| Name | Bodhisatya Ghosh |
|---|---|
| UID no. | 2021700026 |

| Experiment 2 |
|---|

| HONOUR PLEDGE | " I hereby declare that the documentation, code and output attached with this lab experiment has been completed by me in accordance with the highest standards of honesty. I confirm that I have not plagiarized OR used unauthorized materials OR given or received illegitimate help for completing this experiment. I will uphold equity and honesty in the evaluation of my work and if found guilty of plagiarism or dishonesty, will bear the consequences as outlined in the 'integrity' section of the lab rubrics. I am doing so in order to maintain a community built around this code of honour" |
|---|---|
| PROBLEM STATEMENT : | **Data Cleaning and Preprocessing:**<br><br>● Handle missing values by imputing them with mean, median, or mode. Reason out which is more suitable (mean, median or mode) for your dataset and which is not<br><br>● Removing outliers based on a specific threshold. Give reasons for your choice of threshold. What do the outliers in your dataset tell you?<br><br>● Transform variables using log transformation or standardization. What possibly can go wrong when you do not standardize your data? What are the reasons for using log transformation on your variables and when should you definitely use it?<br><br>● Remove duplicate records from a data frame. In case your dataset does not have a exact duplicate rows, can you reason about strategies for identifying and deduplicating your dataset based on a subset of features?<br><br>● Standardize date formats across your dataset Is there a certain date-format that you would prefer? why? |
| THEORY: | |

**PROGRAM:**

**Name** : Bodhisatya Ghosh
**Class** : CSE DS
**UID** : 2021700026
**Subject** : BAP
**Experiment number** : 2

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

from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
```

```python
data = pd.read_csv('../exp 1/penguins_lter.csv')
```

```python
data.drop(columns=['studyName','Individual ID','Comments'], inplace=True)
data.head(1)
```

| | Sample Number | Species | Region | Island | Stage | Clutch Completion | Date Egg | Culmen Length (mm) | Culmen Depth (mm) | Flipper Length (mm) | Body Mass (g) | Sex | Delta 15 N (o/oo) | Delta 13 C (o/oo) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | Adelie Penguin (Pygoscelis adeliae) | Anvers | Torgersen | Adult, 1 Egg Stage | Yes | 11/11/07 | 39.1 | 18.7 | 181.0 | 3750.0 | MALE | NaN | NaN |

```python
data.isna().sum()
```

```
Sample Number        0
Species              0
Region               0
Island               0
Stage                0
Clutch Completion    0
Date Egg             0
Culmen Length (mm)   2
Culmen Depth (mm)    2
Flipper Length (mm)  2
Body Mass (g)        2
Sex                  10
Delta 15 N (o/oo)    14
Delta 13 C (o/oo)    13
dtype: int64
```
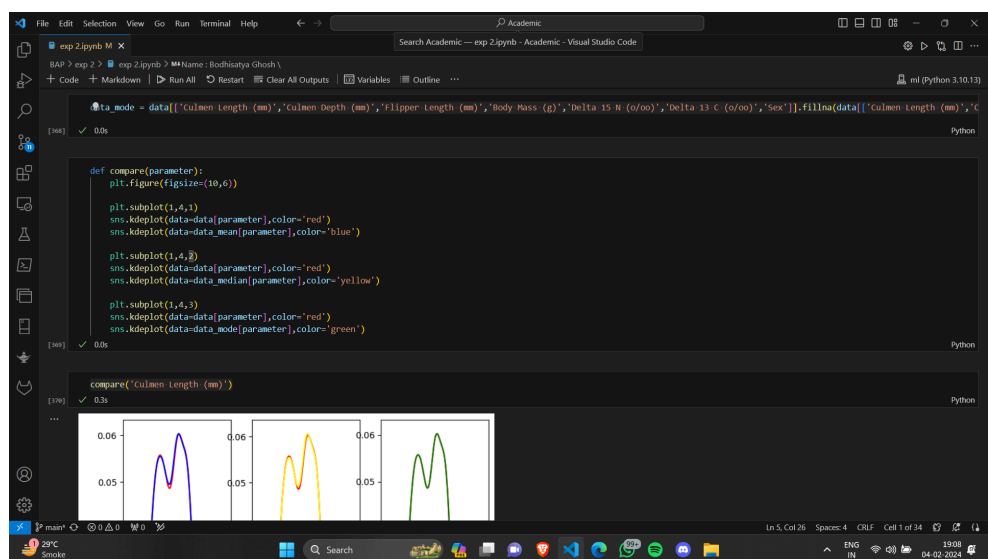
1. Handle missing values by imputing them with mean, median, or mode. Reason out which is more suitable (mean, median or mode) for your dataset and which is not

```python
data_mean = data[['Culmen Length (mm)','Culmen Depth (mm)','Flipper Length (mm)','Body Mass (g)','Delta 15 N (o/oo)','Delta 13 C (o/oo)']].fillna(data[['Culmen Length (mm)','Culmen
```

```python
data_median = data[['Culmen Length (mm)','Culmen Depth (mm)','Flipper Length (mm)','Body Mass (g)','Delta 15 N (o/oo)','Delta 13 C (o/oo)']].fillna(data[['Culmen Length (mm)','Culme
```

