

# Final Project Report

## Introduction

## Project Overviews

Anemia Sense is a machine learning-based diagnostic tool developed to assist in the accurate detection of anemia, a condition marked by a deficiency of red blood cells (RBCs) or hemoglobin. By utilizing key blood test parameters, the system applies various supervised learning algorithms to analyze patterns and predict the likelihood of anemia in patients. This approach improves diagnostic reliability and helps reduce the burden on manual assessment methods.

The main goal of Anemia Sense is to support early diagnosis and enhance clinical decision-making, especially in healthcare settings with limited resources. By offering fast and data-driven insights, the tool can help healthcare providers identify at-risk individuals more efficiently, enabling timely intervention and better patient care.

## Objectives

The primary objectives of this project are:

- To collect and preprocess hematological data relevant to the diagnosis of anemia.
- To explore and analyze patterns within the data that contribute to accurate anemia detection.
- To apply various supervised machine learning algorithms for classifying anemia.
- To evaluate and compare model performance using metrics such as accuracy.
- To select the best-performing model based on a comprehensive evaluation.
- To develop a lightweight, interpretable, and efficient diagnostic tool suitable for real-world deployment.

## Project Initialization and Planning Phase

Date	1August 2025
Skillwallet ID	SWUID20250194750
Project Name	Anemia Sense: Leveraging Machine Learning For Precise Anemia
Maximum Marks	3 Marks

### Problem Statement:

The current anemia diagnosis process poses challenges for both patients and healthcare providers, often impacting the quality of care and timely intervention. Patients, particularly those in resource-limited settings, face obstacles such as delayed test results, limited access to diagnostic tools, and inconsistent assessments. These issues contribute to a suboptimal healthcare experience and can hinder early detection and treatment.

To improve clinical outcomes and enhance the diagnostic journey, *Anemia Sense* aims to address these challenges. By identifying specific pain points in the existing diagnostic workflow and implementing a data-driven, machine learning-based solution, we strive to deliver a faster, more accurate, and user-friendly experience. This approach supports healthcare providers in making informed decisions, ultimately fostering greater trust and improving patient care.



Problem Statement (PS)	I am (Customer)	I'm trying to	But	Because	Which makes me feel
PS-1	A healthcare provider in a rural clinic.	Diagnose anemia accurately and quickly.	I have limited access to advanced diagnostic tools.	Resources are scarce and manual methods are time-consuming.	Concerned about delayed treatment and patient outcomes.

## Project Initialization and Planning Phase

Date	1 August 2025
Skillwallet ID	SWUID20250194750
Project Title	Anemia Sense: Leveraging Machine Learning For Precise Anemia
Maximum Marks	3 Marks

### Project Proposal (Proposed Solution) report

The proposal report aims to transform anemia diagnosis using machine learning, improving both efficiency and accuracy. It addresses diagnostic limitations in current practices, offering a smarter and more accessible solution. Key features include a machine learning-based classification model and near real-time prediction capabilities.

Project Overview	
Objective	The primary objective is to revolutionize anemia detection by implementing advanced supervised machine learning techniques, enabling faster, more accurate, and scalable assessments.
Scope	The project comprehensively analyzes and enhances the diagnostic process for anemia, integrating hematological data and machine learning algorithms to create an efficient, lightweight, and interpretable diagnostic system.
Problem Statement	
Description	Inaccuracies, delays, and inconsistencies in traditional anemia diagnostic methods negatively impact clinical decision-making and patient care, particularly in resource-constrained environments.
Impact	Addressing these issues will lead to improved diagnostic reliability, timely medical intervention, and better patient outcomes—supporting healthcare systems and professionals in delivering quality care more effectively.
Proposed Solution	
Approach	Applying supervised machine learning algorithms to analyze key blood parameters (such as hemoglobin levels, etc.), enabling accurate and early detection of anemia.
Key Features	- Implementation of a machine learning-based anemia classification model.

	<ul style="list-style-type: none"> <li>- Near real-time predictions to support prompt clinical decisions.</li> <li>- Lightweight, interpretable system suitable for deployment in low-resource settings.</li> </ul>
--	---

## Resource Requirements

Resource Type	Description	Specification/Allocation
<b>Hardware</b>		
Computing Resources	CPU/GPU specifications, number of cores	T4 GPU
Memory	RAM specifications	8 GB
Storage	Disk space for data, models, and logs	1 TB SSD
<b>Software</b>		
Frameworks	Python frameworks	Flask
Libraries	Additional libraries	scikit-learn, pandas, numpy, matplotlib, seaborn
Development Environment	IDE	Jupyter Notebook, pycharm
<b>Data</b>		
Data	Source, size, format	Kaggledataset,614,csv

## Initial Project Planning Report

Date	25-01-2024
Skillwallet ID	SWUID20250194750
Project Name	Anemia Sense: Leveraging Machine Learning For Precise Anemia
Maximum Marks	4 Marks

## Product Backlog, Sprint Schedule, and Estimation

Sprint	Functional Requirement (Epic)	User Story Number	User Story / Task	Priority	Team Members	Sprint Start Date	Sprint End Date (Planned)
Sprint-1	Data Collection and Preprocessing	ASLMF LPS-1,2	Data gathering & loading	High	Abhijay yadav	25/07/2025	27/07/2025
Sprint-1	Data Collection and Preprocessing	ASLMF LPS-3,4	Handle missing values	High	Abhijay yadav	25/07/2025	27/07/2025
Sprint-1	Data Collection and Preprocessing	ASLMF LPS-5,6	Encode categorical variables	Medium	Abhijay yadav	25/07/2025	27/07/2025
Sprint-1	Data Collection and Preprocessing	ASLMF LPS-7,8	Dataset balancing	High	Abhijay yadav	25/07/2025	27/07/2025
Sprint-2	Model Development	ASLMF LPS-9-12	Train ML models	High	Abhijay yadav	28/07/2025	30/07/2025
Sprint-2	Model Development	ASLMP FPS-13-15	Evaluate models (accuracy, F1, etc.)	Medium	Abhijay yadav	28/07/2025	30/07/2025
Sprint-3	Model Optimization	ASLMF PS-16-18	Hyperparameter tuning	High	Abhijay yadav	31/07/2025	01/08/2025
Sprint-3	Model Optimization	ASLMF PS-19,20	Final model selection & justification	Medium	Abhijay yadav	31/07/2025	01/08/2025

