

Data Collection and Preprocessing Phase

Date	7 June 2024
Team ID	SWTID1720112707
Project Title	Anemia Sense: Leveraging Machine Learning For Precise Anemia Recognitions
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	<pre>df.head()</pre>

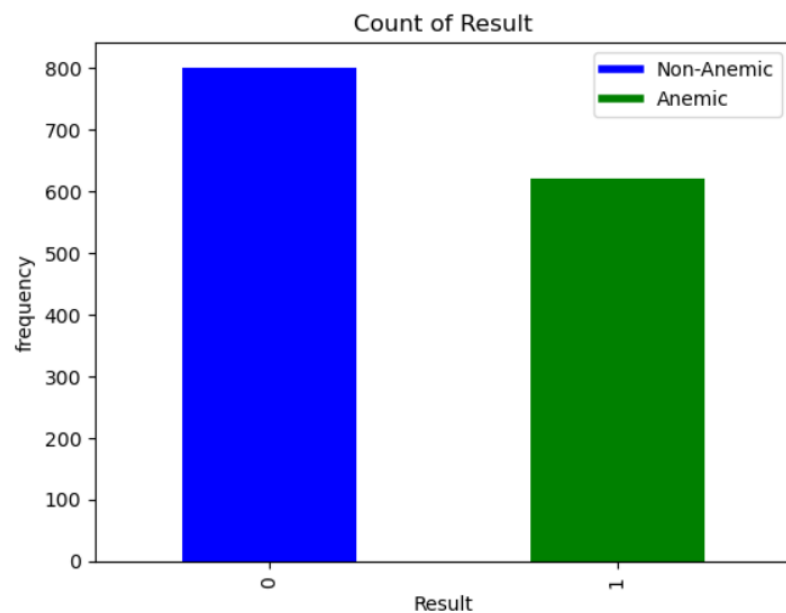
```
df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1421 entries, 0 to 1420
Data columns (total 6 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Gender      1421 non-null   int64
1   Hemoglobin  1421 non-null   float64
2   MCH         1421 non-null   float64
3   MCHC        1421 non-null   float64
4   MCV         1421 non-null   float64
5   Result      1421 non-null   int64
dtypes: float64(4), int64(2)
memory usage: 66.7 KB
  
