# Experiment 1 : Data Cleaning & EDA

```python
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import warnings
```

```python
1 data = pd.read_csv("/content/data.csv")
```

```python
1 data.head(5)
```

|   | Country | Age | Salary | Purchased |
|---|---------|-----|--------|-----------|
| 0 | France | 44.0 | 72000.0 | No |
| 1 | Spain | 27.0 | 48000.0 | Yes |
| 2 | Germany | 30.0 | 54000.0 | No |
| 3 | Spain | 38.0 | 61000.0 | No |
| 4 | Germany | 40.0 | NaN | Yes |

```python
1 data.shape
```

```
(10, 4)
```

```python
1 pd.set_option('display.max_columns', None)
2 pd.set_option('display.max_rows', None)
```

```python
1 data.head(2)
```

|   | Country | Age | Salary | Purchased |
|---|---------|-----|--------|-----------|
| 0 | France | 44.0 | 72000.0 | No |
| 1 | Spain | 27.0 | 48000.0 | Yes |

```python
1 data.tail(2)
```

|   | Country | Age | Salary | Purchased |
|---|---------|-----|--------|-----------|
| 8 | Germany | 50.0 | 83000.0 | No |
| 9 | France | 37.0 | 67000.0 | Yes |

```python
1 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10 entries, 0 to 9
Data columns (total 4 columns):
 #   Column     Non-Null Count  Dtype
---  ------     --------------  -----
