

1 - Introduction and Overview

- Project idea in detail:

Optical character recognition or optical character reader (OCR) is the automated conversion of images of typed, handwritten, or printed text into machine-encoded text, whether from a scanned document, a photo of a document, a scene-photo (e.g., the text on signs and billboards in a landscape photo) or from subtitle text superimposed on an image (e.g., from a television broadcast). The goal of a character recognition system is to transform handwritten text documents on paper into a digital format that can be manipulated by word processor software. The system is required to identify a given input character form by mapping it to a single character in each character set.

- Similar applications in the market:

1 - Evernote:

This app can quickly capture all your handwritten ideas and notes with Evernote's built-in camera on Android and Apple devices. It was already possible to digitize texts in Evernote with the help of the Penultimate app, but an update also allows it to be done from your own application thanks to the Evernote Scannable function. Among its various options, it allows you to digitize the text you find, expanding the possibilities and being a perfect ally for, for example, passing notes between colleagues.

२- Pen to print:

his app can quickly capture all your handwritten ideas and notes with Evernote's built in camera on Android and Apple devices. It was already possible to digitize texts in Evernote with the help of the Penultimate app, but an update also allows it to be done from your own application thanks to the Evernote Scannable function. Among its various options, it allows you to digitize the text you find, expanding the possibilities and being a perfect ally for, for example, passing notes between colleagues.

३- Write for iPad:

If you prefer to write longhand but need to see your text in digital format, consider Write Pad for iPad. You can configure a host of options to recognize input forms and predefined commands, or you can input lettering with your finger or a stylus.

-A Literature Review of Academic publications(papers/books/articles) relevant to the problem

1)

Optical Character Recognition (OCR) is a piece of software that converts printed text and images into digitized form such that it can be manipulated by machine. Unlike the human brain which has the capability to recognize the text/ characters very easily from an image, machines are not intelligent enough to perceive the information available in image. Therefore, many research efforts have been put forward that attempts to transform a document image to format understandable for machine.

OCR is a complex problem because of the variety of languages, fonts and styles in which text can be written, and the complex rules of languages etc. Hence, techniques from different disciplines of computer science (i.e. image processing, pattern classification and natural language processing etc. are employed to address different challenges. This paper introduces the reader

to the problem. It enlightens the reader with the historical perspectives, applications, challenges, and techniques of OCR.

TYPES OF OPTICAL CHARACTER RECOGNITION SYSTEMS

There has been multitude of directions in which research on OCR has been carried out during past years. This section discusses different types of OCR systems that have emerged because of this research. We can categorize these systems based on image acquisition mode, character connectivity, font-restrictions etc. Fig. 1 categorizes the character recognition system.

Based on the type of input, the OCR systems can be categorized as handwriting recognition and machine printed character recognition. The former is a relatively simpler problem because characters are usually of uniform dimensions, and the positions of characters on the page can be predicted [3].

Handwriting character recognition is a very tough job due to the different writing style of the user as well as different pen movements by the user for the same character. These systems can be divided into two sub-categories i.e. on-line and off-line systems. The former is performed in real-time while the users are writing the character. They are less complex as they can capture the temporal or time-based information i.e. speed, velocity, number of strokes made, direction of writing of strokes etc. In addition, there no need for thinning techniques as the trace of the pen is few pixels wide. The offline recognition systems operate on static data i.e. the input is a bitmap. Hence, it is very difficult to perform recognition.

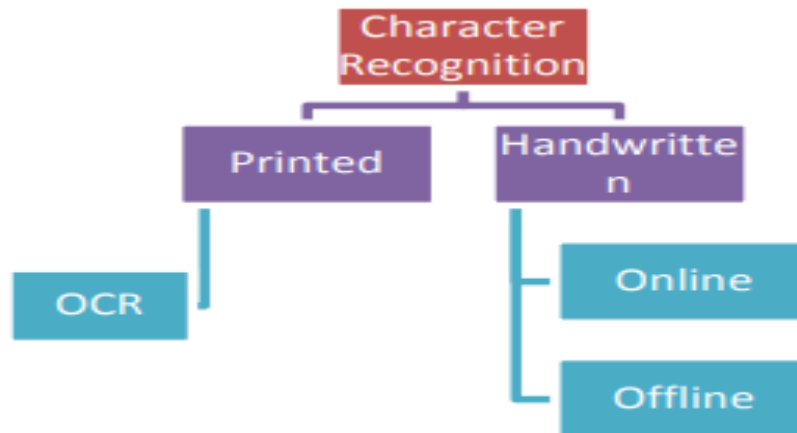


Figure.1: Types of character recognition system

۲)

