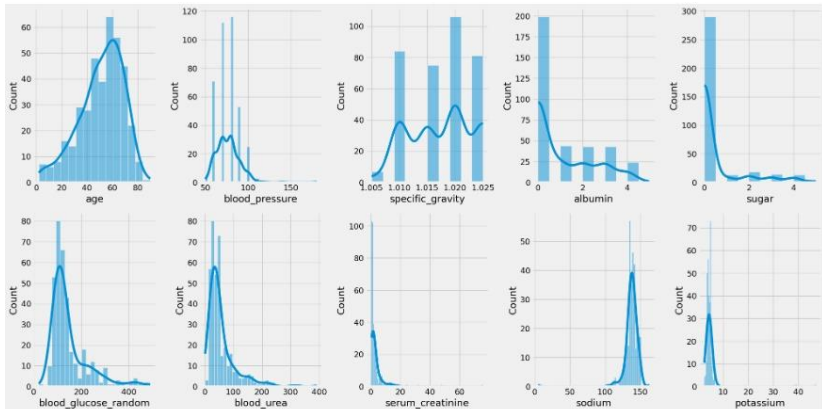


Data Collection and Preprocessing Phase

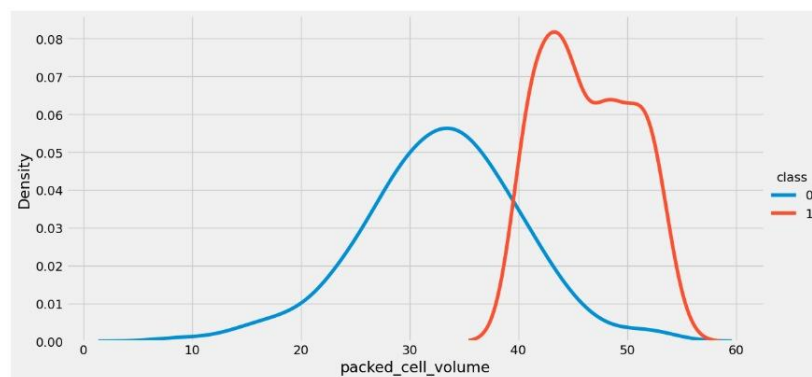
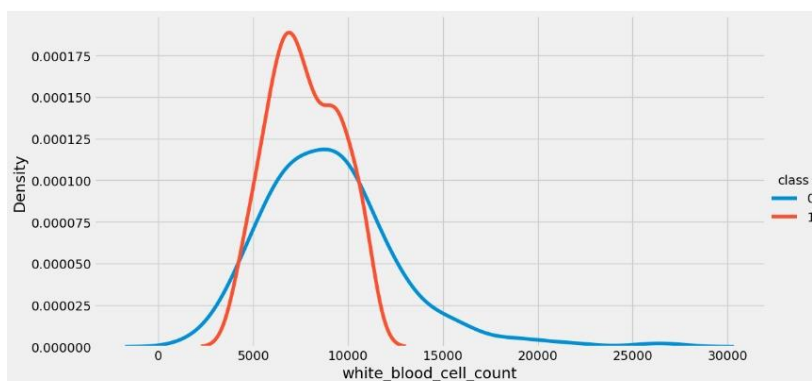
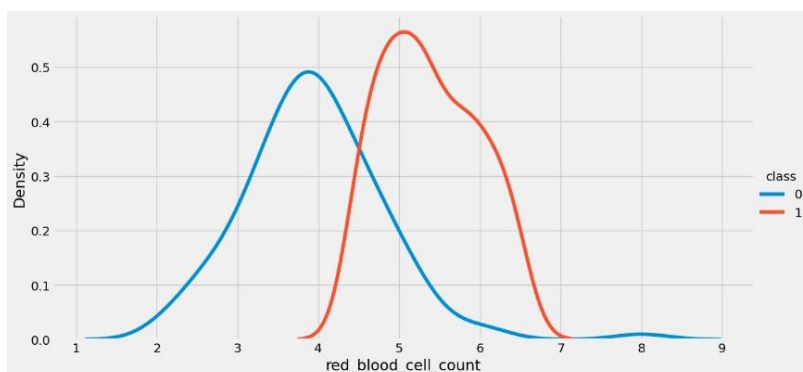
Date	4 th June 2024
Team ID	SWTID1720164961
Project Title	Early Prediction of Chronic Kidney Disease Using Machine Learning
Maximum Marks	6 Marks

Data Exploration and Preprocessing Template

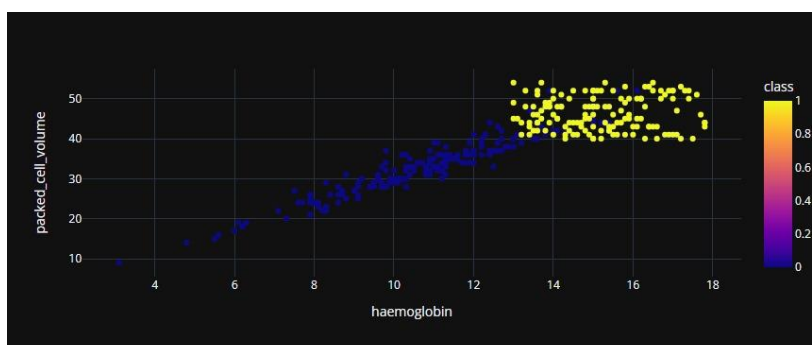
Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