 0   Country    10 non-null     object
 1   Age        9 non-null      float64
 2   Salary     9 non-null      float64
 3   Purchased  10 non-null     object
```

```
dtypes: float64(2), object(2)
memory usage: 448.0+ bytes
```

```
1 data.isnull().sum()
```

```
                    0
    Country         0
    Age             1
    Salary          1
    Purchased       0

    dtype: int64
```
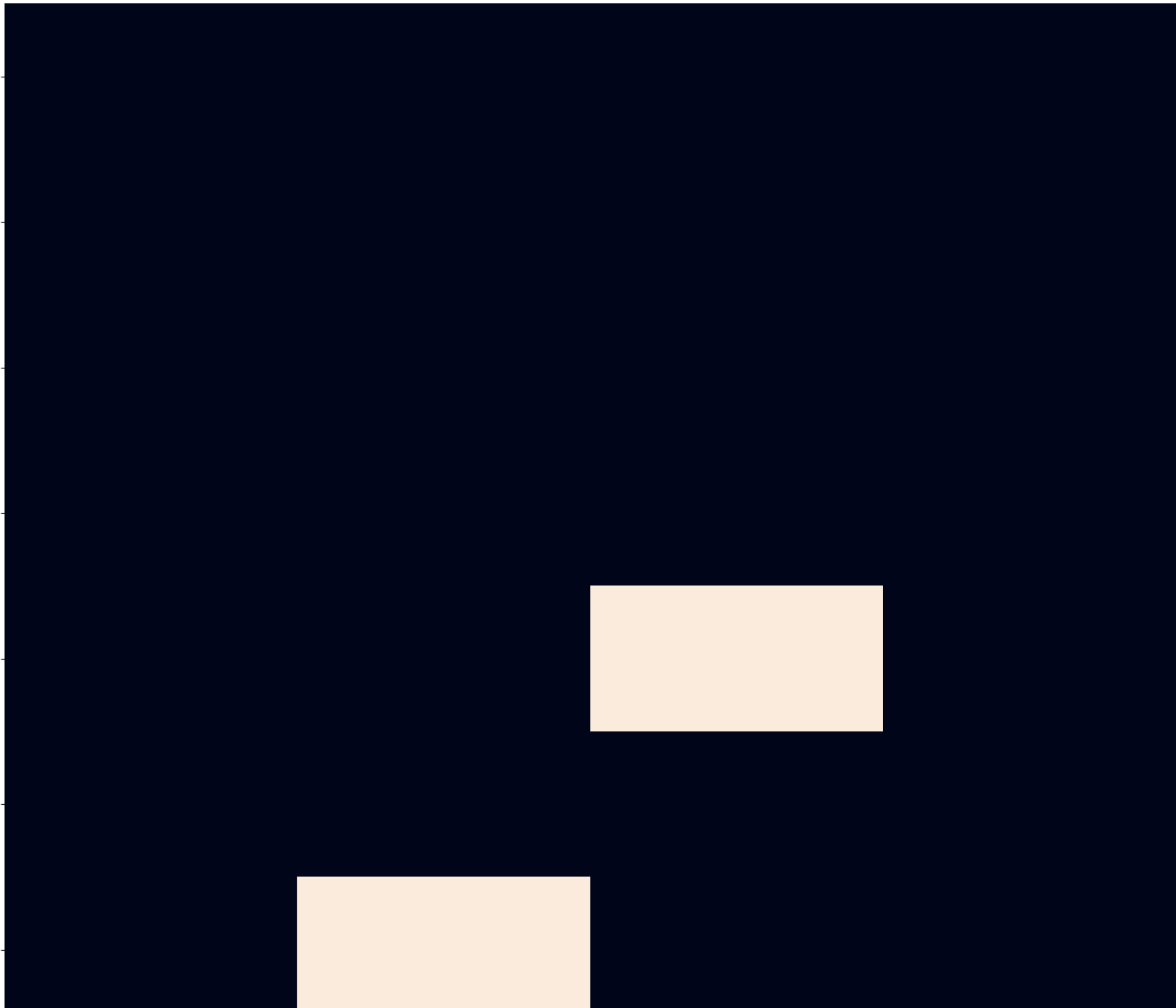
```
1 plt.figure(figsize=(25,25))
2 sns.heatmap(data.isnull())
3 plt.show()
```

```
1 missing_value_percent = data.isnull().sum() / data.shape[0] * 100
2 print(missing_value_percent)
```

```
Country      0.0
Age         10.0
Salary      10.0
Purchased    0.0
dtype: float64
```

```
1 missing_value_column = missing_value_percent[missing_value_percent > 17].keys()
2 print(missing_value_column)
```

```
Index([], dtype='object')
```

```
1 data1 = data.drop(columns = missing_value_column)
```
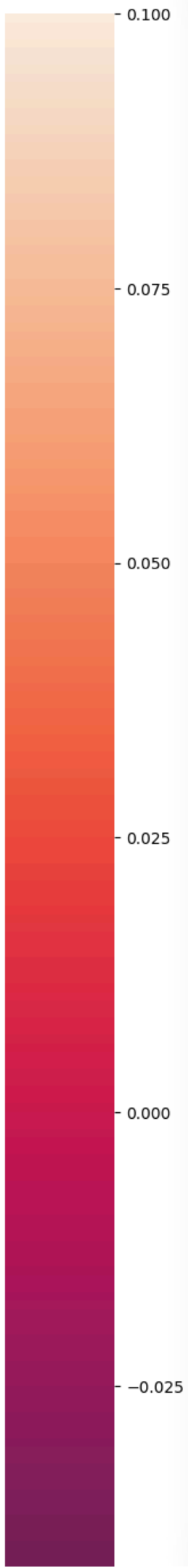
```
1 data1.shape
```

```
(10, 4)
```

```
1 data2 = data1.dropna()
```

```
1  data2.shape
```

```
(8, 4)
```

```
1  plt.figure(figsize=(25,25))
2  sns.heatmap(data1.isnull())
3  plt.show()
```