## Screenshot:

This screenshot shows the Jira Board view for the project "Anemia Sense: Leveraging Machine Learning For Precise Anemia Sensing". The board is organized into three columns: TO DO (6 items), IN PROGRESS (0 items), and DONE (3 items). Each item is a card representing a task, with a status label (e.g., MODEL TRAINING, MODEL EVALUATION) and a progress indicator (a bar chart). The tasks are as follows:

Task	Status	Progress
Train Naive Bayes model	MODEL TRAINING	2/2
Train SVM model	MODEL TRAINING	2/2
Train Gradient Boosting model	MODEL TRAINING	2/2
Calculate accuracy for all models	MODEL EVALUATION	1/1
Generate classification reports	MODEL EVALUATION	1/1
Train Logistic Regression model	MODEL TRAINING	2/2
Train Random Forest Classifier	MODEL TRAINING	2/2
Train Decision Tree Classifier	MODEL TRAINING	2/2

This screenshot shows the Jira Timeline view for the same project. The timeline is organized into Sprints, with a vertical bar indicating the current sprint's progress. The tasks are as follows:

Task	Status	Progress
ASLMLFPS-1 Data Collection	ASLMLFPS-1	2/2
ASLMLFPS-5 Data Preprocessing	ASLMLFPS-5	2/2
ASLMLFPS-9 Model Training	ASLMLFPS-9	2/2
ASLMLFPS-16 Model Evaluation	ASLMLFPS-16	2/2
ASLMLFPS-20 Model Optimization	ASLMLFPS-20	2/2
ASLMLFPS-24 Model Deployment	ASLMLFPS-24	2/2

## Data Collection and Preprocessing Phase

Date	1 August 2025
Skillwallet ID	SWUID20250194750
Project Title	Anemia Sense: Leveraging Machine Learning For Precise Anemia
Maximum Marks	2 Marks

### Data Collection Plan & Raw Data Sources Identification Report:

Elevate your data strategy with a well-structured Data Collection Plan and comprehensive Raw Data Sources report, ensuring meticulous curation and data integrity to support reliable, data-driven anemia diagnosis.

#### Data Collection Plan:

Section	Description
Project Overview	The machine learning project aims to predict the presence of anemia based on patient hematological parameters. Using datasets that include features such as hemoglobin levels, etc. components, the objective is to build a robust model that accurately classifies anemia status—facilitating early detection and better clinical decision-making.
Data Collection Plan	<ul style="list-style-type: none"> <li>● Search for datasets related to anemia diagnosis, including hematological test results and demographic patient data.</li> <li>● Prioritize datasets with labeled outcomes (anemia vs. non-anemia) and diverse population samples.</li> <li>● Ensure inclusion of common clinical features such as Hemoglobin (Hb), MCV, MCH and MCHC.</li> </ul>
Raw Data Sources Identified	The raw data sources for this project include publicly available medical datasets from platform such as <b>Kaggle</b> . These repositories provide anonymized patient blood test records suitable for machine learning analysis. The datasets typically include clinical variables crucial for anemia classification, such as hemoglobin concentration, MCHC, and demographic features like sex .

### Raw Data Sources Report:

Source Name	Description	Location/URL	Format	Size	Access Permissions
Kaggle Dataset	The dataset comprises patient details (gender) and hematological metrics (hemoglobin, MCHC, MCV, MCH), along with anemia diagnosis outcomes. It is used to predict if a patient is likely to suffer from anemia using a binary classification algorithm.	<a href="https://www.kaggle.com/datasets/biswaranjanrao/anemia-dataset">https://www.kaggle.com/datasets/biswaranjanrao/anemia-dataset</a>	CSV	34 kB	Public



## Data Collection and Preprocessing Phase

Date	1 August 2025
Skillwallet ID	SWUID20250194750
Project Title	Anemia Sense: Leveraging Machine Learning For Precise Anemia
Maximum Marks	2 Marks

### Data Quality Report:

The Data Quality Report will summarize data quality issues from the selected source, including severity levels and resolution plans. It will aid in systematically identifying and rectifying data discrepancies.

### Data Quality Report:

Data Source	Data Quality Issue	Severity	Resolution Plan
Kaggle Dataset	Missing values in 'Hemoglobin', 'MCHC', 'MCV', and 'MCH' columns	High	Apply mean or median imputation, or use clinically validated ranges for plausible medical values.
Kaggle Dataset	Categorical data in the 'Gender' and 'Result' columns	Moderate	Apply label encoding or one-hot encoding where appropriate.

