# Data Collection and Preprocessing Phase

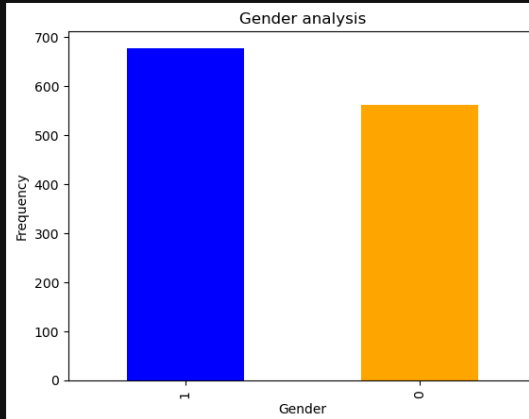| Date | 10 July 2024 |
|---|---|
| Team ID | SWTID1720078683 |
| Project Title | Anemia Sense: Leveraging Machine Learning for Precise Anemia Recognitions |
| Maximum Marks | 6 Marks |

**Data Exploration and Preprocessing**

Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

| Section | Description |
|---|---|
| Data Overview |  |

Univariate Analysis

```
[14]: gender = data['Gender'].value_counts()
      gender.plot(kind = 'bar',color = ['blue','orange'])
      plt.xlabel('Gender')
      plt.ylabel('Frequency')
      plt.title('Gender analysis')
```

```
[14]: Text(0.5, 1.0, 'Gender analysis')
```



```
[15]: sns.displot(data['Hemoglobin'],kde = True)
```

```
<seaborn.axisgrid.FacetGrid at 0x2677de5b190>
```

| | |
|---|---|
| Bivariate Analysis |  |
| Multivariate Analysis |  |

## Data Preprocessing Code Screenshots

| | |
|---|---|
| Loading Data | ```python
data = pd.read_csv('anemia.csv')
``` |
| Handling Missing Data | ```python
data.isnull().any()
```<br><br>```python
data.isnull().sum()
``` |
| Data Transformation | ```python
from sklearn.utils import resample


major = data[data['Result'] == 0]
minor = data[data['Result'] == 1]
undersampling = resample(major,replace = False,n_samples = len(minor),random_state = 47)
data = pd.concat([undersampling,minor])
print(data['Result'].value_counts())
```<br><br>```
Result
0    620
1    620
Name: count, dtype: int64
``` |
| Save Processed Data | ```python
data.to_csv('anemia.csv',index=False)
``` |