# Data Collection and Preprocessing Phase

| Date | 14 July 2024 |
|---|---|
| Team ID | 739872 |
| Project Title | Blood Donation Prediction |
| Maximum Marks | 6 Marks |

**Data Exploration and Preprocessing Template**

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 |  |

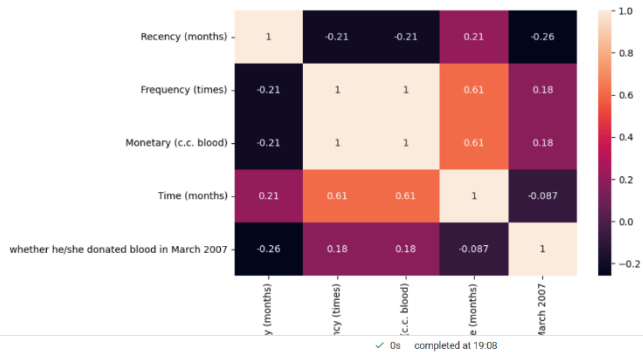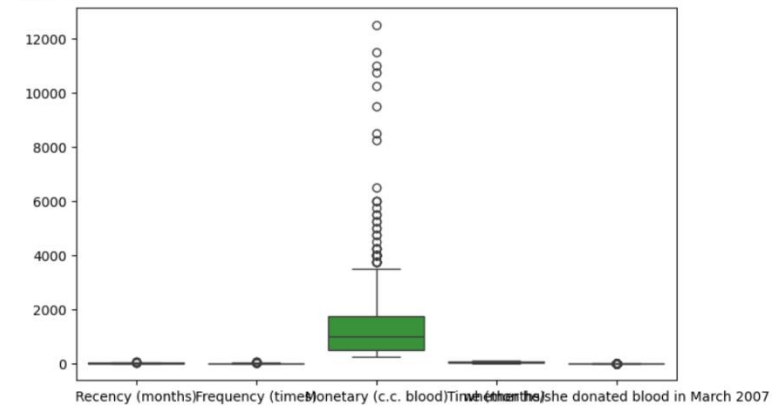| | |
|---|---|
| Bivariate Analysis | <br><br>**BIVARIATE ANALYSIS**<br><br>```python<br>sns.scatterplot(x=df['Monetary (c.c. blood)'],y=df['Time (months)'])<br>```<br>`<Axes: xlabel='Monetary (c.c. blood)', ylabel='Time (months)'>`<br><br> |
| Multivariate Analysis | **MULTIVARIATE ANALYSIS**<br><br>```python<br>[92] sns.heatmap(df.corr(),annot=True)<br>```<br>`<Axes: >`<br><br><br><br>✓ 0s completed at 19:08 |
| Outliers and Anomalies | ```python<br>sns.boxplot(df)<br>```<br>`<Axes: >`<br><br> |

## Data Preprocessing Code Screenshots

| | |
|---|---|
| Loading Data | ```python
[62] df = pd.read_csv("/content/transfusion (2).csv")
df.head()
```<br><br>Recency (months) Frequency (times) Monetary (c.c. blood) Time (months) whether he/she donated blood in March 2007<br>0    2    50    12500    98    1<br>1    0    13    3250    28    1<br>2    1    16    4000    35    1<br>3    2    20    5000    45    1<br>4    1    24    6000    77    0 |
| Handling Missing Data | ```python
df.isnull().sum()
```<br><br>Recency      0<br>Frequency    0<br>Monetary     0<br>Time         0<br>Donated      0<br>dtype: int64 |
| Data Transformation | ```python
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
scaled_features = scaler.fit_transform(df.drop(columns=["whether he/she donated blood in March 2007"]))
scaled_df = pd.DataFrame(scaled_features, columns=df.columns[:-1])

scaled_df["whether he/she donated blood in March 2007"] = df["whether he/she donated blood in March 2007"].valu

# Display the transformed DataFrame
print("\nTransformed DataFrame:")
print(scaled_df.head())
```<br><br>Transformed DataFrame:<br>    Recency (months)  Frequency (times)  Monetary (c.c. blood)  Time (months)  \<br>0    0.000000    0.923077    0.923077    0.270833<br>1    0.054054    0.230769    0.230769    0.020833<br>2    0.027027    0.461538    0.461538    0.125000<br>3    0.013514    0.846154    0.846154    0.343750<br>4    0.027027    0.615385    0.615385    0.208333<br><br>    whether he/she donated blood in March 2007<br>0    1<br>1    0<br>2    1<br>3    0<br>4    1 |
| Feature Engineering | Attached the codes in final submission. |
| Save Processed Data | - |