# Project Report

## 1. INTRODUCTION

1.1 Project Overview

1.2 Purpose

## 2. LITERATURE SURVEY

2.1 Existing problem

2.2 References

2.3 Problem Statement Definition

## 3. IDEATION & PROPOSED SOLUTION

3.1 Empathy Map Canvas

3.2 Ideation & Brainstorming

## 4. REQUIREMENT ANALYSIS

4.1 Functional requirement

4.2 Non-Functional requirements

## 5. PROJECT DESIGN

5.1 Data Flow Diagrams & User Stories

5.2 Solution Architecture

## 6. PROJECT PLANNING & SCHEDULING

6.1 Technical Architecture

6.2 Sprint Planning & Estimation

6.3 Sprint Delivery Schedule

## 7. CODING & SOLUTIONING

(Explain the features added in the project along with code)

7.1 Feature 1

## 8. PERFORMANCE TESTING

8.1 Performance Metrics

## 9. RESULTS

9.1 Output Screenshots

## 10. ADVANTAGES & DISADVANTAGES

## 11. CONCLUSION

## 12. FUTURE SCOPE

## 13. APPENDIX

Source Code

GitHub & Project Demo

## 1.Introduction:

### 1.1 Project Overview:

### Project Name:

Lymphography Classification Using ML

### Project Objective:

The primary objective of this project is to develop an accurate and robust lymphography classification system using the Random Forest algorithm. The goal is to leverage machine learning techniques to analyze medical data and classify lymphography data into distinct categories, such as normal, metastasis, or Fibrosis.

### Key Components and Features:

### 1. Data Collection:

The project starts with the collection of relevant data, including Lymphatics, change in node, extravasates, special forms, dislocation, regeneration and other medical information.

### 2.Data Preprocessing:

Perform thorough preprocessing on the lymphography dataset, including data cleaning, normalization, and feature extraction, to ensure the quality and relevance of input data for the Random Forest algorithm.

### 3.Model Development:

Implement and fine-tune a Random Forest classifier to effectively learn from the preprocessed lymphography data. Optimize hyperparameters to enhance the model's performance in terms of accuracy, precision, recall, and F1 score.

### 4.Feature Analysis:

Conduct an in-depth analysis of feature importance within the Random Forest model to identify the key factors influencing lymphography classification. This analysis can provide valuable insights into the medical relevance of certain features.

**5.Model Evaluation:**

Rigorously evaluate the performance of the developed Random Forest model using appropriate metrics such as accuracy, precision, recall, F1 score, and area under the receiver operating characteristic curve. Ensure the model's generalizability on both training and testing datasets.

**6. User Interface:**

The project may include user-friendly interfaces for healthcare providers. These interfaces can display real-time data, predictions, and recommendations in an easily understandable format.

**7.Integrating with Flask:**

We will be building a web application that is integrated to the model we built. A UI is provided for the uses where he has to enter the values for predictions. The enter values are given to the saved model and prediction is showcased on the UI.

**8. Deploying the Model:**

When our model is ready for prediction, we deploy it using services like AWS.

**Benefits:**

- Improved Diagnostic Accuracy
- Early Detection
- Personalized Treatment
- Improved Patient Experience
- Reduced Healthcare Costs

**Challenges:**

- Data Quality and Quantity
- Interpretability of Results
- Overfitting and Model Complexity
- Computational Resources
- Ethical and Regulatory Considerations

**1.2 Purpose:**

The primary purpose of this project is to improve the accuracy and reliability of lymphography diagnostics. By leveraging machine learning techniques, specifically the Random Forest algorithm, the project aims to develop a robust classification system capable of accurately identifying patterns associated with different lymphatic conditions.

It can be framed in the context of addressing key challenges in medical diagnostics and contributing to advancements in healthcare. The project aims to contribute to early diagnosis by creating a classification model that provides rapid and accurate results.

By understanding the features and patterns indicative of different lymphatic conditions, the project aims to provide valuable information that may inform the diagnosis and treatment of related disorders.

**2.Literature Survey:**

**2.1 Existing Problem:**

lymphography classification faces challenges related to manual analysis, subjectivity, and potential diagnostic errors. Traditional methods struggle with complexity, leading to limitations in accuracy and efficiency.

Lymphography is a crucial diagnostic tool in medical imaging, playing a key role in identifying lymphatic system disorders. Current manual methods of lymphography classification suffer from subjectivity and potential diagnostic errors. To address these challenges, machine learning algorithms, particularly Random Forest, have shown promise in automating the classification process.

**2.2 References:**

Zelikovski A, Manoach M, Giler S, Urca I. Lympha-press, a new pneumatic device for the treatment of lymphedema of the limbs. *Lymphology* 1980; 13: 68–73.

Fyfe NC, Wolfe JH, Kinmonth JB. 'Die-back' in primary lymphedema — lymphographie and clinical correlations. *Lymphology* 1982; 15: 66–69.