```

```
df.shape
```

```
(1421, 6)
```

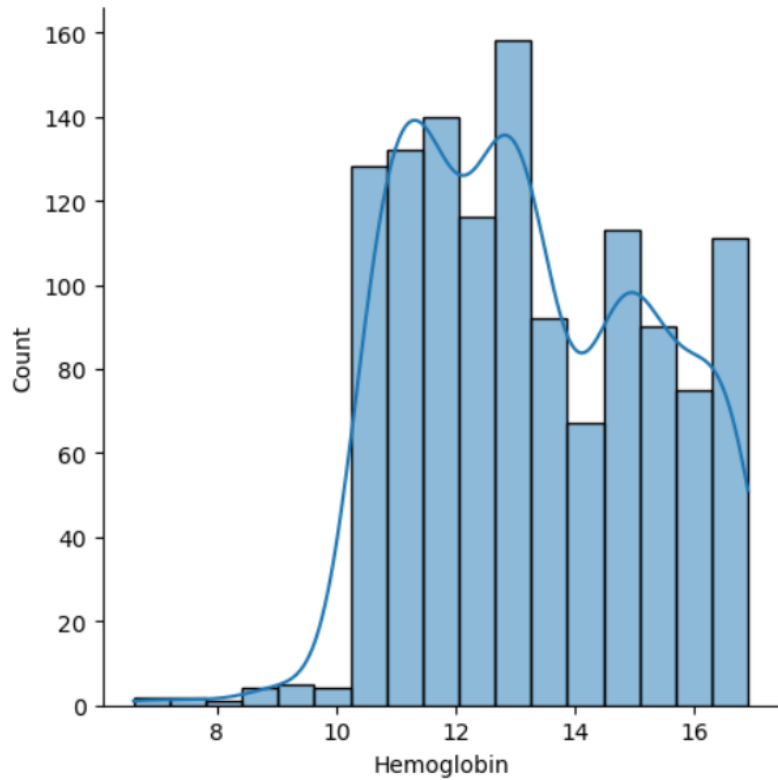


	<div><pre>df.describe()</pre></div> <table><thead><tr><th></th><th>Gender</th><th>Hemoglobin</th><th>MCH</th><th>MCHC</th><th>MCV</th><th>Result</th></tr></thead><tbody><tr><td>count</td><td>1240.000000</td><td>1240.000000</td><td>1240.000000</td><td>1240.000000</td><td>1240.000000</td><td>1240.000000</td></tr><tr><td>mean</td><td>0.540323</td><td>13.218145</td><td>22.903952</td><td>30.277984</td><td>85.620968</td><td>0.500000</td></tr><tr><td>std</td><td>0.498573</td><td>1.976190</td><td>3.993624</td><td>1.394515</td><td>9.673794</td><td>0.500202</td></tr><tr><td>min</td><td>0.000000</td><td>6.600000</td><td>16.000000</td><td>27.800000</td><td>69.400000</td><td>0.000000</td></tr><tr><td>25%</td><td>0.000000</td><td>11.500000</td><td>19.400000</td><td>29.100000</td><td>77.300000</td><td>0.000000</td></tr><tr><td>50%</td><td>1.000000</td><td>13.000000</td><td>22.700000</td><td>30.400000</td><td>85.300000</td><td>0.500000</td></tr><tr><td>75%</td><td>1.000000</td><td>14.900000</td><td>26.200000</td><td>31.500000</td><td>94.225000</td><td>1.000000</td></tr><tr><td>max</td><td>1.000000</td><td>16.900000</td><td>30.000000</td><td>32.500000</td><td>101.600000</td><td>1.000000</td></tr></tbody></table>		Gender	Hemoglobin	MCH	MCHC	MCV	Result	count	1240.000000	1240.000000	1240.000000	1240.000000	1240.000000	1240.000000	mean	0.540323	13.218145	22.903952	30.277984	85.620968	0.500000	std	0.498573	1.976190	3.993624	1.394515	9.673794	0.500202	min	0.000000	6.600000	16.000000	27.800000	69.400000	0.000000	25%	0.000000	11.500000	19.400000	29.100000	77.300000	0.000000	50%	1.000000	13.000000	22.700000	30.400000	85.300000	0.500000	75%	1.000000	14.900000	26.200000	31.500000	94.225000	1.000000	max	1.000000	16.900000	30.000000	32.500000	101.600000	1.000000
	Gender	Hemoglobin	MCH	MCHC	MCV	Result																																																										
count	1240.000000	1240.000000	1240.000000	1240.000000	1240.000000	1240.000000																																																										
mean	0.540323	13.218145	22.903952	30.277984	85.620968	0.500000																																																										
std	0.498573	1.976190	3.993624	1.394515	9.673794	0.500202																																																										
min	0.000000	6.600000	16.000000	27.800000	69.400000	0.000000																																																										
25%	0.000000	11.500000	19.400000	29.100000	77.300000	0.000000																																																										
50%	1.000000	13.000000	22.700000	30.400000	85.300000	0.500000																																																										
75%	1.000000	14.900000	26.200000	31.500000	94.225000	1.000000																																																										
max	1.000000	16.900000	30.000000	32.500000	101.600000	1.000000																																																										
Univariate Analysis	<div><pre># Univariate Analysis output = df['Gender'].value_counts() output.plot(kind = 'bar', color=['orange', 'green']) plt.xlabel('Gender') plt.ylabel('Frequency') plt.title('Gender count') plt.show()</pre></div> <div></div>																																																															

```
sns.displot(df['Hemoglobin'], kde = True)
```

```
D:\Anaconda\Lib\site-packages\seaborn\_oldcore.py:1119: FutureWarning
ert inf values to NaN before operating instead.
with pd.option_context('mode.use_inf_as_na', True):
```

```
<seaborn.axisgrid.FacetGrid at 0x24dbfc157d0>
```



