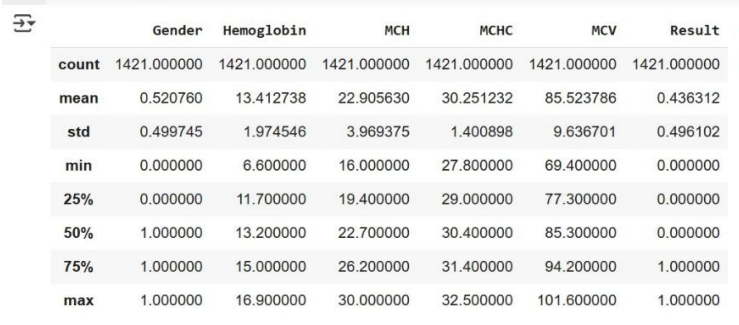


Data Collection and Preprocessing Phase

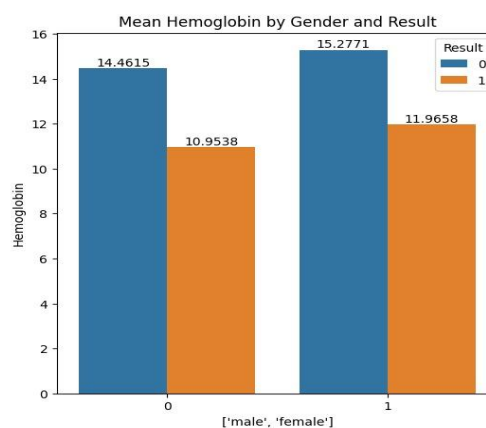
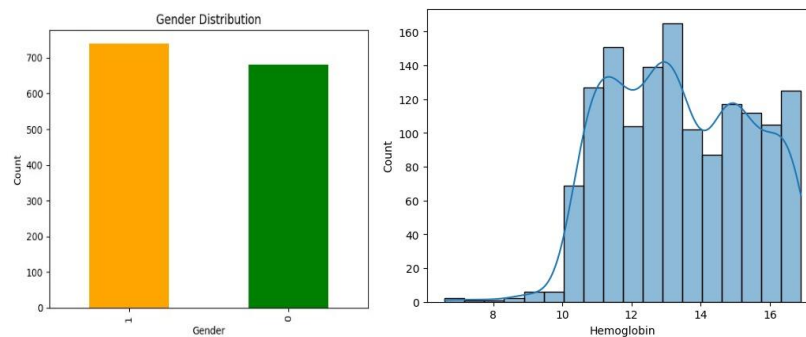
Date	03 August 2025
Project Title	Anemia Sense - Machine Learning for Precise Anemia Recognition
Maximum Marks	6 Marks

Data Exploration and Preprocessing Report

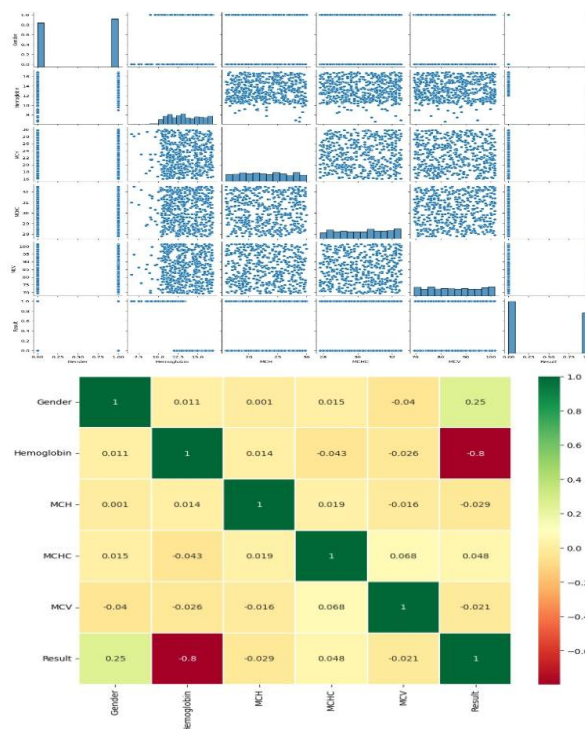
Dataset variables will be statistically analysed to identify patterns and outliers, with Python employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modelling, and forming a strong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