O'Brien BMcC, Shafiroff BB. Microlymphaticovenous and resectional surgery in obstructive lymphoedema. *World J Surg* 1979; 3: 3–15.

https://www.semanticscholar.org/paper/Lymph-Diseases-Prediction-Using-Random-Forest-and-Almayyan/e864236c21839cc26b4a0a2ef37cc03b8179a18e?p2df
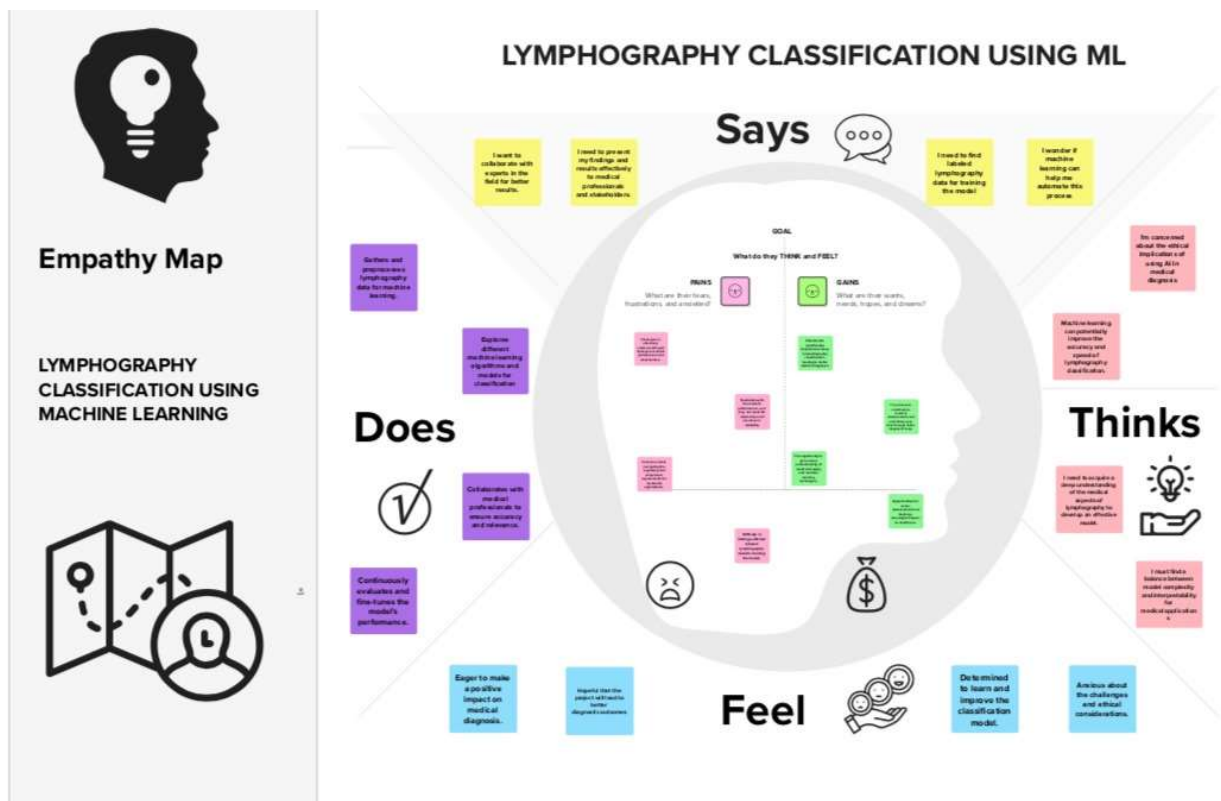
## 2.3 Problem Statement Definition

Lymphography, a vital diagnostic technique for assessing the lymphatic system, faces challenges in accurate and efficient classification of lymphatic disorders. Manual analysis is prone to subjectivity and potential diagnostic errors, and existing automated methods often lack the necessary accuracy and interpretability. The objective is to develop a robust lymphography classification system that overcomes these challenges using the Random Forest algorithm.

The project aims to address the limitations of current lymphography classification methods by leveraging the capabilities of the Random Forest algorithm. The specific challenges include the need for accurate and interpretable classification, handling the complexity of lymphography data, and providing a system that can generalize well across diverse datasets. The project seeks to design, implement, and optimize a Random Forest-based classification system for lymphographic data, with a focus on achieving high accuracy, interpretability, and adaptability to varying clinical scenarios.

## 3.Ideation & Proposed Solution:

## 3.1 Empathy Map Canvas:

## 3.2 Ideation & Brainstorming:

## Brainstrom:

### VASHNAVI

Develop a ml algorithm that can be trained on lymphography data to classify them into different categories, such as normal, metastases, malignant lymph, and fibrosis

Collect, clean, and split the lymphography dataset from a CSV file for machine learning.

Design the system to be scalable, so it can handle larger datasets and potentially be used in a broader range of healthcare settings

Consider tracking and analyzing patient outcomes to assess the impact of the system on diagnosis and treatment decisions

### SRAVANI

Deploy the model for real-time lymphography classification.

Establish a feedback loop with healthcare providers to continually improve the system based on their real-world experiences and feedback.

Involve medical experts in the model development process to provide domain knowledge, validate results, and ensure that the model aligns with clinical practices

Evaluate the model's performance on the test set to get an unbiased estimate of its accuracy.Compute metrics such as accuracy, precision

### ARUN KUMAR

develop a model that can classify lymphography data based on their texture and shape. This model could be trained on a large dataset.

.Use a multi-task learning approach to train a machine learning model to perform multiple tasks simultaneously

Train the Random Forest model with the prepared lymphography data.

Develop a ml algorithm that can be used to track changes in lymphography over time for monitoring disease progression and treatment response.

### SAIKIRAN

Use a data augmentation approach to increase the size and diversity of the lymphography dataset

Group tasks into data preparation, model selection and training, and model evaluation and deployment for efficient project management

Assess the model's performance using accuracy, precision, recall, and F1 score.