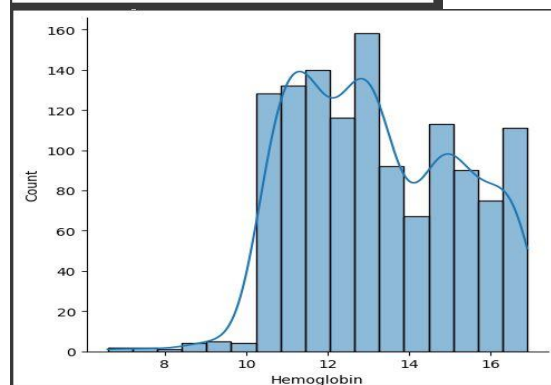
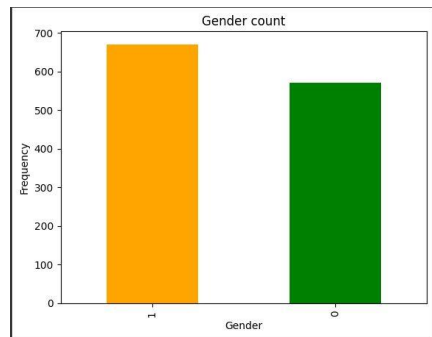
## Data Collection and Preprocessing Phase

Date	1 August 2025
Skillwallet ID	SWUID20250194750
Project Title	Anemia Sense: Leveraging Machine Learning For Precise Anemia
Maximum Marks	6 Marks

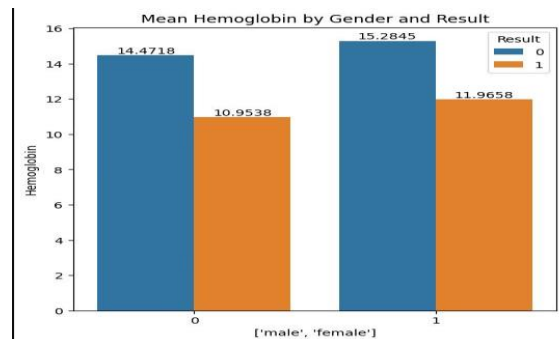
### Data Exploration and Preprocessing Report

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.

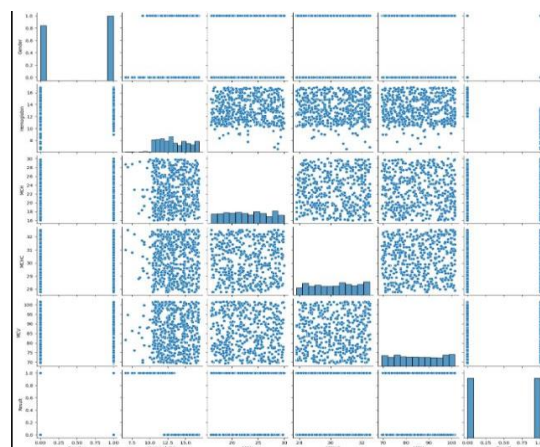
Section	Description																																																															
Data Overview	<u>Dimension:</u> 614rows×13columns <u>Descriptive statistics:</u>																																																															
	<table><tr><th></th><th>Gender</th><th>Hemoglobin</th><th>MCH</th><th>MCHC</th><th>MCV</th><th>Result</th></tr><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></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																																																																

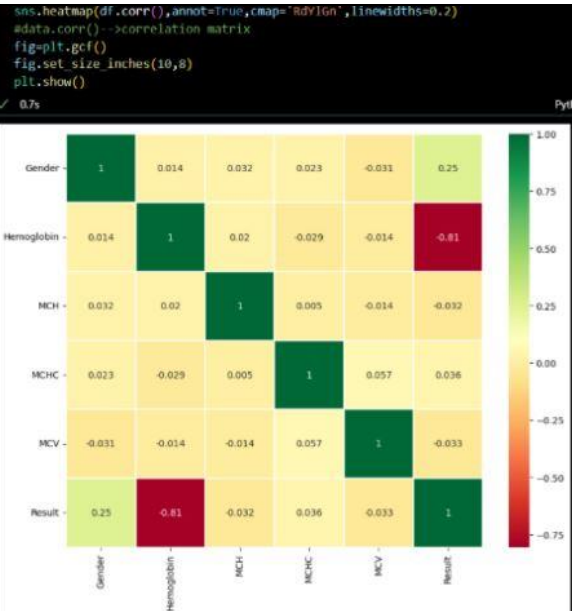


## Bivariate Analysis



## Multivariate Analysis





## Splitting Data into Train and Test

```
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test = train_test_split(X,Y, test_size=0.2, random_state=20)

print(x_train.shape)
print(x_test.shape)
print(y_train.shape)
print(y_test.shape)
```

```
(992, 5)
(248, 5)
(992,)
(248,)
```

## Data Preprocessing Code Screenshots

### Loading Data

```
df = pd.read_csv('anemia.csv')
df
```

	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
...	...	...	...	...	...	...
1416	0	10.6	25.4	28.2	82.9	1
1417	1	12.1	28.3	30.4	86.9	1
1418	1	13.1	17.7	28.1	80.7	1
1419	0	14.3	16.2	29.5	95.2	0
1420	0	11.8	21.2	28.4	98.1	1

1421 rows x 6 columns

## Handling Missing Data

```
df = pd.read_csv('anemia.csv')

df.info()
df.shape
```

```
<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
(1421, 6)
```

