# Data Collection and Preprocessing Phase
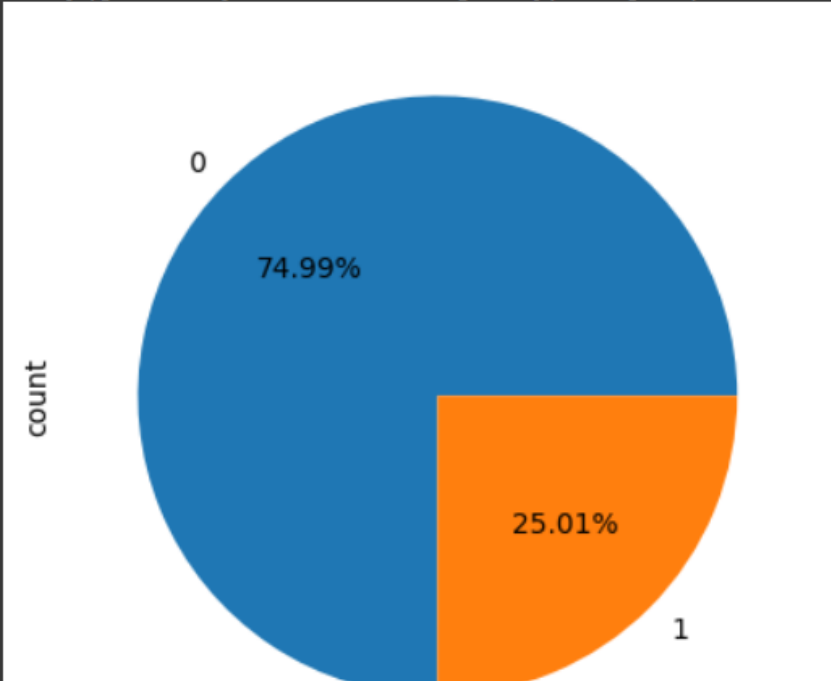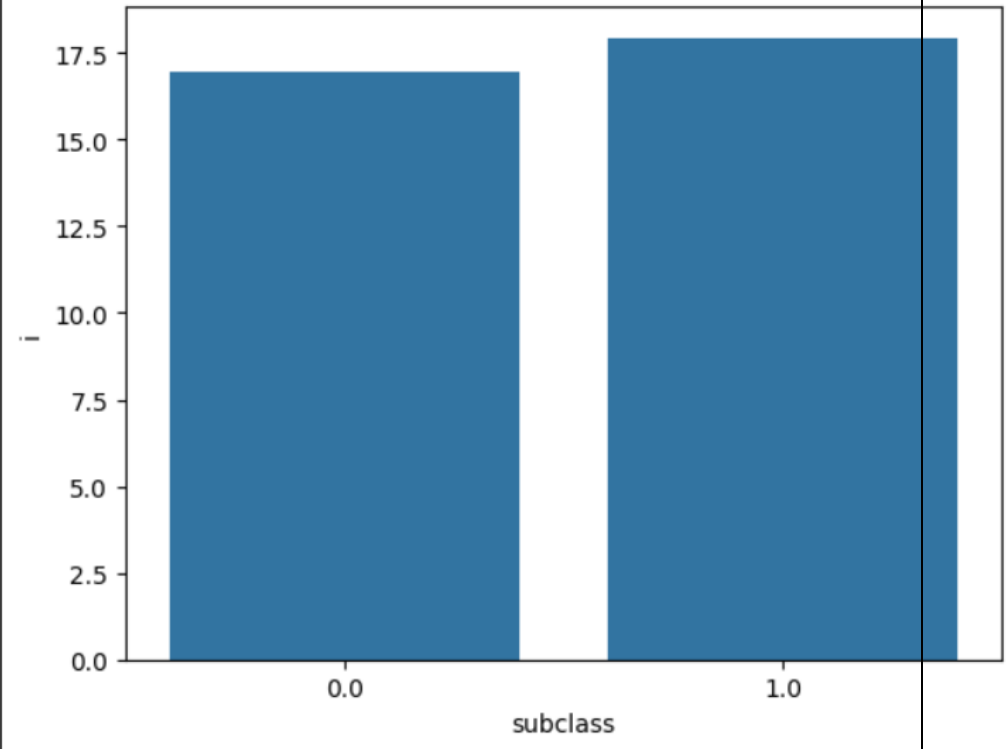
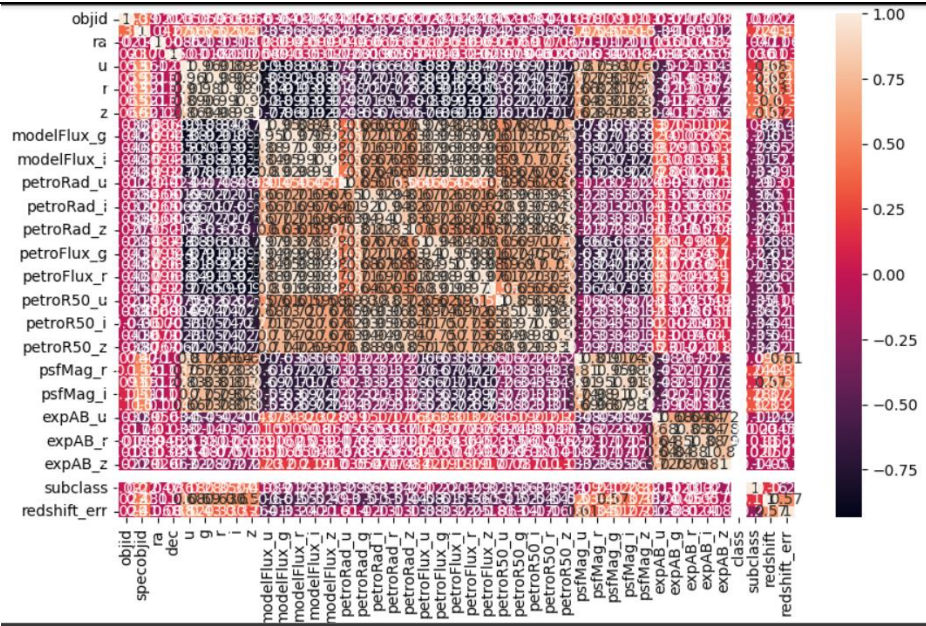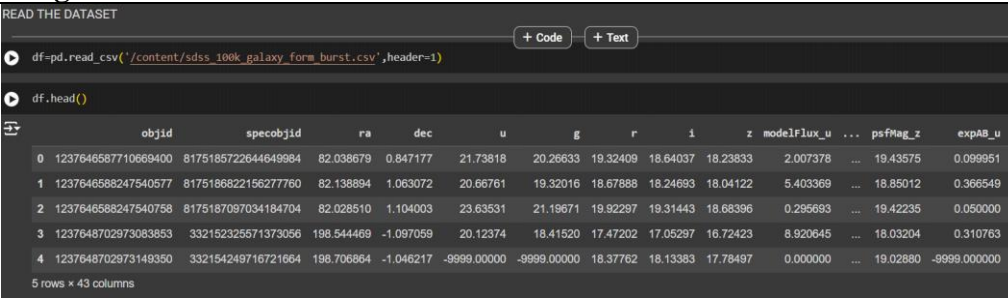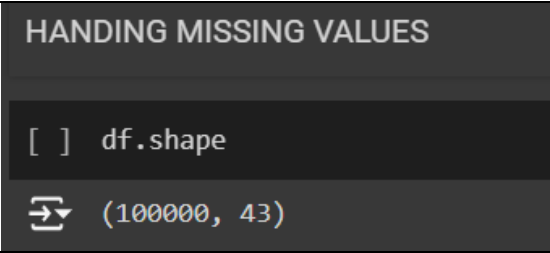| Date | 15 July 2024 |
|---|---|
| Team ID | 740051 |
| Project Title | **SDSS galaxy classification using Machine Learning** |
| Maximum Marks | 6 Marks |

