# Data Collection and Preprocessing Phase

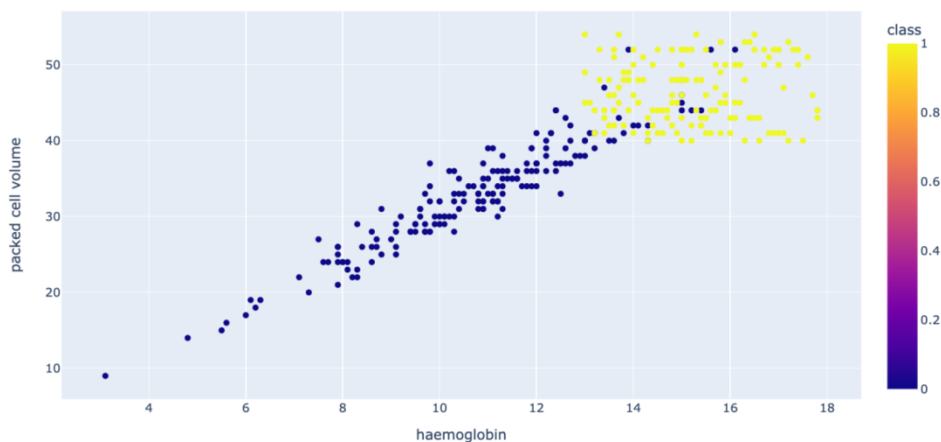| | |
|---|---|
| Date | 08 July 2024 |
| Team ID | SWTID1720174514 |
| Project Title | Early Prediction Of Chronic Kidney Disease Using Machine Learning |
| Maximum Marks | 6 Marks |

**Data Exploration and Preprocessing Template**

The variables of the dataset will be statistically examined to find general trends and extremes, and for this, a tool such as Python used for preprocessing like normalization and feature engineering activities. Data cleaning will find missing value analysis it determines the ways of handling missing values and outliers to improve the quality of the data in the upcoming analysis or modeling process.

| Section | Description |
|---|---|
| Data Overview |  |
| Univariate Analysis |  |

| | |
|---|---|
| Bivariate Analysis |  |
| Multivariate Analysis | `<Axes: xlabel='age', ylabel='blood pressure'>`<br> |

| Outliers and Anomalies |  |

**Data Preprocessing Code Screenshots**

| Loading Data |  |

| Handling Missing Data |  |

```
df['diabetes mellitus'].replace(to_replace = {'\tno':'no','\tyes':'yes',' yes':'yes'},inplace=True)

df['coronary artery disease'] = df['coronary artery disease'].replace(to_replace = '\tno', value='no')

df['class'] = df['class'].replace(to_replace = 'ckd\t', value = 'ckd')


for col in cat_col:
    print('{} has {} values  '.format(col, df[col].unique()))
    print('\n')
```

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