Bivariate Analysis

```
# Bivariate Analysis
```

```
# Calculate mean hemoglobin levels grouped by gender and result
```

```
mean_hemoglobin = df.groupby(['Gender', 'Result'])['Hemoglobin'].mean().reset_index()
```

```
# Pivot the data to get a format suitable for plotting
```

```
mean_hemoglobin_pivot = mean_hemoglobin.pivot(index='Gender', columns='Result', values='Hemoglobin')
```

```
# Plot the histogram
```

```
mean_hemoglobin_pivot.plot(kind='bar', color=['blue', 'green'], edgecolor='black')
```

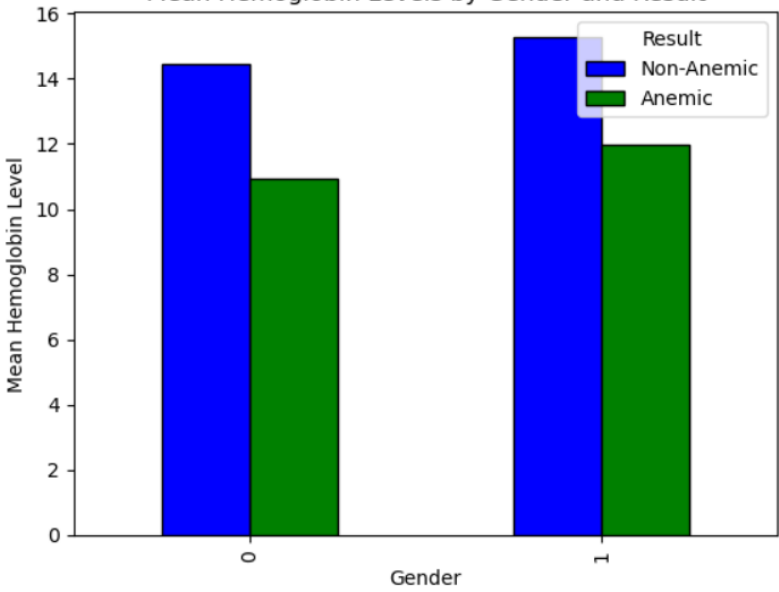
```
plt.xlabel('Gender')
```

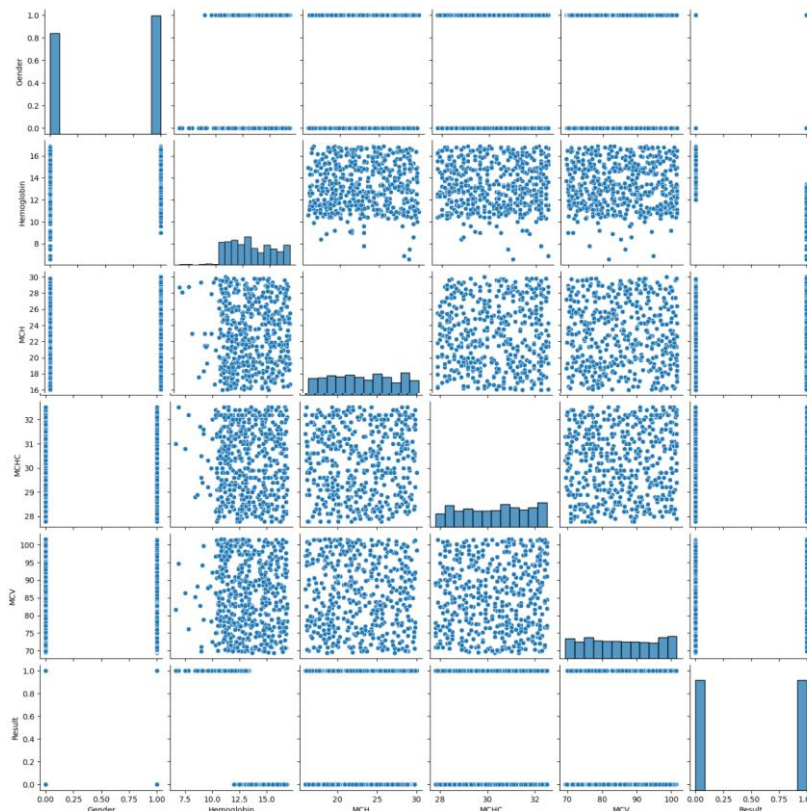
```
plt.ylabel('Mean Hemoglobin Level')
```

```
plt.title('Mean Hemoglobin Levels by Gender and Result')
```

```
plt.legend(title='Result', labels=['Non-Anemic', 'Anemic'])
```

```
plt.show()
```