546
552
558
564
570
576
582
588
594
600
606
612
618
624
630

−0.050

```
1 data2.isnull().sum().sum()
```

0

Experiment 2: Outliers estimation & Plot it through violen/whisker

```
1 # Importing
2 import sklearn
3 from sklearn.datasets import load_diabetes
4 import pandas as pd
5 import matplotlib.pyplot as plt
6
7 # Load the dataset
8 diabetics = load_diabetes()
9
10 # Create the dataframe
11 column_name = diabetics.feature_names
12 df_diabetics = pd.DataFrame(diabetics.data)
13 df_diabetics.columns = column_name
14 print(df_diabetics.head())
15
```

```
        age       sex       bmi        bp        s1        s2        s3  \
0  0.038076  0.050680  0.061696  0.021872 -0.044223 -0.034821 -0.043401
1 -0.001882 -0.044642 -0.051474 -0.026328 -0.008449 -0.019163  0.074412
2  0.085299  0.050680  0.044451 -0.005670 -0.045599 -0.034194 -0.032356
3 -0.089063 -0.044642 -0.011595 -0.036656  0.012191  0.024991 -0.036038
4  0.005383 -0.044642 -0.036385  0.021872  0.003935  0.015596  0.008142

        s4        s5        s6
0 -0.002592  0.019907 -0.017646
1 -0.039493 -0.068332 -0.092204
2 -0.002592  0.002861 -0.025930
3  0.034309  0.022688 -0.009362
4 -0.002592 -0.031988 -0.046641
```
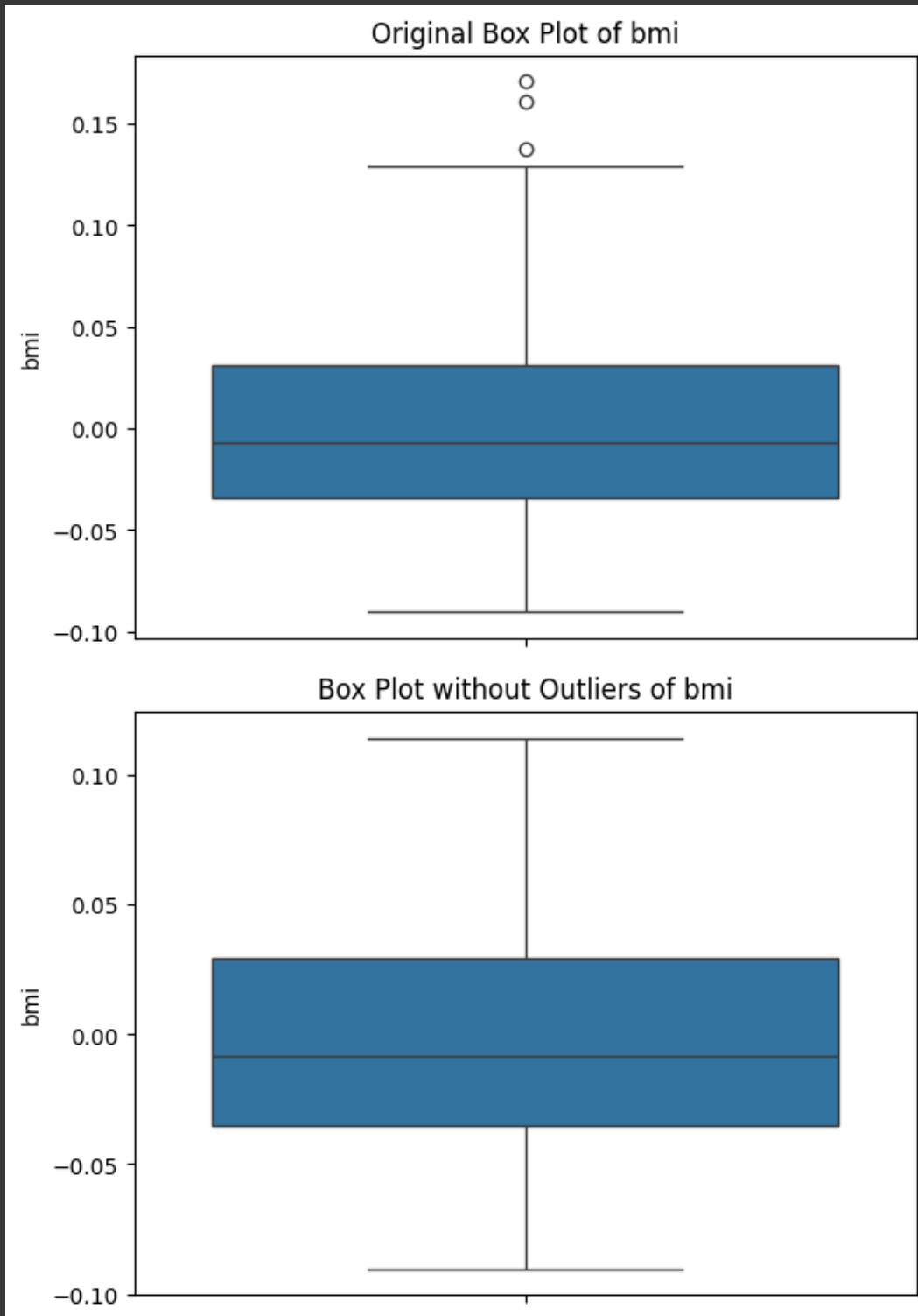
```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4
5 def removal_box_plot(df, column, threshold):
6     sns.boxplot(df[column])
7     plt.title(f'Original Box Plot of {column}')
8     plt.show()
9
10     removed_outliers = df[df[column] <= threshold]
11
12     sns.boxplot(removed_outliers[column])
13     plt.title(f'Box Plot without Outliers of {column}')
14     plt.show()
15     return removed_outliers
```

```
16
17
18  threshold_value = 0.12
19
20  no_outliers = removal_box_plot(df_diabetics, 'bmi', threshold_value)
21
```