Develop a ml-based lymphography classification system that is integrated with a telemedicine platform

**Group Ideas:**

**Group ideas**

Take turns sharing your ideas while clustering similar or related notes as you go. Once all sticky notes have been grouped, give each cluster a sentence-like label. If a cluster is bigger than six sticky notes, try and see if you and break it up into smaller sub-groups.

Stratified split the CSV file into two datasets based on the lymphography classification labels.
This will ensure that the training and test sets have a similar distribution of labels, which can improve the model's performance.

Choose a machine learning model that is appropriate for the classification task. Consider using a random forest approach to fine-tune a pre-trained model on the lymphography dataset.

Train the model on the training set using the selected hyperparameters.Monitor the model's performance on the validation set to prevent overfitting.

Evaluate the model's performance on the test set to get an unbiased estimate of its accuracy.Compute metrics such as accuracy, precision
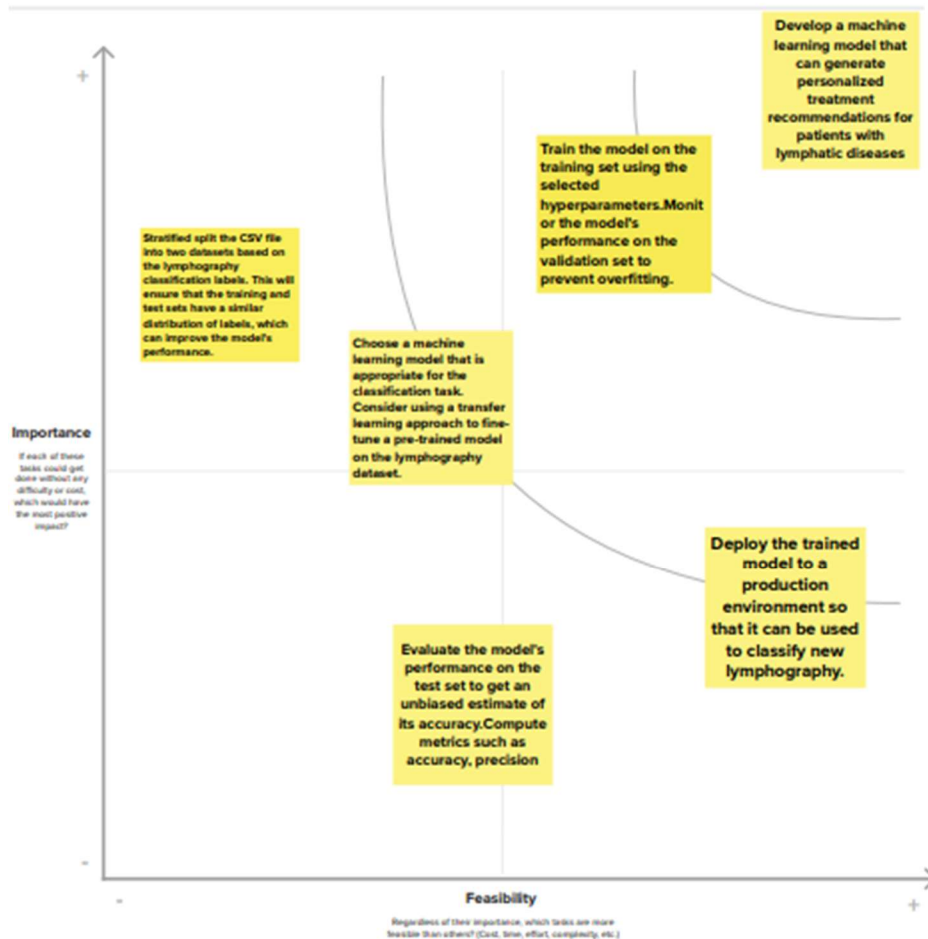
Deploy the trained model to a production environment so that it can be used to classify new lymphography.

Develop a machine learning model that can generate personalized treatment recommendations for patients with lymphatic diseases

④
**Prioritize**

Your team should all be on the same page about what's important moving forward. Place your ideas on this grid to determine which ideas are important and which are feasible.

**Importance**

If each of these tasks could get done without any difficulty or cost, which would have the most positive impact?

Develop a machine learning model that can generate personalized treatment recommendations for patients with lymphatic diseases

Train the model on the training set using the selected hyperparameters.Monitor the model's performance on the validation set to prevent overfitting.

Stratified split the CSV file into two datasets based on the lymphography classification labels. This will ensure that the training and test sets have a similar distribution of labels, which can improve the model's performance.

Choose a machine learning model that is appropriate for the classification task. Consider using a transfer learning approach to fine-tune a pre-trained model on the lymphography dataset.

Deploy the trained model to a production environment so that it can be used to classify new lymphography.

Evaluate the model's performance on the test set to get an unbiased estimate of its accuracy.Compute metrics such as accuracy, precision

**Feasibility**

Regardless of their importance, which tasks are more feasible than others? (Cost, time, effort, complexity, etc.)

---

## 4.Requirement Analysis:

### 4.1 Functional Requirement:

The system should perform data preprocessing tasks, including image normalization, feature extraction, and handling missing data.Implement a Random Forest classification model capable of learning from preprocessed lymphography data. Optimize hyperparameters for the Random Forest algorithm to enhance classification performance.