	<p>Mean Hemoglobin Levels by Gender and Result</p>  <table><thead><tr><th>Gender</th><th>Non-Anemic</th><th>Anemic</th></tr></thead><tbody><tr><td>0</td><td>~14.5</td><td>~11.0</td></tr><tr><td>1</td><td>~15.5</td><td>~12.0</td></tr></tbody></table>	Gender	Non-Anemic	Anemic	0	~14.5	~11.0	1	~15.5	~12.0
Gender	Non-Anemic	Anemic								
0	~14.5	~11.0								
1	~15.5	~12.0								
Multivariate Analysis	<pre># Multivariate Analysis # Convert infinite values to NaN df.replace([np.inf, -np.inf], np.nan, inplace=True) df.dropna(inplace=True) sns.pairplot(df) plt.show()</pre>									



Multivariate Analysis

```
print(df.dtypes)

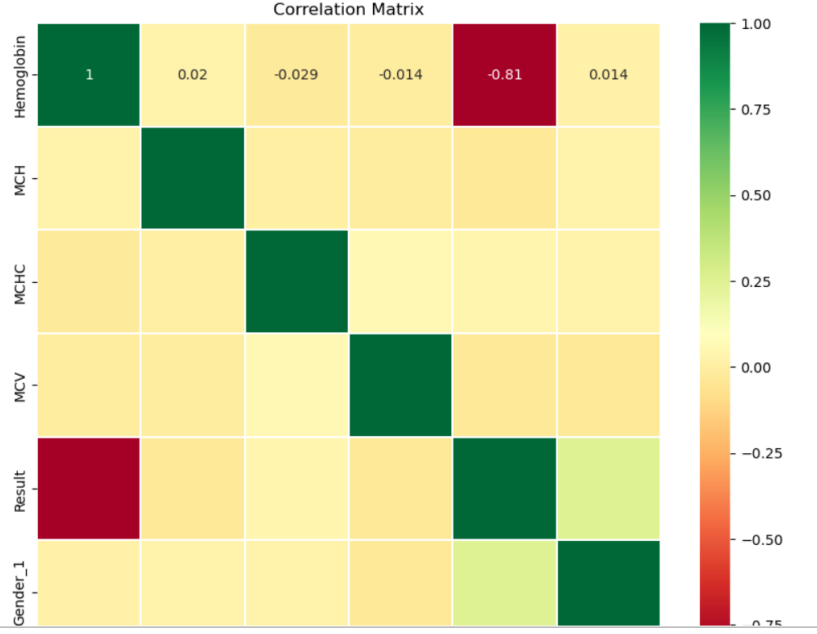
df = df.apply(pd.to_numeric, errors='coerce')

print(df.isnull().sum())
df = df.dropna() # or use an appropriate imputation method

df = pd.get_dummies(df, columns=['Gender'], drop_first=True)

plt.figure(figsize=(10, 8))
sns.heatmap(df.corr(), annot=True, cmap='RdYlGn', linewidths=0.2)
plt.title('Correlation Matrix')
plt.show()
```

```
Gender          int64
Hemoglobin      float64
MCH             float64
MCHC           float64
MCV            float64
Result         int64
dtype: object
Gender          0
Hemoglobin      0
MCH             0
MCHC           0
MCV            0
Result         0
dtype: int64
```