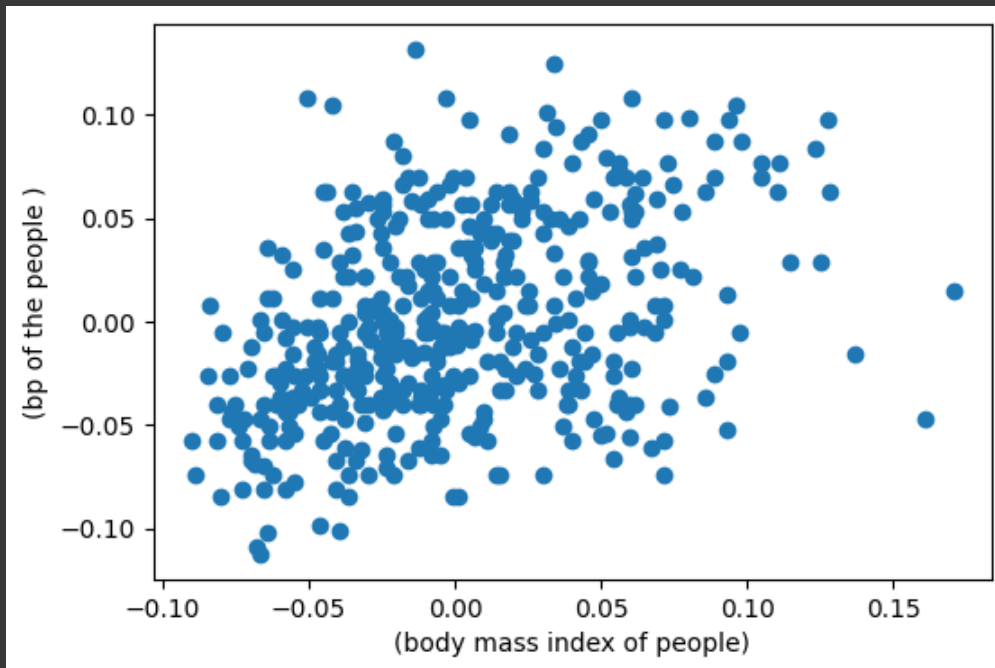Original Box Plot of bmi

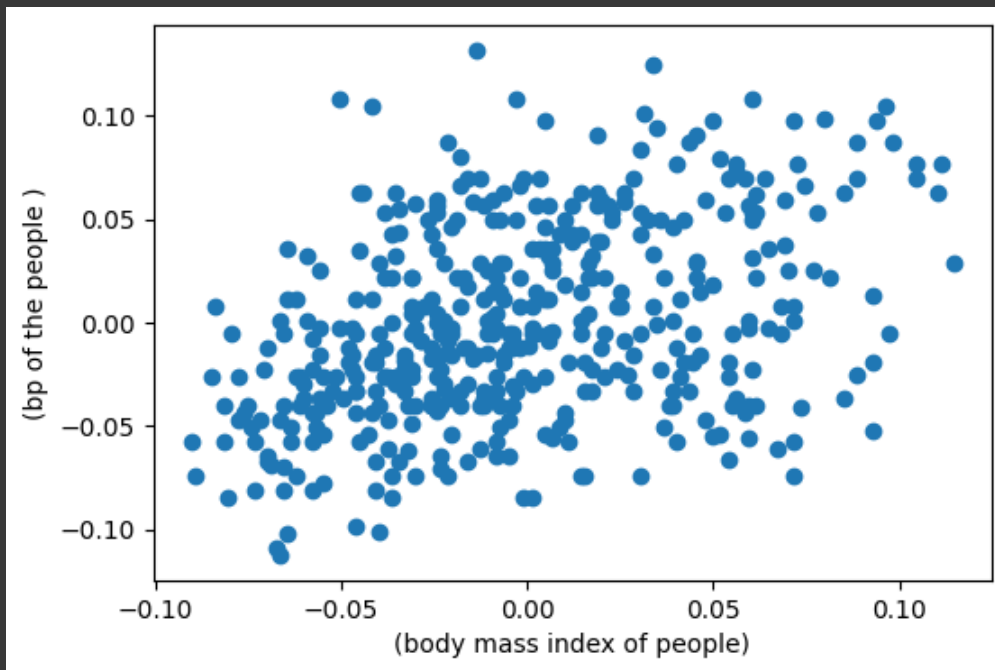Box Plot without Outliers of bmi

```
1   fig, ax = plt.subplots(figsize=(6, 4))
2   ax.scatter(df_diabetics['bmi'], df_diabetics['bp'])
3   ax.set_xlabel('(body mass index of people)')
4   ax.set_ylabel('(bp of the people )')
5   plt.show()
6
```

```python
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

outlier_indices = np.where((df_diabetics['bmi'] > 0.12) &
  (df_diabetics['bp'] < 0.8))

no_outliers = df_diabetics.drop(outlier_indices[0])

# Scatter plot without outliers
fig, ax_no_outliers = plt.subplots(figsize=(6, 4))
ax_no_outliers.scatter(no_outliers['bmi'], no_outliers['bp'])
ax_no_outliers.set_xlabel('(body mass index of people)')
ax_no_outliers.set_ylabel('(bp of the people )')
plt.show()
```



```python
import numpy as np

threshold_z = 2
```

```
 4
 5 outlier_indices = np.where(z > threshold_z)[0]
 6 no_outliers = df_diabetics.drop(outlier_indices)
 7 print("Original DataFrame Shape:", df_diabetics.shape)
 8 print("DataFrame Shape after Removing Outliers:", no_outliers.shape)
 9
```

```
Original DataFrame Shape: (442, 10)
DataFrame Shape after Removing Outliers: (426, 10)
```

```
 1 # IQR
 2 Q1 = np.percentile(df_diabetics['bmi'], 25, method='midpoint')
 3 Q3 = np.percentile(df_diabetics['bmi'], 75, method='midpoint')
 4 IQR = Q3 - Q1
 5 print(IQR)
 6
```

```
0.06520763046978838
```

```
 1 # Above Upper bound
 2 upper = Q3+1.5*IQR