Provide functionality for analyzing and interpreting feature importance within the Random Forest model to identify key factors influencing lymphography classification. Conduct rigorous model evaluation using metrics such as accuracy, precision, recall, F1 score. Ensure the model's performance on both training and testing datasets.
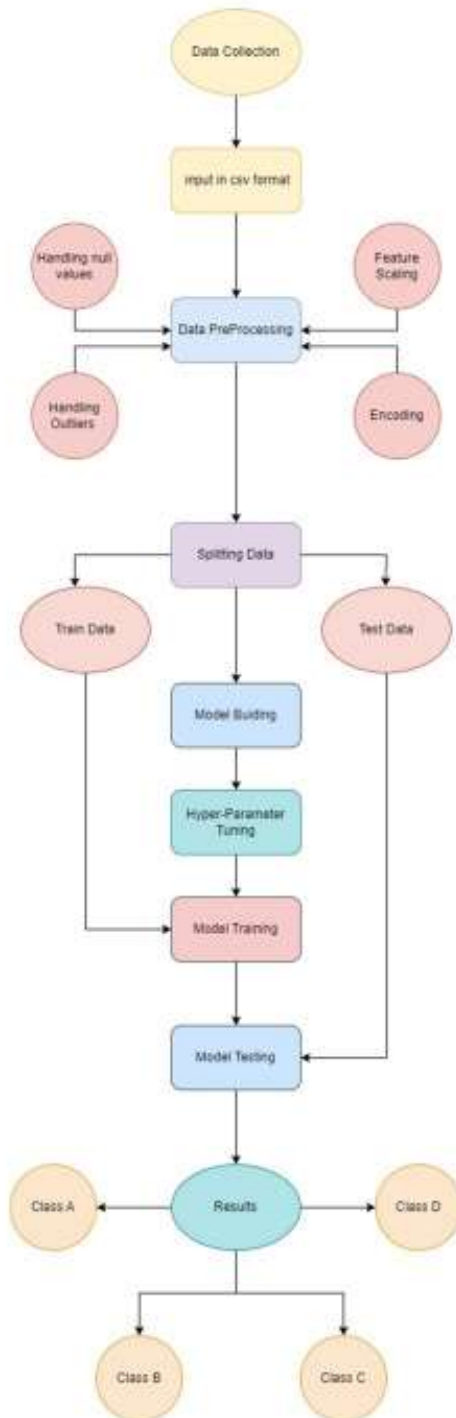
## 4.2 Non-Functional Requirement:

The system should be able to process and classify lymphography data within a reasonable timeframe to meet real-time clinical requirements. The system should be designed to handle an increasing volume of lymphography data as the dataset grows over time. Ensure compatibility with standard medical diagnosis and integrate seamlessly with existing healthcare information systems.

The classification system should be reliable, providing consistent and accurate results across different datasets and under varying conditions. Implement robust security measures to protect patient data and ensure compliance with healthcare privacy regulations. The system should be designed with modularity and code maintainability in mind to facilitate future updates and improvements.

**5.Project Design:**

**5.1 DataFlow Diagrams & User Stories:**

| User Type | Functional Requirement (Epic) | User Story Number | User Story / Task | Acceptance criteria | Priority | Release |
|---|---|---|---|---|---|---|
| Patient | Interface | US1 | Need a friendly interface with proper labels | Easily navigable interface | High | Sprint-1 |
| | | US2 | Expects fields to enter information | Distinctly visible fields | High | Sprint-1 |
| | Prediction | US3 | Needs prediction any number of times in a single page session | Results based on varied inputs | High | Sprint-1 |
| | | US4 | Expects a clear result on the type of lymphography Disease, if any | Single,most probable disease category | High | Sprint-2 |
| | | US5 | Needs a clear description of the predicted disease | Elongated description of the prediction | High | Sprint-2 |

## 5.2 Solution Architecture:

The basic architecture of the proposed solution revolves around the fundamental machine building using Machine Learning Algorithm, which is Random Forest in our project.
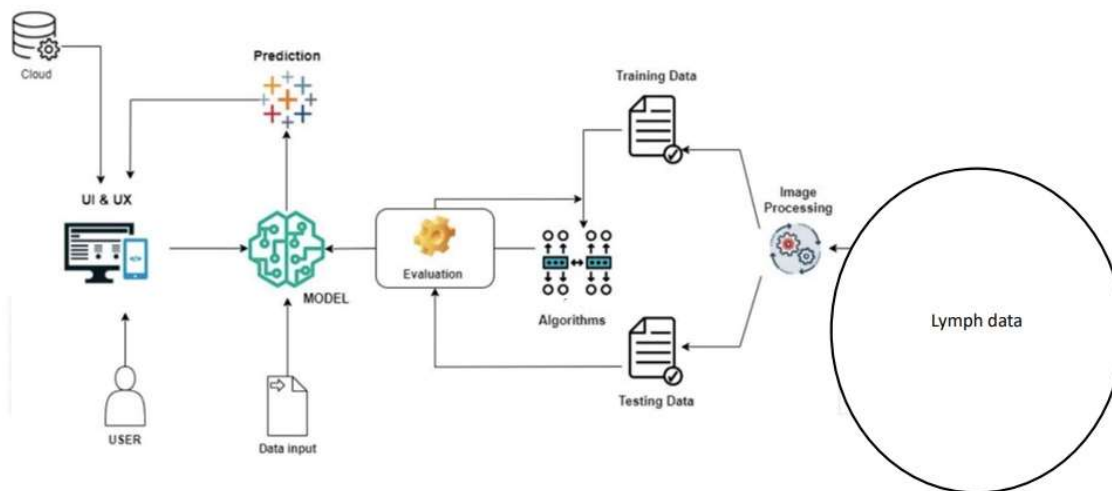
Building blocks