	<div><p>Correlation Matrix</p><table><tr><th></th><th>Hemoglobin</th><th>MCH</th><th>MCHC</th><th>MCV</th><th>Result</th><th>Gender_1</th></tr><tr><th>Hemoglobin</th><td>1</td><td>0.02</td><td>-0.029</td><td>-0.014</td><td>-0.81</td><td>0.014</td></tr><tr><th>MCH</th><td>0.02</td><td>1</td><td>0.01</td><td>0.01</td><td>0.01</td><td>0.01</td></tr><tr><th>MCHC</th><td>-0.029</td><td>0.01</td><td>1</td><td>0.01</td><td>0.01</td><td>0.01</td></tr><tr><th>MCV</th><td>-0.014</td><td>0.01</td><td>0.01</td><td>1</td><td>0.01</td><td>0.01</td></tr><tr><th>Result</th><td>-0.81</td><td>0.01</td><td>0.01</td><td>0.01</td><td>1</td><td>0.01</td></tr><tr><th>Gender_1</th><td>0.014</td><td>0.01</td><td>0.01</td><td>0.01</td><td>0.01</td><td>1</td></tr></table></div>		Hemoglobin	MCH	MCHC	MCV	Result	Gender_1	Hemoglobin	1	0.02	-0.029	-0.014	-0.81	0.014	MCH	0.02	1	0.01	0.01	0.01	0.01	MCHC	-0.029	0.01	1	0.01	0.01	0.01	MCV	-0.014	0.01	0.01	1	0.01	0.01	Result	-0.81	0.01	0.01	0.01	1	0.01	Gender_1	0.014	0.01	0.01	0.01	0.01	1
	Hemoglobin	MCH	MCHC	MCV	Result	Gender_1																																												
Hemoglobin	1	0.02	-0.029	-0.014	-0.81	0.014																																												
MCH	0.02	1	0.01	0.01	0.01	0.01																																												
MCHC	-0.029	0.01	1	0.01	0.01	0.01																																												
MCV	-0.014	0.01	0.01	1	0.01	0.01																																												
Result	-0.81	0.01	0.01	0.01	1	0.01																																												
Gender_1	0.014	0.01	0.01	0.01	0.01	1																																												
Data Preprocessing Code Screenshots																																																		
Loading Data	<div><pre>import pandas as pd import numpy as np import matplotlib.pyplot as plt import seaborn as sns df = pd.read_csv('D:/anemia.csv')</pre></div>																																																	
Handling Missing Data	<div><pre>print("Initial Count of Missing Values in the Dataset\n") print(df.isnull().sum()) # Handling missing values # Filling missing numerical values with the median for column in df.select_dtypes(include=[np.number]).columns: df[column].fillna(df[column].median(), inplace=True) for column in df.select_dtypes(include=[object]).columns: df[column].fillna(df[column].mode()[0], inplace=True) print("\nFinal Count of Missing Values in the Dataset\n") print(df.isnull().sum())</pre></div>																																																	

Data Transformation	<pre> #we can see that the female count is more than the male so, # we can balance it using the undersampling from sklearn.utils import resample majorclass = df[df['Result'] == 0] minorclass = df[df['Result'] == 1] major_downsample = resample(majorclass, replace=False, n_samples=len(minorclass), random_state=42) df = pd.concat([major_downsample, minorclass]) print(df['Result'].value_counts()) </pre>
Feature Engineering	<pre> for column in df.select_dtypes(include=[np.number]).columns: df[column].fillna(df[column].median(), inplace=True) for column in df.select_dtypes(include=[object]).columns: df[column].fillna(df[column].mode()[0], inplace=True) </pre>
Save Processed Data	<pre> df = pd.concat([major_downsample, minorclass]) </pre>