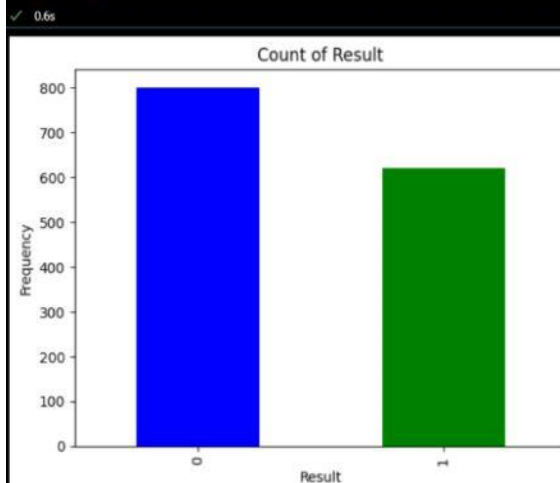
```
df.isnull().sum()
```

	0
Gender	0
Hemoglobin	0
MCH	0
MCHC	0
MCV	0
Result	0

```
dtype: int64
```

## Handling Imbalanced Values

```
results = df['Result'].value_counts()
results.plot(kind = 'bar', color=['blue', 'green'])
plt.xlabel('Result')
plt.ylabel('Frequency')
plt.title('Count of Result')
plt.show()
```



```
from sklearn.utils import resample
majorclass = df[df['Result'] == 0]
minorclass = df[df['Result'] == 1]

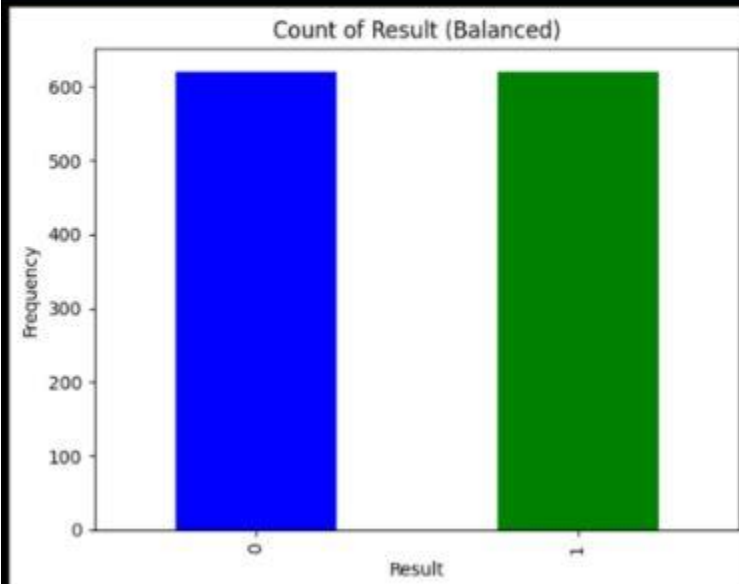
major_downsample = resample(majorclass, replace=False, n_samples=len(minorclass), random_state=42)
df = pd.concat([major_downsample, minorclass])

print(df['Result'].value_counts())
```

```
Result
0    620
1    620
Name: count, dtype: int64
```

```
result_balanced = df['Result'].value_counts()  
result_balanced.plot(kind='bar', color=['blue', 'green'])  
plt.xlabel('Result')  
plt.ylabel('Frequency')  
plt.title('Count of Result (Balanced)')  
plt.show()
```

✓ 0.4s



## Model Development Phase

Date	1 August 2025
Skillwallet ID	SWUID20250194750
Project Title	Anemia Sense: Leveraging Machine Learning For Precise Anemia
Maximum Marks	5 Marks

### Feature Selection Report

In the forthcoming update, each feature will be accompanied by a brief description. Users will indicate whether it's selected or not, providing reasoning for their decision. This process will streamline decision-making and enhance transparency in feature selection.

Feature	Description	Selected (Yes/No)	Reasoning
Gender	Patient's gender (0 = male, 1 = female)	Yes	Gender may influence anemia risk due to physiological differences (e.g., menstruation, pregnancy).
Hemoglobin	Oxygen-carrying protein in red blood cells	Yes	A direct and critical indicator of anemia.
MCH	Average amount of hemoglobin per red blood cell	Yes	Helps classify the type of anemia and assess severity.
MCHC	Average concentration of hemoglobin in red blood cells	Yes	Important for diagnosing anemia subtypes.
MCV	Average size of red blood cells	Yes	An essential metric in anemia classification (e.g., microcytic, normocytic, macrocytic anemia).

Result	Anemia diagnosis outcome (0 = not anemic, 1 = anemic)	Yes	Target variable for the classification model; crucial for training and evaluation.
--------	---	-----	--



## Model Development Phase

Date	1 August 2025
Skillwallet ID	SWUID20250194750
Project Title	Anemia Sense: Leveraging Machine Learning For Precise Anemia
Maximum Marks	6 Marks

### Model Selection Report

In the forthcoming Model Selection Report, various models will be outlined, detailing their descriptions, hyperparameters, and performance metrics, including Accuracy or F1 Score. This comprehensive report will provide insights into the chosen models and their effectiveness in predicting anemia.

Model	Description	Hyperparameters	Performance Metric (e.g., Accuracy, F1 Score)
Linear Regression	Statistical method adapted for classification; models linear relationship between features and anemia outcome.	-	Accuracy score = 99.19%
Decision Tree Classifier	Tree-based model; easy to interpret, captures non-linear relationships, useful for early insights.	-	Accuracy score = 100.00%
Random Forest Classifier	Ensemble of decision trees; reduces overfitting, improves generalization, and ranks features effectively.	-	Accuracy score = 100.00%
Gaussian Naive Bayes	Probabilistic model; assumes feature independence, efficient with small datasets and performs well in practice.	-	Accuracy score = 97.98%

Support Vector Classifier	Finds optimal hyperplane for classification; effective in high-dimensional spaces and robust to overfitting.	-	Accuracy score = 93.95%
Gradient Boost Classifier	Sequential ensemble method; minimizes prediction error, strong performance on complex datasets.	-	Accuracy score = 100.00%

## Model Development Phase

Date	1 August 2025
Skillwallet ID	SWUID20250194750
Project Title	Anemia Sense: Leveraging Machine Learning For Precise Anemia
Maximum Marks	4 Marks

### Initial Model Training Code, Model Validation and Evaluation Report

The initial model training code will be showcased in the future through a screenshot. The model validation and evaluation report will include classification reports, accuracy, and confusion matrices for multiple models, presented through respective screenshots.