Section	Description																																																															
Data Overview	<p><u>Dimension:</u> 1421 rows × 6 columns</p> <p><u>Descriptive statistics:</u></p> <pre>df.describe()</pre>  <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>1421.000000</td><td>1421.000000</td><td>1421.000000</td><td>1421.000000</td><td>1421.000000</td><td>1421.000000</td></tr><tr><td>mean</td><td>0.520760</td><td>13.412738</td><td>22.905630</td><td>30.251232</td><td>85.523786</td><td>0.436312</td></tr><tr><td>std</td><td>0.499745</td><td>1.974546</td><td>3.969375</td><td>1.400898</td><td>9.636701</td><td>0.496102</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>11.700000</td><td>19.400000</td><td>29.000000</td><td>77.300000</td><td>0.000000</td></tr><tr><td>50%</td><td>1.000000</td><td>13.200000</td><td>22.700000</td><td>30.400000</td><td>85.300000</td><td>0.000000</td></tr><tr><td>75%</td><td>1.000000</td><td>15.000000</td><td>26.200000</td><td>31.400000</td><td>94.200000</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	1421.000000	1421.000000	1421.000000	1421.000000	1421.000000	1421.000000	mean	0.520760	13.412738	22.905630	30.251232	85.523786	0.436312	std	0.499745	1.974546	3.969375	1.400898	9.636701	0.496102	min	0.000000	6.600000	16.000000	27.800000	69.400000	0.000000	25%	0.000000	11.700000	19.400000	29.000000	77.300000	0.000000	50%	1.000000	13.200000	22.700000	30.400000	85.300000	0.000000	75%	1.000000	15.000000	26.200000	31.400000	94.200000	1.000000	max	1.000000	16.900000	30.000000	32.500000	101.600000	1.000000
		Gender	Hemoglobin	MCH	MCHC	MCV	Result																																																									
	count	1421.000000	1421.000000	1421.000000	1421.000000	1421.000000	1421.000000																																																									
	mean	0.520760	13.412738	22.905630	30.251232	85.523786	0.436312																																																									
	std	0.499745	1.974546	3.969375	1.400898	9.636701	0.496102																																																									
	min	0.000000	6.600000	16.000000	27.800000	69.400000	0.000000																																																									
	25%	0.000000	11.700000	19.400000	29.000000	77.300000	0.000000																																																									
	50%	1.000000	13.200000	22.700000	30.400000	85.300000	0.000000																																																									
	75%	1.000000	15.000000	26.200000	31.400000	94.200000	1.000000																																																									
	max	1.000000	16.900000	30.000000	32.500000	101.600000	1.000000																																																									
Univariate Analysis																																																																