- Data Set
- Model (Built using scikit Learn library (Python))
- Front-End interface
- Back-End support (To host the application )

**Work Flow**

- Collect the data
- Data Preprocessing
- Splitting the data into
- Train Data
- Test Data
- Validation Data
- Initializing the model
- Training the model
- Testing the model
- Saving the model
- Integrating Flask with the ML model
- Hosting the application

**Solution Architecture Diagram:**

# 6. Project Planning & Scheduling
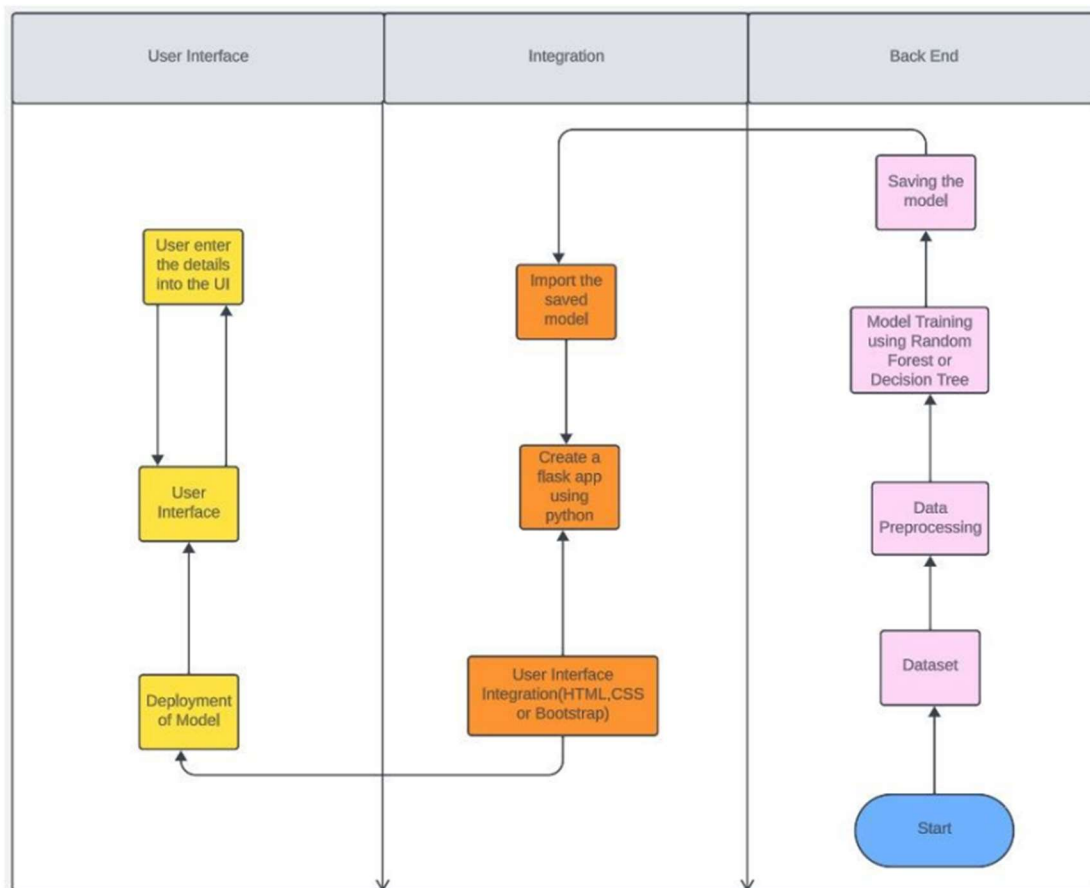
## 6.1 Technical Architecture:

**Table-1: Components & Technologies**

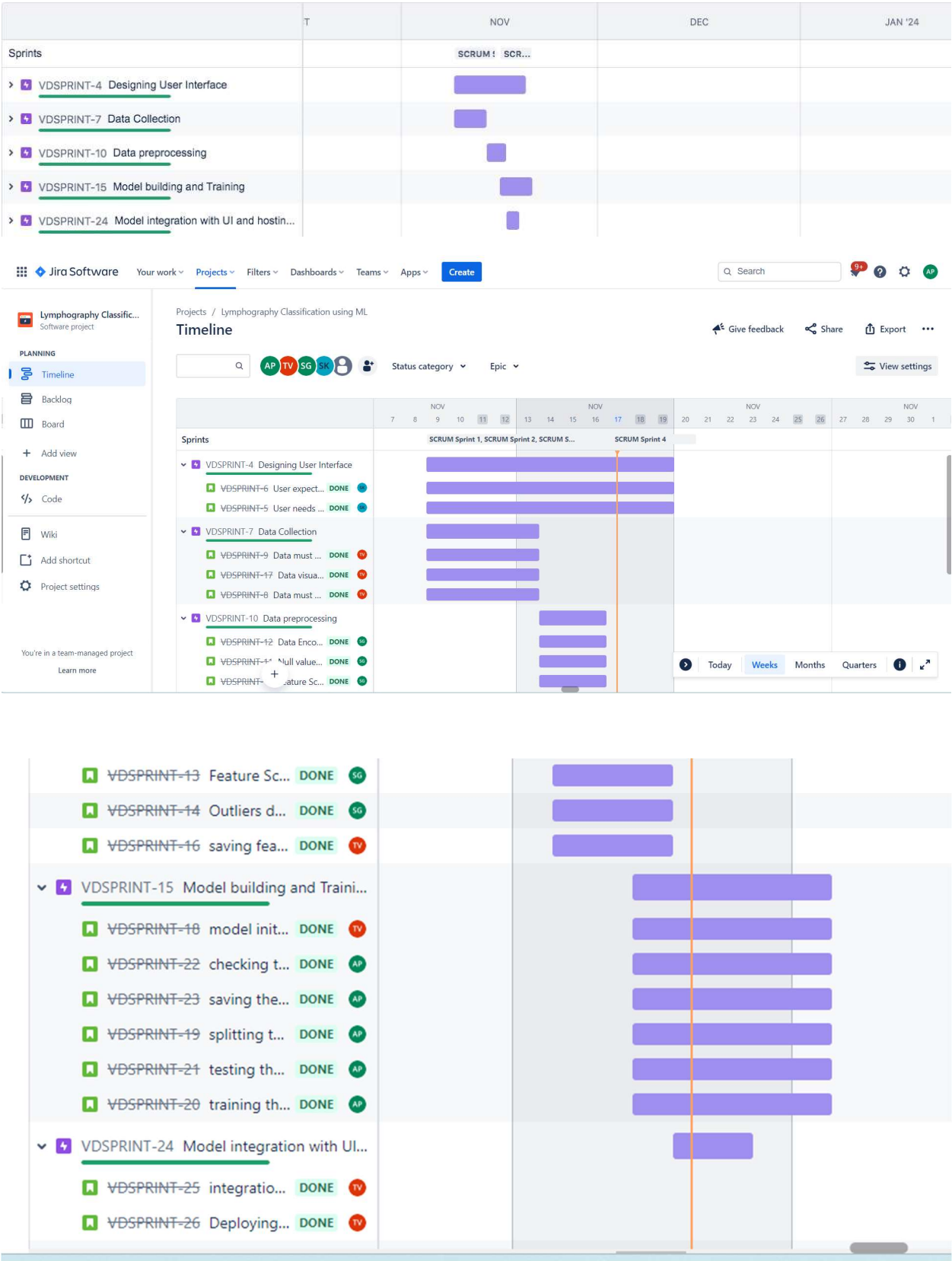| SNO | Component | Description | Technology |
|-----|-----------|-------------|------------|
| 1 | User Interface | Web UI | HTML, CSS, JavaScript |
| 2 | Application Logic-1 | Data Preprocessing | Python, Numpy |
| 3 | Application Logic-2 | Creating ML model | Necessary Python Libraries |
| 4 | Application Logic-3 | Web application | Flask |
| 5 | Machine Learning Model | ML model using Random Forest | Machine learning algorithm (Random Forest) from scikit learn |
| 6 | Infrastructure (Server / Cloud) | Application Deployment on Cloud Server | AWS EC2 |

**Table-2: Application Characteristics:**

| SNO | Characteristics | Description | Technology |
|-----|-----------------|-------------|------------|
| 1 | Open-Source Frameworks | Flask | Technology of Open Source framework |
| 2 | Security Implementations | CSRF Protection, Secure Flag For Cookies | SHA-256, Encryptions, IAM Controls, OWASP etc. |