 3 upper_array = np.array(df_diabetics['bmi'] >= upper)
 4 print("Upper Bound:", upper)
 5 print(upper_array.sum())
 6
 7 # Below Lower bound
 8 lower = Q1-1.5*IQR
 9 lower_array = np.array(df_diabetics['bmi'] <= lower)
10 print("Lower Bound:", lower)
11 print(lower_array.sum())
12
```

```
Upper Bound: 0.12879000811776306
3
Lower Bound: -0.13204051376139045
0
```

```
 1  # Importing
 2  import sklearn
 3  from sklearn.datasets import load_diabetes
 4  import pandas as pd
 5
 6  # Load the dataset
 7  diabetes = load_diabetes()
 8
 9  # Create the dataframe
10  column_name = diabetes.feature_names
11  df_diabetes = pd.DataFrame(diabetes.data)
12  df_diabetes .columns = column_name
13  df_diabetes .head()
14  print("Old Shape: ", df_diabetes.shape)
15
16  ''' Detection '''
17  # IQR
18  # Calculate the upper and lower limits
19  Q1 = df_diabetes['bmi'].quantile(0.25)
20  Q3 = df_diabetes['bmi'].quantile(0.75)
21  IQR = Q3 - Q1
22  lower = Q1 - 1.5*IQR
23  upper = Q3 + 1.5*IQR
24
25  # Create arrays of Boolean values indicating the outlier rows
26  upper_array = np.where(df_diabetes['bmi'] >= upper)[0]
```

```
27    lower_array = np.where(df_diabetes['bmi'] <= lower)[0]
28
29    # Removing the outliers
30    df_diabetes.drop(index=upper_array, inplace=True)
31    df_diabetes.drop(index=lower_array, inplace=True)
32
33    # Print the new shape of the DataFrame
34    print("New Shape: ", df_diabetes.shape)
35
```

```
Old Shape:  (442, 10)
New Shape:  (439, 10)
```

