_layout()
         21 plt.show()
         22

```



Exp 7

Continue exploring the data, and provide a brief summary of what you discover.

Descriptive Statistics: Calculate summary statistics such as mean, median, standard deviation, minimum, maximum, and quartiles for quantitative variables. For categorical variables, count the frequency of each category. **Data Visualization:** Create visualizations such as histograms, box plots, scatter plots, and bar plots to understand the distribution, relationships, and patterns in the data. **Correlation Analysis:** Compute correlations between pairs of quantitative variables to identify any linear relationships. Visualize the correlation matrix using a heatmap. **Outlier Detection:** Identify any outliers in the data using statistical methods or visualization techniques. **Missing Values:** Check for missing values in the dataset and decide on appropriate strategies for handling them, such as imputation or removal. **Feature Engineering:** Create new features or transform existing features to better represent the underlying patterns in the data. **Dimensionality Reduction:** Apply techniques such as principal component analysis (PCA) or t-distributed stochastic neighbor embedding (t-SNE) to reduce the dimensionality of the data and visualize high-dimensional data in lower dimensions. **Cluster Analysis:** Explore whether there are natural groupings or clusters in the data using clustering algorithms such as K-means or hierarchical clustering. Based on

the results of these analyses, we can provide a summary of the key findings, insights, and patterns observed in the data. This summary can help stakeholders better understand the

Exp 8

Carry out the pearson product moment correlation, spearman rho correlation and kendall's tau

In [24]:

```
1 import numpy as np
2 from scipy.stats import pearsonr, spearmanr, kendalltau
3
4 x = np.array([1, 2, 3, 4, 5])
5 y = np.array([2, 3, 4, 5, 6])
6
7 pearson_corr, pearson_p_value = pearsonr(x, y)
8 print("Pearson correlation coefficient:", pearson_corr)
9 print("Pearson p-value:", pearson_p_value)
10
11 spearman_corr, spearman_p_value = spearmanr(x, y)
12 print("Spearman's rank correlation coefficient (Spearman rho):", spearman_corr)
13 print("Spearman p-value:", spearman_p_value)
14
15 kendall_corr, kendall_p_value = kendalltau(x, y)
16 print("Kendall's tau correlation coefficient:", kendall_corr)
17 print("Kendall p-value:", kendall_p_value)
18
```

Pearson correlation coefficient: 1.0

Pearson p-value: 0.0

Spearman's rank correlation coefficient (Spearman rho): 0.9999999999999999

Spearman p-value: 1.4042654220543672e-24

Kendall's tau correlation coefficient: 0.9999999999999999

Kendall p-value: 0.016666666666666666

Exp 9

Perform Simple Hypothesis testing, student's t-test, paired t and u test, correlation and covariance, tests for association

```
In [26]: 1 import numpy as np
2 from scipy.stats import ttest_ind, ttest_rel, wilcoxon, chi2_contingency
3 import pandas as pd
4
5 # Example data
6 data = {
7     'X': [1, 2, 3, 4, 5],
8     'Y': [2, 3, 4, 5, 6],
9     'Category': ['A', 'B', 'A', 'B', 'A']
10 }
11 df = pd.DataFrame(data)
12
13 # Simple Hypothesis Testing
14 mean_value = 3 # Known value for comparison
15 sample_mean = np.mean(df['X'])
16 t_statistic, p_value = ttest_ind(df['X'], [mean_value])
17 print("Simple Hypothesis Testing:")
18 print("Sample mean:", sample_mean)
19 print("t-statistic:", t_statistic)
20 print("p-value:", p_value)
21
22 # Student's t-test
23 t_statistic, p_value = ttest_ind(df[df['Category'] == 'A']['Y'], df[df['Category'] == 'B']['Y'])
24 print("\nStudent's t-test:")
25 print("t-statistic:", t_statistic)
26 print("p-value:", p_value)
27
28 # Paired t-test
29 t_statistic, p_value = ttest_rel(df['X'], df['Y'])
30 print("\nPaired t-test:")
31 print("t-statistic:", t_statistic)
32 print("p-value:", p_value)
33
34 # Wilcoxon Signed-Rank Test (U Test)
35 statistic, p_value = wilcoxon(df['X'], df['Y'])
36 print("\nWilcoxon Signed-Rank Test (U Test):")
37 print("Statistic:", statistic)
38 print("p-value:", p_value)
39
40 # Correlation and Covariance
41 correlation = df['X'].corr(df['Y'])
42 covariance = df['X'].cov(df['Y'])
43 print("\nCorrelation and Covariance:")
44 print("Correlation coefficient:", correlation)
45 print("Covariance:", covariance)
46
47 # Tests for Association
48 contingency_table = pd.crosstab(df['X'], df['Category'])
49 chi2, p_value, _, _ = chi2_contingency(contingency_table)
50 print("\nTests for Association:")
51 print("Chi-square statistic:", chi2)
52 print("p-value:", p_value)
53
```


Simple Hypothesis Testing:

Sample mean: 3.0

t-statistic: 0.0

p-value: 1.0

Student's t-test:

t-statistic: 0.0

p-value: 1.0

Paired t-test:

t-statistic: -inf

p-value: 0.0

Wilcoxon Signed-Rank Test (U Test):

Statistic: 0.0

p-value: 0.0625

Correlation and Covariance:

Correlation coefficient: 0.9999999999999999

Covariance: 2.5

Tests for Association:

Chi-square statistic: 5.000000000000001

p-value: 0.2872974951836456

C:\Users\hites\anaconda3\Lib\site-packages\scipy\stats_axis_nan_policy.py:523: RuntimeWarning: Precision loss occurred in moment calculation due to catastrophic cancellation. This occurs when the data are nearly identical. Results may be unreliable.

```
res = hypotest_fun_out(*samples, **kws)
```

Exp 10

Perform Programming for Eigen values and Eigen vectors

In [28]:

```
1 import numpy as np
2
3 # Example matrix
4 A = np.array([[1, 2],
5               [2, 1]])
6
7 # Compute eigenvalues and eigenvectors
8 eigenvalues, eigenvectors = np.linalg.eig(A)
9
10 # Print eigenvalues and eigenvectors
11 print("Eigenvalues:")
12 print(eigenvalues)
13 print("\nEigenvectors:")
14 print(eigenvectors)
15
```

Eigenvalues:

```
[ 3. -1.]
```

Eigenvectors:

```
[[ 0.70710678 -0.70710678]
 [ 0.70710678  0.70710678]]
```