### Initial Model Training Code:

```
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
from sklearn.metrics import classification_report

logistic_regression = LogisticRegression()
logistic_regression.fit(x_train,y_train)
y_pred = logistic_regression.predict(x_test)
acc_lr = accuracy_score(y_test,y_pred)
c_lr = classification_report(y_test,y_pred)
print('Accuracy Score: ',acc_lr)
print(c_lr)
```

```
from sklearn.naive_bayes import GaussianNB
NB = GaussianNB()
NB.fit(x_train,y_train)
y_pred = NB.predict(x_test)
acc_nb = accuracy_score(y_test,y_pred)
c_nb = classification_report(y_test,y_pred)
print('Accuracy Score: ',acc_nb)
print(c_nb)
```

```
from sklearn.ensemble import GradientBoostingClassifier
GBC = GradientBoostingClassifier()
GBC.fit(x_train,y_train)
y_pred = GBC.predict(x_test)
acc_gbc = accuracy_score(y_test,y_pred)
c_gbc = classification_report(y_test,y_pred)
print('Accuracy Score: ',acc_gbc)
print(c_gbc)
```

```
from sklearn.tree import DecisionTreeClassifier
decision_tree_model = DecisionTreeClassifier()
decision_tree_model.fit(x_train,y_train)
y_pred = decision_tree_model.predict(x_test)
acc_dt = accuracy_score(y_test,y_pred)
c_dt = classification_report(y_test,y_pred)
print('Accuracy Score: ',acc_dt)
print(c_dt)

from sklearn.ensemble import RandomForestClassifier
random_forest = RandomForestClassifier()
random_forest.fit(x_train,y_train)
y_pred = random_forest.predict(x_test)
acc_rf = accuracy_score(y_test,y_pred)
c_rf = classification_report(y_test,y_pred)
print('Accuracy Score: ',acc_rf)
print(c_rf)

from sklearn.svm import SVC
support_vector = SVC()
support_vector.fit(x_train,y_train)
y_pred = support_vector.predict(x_test)
acc_svc = accuracy_score(y_test,y_pred)
c_svc = classification_report(y_test,y_pred)
print('Accuracy Score: ',acc_svc)
print(c_svc)
```

### Model Validation and Evaluation Report:

Model	Classification Report	F1 Score	Confusion Matrix
Linear Regression	<pre> precision    recall  f1-score   support  0           1.00      0.98      0.99         113 1           0.99      1.00      0.99         135  accuracy          0.99 macro avg          0.99 weighted avg       0.99</pre>	99%	<pre> con_lr = confusion_matrix(y_test, y_pred) print(con_lr)  [[111  2]  [ 0 135]]</pre>

Decision Tree	<pre>print('Accuracy Score: ',acc_dt) print(c_dt)</pre> <table><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr><tr><td>0</td><td>1.00</td><td>1.00</td><td>1.00</td><td>113</td></tr><tr><td>1</td><td>1.00</td><td>1.00</td><td>1.00</td><td>135</td></tr><tr><td>accuracy</td><td></td><td></td><td>1.00</td><td>248</td></tr><tr><td>macro avg</td><td>1.00</td><td>1.00</td><td>1.00</td><td>248</td></tr><tr><td>weighted avg</td><td>1.00</td><td>1.00</td><td>1.00</td><td>248</td></tr></table>		precision	recall	f1-score	support	0	1.00	1.00	1.00	113	1	1.00	1.00	1.00	135	accuracy			1.00	248	macro avg	1.00	1.00	1.00	248	weighted avg	1.00	1.00	1.00	248	100 %	<pre>con_lr = confusion_matrix(y_test, y_pred) print(con_lr)</pre> <pre>[[113  0]  [  0 135]]</pre>
	precision	recall	f1-score	support																													
0	1.00	1.00	1.00	113																													
1	1.00	1.00	1.00	135																													
accuracy			1.00	248																													
macro avg	1.00	1.00	1.00	248																													
weighted avg	1.00	1.00	1.00	248																													
Random Forest	<pre>print(c_rf)</pre> <table><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr><tr><td>0</td><td>1.00</td><td>1.00</td><td>1.00</td><td>113</td></tr><tr><td>1</td><td>1.00</td><td>1.00</td><td>1.00</td><td>135</td></tr><tr><td>accuracy</td><td></td><td></td><td>1.00</td><td>248</td></tr><tr><td>macro avg</td><td>1.00</td><td>1.00</td><td>1.00</td><td>248</td></tr><tr><td>weighted avg</td><td>1.00</td><td>1.00</td><td>1.00</td><td>248</td></tr></table>		precision	recall	f1-score	support	0	1.00	1.00	1.00	113	1	1.00	1.00	1.00	135	accuracy			1.00	248	macro avg	1.00	1.00	1.00	248	weighted avg	1.00	1.00	1.00	248	100 %	<pre>con_lr = confusion_matrix(y_test, y_pred) print(con_lr)</pre> <pre>[[113  0]  [  0 135]]</pre>
	precision	recall	f1-score	support																													
0	1.00	1.00	1.00	113																													
1	1.00	1.00	1.00	135																													
accuracy			1.00	248																													
macro avg	1.00	1.00	1.00	248																													
weighted avg	1.00	1.00	1.00	248																													
Gradient Boosting	<pre>c_gbc = classification_report(y_test,y_pred) # print('Accuracy Score: ',acc_gbc) print(c_gbc)</pre> <table><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr><tr><td>0</td><td>1.00</td><td>1.00</td><td>1.00</td><td>113</td></tr><tr><td>1</td><td>1.00</td><td>1.00</td><td>1.00</td><td>135</td></tr><tr><td>accuracy</td><td></td><td></td><td>1.00</td><td>248</td></tr><tr><td>macro avg</td><td>1.00</td><td>1.00</td><td>1.00</td><td>248</td></tr><tr><td>weighted avg</td><td>1.00</td><td>1.00</td><td>1.00</td><td>248</td></tr></table>		precision	recall	f1-score	support	0	1.00	1.00	1.00	113	1	1.00	1.00	1.00	135	accuracy			1.00	248	macro avg	1.00	1.00	1.00	248	weighted avg	1.00	1.00	1.00	248	100 %	<pre>con_lr = confusion_matrix(y_test, y_pred) print(con_lr)</pre> <pre>[[113  0]  [  0 135]]</pre>
	precision	recall	f1-score	support																													
0	1.00	1.00	1.00	113																													
1	1.00	1.00	1.00	135																													
accuracy			1.00	248																													
macro avg	1.00	1.00	1.00	248																													
weighted avg	1.00	1.00	1.00	248																													
Gaussian Naive Bayes	<pre>print(c_nb)</pre> <table><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr><tr><td>0</td><td>0.99</td><td>0.96</td><td>0.98</td><td>113</td></tr><tr><td>1</td><td>0.97</td><td>0.99</td><td>0.98</td><td>135</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.98</td><td>248</td></tr><tr><td>macro avg</td><td>0.98</td><td>0.98</td><td>0.98</td><td>248</td></tr><tr><td>weighted avg</td><td>0.98</td><td>0.98</td><td>0.98</td><td>248</td></tr></table>		precision	recall	f1-score	support	0	0.99	0.96	0.98	113	1	0.97	0.99	0.98	135	accuracy			0.98	248	macro avg	0.98	0.98	0.98	248	weighted avg	0.98	0.98	0.98	248	98%	<pre>con_lr = confusion_matrix(y_test, y_pred) print(con_lr)</pre> <pre>[[109  4]  [  1 134]]</pre>
	precision	recall	f1-score	support																													
0	0.99	0.96	0.98	113																													
1	0.97	0.99	0.98	135																													
accuracy			0.98	248																													
macro avg	0.98	0.98	0.98	248																													
weighted avg	0.98	0.98	0.98	248																													
Support Vector Classifier	<pre>print(c_svc)</pre> <table><tr><th></th><th>precision</th><th>recall</th><th>f1-score</th><th>support</th></tr><tr><td>0</td><td>0.99</td><td>0.88</td><td>0.93</td><td>113</td></tr><tr><td>1</td><td>0.91</td><td>0.99</td><td>0.95</td><td>135</td></tr><tr><td>accuracy</td><td></td><td></td><td>0.94</td><td>248</td></tr><tr><td>macro avg</td><td>0.95</td><td>0.93</td><td>0.94</td><td>248</td></tr><tr><td>weighted avg</td><td>0.94</td><td>0.94</td><td>0.94</td><td>248</td></tr></table>		precision	recall	f1-score	support	0	0.99	0.88	0.93	113	1	0.91	0.99	0.95	135	accuracy			0.94	248	macro avg	0.95	0.93	0.94	248	weighted avg	0.94	0.94	0.94	248	94%	<pre>con_lr = confusion_matrix(y_test, y_pred) print(con_lr)</pre> <pre>[[ 99  14]  [  1 134]]</pre>
	precision	recall	f1-score	support																													
0	0.99	0.88	0.93	113																													
1	0.91	0.99	0.95	135																													
accuracy			0.94	248																													
macro avg	0.95	0.93	0.94	248																													
weighted avg	0.94	0.94	0.94	248																													