Experiment 3: Imputation of Missing Values KNN, Interpolation,Min-Max/Z-score.

```
 1 # import necessary libraries
 2 import numpy as np
 3 import pandas as pd
 4
 5 # import the KNNimputer class
 6 from sklearn.impute import KNNImputer
 7
 8
 9 # create dataset for marks of a student
10 dict = {'Maths': [80, 90, np.nan, 95],
11         'Chemistry': [60, 65, 56, np.nan],
12         'Physics': [np.nan, 57, 80, 78],
13         'Biology': [78, 83, 67, np.nan]}
14
15 # creating a data frame from the list
16 Before_imputation = pd.DataFrame(dict)
17 # print dataset before imputation
18 print("Data Before performing imputation\n", Before_imputation)
19
20 # create an object for KNNImputer
21 imputer = KNNImputer(n_neighbors=2)
22 After_imputation = imputer.fit_transform(Before_imputation)
23 # print dataset after performing the operation
24 print("\n\nAfter performing imputation\n", After_imputation)
25
```

```
Data Before performing imputation
     Maths  Chemistry  Physics  Biology
0    80.0       60.0      NaN     78.0
1    90.0       65.0     57.0     83.0
2     NaN       56.0     80.0     67.0
3    95.0        NaN     78.0      NaN


After performing imputation
 [[80.   60.   68.5 78. ]
 [90.   65.   57.   83. ]
 [87.5 56.   80.   67. ]
 [95.   58.   78.   72.5]]
```
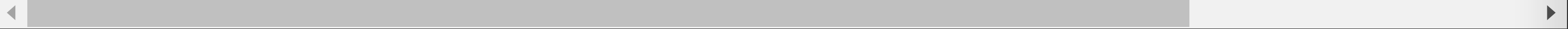
```
1 import pandas as pd
2 import numpy as np
3 a=pd.Series([0, 1, np.nan, 3,4,5,7])
4 a.interpolate()
```

|   | 0 |
|---|---|
| **0** | 0.0 |
| 1 | 1.0 |
| **2** | 2.0 |
| 3 | 3.0 |
| **4** | 4.0 |
| 5 | 5.0 |
| **6** | 7.0 |

dtype: float64

Experiment 4: Pivot table, Melt function.

```python
1 import pandas as pd
2 import numpy as np
3 from sklearn.datasets import fetch_openml
4
5 X,y = fetch_openml("autos", version=1, as_frame=True, return_X_y=True)
6 data = X
7 data['target'] = y
```

```python
1 pivot = np.round(pd.pivot_table(data, values='price',
2                                 index='num-of-doors',
3                                 columns='fuel-type',
4                                 aggfunc=np.mean),2)
5 pivot
```

