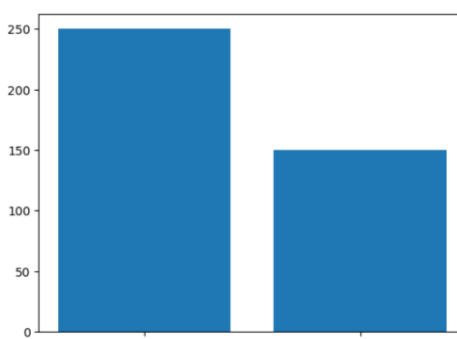
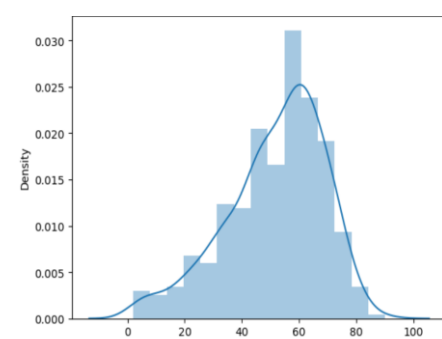


Data Collection and Preprocessing Phase

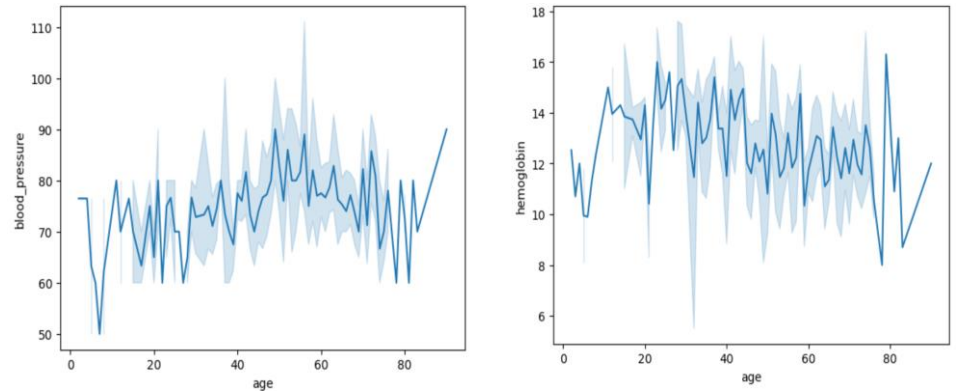
Date	12 July 2024
Team ID	SWTID1720083491
Project Title	Early Prediction of Chronic Kidney Disease Using Machine Learning
Maximum Marks	6 Marks

Data Exploration and Preprocessing Template

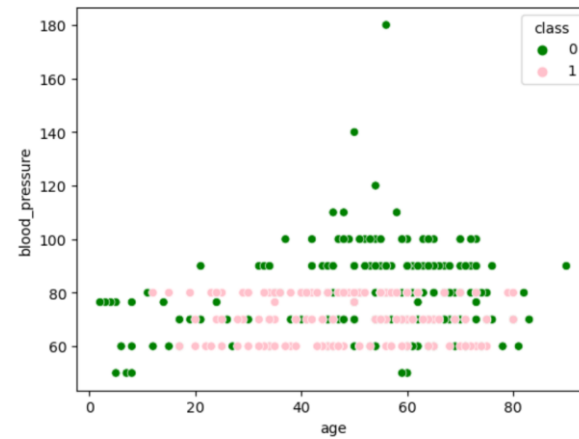
Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

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

Bivariate Analysis



Multivariate Analysis



Outliers and Anomalies

```
import pandas as pd
import numpy as np

def remove_outliers_with_median(fd, cols_to_check, threshold=3.5):

    for col in cols_to_check:
        if pd.api.types.is_numeric_dtype(fd[col]): # Check for numeric data type
            # Calculate z-scores efficiently using vectorized operations
            z_scores = np.abs(stats.zscore(fd[col]))

            # Identify outliers based on z-score threshold
            outlier_indices = np.where(z_scores > threshold)[0]

            # Replace outliers with median in one step
            fd.loc[outlier_indices, col] = fd[col].median()

    return fd

cols_to_check = ["bgr", "bu", "wc", "sc", "bp"] # Your continuous columns
df_without_outliers = remove_outliers_with_median(df1.copy(), cols_to_check)

df_without_outliers.isna().sum()
```