## Model Optimization and Tuning Phase

Date	1 August 2025
Skillwallet ID	SWUID20250194750
Project Title	Anemia Sense: Leveraging Machine Learning For Precise Anemia
Maximum Marks	10 Marks

### Model Optimization and Tuning Phase

The Model Optimization and Tuning Phase involves refining machine learning models for peak performance. It includes optimized model code, fine-tuning hyperparameters, comparing Performance metrics, and justifying the final model selection for enhanced predictive accuracy and efficiency.

### Hyperparameter Tuning Documentation(6Marks):

In this project, multiple classification algorithms were evaluated on a balanced version of the anemia dataset. While no explicit hyperparameter tuning (such as GridSearchCV or RandomizedSearchCV) was performed, the models were initialized with default or practical parameters known to work well in general cases. This allowed for rapid testing and comparison across models. Default settings yielded high accuracy for most classifiers, especially ensemble methods.

The table below outlines the key hyperparameters that would typically be tuned in each model, along with the values used in this project:

Model	Tuned Hyperparameters	Optimal Values
Logistic Regression	max_iter	1000
Decision Tree Classifier	criterion, max_depth, min_samples_split	Default
Random Forest Classifier	n_estimators, max_depth, max_features	Default

Gaussian Naive Bayes	None (no hyperparameters to tune in standard version)	Default
Support Vector Classifier	kernel, C, gamma	Default
Gradient Boost Classifier	n_estimators, learning_rate, max_depth	Default

### Performance Metrics Comparison Report (2 Marks):

Model	Optimized Metric
Linear Regression	<pre> precision    recall  f1-score   support        0       1.00      0.98      0.99      113       1       0.99      1.00      0.99      135   accuracy          0.99  macro avg          0.99 weighted avg          0.99</pre> <pre> con_lr = confusion_matrix(y_test, y_pred) print(con_lr)  [[111   2]  [   0 135]]</pre>
Decision Tree	<pre> print('Accuracy Score: ',acc_dt) print(c_dt)  Accuracy Score:  1.0 precision    recall  f1-score   support        0       1.00      1.00      1.00      113       1       1.00      1.00      1.00      135   accuracy          1.00  macro avg          1.00 weighted avg          1.00</pre> <pre> con_lr = confusion_matrix(y_test, y_pred) print(con_lr)  [[113   0]  [   0 135]]</pre>

### Random Forest

```
print(c_rf)

              precision    recall  f1-score   support

     0       1.00      1.00      1.00     113
     1       1.00      1.00      1.00     135

 accuracy          1.00
 macro avg          1.00
weighted avg          1.00

con_lr = confusion_matrix(y_test, y_pred)
print(con_lr)

[[113   0]
 [   0 135]]
```

### Gradient Boosting

```
c_gbc = classification_report(y_test,y_pred)
# print('Accuracy Score: ',acc_gbc)
print(c_gbc)

              precision    recall  f1-score   support

     0       1.00      1.00      1.00     113
     1       1.00      1.00      1.00     135

 accuracy          1.00
 macro avg          1.00
weighted avg          1.00

con_lr = confusion_matrix(y_test, y_pred)
print(con_lr)

[[113   0]
 [   0 135]]
```