| 3 | Scalable Architecture | 3 – tier, Micro-services | Micro web applications using Flask |
| 4 | Availability | U se of load balancers (ALB), distributed servers etc, | Application Load balancer Werkzeug,Jinja2,S inatra RubyFramework |

| | | Orm-Agnostic, Web Framework,Wsgi 1.0Compliant, Http Request Handling Functionality High Flexibility | SQLAlchemy,Extensions, Werkzeug,Jinja2,Sinatra RubyFramework |
|---|---|---|---|
| 5 | Performance | | |

## 6.2 Sprint Planning and Estimation:

| Sprint | Functional Requirement (Epic) | User Story Number | User Story / Task | Story Points | Priority | Team Memebers |
|---|---|---|---|---|---|---|
| Sprint-1 | Interface | US1 | Need a friendly interface with proper labels | 2 | High | Sai Kiran |
| | | US2 | Expects fields to enter information | 1 | High | Sai Kiran |
| | Prediction | US3 | Needs prediction any number of times in a single page session | 1 | High | Vaishnavi |
| Sprint 2 | | US4 | Expects a clear result on the type of lymphography Disease, if any | 1 | High | Sravani |
| | | US5 | Needs a clear description of the predicted disease | 1 | High | Arun |

## 6.3 Sprint Delivery Schedule:



| | | T | NOV | DEC | JAN '24 |
|---|---|---|---|---|---|
| Sprints | | | SCRUM ! SCR... | | |
| > VDSPRINT-4 Designing User Interface | | | | | |
| > VDSPRINT-7 Data Collection | | | | | |
| > VDSPRINT-10 Data preprocessing | | | | | |
| > VDSPRINT-15 Model building and Training | | | | | |
| > VDSPRINT-24 Model integration with UI and hostin... | | | | | |