Data Preprocessing Code Screenshots

Loading Data	<pre>#2. Read the Dataset df = pd.read_csv('chronickidneydisease.csv') df</pre> <table><thead><tr><th></th><th>id</th><th>age</th><th>bp</th><th>sg</th><th>al</th><th>su</th><th>rbc</th><th>pc</th><th>pcc</th><th>ba</th><th>...</th><th>pcv</th><th>wc</th><th>rc</th><th>htn</th><th>dm</th><th>cad</th><th>appet</th><th>pe</th><th>ane</th><th>classification</th></tr></thead><tbody><tr><td>0</td><td>0</td><td>48.0</td><td>80.0</td><td>1.020</td><td>1.0</td><td>0.0</td><td>NaN</td><td>normal</td><td>notpresent</td><td>notpresent</td><td>...</td><td>44</td><td>7800</td><td>5.2</td><td>yes</td><td>yes</td><td>no</td><td>good</td><td>no</td><td>no</td><td>ckd</td></tr><tr><td>1</td><td>1</td><td>7.0</td><td>50.0</td><td>1.020</td><td>4.0</td><td>0.0</td><td>NaN</td><td>normal</td><td>notpresent</td><td>notpresent</td><td>...</td><td>38</td><td>6000</td><td>NaN</td><td>no</td><td>no</td><td>no</td><td>good</td><td>no</td><td>no</td><td>ckd</td></tr><tr><td>2</td><td>2</td><td>62.0</td><td>80.0</td><td>1.010</td><td>2.0</td><td>3.0</td><td>normal</td><td>normal</td><td>notpresent</td><td>notpresent</td><td>...</td><td>31</td><td>7500</td><td>NaN</td><td>no</td><td>yes</td><td>no</td><td>poor</td><td>no</td><td>yes</td><td>ckd</td></tr><tr><td>3</td><td>3</td><td>48.0</td><td>70.0</td><td>1.005</td><td>4.0</td><td>0.0</td><td>normal</td><td>abnormal</td><td>present</td><td>notpresent</td><td>...</td><td>32</td><td>6700</td><td>3.9</td><td>yes</td><td>no</td><td>no</td><td>poor</td><td>yes</td><td>yes</td><td>ckd</td></tr><tr><td>4</td><td>4</td><td>51.0</td><td>80.0</td><td>1.010</td><td>2.0</td><td>0.0</td><td>normal</td><td>normal</td><td>notpresent</td><td>notpresent</td><td>...</td><td>35</td><td>7300</td><td>4.6</td><td>no</td><td>no</td><td>no</td><td>good</td><td>no</td><td>no</td><td>ckd</td></tr></tbody></table>		id	age	bp	sg	al	su	rbc	pc	pcc	ba	...	pcv	wc	rc	htn	dm	cad	appet	pe	ane	classification	0	0	48.0	80.0	1.020	1.0	0.0	NaN	normal	notpresent	notpresent	...	44	7800	5.2	yes	yes	no	good	no	no	ckd	1	1	7.0	50.0	1.020	4.0	0.0	NaN	normal	notpresent	notpresent	...	38	6000	NaN	no	no	no	good	no	no	ckd	2	2	62.0	80.0	1.010	2.0	3.0	normal	normal	notpresent	notpresent	...	31	7500	NaN	no	yes	no	poor	no	yes	ckd	3	3	48.0	70.0	1.005	4.0	0.0	normal	abnormal	present	notpresent	...	32	6700	3.9	yes	no	no	poor	yes	yes	ckd	4	4	51.0	80.0	1.010	2.0	0.0	normal	normal	notpresent	notpresent	...	35	7300	4.6	no	no	no	good	no	no	ckd
	id	age	bp	sg	al	su	rbc	pc	pcc	ba	...	pcv	wc	rc	htn	dm	cad	appet	pe	ane	classification																																																																																																																
0	0	48.0	80.0	1.020	1.0	0.0	NaN	normal	notpresent	notpresent	...	44	7800	5.2	yes	yes	no	good	no	no	ckd																																																																																																																
1	1	7.0	50.0	1.020	4.0	0.0	NaN	normal	notpresent	notpresent	...	38	6000	NaN	no	no	no	good	no	no	ckd																																																																																																																
2	2	62.0	80.0	1.010	2.0	3.0	normal	normal	notpresent	notpresent	...	31	7500	NaN	no	yes	no	poor	no	yes	ckd																																																																																																																
3	3	48.0	70.0	1.005	4.0	0.0	normal	abnormal	present	notpresent	...	32	6700	3.9	yes	no	no	poor	yes	yes	ckd																																																																																																																
4	4	51.0	80.0	1.010	2.0	0.0	normal	normal	notpresent	notpresent	...	35	7300	4.6	no	no	no	good	no	no	ckd																																																																																																																
Handling Missing Data	<pre># Handling the Null Values # with mean() df['blood_glucose_random'].fillna(df['blood_glucose_random'].mean(),inplace=True) df['blood_pressure'].fillna(df['blood_pressure'].mean(),inplace=True) df['blood_urea'].fillna(df['blood_urea'].mean(),inplace=True) df['hemoglobin'].fillna(df['hemoglobin'].mean(),inplace=True) df['packed_cell_volume'].fillna(df['packed_cell_volume'].mean(),inplace=True) df['potassium'].fillna(df['potassium'].mean(),inplace=True) df['red_blood_cell_count'].fillna(df['red_blood_cell_count'].mean(),inplace=True) df['serum_creatinine'].fillna(df['serum_creatinine'].mean(),inplace=True) df['sodium'].fillna(df['sodium'].mean(),inplace=True) df['white_blood_cell_count'].fillna(df['white_blood_cell_count'].mean(),inplace=True) # with mode() df['age'].fillna(df['age'].mode()[0],inplace=True) df['hypertension'].fillna(df['hypertension'].mode()[0],inplace=True) df['pus_cell_clumps'].fillna(df['pus_cell_clumps'].mode()[0],inplace=True) df['appetite'].fillna(df['appetite'].mode()[0],inplace=True) df['albumin'].fillna(df['albumin'].mode()[0],inplace=True) df['pus_cell'].fillna(df['pus_cell'].mode()[0],inplace=True) df['red_blood_cell'].fillna(df['red_blood_cell'].mode()[0],inplace=True) df['coronary_artery_disease'].fillna(df['coronary_artery_disease'].mode()[0],inplace=True) df['bacteria'].fillna(df['bacteria'].mode()[0],inplace=True) df['anemia'].fillna(df['anemia'].mode()[0],inplace=True) df['sugar'].fillna(df['sugar'].mode()[0],inplace=True) df['diabetesmellitus'].fillna(df['diabetesmellitus'].mode()[0],inplace=True) df['pedal_edema'].fillna(df['pedal_edema'].mode()[0],inplace=True) df['specific_gravity'].fillna(df['specific_gravity'].mode()[0],inplace=True)</pre>																																																																																																																																				
Data Transformation	<pre># Label Encoding for i in catcols: print("Label Encoding of:",i) le= LabelEncoder() print(c(df[i])) df[i]=le.fit_transform(df[i]) print(c(df[i])) print("*****100)</pre>																																																																																																																																				