Three classifications of optical reader machines are normally used:

- Mark Readers
- Bar Code Readers
- Character Readers

In some usage the term OCR is applied to systems of all three types. However, OCR is more correctly used for character readers only, though in some instances, systems possess all three capabilities.

How OCR Works:

۱. Image Acquisition:

- OCR starts with the acquisition of an image, which can be a scanned document, a photo of a document, or a PDF file.

۲. Preprocessing:

- The image may undergo preprocessing steps, such as noise reduction, binarization (converting to black and white).

۳. Text Detection:

- OCR algorithms locate areas in the image where text is present. This involves identifying lines, paragraphs, and individual characters.

ξ. Character Recognition:

- OCR recognizes individual characters within the detected text regions. This process involves pattern recognition and machine learning algorithms.

ο. Text Post-processing:

- The recognized characters are often post-processed to correct errors and improve accuracy. This may involve dictionary-based corrections, language modeling, or contextual analysis.

Ϛ. Output:

- The final output is usually in a text format (e.g., TXT, DOC, PDF with selectable text), allowing users to edit, search, and manipulate the content.

Resources

ϛ. First s:

- Journal of Information & Communication Technology-JICT Vol. 10 Issue. 2, December 2016 (A Survey on Optical Character Recognition System) Noman Islam, Zeeshan Islam, Nazia Noor. <https://arxiv.org/abs/1710.05703>

Ϝ. Second s:

- An Overview of Optical Character Recognition (OCR) Theology and Techniques. <https://apps.dtic.mil/sti/citations/ADA131341>

ϝ. Third s:

- Google scholar. https://scholar.google.com/scholar?hl=ar&as_sdt=0%2C0&q=Over+view+of+OCR&btnG=

ξ. Fourth s:

- Using Random Forests for Handwritten Digit Recognition. https://www.researchgate.net/publication/4288221_Using_Random_Forests_for_Handwritten_Digit_Recognition

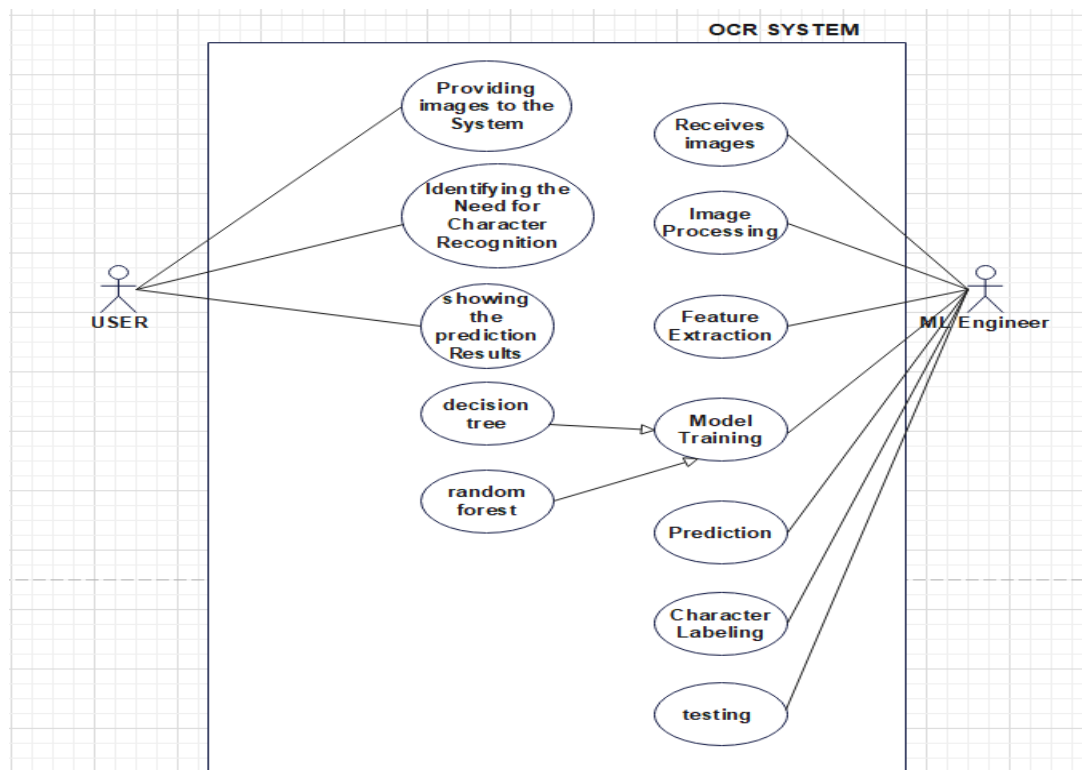
ο. Fifth s:

- Preprocessing The image resource: [Neural network EMNIST | Kaggle](#)

٢- Proposed solution & dataset:

- Main Functionalities

- ١- load image from dataset.
- ٢- Image Processing.
- ٣- Training Model.
- ٤- testing Model.
- ٥- Show the results and the accuracy of testing model.
- ٦- predict an image from dataset.



- Dataset:

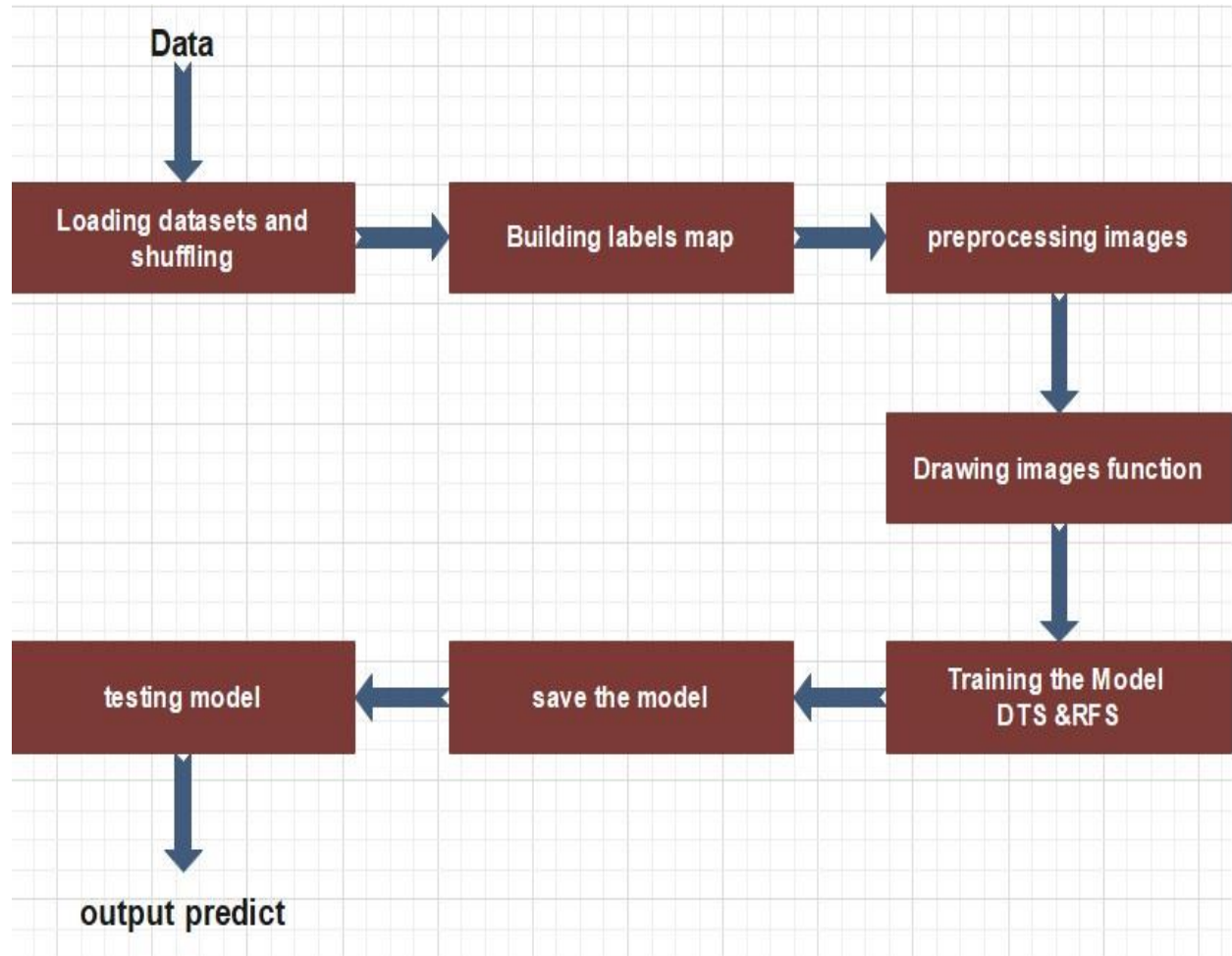
In Kaggle: <https://www.kaggle.com/crawford/emnist>