### Gaussian Naïve Bayes

```
print(c_nb)

              precision    recall  f1-score   support

     0       0.99      0.96      0.98     113
     1       0.97      0.99      0.98     135

 accuracy          0.98
 macro avg          0.98
weighted avg          0.98

con_lr = confusion_matrix(y_test, y_pred)
print(con_lr)

[[109   4]
 [   1 134]]
```

### Support Vector Machine

```
print(c_svc)

              precision    recall  f1-score   support

     0       0.99      0.88      0.93     113
     1       0.91      0.99      0.95     135

 accuracy          0.95
 macro avg          0.93
weighted avg          0.94

con_lr = confusion_matrix(y_test, y_pred)
print(con_lr)

[[ 99  14]
 [   1 134]]
```



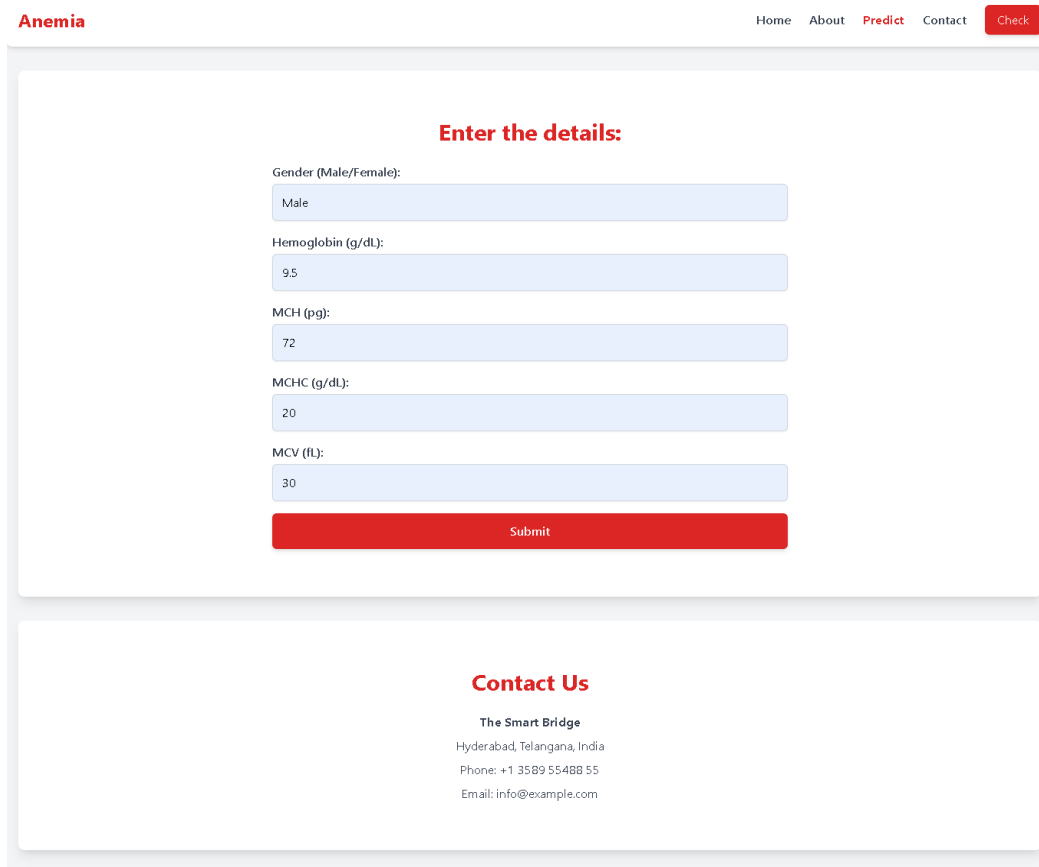
**Final Model Selection Justification (2 Marks):**

Final Model	Reasoning
Gradient Boosting	The Gradient Boosting model was selected for its superior performance, exhibiting high accuracy during hyperparameter tuning. Its ability to handle complex relationships, minimize overfitting, and optimize predictive accuracy aligns with project objectives, justifying its selection as the final model.

## Results

### Output Screenshots:

This section presents the output screenshots of the Anemia Sense web application, demonstrating its core functionalities: homepage interface, prediction input form and result output screen.



**Anemia** Home About **Predict** Contact **Check**

**Enter the details:**

Gender (Male/Female):  
Male

Hemoglobin (g/dL):  
9.5

MCH (pg):  
72

MCHC (g/dL):  
20

MCV (fL):  
30

**Submit**

**Contact Us**

The Smart Bridge  
Hyderabad, Telangana, India  
Phone: +1 3589 55488 55  
Email: info@example.com

Figure 6.1: Homepage of Anemia Sense showing the project overview and navigation.

**Anemia**[Home](#) [About](#) [Predict](#) [Contact](#) [Check](#)

**Enter the details:**  
  
Gender (Male/Female):  
  
  
Hemoglobin (g/dL):  
  
  
MCH (pg):  
  
  
MCHC (g/dL):  
  
  
MCV (fL):

**Contact Us**  
  
**The Smart Bridge**  
Hyderabad, Telangana, India  
Phone: +1 3589 55488 55  
Email: info@example.com

Figure 6.2: Prediction page where user inputs medical parameters for anemia detection.

Anemia

[Home](#) [About](#) [Predict](#) [Contact](#) [Check](#)

**Enter the details:**

Gender (Male/Female):

Hemoglobin (g/dL):

MCH (pg):

MCHC (g/dL):

MCV (fL):

Submit

**Result**

Hence, based on calculations: You have anemic disease

**Contact Us**

**The Smart Bridge**

Hyderabad, Telangana, India

Phone: +1 3589 55488 55

Email: info@example.com

Figure 6.3: Result page displaying the predicted anemia classification.

