

Data Collection and Preprocessing Phase

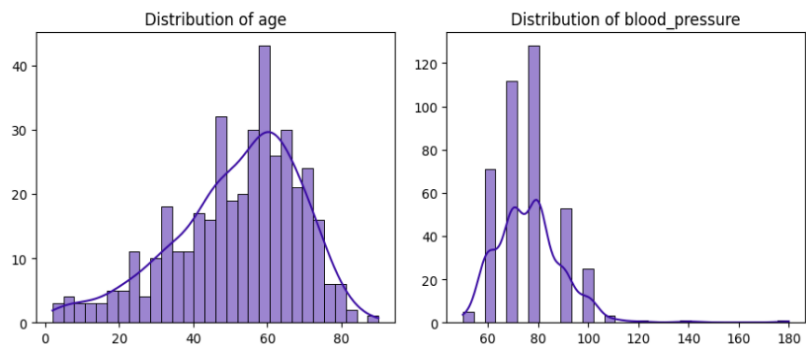
Date	22 June 2025
Team ID	SWTID1749634408
Project Title	Early Prediction for Chronic Kidney Disease Detection: A Progressive Approach to Health Management
Maximum Marks	6 Marks

Data Exploration and Preprocessing Report

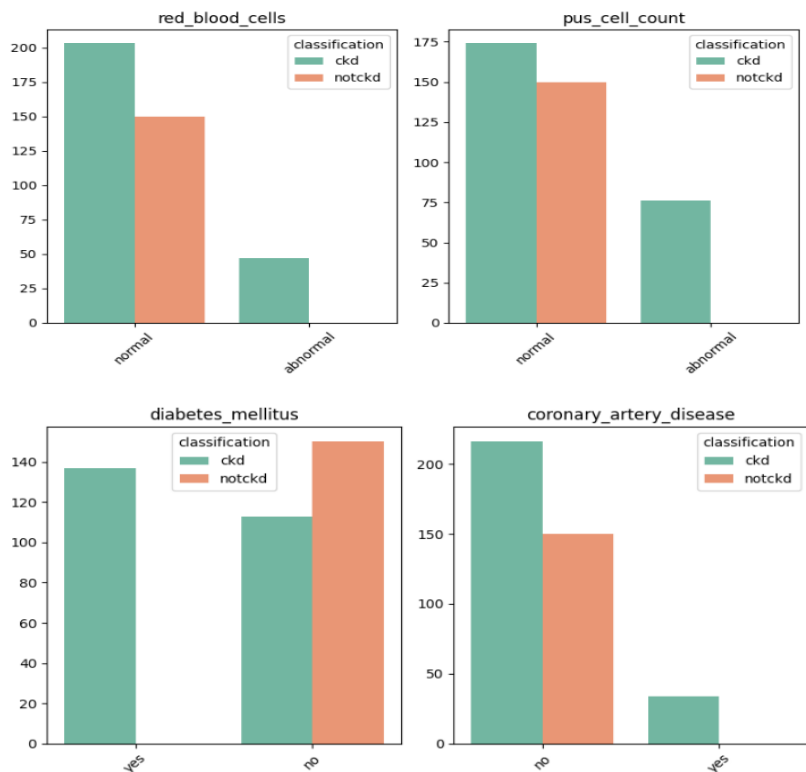
Dataset variables will be statistically analyzed to identify patterns and outliers, with Python employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modeling, and forming a strong foundation for insights and predictions.

Section	Description
Data Overview	<u>Dimension:</u> 614 rows × 13 columns
	<u>Descriptive statistics:</u>