**7. Coding & Solutioning**

  **7.1** A classification Model based on Random forest with Highly desirable Performance.

```
[34]   1 grid_search.fit(x_train,y_train)

       Fitting 2 folds for each of 96 candidates, totalling 192 fits
       ┌ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ┐
       │  ▸          GridSearchCV                         │
       │  ▸ estimator: RandomForestClassifier             │
       │       ┌──────────────────────────────┐          │
       │       │ ▸ RandomForestClassifier       │          │
       │       └──────────────────────────────┘          │
       └ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ┘
```

```
[35]   1 grid_search.best_score_

       0.788135593220339
```

```
[36]   1 rf_best=grid_search.best_estimator_
       2 rf_best
```

```
                        RandomForestClassifier
RandomForestClassifier(max_depth=10, max_features=2, min_samples_leaf=3)
```

```
   ●   1 rf_classify=RandomForestClassifier(random_state=42,
       2                                     n_jobs=-1,
       3                                     max_depth=9,
       4                                     min_samples_split=2,
       5                                     max_features='sqrt',
       6                                     n_estimators=90,
       7                                     bootstrap=True)
```

```
[38]   1 rf_classify.fit(x_train,y_train)
```

**7.2** Highly Scalable application deployed on Cloud.

## 8. Performance Testing

### 8.1 Performance Metrics

### Confusion Matrix:

The function is called and executed the confmatrix function block.The output is displayed as below:

```
[39]  1 from sklearn.metrics import accuracy_score
```
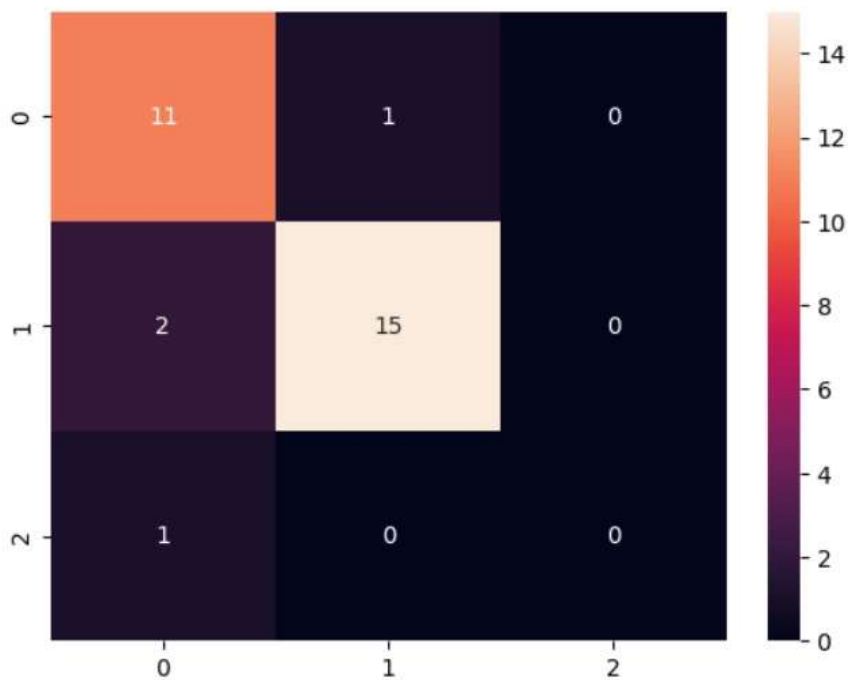
```
[40]  1 prediction=rf_classify.predict(x_test)
```

```
[41]  1 from sklearn.metrics import accuracy_score,f1_score,confusion_matrix,classification_report
```

```
1 confusion_matrix(y_test,prediction)
```

```
1 import seaborn as sns
2 sns.heatmap(cf, annot=True)
```

<Axes: >

**Model Accuracy:**

Model accuracy is defined as the number of classifications a model correctly predicts divided by the total number of predictions made. It's a way of assessing the performance of a model, but certainly not the only way.

\

```
[43]   1 accuracy_score(y_test,prediction)
```

```
0.8666666666666667
```

**Classification Report:**

```
[44]   1 print(classification_report(y_test,prediction))
```

```
                precision    recall  f1-score   support

           2        0.79      0.92      0.85        12
           3        0.94      0.88      0.91        17
           4        0.00      0.00      0.00         1

    accuracy                            0.87        30
   macro avg        0.57      0.60      0.59        30
weighted avg        0.85      0.87      0.85        30
```

# 9.Results

## 9.1 Output ScreenShots

```
1 print(rf_classify.predict(ms.transform([[4,2,1,1,1,1,1,2,1,2,2,2,4,8,1,1,2,2]])))
```

```
[3]
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature
  warnings.warn(
```

```
1 print(rf_classify.predict(ms.transform([[3,1,1,1,1,1,1,1,1,2,2,4,3,5,1,2,2,1]])))
```

```
[2]
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature
  warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature
  warnings.warn(
```

**Lymphography Prediction**

Enter the Lymphatics :

`4`

Enter lym.nodes enlar :

`2`

Enter the block of affere :

`2`

Enter changes in lym :

`2`

Enter the bl. of lymph. c :

`1`

Enter defect in node :

`2`

Enter the bl. of lymph. s :

`1`

Enter changes in node :

`4`

Enter the bypass :

`1`

Enter changes in stru :

`8`

Enter the extravasates :

`1`

Enter special forms :

`1`

Enter regeneration of :

`1`

Enter dislocation of :

`1`

Enter the block of affere :

2

Enter changes in lym :

2

Enter the bl. of lymph. c :

1

Enter defect in node :

2

Enter the bl. of lymph. s :

1

Enter changes in node :

4

Enter the bypass :

1

Enter changes in stru :

8

Enter the extravasates :

1

Enter special forms :

1

Enter regeneration of :

1

Enter dislocation of :

1

Enter the early uptake in :

2

Enter exclusion of no :

2

Enter the lym.nodes dimin :

1

Enter no. of nodes in :

2

submit

Output:

**Prediction Result**

**MALIGN LYMPH**

When people talk about malignancy in the context of the lymphatic system, they often refer to cancer that has spread to the lymph nodes or originated in the lymphatic system. Lymph nodes are small, bean-shaped structures that produce and store cells that help fight infection. If cancer cells break away from a tumor, they can travel through the lymphatic system and form new tumors in other parts of the body.Common cancers that can involve lymph nodes include lymphomas (cancers of the lymphatic system) and metastatic cancers (cancers that have spread from their original site to other parts of the body).