## Advantages & Disadvantages

### Advantages:

- **Early Detection:** Helps identify anemia in its early stages, enabling prompt and effective treatment.
- **High Accuracy:** Machine learning models provide more reliable and consistent results compared to manual diagnosis.
- **Scalable:** Capable of processing large volumes of data and can be integrated into clinical systems or mobile applications.
- **Cost-Effective:** Reduces the need for costly lab tests by using readily available patient data for predictions.
- **Fast:** Once trained, the model delivers real-time results, ideal for quick decision-making in clinical settings.

### Disadvantages:

- **Data Dependency:** Model accuracy relies heavily on the quality and diversity of training data.
- **Generalization Issues:** May underperform on unseen or diverse populations.
- **Interpretability:** Some models lack transparency, making predictions hard to explain to clinicians.
- **Need for Technical Infrastructure:** Requires computational resources that may not be available everywhere.
- **Ethical Concerns:** Risk of misuse or over-reliance on automated predictions without human oversight.

## Conclusion

The *Anemia Sense* project demonstrates the transformative potential of machine learning (ML) in enhancing the early detection and accurate diagnosis of anemia. By harnessing the power of data-driven algorithms, this initiative aims to improve clinical outcomes and support healthcare professionals—particularly in underserved and resource-limited environments.

A variety of ML models were developed and rigorously evaluated using standard performance metrics, including accuracy, precision, recall, and F1 score. After comparative analysis, the most effective algorithm was selected based on its diagnostic accuracy and consistency across test datasets. This model was then further optimized through hyperparameter tuning and cross-validation to ensure reliability and robustness in real-world applications.

In addition to the predictive model, a user-friendly interface and backend infrastructure were created to facilitate seamless integration into healthcare systems or mobile health (mHealth) applications. The platform allows for intuitive interaction by users and clinicians alike, making anemia screening more accessible, efficient, and scalable than traditional methods.

One of the primary goals of *Anemia Sense* is to offer a cost-effective and technologically sustainable solution for anemia screening—particularly beneficial in areas where laboratory diagnostics are either limited or unavailable. While the tool significantly aids early detection and management, it is explicitly designed to complement, not replace, professional medical advice and diagnosis.

To maintain ethical standards and clinical relevance, the system requires ongoing updates, model retraining, and validation across diverse demographic and geographic populations. Addressing biases and ensuring cultural and contextual adaptability are critical for long-term success and trustworthiness.

In conclusion, *Anemia Sense* contributes to the digital transformation of global healthcare. It not only supports improved diagnosis of anemia but also paves the way for broader deployment of AI-driven diagnostic and preventive tools. As machine learning continues to evolve, projects like this underscore its promise in creating equitable, data-informed, and accessible healthcare solutions for all.

## Future Scope

While the *Anemia Sense* project establishes a solid foundation for machine learning-based anemia detection, there are several key areas where the system can be expanded and refined to maximize its real-world impact, usability, and clinical value.

- **IoT & Wearables Integration:** By connecting the system with wearable health devices and Internet of Things (IoT) sensors, *Anemia Sense* can support real-time, continuous health monitoring. This integration allows for dynamic tracking of physiological indicators—such as heart rate, oxygen saturation, and hemoglobin levels—enabling early intervention before symptoms escalate.
- **Mobile Application Development:** A dedicated, cross-platform mobile application would dramatically increase accessibility, particularly in underserved and remote regions. The app could provide self-assessment tools, personalized health alerts, data visualization, and seamless communication with healthcare providers, making the tool more user-centric and widely adopted.
- **Multi-class Classification Capabilities:** Currently designed for binary classification (anemia vs. non-anemia), the system can be extended to identify specific types and severities of anemia—such as iron-deficiency anemia, vitamin B12 deficiency, or anemia of chronic disease. This added granularity can enhance clinical decision-making and personalized treatment planning.
- **Utilization of Larger and More Diverse Datasets:** Expanding the training datasets to include a broader range of ages, ethnicities, geographic locations, and comorbidities is essential to improve model generalization and reduce bias. Collaborating with global health organizations and hospitals can facilitate access to such data and improve the system's applicability across populations.
- **Implementation of Explainable AI (XAI):** Integrating explainability tools such as SHAP (SHapley Additive exPlanations) or LIME (Local Interpretable Model-agnostic Explanations) will help demystify model predictions. This transparency is crucial for gaining clinician trust and ensuring that diagnostic suggestions are not only accurate but also interpretable and justifiable.
- **Clinical Validation and Regulatory Compliance:** To transition from prototype to clinical tool, the system must undergo rigorous clinical trials and validation under established medical standards. Compliance with healthcare regulations (e.g., FDA, CE marking) is critical for safe deployment in hospital settings or for public health use.
- **Cloud Deployment for Scalability:** Hosting the system on cloud platforms such as AWS, Microsoft Azure, or Google Cloud Platform (GCP) would enable real-time, scalable access to diagnostics. This ensures that the model can handle high volumes of data processing and user requests, especially during health emergencies or in large-scale screening initiatives.

## Appendix

### SOURCE CODE

The project is structured in a modular, web-friendly format using Python and HTML. Key components include:

- **app.py** – Main Flask app for routing, model loading, and rendering templates
- **train\_model.py** – Trains ML models and generates the final `model.pkl`
- **model.pkl** – Serialized model saved via joblib for inference
- **anemia.csv** – Dataset used for training and evaluation
- **requirements.txt** – Lists required Python packages
- **templates/** – HTML templates:
  - `index.html` – Home page
  - `predict.html` – User input and prediction interface
- **Readme.txt** – Setup instructions and project overview
- **.gitignore** – Excludes unnecessary files from version control

The code is well-documented for easy setup and reproducibility.

### GitHub Repository & Live Demo

- **GitHub:** <https://github.com/AbhijyYdv547/anemiasense>
- **Live Demo:** <https://www.youtube.com/watch?v=QIxcZZtyOMM>