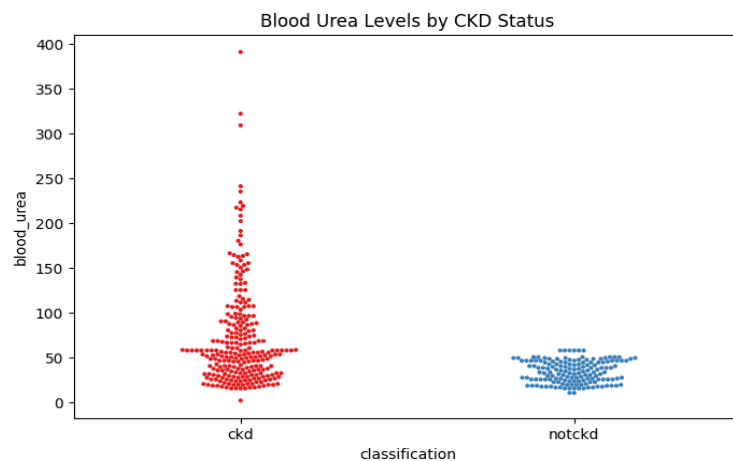
Univariate Analysis



Bivariate Analysis



Multivariate Analysis



Outliers and Anomalies	-																																																																		
Data Preprocessing Code Screenshots																																																																			
Loading Data	<div>df.head()</div> <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></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></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></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></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></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></tr></tbody></table>		id	age	bp	sg	al	su	rbc	pc	pcc	ba	0	0	48.0	80.0	1.020	1.0	0.0	NaN	normal	notpresent	notpresent	1	1	7.0	50.0	1.020	4.0	0.0	NaN	normal	notpresent	notpresent	2	2	62.0	80.0	1.010	2.0	3.0	normal	normal	notpresent	notpresent	3	3	48.0	70.0	1.005	4.0	0.0	normal	abnormal	present	notpresent	4	4	51.0	80.0	1.010	2.0	0.0	normal	normal	notpresent	notpresent
	id	age	bp	sg	al	su	rbc	pc	pcc	ba																																																									
0	0	48.0	80.0	1.020	1.0	0.0	NaN	normal	notpresent	notpresent																																																									
1	1	7.0	50.0	1.020	4.0	0.0	NaN	normal	notpresent	notpresent																																																									
2	2	62.0	80.0	1.010	2.0	3.0	normal	normal	notpresent	notpresent																																																									
3	3	48.0	70.0	1.005	4.0	0.0	normal	abnormal	present	notpresent																																																									
4	4	51.0	80.0	1.010	2.0	0.0	normal	normal	notpresent	notpresent																																																									
Handling Missing Data	<pre>df['age'] = df['age'].fillna(df['age'].mode()[0]) df['hypertension'] = df['hypertension'].fillna(df['hypertension'].mode()[0]) df['pus_cell_clumps'] = df['pus_cell_clumps'].fillna(df['pus_cell_clumps'].mode()[0]) df['appetite'] = df['appetite'].fillna(df['appetite'].mode()[0]) df['albumin'] = df['albumin'].fillna(df['albumin'].mode()[0]) df['pus_cell_count'] = df['pus_cell_count'].fillna(df['pus_cell_count'].mode()[0]) df['red_blood_cells'] = df['red_blood_cells'].fillna(df['red_blood_cells'].mode()[0]) df['coronary_artery_disease'] = df['coronary_artery_disease'].fillna(df['coronary_artery_disease'].mode()[0]) df['bacteria'] = df['bacteria'].fillna(df['bacteria'].mode()[0]) df['anemia'] = df['anemia'].fillna(df['anemia'].mode()[0]) df['sugar'] = df['sugar'].fillna(df['sugar'].mode()[0]) df['diabetes_mellitus'] = df['diabetes_mellitus'].fillna(df['diabetes_mellitus'].mode()[0]) df['pedal_edema'] = df['pedal_edema'].fillna(df['pedal_edema'].mode()[0]) df['specific_gravity'] = df['specific_gravity'].fillna(df['specific_gravity'].mode()[0])</pre>																																																																		
Data Transformation	<pre>df['diabetes_mellitus'] = df['diabetes_mellitus'].replace(to_replace = {'\tno' : 'no', '\tyes' : 'yes', 'yes' : 'yes'}) df['coronary_artery_disease'] = df['coronary_artery_disease'].replace(to_replace = {'\tno' : 'no'}) df['classification'] = df['classification'].replace(to_replace = {'ckd\t' : 'ckd'}) df['packed_cell_volume'] = pd.to_numeric(df['packed_cell_volume'], errors = 'coerce') df['white_blood_cell_count'] = pd.to_numeric(df['white_blood_cell_count'], errors = 'coerce') df['red_blood_cell_count'] = pd.to_numeric(df['red_blood_cell_count'], errors = 'coerce')</pre>																																																																		
Feature Engineering	Attached the codes in final submission.																																																																		
Save Processed Data	-																																																																		