	<pre># Scaling - MinMax Scaler from sklearn.preprocessing import MinMaxScaler scale=MinMaxScaler() X_scaled=scale.fit_transform(x) X_scaled array([[1. , 1. , 0.21153846, ..., 0. , 1. , 0.], [1. , 1. , 0.2693088 , ..., 0. , 0. , 0.], [1. , 1. , 0.85683761, ..., 1. , 1. , 0.], ..., [1. , 1. , 0.16666667, ..., 0. , 0. , 0.], [1. , 1. , 0.1965812 , ..., 0. , 0. , 0.], [1. , 1. , 0.23290598, ..., 0. , 0. , 0.]])</pre>
Feature Engineering	<pre>#Rectifying the unknown class "ckd\t" df['class']=df['class'].replace("ckd\t","ckd") df['class'].unique() # Removing the columns which are non Numerical contcols.remove('specific_gravity') contcols.remove('albumin') contcols.remove('sugar') print(contcols) # Adding columns which we found Continuous contcols.add('red_blood_cell_count') contcols.add('packed_cell_volume') contcols.add('white_blood_cell_count') print(contcols) # Rectifying the Categorical columns Classes df['coronary_artery_disease'] = df.coronary_artery_disease.replace('\tno','no') c(df['coronary_artery_disease']) Counter({'no': 364, 'yes': 34, nan: 2}) df['diabetesmellitus'] = df.diabetesmellitus.replace(to_replace={'\tno':'no','\tyes':'yes',' yes':'yes'}) c(df['diabetesmellitus']) Counter({'no': 261, 'yes': 137, nan: 2})</pre>

Save Processed Data

```
# Choose a filename for the CSV (replace 'your_file.csv' with your desired name)
filename = 'Final_Data.csv'

# Export the DataFrame to a CSV file
df_filtered.to_csv(filename, index=False) # Optional: set index=False to exclude the row index from the CSV
print(f'DataFrame exported to CSV file: {filename}')
```