# Data Collection and Preprocessing Phase

| | |
|---|---|
| Date | 06 July 2024 |
| Team ID | SWTID1720082372 |
| Project Title | Early Prediction of Chronic Kidney Disease Using Machine Learning |
| Maximum Marks | 6 Marks |

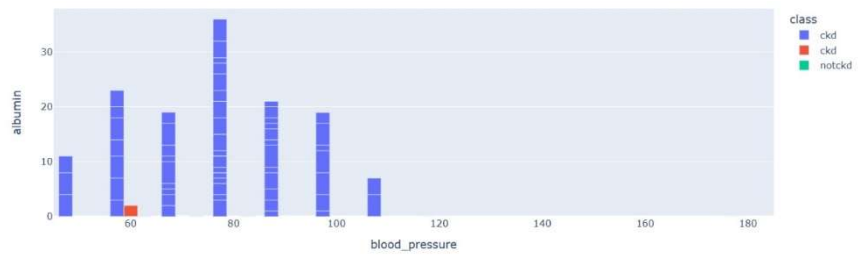**Data Exploration and Preprocessing Template**

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 | <u>Dimension</u>:<br>400 rows X 26 columns<br><u>Descriptive analysis:</u><br><br>|  | age | blood_pressure | specific_gravity | albumin | sugar | blood glucose random | blood_urea | serum_creatinine | sodium | potassium | hemoglobin |<br>|---|---|---|---|---|---|---|---|---|---|---|---|<br>| count | 391.000000 | 388.000000 | 353.000000 | 354.000000 | 351.000000 | 356.000000 | 381.000000 | 383.000000 | 313.000000 | 312.000000 | 348.000000 |<br>| mean | 51.483376 | 76.469072 | 1.017408 | 1.016949 | 0.450142 | 148.036517 | 57.425722 | 3.072454 | 137.528754 | 4.627244 | 12.526437 |<br>| std | 17.169714 | 13.683637 | 0.005717 | 1.352679 | 1.099191 | 79.281714 | 50.503006 | 5.741126 | 10.408752 | 3.193904 | 2.912587 |<br>| min | 2.000000 | 50.000000 | 1.005000 | 0.000000 | 0.000000 | 22.000000 | 1.500000 | 0.400000 | 4.500000 | 2.500000 | 3.100000 |<br>| 25% | 42.000000 | 70.000000 | 1.010000 | 0.000000 | 0.000000 | 99.000000 | 27.000000 | 0.900000 | 135.000000 | 3.800000 | 10.300000 |<br>| 50% | 55.000000 | 80.000000 | 1.020000 | 0.000000 | 0.000000 | 121.000000 | 42.000000 | 1.300000 | 138.000000 | 4.400000 | 12.650000 |<br>| 75% | 64.500000 | 80.000000 | 1.020000 | 2.000000 | 0.000000 | 163.000000 | 66.000000 | 2.800000 | 142.000000 | 4.900000 | 15.000000 |<br>| max | 90.000000 | 180.000000 | 1.025000 | 5.000000 | 5.000000 | 490.000000 | 391.000000 | 76.000000 | 163.000000 | 47.000000 | 17.800000 | |
| Univariate Analysis |  |

| | |
|---|---|
| |  Histogram of age / Boxplot of age |
| Bivariate Analysis |  Relationship between Blood Age and Hemoglobin<br><br> Correlation Matrix |

| | |
|---|---|
| Multivariate Analysis |  |
| Outliers and Anomalies | - |

## Data Preprocessing Code Screenshots

| | |
|---|---|
| Loading Data | **Loading Datatset**<br> |
| Handling Missing Data | **Handling missing values of red blood cell and pus cell**<br> |

| | |
|---|---|
| Data Transformation | **Drop id Column**<br><br>`[9]: data.drop(["id"],axis=1, inplace=True)`<br><br>**Renaming Columns**<br><br>`[10]: data.columns`<br><br>`[10]: Index(['age', 'bp', 'sg', 'al', 'su', 'rbc', 'pc', 'pcc', 'ba', 'bgr', 'bu',`<br>`        'sc', 'sod', 'pot', 'hemo', 'pcv', 'wc', 'rc', 'htn', 'dm', 'cad',`<br>`        'appet', 'pe', 'ane', 'classification'],`<br>`        dtype='object')`<br><br>`[11]: data.columns=['age','blood_pressure','specific_gravity','albumin',`<br>`       'sugar','red_blood_cells','pus_cell','pus_cell_clumps','bacteria', 'blood glucose random','blood_urea','serum_creatinine','sodium','potassium',`<br>`       'hemoglobin','packed_cell_volume','white_blood_cell_count','red_blood_cell_count',`<br>`       'hypertension','diabetesmellitus','coronary_artery_disease', 'appetite', 'pedal_edema','anemia','class']`<br><br>`[12]: data.columns`<br><br>**Correcting Data type**<br><br>`[48]: features=['red_blood_cells','packed_cell_volume','white_blood_cell_count']`<br>`      def convert_dtype(data,feature):`<br>`          data[feature] = pd.to_numeric(data[feature], errors='coerce')`<br><br>`      for feature in features:`<br>`          convert_dtype(data,feature)`<br><br>`      data.dtypes`<br><br>**Cleaning categorical columns**<br><br>`[49]: cat_col=[col for col in data.columns if data[col].dtype=='object']`<br>`      for col in cat_col:`<br>`          print('{} has {} values '.format(col,data[col].unique()))`<br>`          print('\n')` |
| Feature Engineering | Attached the codes in final submission. |
| Save Processed Data | - |