<ipython-input-2-3c77819066fb>:1: FutureWarning: The default value of observed=False is deprecated and will change to observed=True in a future version of pandas. Specify observed=False to silence this warning and r
    pivot = np.round(pd.pivot_table(data, values='price',
<ipython-input-2-3c77819066fb>:1: FutureWarning: The provided callable <function mean at 0x7f58e7b6c280> is currently using DataFrameGroupBy.mean. In a future version of pandas, the provided callable will be used di
    pivot = np.round(pd.pivot_table(data, values='price',

| fuel-type | diesel | gas |
|---|---|---|
| num-of-doors | | |
| **four** | 16432.38 | 13092.81 |
| two | 14350.00 | 12762.76 |

```python
1 pivot = np.round(pd.pivot_table(data, values='price',
2                                 index=['num-of-doors', 'body-style'],
3                                 columns=['fuel-type', 'fuel-system'],
4                                 aggfunc=np.mean,
5                                 fill_value=0),2)
6 pivot
```

| | | fuel-type | diesel | gas | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | fuel-system | idi | 1bbl | 2bbl | 4bbl | mfi | mpfi | spdi | spfi |
| num-of-doors | body-style | | | | | | | | | |
| four | hatchback | | 7788.00 | 0.00 | 7813.71 | 0.0 | 0.0 | 10618.00 | 0.00 | 0.0 |
| | sedan | | 16328.92 | 8811.67 | 7711.19 | 0.0 | 0.0 | 18425.68 | 9279.00 | 0.0 |
| | wagon | | 19727.67 | 7295.00 | 8028.89 | 0.0 | 0.0 | 14213.42 | 0.00 | 0.0 |
| two | convertible | | 0.00 | 0.00 | 0.00 | 0.0 | 0.0 | 21890.50 | 0.00 | 0.0 |
| | hardtop | | 28176.00 | 0.00 | 8249.00 | 0.0 | 0.0 | 23540.50 | 0.00 | 0.0 |
| | hatchback | | 0.00 | 7054.43 | 6701.67 | 12145.0 | 12964.0 | 14581.50 | 11479.43 | 11048.0 |
| | sedan | | 7437.00 | 0.00 | 7570.00 | 0.0 | 0.0 | 21034.00 | 0.00 | 0.0 |

```
1  np.round(pd.pivot_table(data, values='price',
2                          index=['body-style'],
3                          columns=['num-of-doors'],
4                          aggfunc=[np.mean, np.median],
5                          fill_value=0),2)
```
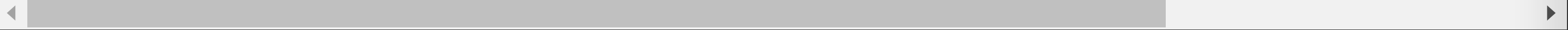
| | mean | | median | |
|---|---|---|---|---|
| num-of-doors | four | two | four | two |
| body-style | | | | |
| convertible | 0.00 | 21890.50 | 0.0 | 17084.5 |
| hardtop | 0.00 | 22208.50 | 0.0 | 19687.5 |
| hatchback | 8372.00 | 10230.79 | 8073.0 | 8970.0 |
| sedan | 14614.13 | 14283.00 | 12555.0 | 8678.0 |
| wagon | 12371.96 | 0.00 | 11694.0 | 0.0 |

```
1 import pandas as pd
2 d1 = {"Name": ["Tom", "Jerry", "Spike"], "ID": [1, 2, 3],
3     "Role": ["Cat", "Mouse", "Dog"]}
4 df = pd.DataFrame(d1)
5 print(df)
6
7 df_melted = pd.melt(df, id_vars=["ID"], value_vars=["Name", "Role"])
8 print(df_melted)
```

```
   Name  ID  Role
0   Tom   1   Cat
```

```
1 Jerry   2  Mouse
2 Spike   3   Dog
  ID variable  value
0  1   Name    Tom
1  2   Name  Jerry
2  3   Name  Spike
3  1   Role    Cat
4  2   Role  Mouse
5  3   Role    Dog
```

```
1 #multiple columns as id_vars
2 df_melted = pd.melt(df, id_vars=["ID", "Name"], value_vars=["Role"])
3 print(df_melted)
```

```
   ID   Name variable  value
0  1    Tom     Role    Cat
1  2  Jerry     Role  Mouse
2  3  Spike     Role    Dog
```

```
1 #skipping columns in melt function
2 df_melted = pd.melt(df, id_vars=["Name"], value_vars=["Role"])
3 print(df_melted)
```

```
    Name variable  value
0    Tom     Role    Cat
1  Jerry     Role  Mouse
2  Spike     Role    Dog
```