Section	Description																																																																																																																																				
Data Overview	<div>400 rows 26 columns</div> <div><pre>[5]: data.head()</pre></div> <table><tr><th>[5]:</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>cla</th></tr><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></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></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></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></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></td></tr></table> <div>5 rows x 26 columns</div>	[5]:	id	age	bp	sg	al	su	rbc	pc	pcc	ba	...	pcv	wc	rc	htn	dm	cad	appet	pe	ane	cla	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		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		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		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		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	
[5]:	id	age	bp	sg	al	su	rbc	pc	pcc	ba	...	pcv	wc	rc	htn	dm	cad	appet	pe	ane	cla																																																																																																																
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																																																																																																																	
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																																																																																																																	
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																																																																																																																	
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																																																																																																																	
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																																																																																																																	
Univariate Analysis	<div></div>																																																																																																																																				

Bivariate Analysis



Multivariate Analysis



Outliers and Anomalies	NA																																																																																																																															
Data Preprocessing Code Screenshots																																																																																																																																
Loading Data	<div><pre>[5]: data.head()</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>cla</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></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></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></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></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></tr></tbody></table><p>5 rows × 26 columns</p></div>		id	age	bp	sg	al	su	rbc	pc	pcc	ba	...	pcv	wc	rc	htn	dm	cad	appet	pe	ane	cla	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	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	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	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	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
	id	age	bp	sg	al	su	rbc	pc	pcc	ba	...	pcv	wc	rc	htn	dm	cad	appet	pe	ane	cla																																																																																																											
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																																																																																																												
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																																																																																																												
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																																																																																																												
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																																																																																																												
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																																																																																																												
Handling Missing Data	<div><pre>[66]: # filling null values, we will use two methods, random sampling for higher null values and # mean/mode sampling for lower null values def random_value_imputation(feature): random_sample = data[feature].dropna().sample(data[feature].isna().sum()) random_sample.index = data[data[feature].isnull()].index data.loc[data[feature].isnull(), feature] = random_sample def impute_mode(feature): mode = data[feature].mode()[0] data[feature] = data[feature].fillna(mode)</pre><pre>[67]: # filling num_cols null values using random sampling method for col in num_cols: random_value_imputation(col)</pre><div><pre>[68]: data[num_cols].isnull().sum()</pre><table><tbody><tr><td>age</td><td>0</td></tr><tr><td>blood_pressure</td><td>0</td></tr><tr><td>specific_gravity</td><td>0</td></tr><tr><td>albumin</td><td>0</td></tr><tr><td>sugar</td><td>0</td></tr><tr><td>blood_glucose_random</td><td>0</td></tr><tr><td>blood_urea</td><td>0</td></tr><tr><td>serum_creatinine</td><td>0</td></tr><tr><td>sodium</td><td>0</td></tr><tr><td>potassium</td><td>0</td></tr><tr><td>haemoglobin</td><td>0</td></tr><tr><td>packed_cell_volume</td><td>0</td></tr><tr><td>white_blood_cell_count</td><td>0</td></tr><tr><td>red_blood_cell_count</td><td>0</td></tr><tr><td>dtype:</td><td>int64</td></tr></tbody></table></div><pre>[69]: # filling "red_blood_cells" and "pus_cell" using random sampling method and rest of cat_cols using mode imputation random_value_imputation('red_blood_cells') random_value_imputation('pus_cell') for col in cat_cols: impute_mode(col)</pre><pre>[71]: data[cat_cols].isnull().sum()</pre></div>	age	0	blood_pressure	0	specific_gravity	0	albumin	0	sugar	0	blood_glucose_random	0	blood_urea	0	serum_creatinine	0	sodium	0	potassium	0	haemoglobin	0	packed_cell_volume	0	white_blood_cell_count	0	red_blood_cell_count	0	dtype:	int64																																																																																																	
age	0																																																																																																																															
blood_pressure	0																																																																																																																															
specific_gravity	0																																																																																																																															
albumin	0																																																																																																																															
sugar	0																																																																																																																															
blood_glucose_random	0																																																																																																																															
blood_urea	0																																																																																																																															
serum_creatinine	0																																																																																																																															
sodium	0																																																																																																																															
potassium	0																																																																																																																															
haemoglobin	0																																																																																																																															
packed_cell_volume	0																																																																																																																															
white_blood_cell_count	0																																																																																																																															
red_blood_cell_count	0																																																																																																																															
dtype:	int64																																																																																																																															
Data Transformation	NA																																																																																																																															
Feature Engineering	NA																																																																																																																															
Save Processed Data	NA																																																																																																																															