Bivariate Analysis



Multivariate Analysis



Outliers and Anomalies	-																																																																								
Data Preprocessing Code Screenshots																																																																									
Loading Data	<div><pre>df.head()</pre><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>0</td><td>1</td><td>14.9</td><td>22.7</td><td>29.1</td><td>83.7</td><td>0</td></tr><tr><td>1</td><td>0</td><td>15.9</td><td>25.4</td><td>28.3</td><td>72.0</td><td>0</td></tr><tr><td>2</td><td>0</td><td>9.0</td><td>21.5</td><td>29.6</td><td>71.2</td><td>1</td></tr><tr><td>3</td><td>0</td><td>14.9</td><td>16.0</td><td>31.4</td><td>87.5</td><td>0</td></tr><tr><td>4</td><td>1</td><td>14.7</td><td>22.0</td><td>28.2</td><td>99.5</td><td>0</td></tr></tbody></table></div>		Gender	Hemoglobin	MCH	MCHC	MCV	Result	0	1	14.9	22.7	29.1	83.7	0	1	0	15.9	25.4	28.3	72.0	0	2	0	9.0	21.5	29.6	71.2	1	3	0	14.9	16.0	31.4	87.5	0	4	1	14.7	22.0	28.2	99.5	0																														
	Gender	Hemoglobin	MCH	MCHC	MCV	Result																																																																			
0	1	14.9	22.7	29.1	83.7	0																																																																			
1	0	15.9	25.4	28.3	72.0	0																																																																			
2	0	9.0	21.5	29.6	71.2	1																																																																			
3	0	14.9	16.0	31.4	87.5	0																																																																			
4	1	14.7	22.0	28.2	99.5	0																																																																			
Handling Missing Data	<div><div><pre>df.info()</pre><pre><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</pre></div><div><pre>df.isnull().sum()</pre><table><tbody><tr><td></td><td>0</td></tr><tr><td>Gender</td><td>0</td></tr><tr><td>Hemoglobin</td><td>0</td></tr><tr><td>MCH</td><td>0</td></tr><tr><td>MCHC</td><td>0</td></tr><tr><td>MCV</td><td>0</td></tr><tr><td>Result</td><td>0</td></tr></tbody></table><p>dtype: int64</p></div></div>		0	Gender	0	Hemoglobin	0	MCH	0	MCHC	0	MCV	0	Result	0																																																										
	0																																																																								
Gender	0																																																																								
Hemoglobin	0																																																																								
MCH	0																																																																								
MCHC	0																																																																								
MCV	0																																																																								
Result	0																																																																								
Data Transformation	<div><pre>x = df.drop('Result', axis = 1)</pre><pre>x</pre><table><thead><tr><th></th><th>Gender</th><th>Hemoglobin</th><th>MCH</th><th>MCHC</th><th>MCV</th></tr></thead><tbody><tr><td>1234</td><td>1</td><td>16.6</td><td>18.8</td><td>28.1</td><td>70.9</td></tr><tr><td>1188</td><td>0</td><td>15.3</td><td>18.3</td><td>30.4</td><td>93.4</td></tr><tr><td>106</td><td>0</td><td>14.8</td><td>20.4</td><td>28.5</td><td>91.1</td></tr><tr><td>954</td><td>0</td><td>14.6</td><td>16.9</td><td>31.9</td><td>78.1</td></tr><tr><td>112</td><td>0</td><td>15.9</td><td>28.7</td><td>31.0</td><td>81.6</td></tr><tr><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td></tr><tr><td>1415</td><td>1</td><td>13.2</td><td>20.1</td><td>28.8</td><td>91.2</td></tr><tr><td>1416</td><td>0</td><td>10.6</td><td>25.4</td><td>28.2</td><td>82.9</td></tr><tr><td>1417</td><td>1</td><td>12.1</td><td>28.3</td><td>30.4</td><td>86.9</td></tr><tr><td>1418</td><td>1</td><td>13.1</td><td>17.7</td><td>28.1</td><td>80.7</td></tr><tr><td>1420</td><td>0</td><td>11.8</td><td>21.2</td><td>28.4</td><td>98.1</td></tr></tbody></table></div>		Gender	Hemoglobin	MCH	MCHC	MCV	1234	1	16.6	18.8	28.1	70.9	1188	0	15.3	18.3	30.4	93.4	106	0	14.8	20.4	28.5	91.1	954	0	14.6	16.9	31.9	78.1	112	0	15.9	28.7	31.0	81.6	1415	1	13.2	20.1	28.8	91.2	1416	0	10.6	25.4	28.2	82.9	1417	1	12.1	28.3	30.4	86.9	1418	1	13.1	17.7	28.1	80.7	1420	0	11.8	21.2	28.4	98.1
	Gender	Hemoglobin	MCH	MCHC	MCV																																																																				
1234	1	16.6	18.8	28.1	70.9																																																																				
1188	0	15.3	18.3	30.4	93.4																																																																				
106	0	14.8	20.4	28.5	91.1																																																																				
954	0	14.6	16.9	31.9	78.1																																																																				
112	0	15.9	28.7	31.0	81.6																																																																				
...																																																																				
1415	1	13.2	20.1	28.8	91.2																																																																				
1416	0	10.6	25.4	28.2	82.9																																																																				
1417	1	12.1	28.3	30.4	86.9																																																																				
1418	1	13.1	17.7	28.1	80.7																																																																				
1420	0	11.8	21.2	28.4	98.1																																																																				
Feature Engineering	Attached the codes in final submission.																																																																								
Save Processed Data	-																																																																								