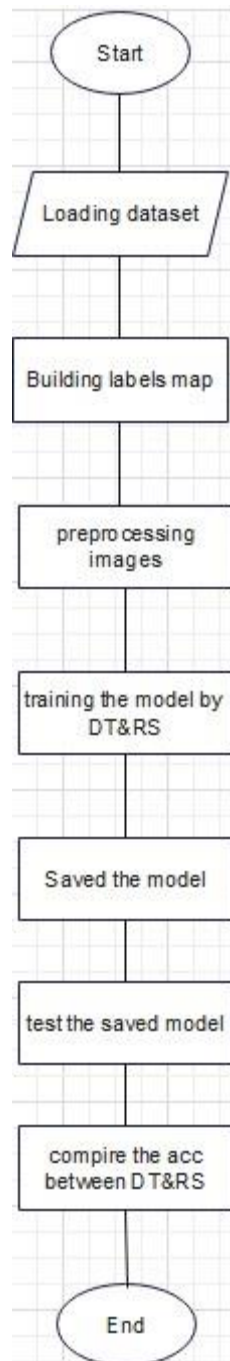
- Github Repo:

[SeifMohamed00/AI ML DecisionTree RF: Decison Tree and Random Forest models done on MINST English letters Dataset \(github.com\)](#)

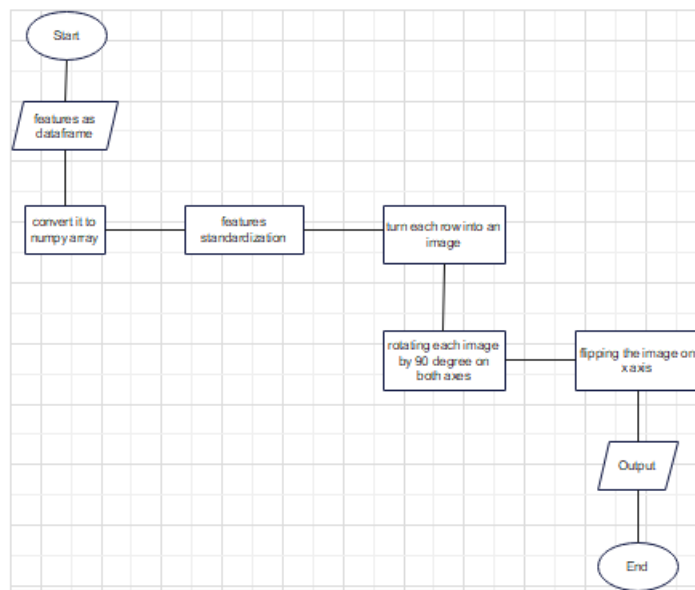
३- Applied Algorithms:

-Block diagram





Flowchart Of Our Project



4 - Experiments & Results:

Experiments: We tried the decision tree algorithm and random forests algorithm. The decision tree algorithm was 70,71% accuracy in the model, and the random forests algorithm was 82,2% accuracy.