**Data Exploration and Preprocessing Template**

Exploration and Preprocessing Template for SDSS galaxy classification for Machine Learning: Load data, handle missing values, explore basic statistics, visualize distributions, encode categorical variables, normalize/scale features, identify outliers, and prepare for modeling

| Section | Description |
|---|---|
| Data Overview | Summary of the dataset, including number of rows and columns, data types of each column, and brief descriptions of each column. |
| Univariate Analysis | Distribution analysis of individual variables using histograms, bar charts, and descriptive statistics (mean, median, mode, standard deviation). **#Univariate Analysis** |

```
[ ]  sub = df["subclass"].value_counts()
     sub
```

```
subclass
0    74993
1    25007
Name: count, dtype: int64
```

```
[ ]  sub.plot(kind="pie",subplots=True,autopct="%1.2f%%")
```

```
array([<Axes: ylabel='count'>], dtype=object)
```



| Bivariate Analysis | Examination of relationships between pairs of variables using scatter plots, correlation matrices, and pairwise plots to identify patterns and trends. **#Bivariate Analysis** |
| --- | --- |

| | |
|---|---|
| | **BIVARIATE ANALYSIS**<br><br>`[ ]  sns.barplot(x='subclass',y='i',data=df)`<br><br>`<Axes: xlabel='subclass', ylabel='i'>`<br><br> |
| Multivariate Analysis | Investigation of interactions between multiple variables using heatmaps, PCA (Principal Component Analysis), and clustering to understand data structure.<br><br>**MULTIVARIATE ANALYSIS**<br><br>```python<br>plt.figure(figsize=(10,6))<br>sns.heatmap(df.corr(),annot=True)<br>plt.show()<br>``` |

| | |
|---|---|
| Outliers and Anomalies | Identification and description of outliers and anomalies, summarized in a table with details on detection method, number of outliers, description, and potential impact. |

**Data Preprocessing Code Screenshots**

| | |
|---|---|
| Loading Data |  |
| Handling Missing Data |  |

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100000 entries, 0 to 99999
Data columns (total 43 columns):
 #   Column      Non-Null Count   Dtype
---  ------      --------------   -----
 0   objid       100000 non-null  int64
 1   specobjid   100000 non-null  uint64
 2   ra          100000 non-null  float64
 3   dec         100000 non-null  float64
 4   u           100000 non-null  float64
 5   g           100000 non-null  float64
 6   r           100000 non-null  float64
 7   i           100000 non-null  float64
 8   z           100000 non-null  float64
 9   modelFlux_u 100000 non-null  float64
 10  modelFlux_g 100000 non-null  float64
 11  modelFlux_r 100000 non-null  float64
 12  modelFlux_i 100000 non-null  float64
 13  modelFlux_z 100000 non-null  float64
 14  petroRad_u  100000 non-null  float64
 15  petroRad_g  100000 non-null  float64
 16  petroRad_i  100000 non-null  float64
 17  petroRad_r  100000 non-null  float64
 18  petroRad_z  100000 non-null  float64
```

```
 19  petroFlux_u    100000 non-null   float64
 20  petroFlux_g    100000 non-null   float64
 21  petroFlux_i    100000 non-null   float64
 22  petroFlux_r    100000 non-null   float64
 23  petroFlux_z    100000 non-null   float64
 24  petroR50_u     100000 non-null   float64
 25  petroR50_g     100000 non-null   float64
 26  petroR50_i     100000 non-null   float64
 27  petroR50_r     100000 non-null   float64
 28  petroR50_z     100000 non-null   float64
 29  psfMag_u       100000 non-null   float64
 30  psfMag_r       100000 non-null   float64
 31  psfMag_g       100000 non-null   float64
 32  psfMag_i       100000 non-null   float64
 33  psfMag_z       100000 non-null   float64
 34  expAB_u        100000 non-null   float64
 35  expAB_g        100000 non-null   float64
 36  expAB_r        100000 non-null   float64
 37  expAB_i        100000 non-null   float64
 38  expAB_z        100000 non-null   float64
 39  class          100000 non-null   object
 40  subclass       100000 non-null   object
 41  redshift       100000 non-null   float64
 42  redshift_err   100000 non-null   float64
dtypes: float64(39), int64(1), object(2), uint64(1)
memory usage: 32.8+ MB
```

For checking the null values, . isnull() function is used. To sum those null values we use . sum() function. From the above image we found that there are no null values present in our dataset. So we can skip handling the missing values step.

| | |
|---|---|
| Data Transformation | - |
| Feature Engineering | - |
| Save Processed Data | - |