## 10. Advantages and Disadvantages

### 10.1 Advantages:

**High Accuracy:**

Random Forest is known for its ability to provide high accuracy in classification tasks. It can effectively handle complex patterns in medical data, contributing to more reliable diagnoses.

**Ensemble Learning:**

Random Forest is an ensemble learning method, combining the predictions of multiple decision trees. This ensemble approach often leads to improved generalization and robustness, reducing the risk of overfitting.

**Feature Importance Analysis:**

Random Forest provides a built-in mechanism for assessing feature importance. This is valuable in the medical domain, as it can offer insights into the relevance of different imaging features for lymphography classification.

**Handle Nonlinear Relationships:**

Random Forest is capable of capturing nonlinear relationships within the data, making it suitable for complex medical classification tasks where features may exhibit intricate interactions.

**Reduced Sensitivity to Noise:**

The ensemble nature of Random Forest makes it less sensitive to noisy data compared to individual decision trees. This is beneficial when working with medical imaging datasets that may have inherent noise or variability.

**Interpretability:**

While Random Forest is an ensemble model, it still provides a degree of interpretability. Feature importance analysis and visualization tools can help medical professionals understand the factors influencing classification decisions.

**Versatility:**

Random Forest can handle both classification and regression tasks, providing versatility in application. This allows for potential extensions of the project to address related medical analysis challenges.

**10.2 Disadvantages:**

**Computational Intensity:**

Training a Random Forest model can be computationally intensive, especially with large datasets and numerous decision trees. This might require substantial computational resources.

**Black-Box Nature:**

Despite providing some interpretability, Random Forest is considered a "black-box" model. Understanding the decision-making process for individual predictions may be challenging, which can be a concern in critical medical applications.

**Overfitting Risk:**

Random Forests are susceptible to overfitting, especially if not properly tuned. Careful hyperparameter tuning and validation are necessary to mitigate this risk and ensure the model generalizes well to new data.

**Training Time:**

The training time for Random Forests can be longer compared to simpler models. This may be a consideration in situations where real-time processing is crucial.

**Memory Usage:**

Random Forests can be memory-intensive, particularly as the number of trees in the ensemble increases. Memory constraints may impact the scalability of the model.

**Limited Performance Gain with Small Datasets:**

Random Forests may not provide a significant performance improvement over simpler models when working with small datasets. This could be a consideration if the available lymphography dataset is limited.

**Difficulty in Handling Imbalanced Data:**

Random Forests may struggle to perform well with highly imbalanced datasets. If the distribution of classes in the lymphography dataset is uneven, this imbalance may affect the model's ability to accurately classify the minority class.

**11.Conclusion**

   In conclusion, this project aimed to develop a robust lymphography classification system using the Random Forest algorithm, addressing challenges in accuracy and interpretability associated with current methods. Through comprehensive data preprocessing, model development, and feature importance analysis, the Random Forest classifier demonstrated its efficacy in accurately classifying lymphatic system disorders. The system's interpretability was enhanced through insightful feature importance analysis, providing valuable insights for medical professionals.

The advantages of Random Forest, including its ability to handle complex patterns, ensemble learning for improved generalization, and feature importance analysis, were leveraged to achieve high accuracy in lymphography classification. The system's usability was emphasized through a user-friendly interface, facilitating seamless integration into clinical workflows.

**12.Future Scope**

Despite the success of the project, challenges such as computational intensity during training, the black-box nature of the model, and potential overfitting risks were acknowledged. Ongoing efforts to optimize these aspects should be considered for further refinement.

This project contributes to the field of medical image analysis by showcasing the potential of machine learning, particularly Random Forest, in improving lymphography diagnostics. The developed classification system has the potential to enhance early detection, support medical professionals, and contribute to personalized treatment plans for patients with lymphatic system disorders.

**Future Directions:**

While this project addressed key challenges, there are avenues for further research and improvement. Future work could focus on:

**Model Optimization:**

Fine-tune hyperparameters and explore advanced techniques to mitigate potential overfitting, reducing computational intensity without compromising accuracy.

**Interpretability Enhancement:**

Investigate methods to enhance the interpretability of the Random Forest model, providing clearer insights into the decision-making process for individual classifications.

**Real-Time Deployment:**

Develop strategies for real-time deployment, ensuring the system's efficiency in clinical settings without compromising accuracy.

**Collaboration with Medical Professionals:**

Collaborate closely with medical professionals to incorporate domain-specific knowledge and ensure the system aligns with the practical needs of healthcare practitioners.

**13.Appendix**

**13.1Source Code :** smartinternz02/SI-GuidedProject-611862-1698429144

**13.2 Github :** smartinternz02/SI-GuidedProject-611862-1698429144

**13.3DemoVideo:**

**https://drive.google.com/file/d/1qb8fbPzUWAEYUfz4XdTT7Q9_P1pP6Z0w/view?usp=sharing**