```python
data_mode = data[['Culmen Length (mm)','Culmen Depth (mm)','Flipper Length (mm)','Body Mass (g)','Delta 15 N (o/oo)','Delta 13 C (o/oo)','Sex']].fillna(data[['Culmen Length (mm)','C
```

```python
def compare(parameter):
    plt.figure(figsize=(10,6))

    plt.subplot(1,4,1)
    sns.kdeplot(data=data[parameter],color='red')
    sns.kdeplot(data=data_mean[parameter],color='blue')

    plt.subplot(1,4,2)
    sns.kdeplot(data=data[parameter],color='red')
    sns.kdeplot(data=data_median[parameter],color='yellow')

    plt.subplot(1,4,3)
    sns.kdeplot(data=data[parameter],color='red')
    sns.kdeplot(data=data_mode[parameter],color='green')
```

```python
compare('Culmen Length (mm)')
```

As we can see, using Mode to replace missing values produces least deviation from original data, hence, mode will be used as .fillna() parameter

```python
data.fillna(data[['Culmen Length (mm)','Culmen Depth (mm)','Flipper Length (mm)','Body Mass (g)','Delta 15 N (o/oo)','Delta 13 C (o/oo)']].mode(),inplace=True)
data.fillna(data['Delta 13 C (o/oo)'].mode(),inplace=True)
data.dropna(inplace=True)
```
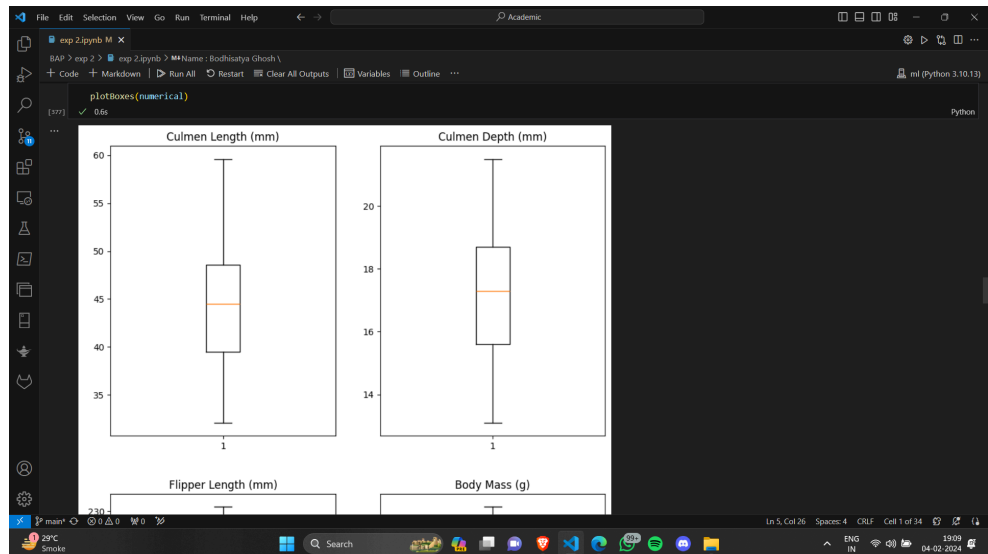
2. Removing outliers based on a specific threshold. Give reasons for your choice of threshold. What do the outliers in your dataset tell you?

```python
def plotBoxes(numerical):
    plt.figure(figsize=(10,20))
    cnt = 1
    for column in numerical:
        if(cnt>len(numerical)):
            break
        plt.subplot(3,2,cnt)
        plt.boxplot(numerical[column])
        plt.title(column)
        cnt+=1
```

```python
numerical = data[['Culmen Length (mm)','Culmen Depth (mm)','Flipper Length (mm)','Body Mass (g)','Delta 15 N (o/oo)','Delta 13 C (o/oo)']]
```

```python
plotBoxes(numerical)
```

As no outliers were found in the boxplots, no removal of outliers was needed

3. Transform variables using log transformation or standardization. What possibly can go wrong when you do not standardize your data? What are the reasons for using log transformation on your variables and when should you definitely use it?
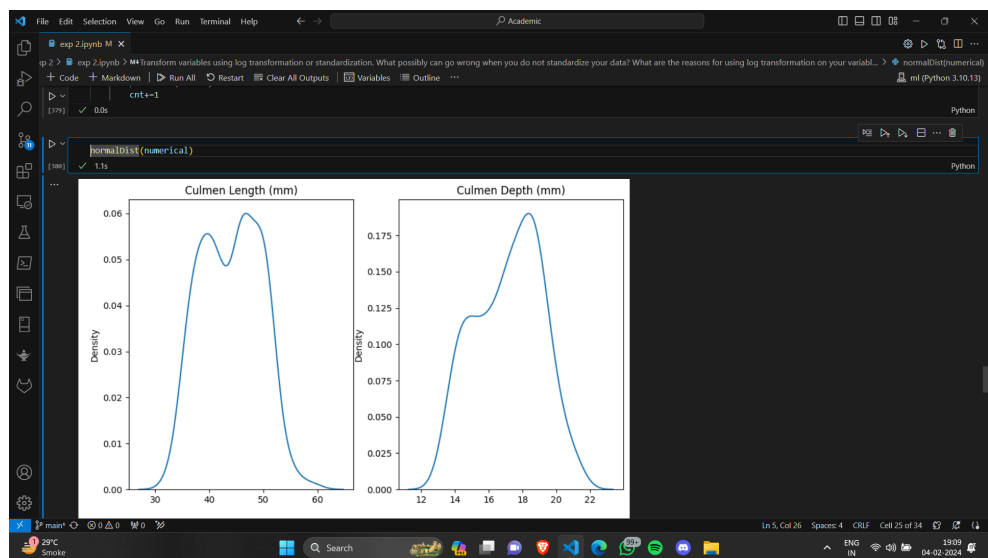
```python
from scipy.stats import norm
```