Experiment 5: Image normalization (rescale -1,1).

```
1 import numpy as np
2 from PIL import Image
3 import matplotlib.pyplot as plt
4
5 # Load an image
6 image_path = r"/content/mascot-logo-design_P1_900x420.jpg"
7
8 try:
9     image = Image.open(image_path)
10     image = image.convert('RGB')  # Ensure image is in RGB format
11 except FileNotFoundError:
12     print(f"File not found: {image_path}")
13     raise
14 except Exception as e:
15     print(f"Error occurred: {e}")
16     raise
17
18 # Convert image to numpy array
19 image_array = np.array(image).astype('float32')
20
21 # Normalize the image to the range [-1, 1]
22 normalized_image_array = (image_array / 127.5) - 1
23
24 # Verify normalization
25 print(f"Min value in normalized image: {np.min(normalized_image_array)}")
26 print(f"Max value in normalized image: {np.max(normalized_image_array)}")
27
28 # Convert normalized image back to the range [0, 255] for visualization
29 denormalized_image_array = (normalized_image_array + 1) * 127.5
30 denormalized_image_array = denormalized_image_array.astype('uint8')
31
32 # Display the original and normalized images
```

```
33 fig, axes = plt.subplots(1, 2, figsize=(12, 6))
34
35 axes[0].imshow(image)
36 axes[0].set_title("Original Image")
37 axes[0].axis('off')
38
39 axes[1].imshow(denormalized_image_array)
40 axes[1].set_title("Normalized Image (rescaled to [0, 255] for display)")
41 axes[1].axis('off')
42
43 plt.show()
44
```

Min value in normalized image: -1.0
Max value in normalized image: 1.0



Experiment 6: Linear regression and Logistic regression.

```
 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.linear_model import LinearRegression
 6 from sklearn.metrics import mean_squared_error, r2_score
 7
 8 # Create a synthetic dataset for linear regression
 9 np.random.seed(0)
10 X = 2 * np.random.rand(100, 1)
11 y = 4 + 3 * X + np.random.randn(100, 1)
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                                               random_state=0)
16
17 # Create and train the linear regression model
18 lin_reg = LinearRegression()
19 lin_reg.fit(X_train, y_train)
20
21 # Predict on the test set
22 y_pred = lin_reg.predict(X_test)
23
24 # Evaluate the model
25 mse = mean_squared_error(y_test, y_pred)
26 r2 = r2_score(y_test, y_pred)
27
28 print(f"Mean Squared Error: {mse}")
29 print(f"R^2 Score: {r2}")
30
31 # Plot the results
```

```python
32 plt.scatter(X_test, y_test, color='blue', label='Actual')
33 plt.plot(X_test, y_pred, color='red', label='Predicted')
34 plt.xlabel("X")
35 plt.ylabel("y")
36 plt.title("Linear Regression")
37 plt.legend()
38 plt.show()
```

Mean Squared Error: 1.0434333815695171