```python
def normalDist(data):
    plt.figure(figsize=(10,20))
    cnt = 1
    for column in data.columns:
        x_axis = np.array(data.loc[:,column])

        plt.subplot(3,2,cnt)

        sns.kdeplot(x_axis)
        plt.title(column)
        cnt+=1
```

```python
# numericalLog = np.log10(numerical)
normalDist(numerical)
```



Culmen Length (mm)    Culmen Depth (mm)

---

```python
        cnt+=1
```

```python
normalDist(numerical)
```



---

## Log Transformation Use Cases:

Skewed Data: When your data has a long tail or is positively skewed, log transformation can help normalize it. Multiplicative Relationships: Log transformation is also useful when there are multiplicative relationships between variables, convertin them into additive relationships.

As we can see 'Culmen Length (mm)','Flipper Length (mm)','Body Mass (g)','Delta 15 N (o/oo)','Delta 13 C (o/oo)' are all right skewed, hence we will use log transform

```python
numericalLog = np.log10(data[['Culmen Length (mm)','Flipper Length (mm)','Body Mass (g)','Delta 15 N (o/oo)']])
data[['Culmen Length (mm)','Flipper Length (mm)','Body Mass (g)','Delta 15 N (o/oo)']] = numericalLog
```
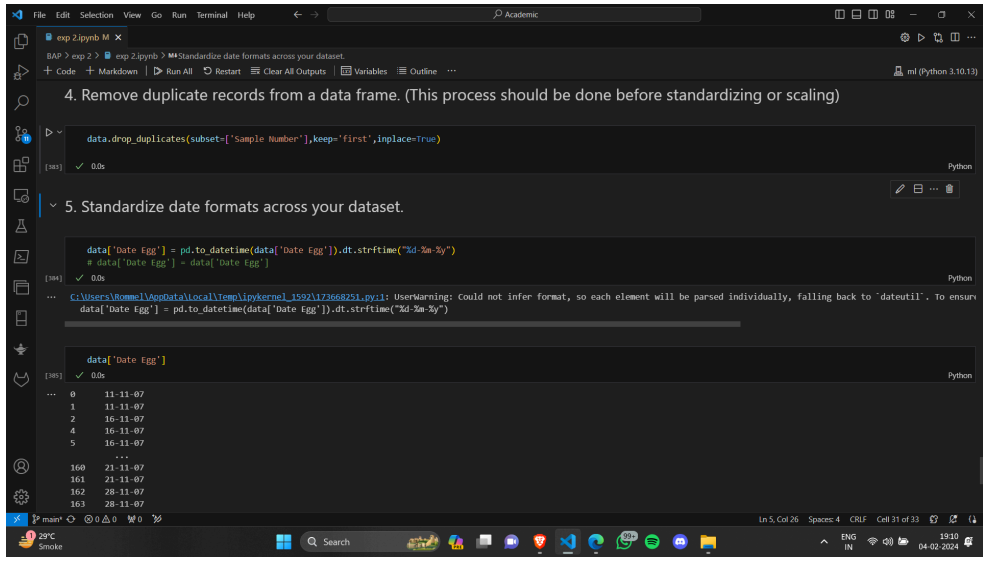
4. Remove duplicate records from a data frame. (This process should be done before standardizing or scaling)

```python
data.drop_duplicates(subset=['Sample Number'],keep='first',inplace=True)
```

```markdown
### 5. Standardize date formats across your dataset.
```

| **RESULT:** | 1. From the python notebook, we can see that using mode to fill null values gives us minimum deviation, hence we will be using mode in fillna().<br>2. From the boxplots, we could see that there are no outliers in the given data, hence we will not remove any data<br>3. If we do not standardize our data, the numerically larger values may end up dominating the model parameters and hinder in accurate calculations.<br>Log transforms are usually used for right skewed data.<br>4. Duplicates were removed using the Sample number as parameters<br>5. Dates were standardized to dd-mm-yy format. |
|---|---|
| **References:** | [How to change the Pandas datetime format in Python? - GeeksforGeeks](#) |

**CONCLUSION:** In this experiment, I have learnt how to handle missing data, fill in out of format data, check for outliers and how to deal with them.