Time: decision tree model took about 4 hours of training because of the grid_search done on the parameters then we took those parameters and used them at Random forest. Random Forest with n_estimators = 100 took about an hour and a half with cross validation folds n_splits = 5 and with n_estimators = 200 took about two hours and so little difference in cross validated accuracy :

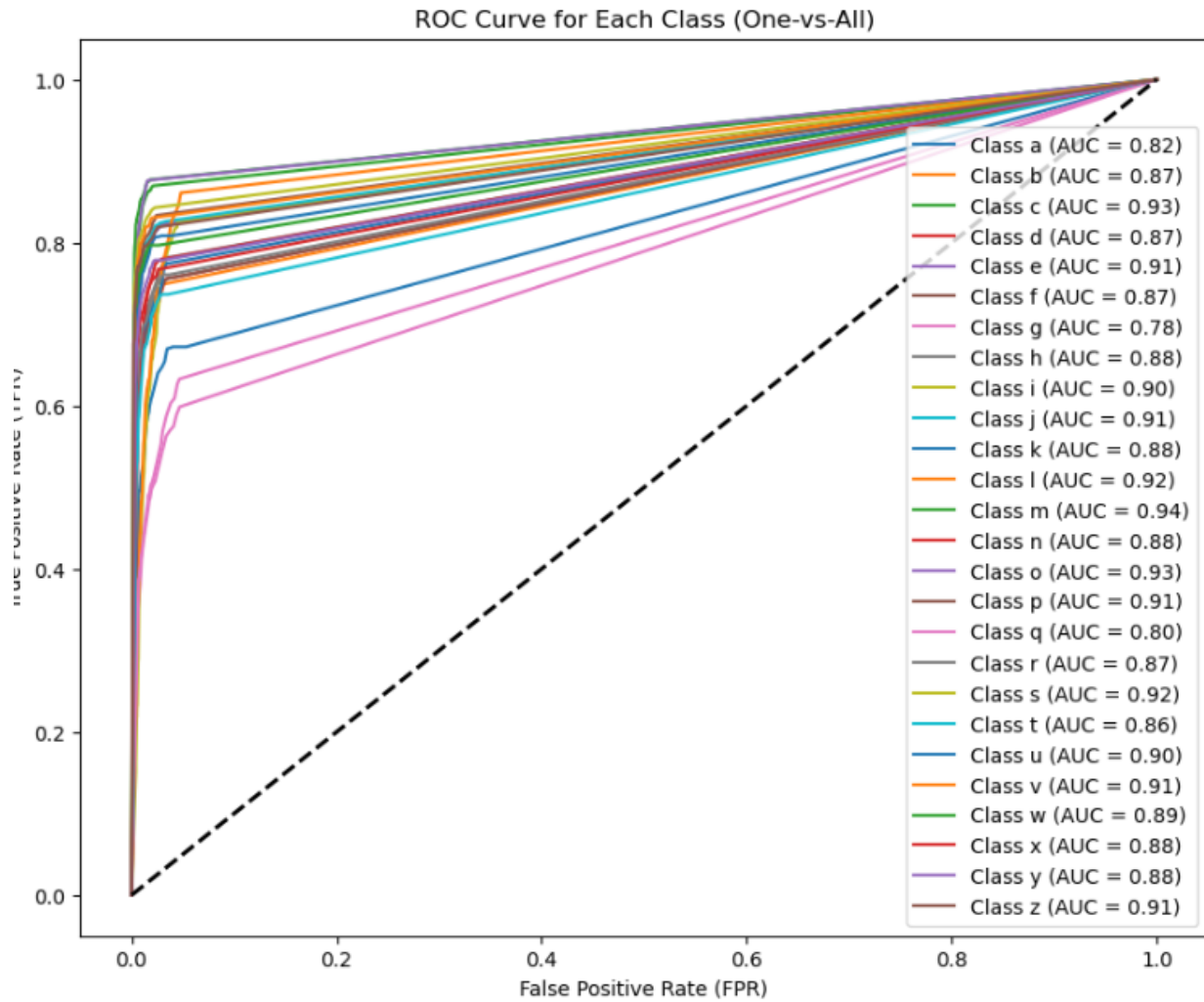
RF_100 accuracy = 81,8 % , RF_200 accuracy = 82,2 %

Params used to search for the best DecisionTree hyperparametr

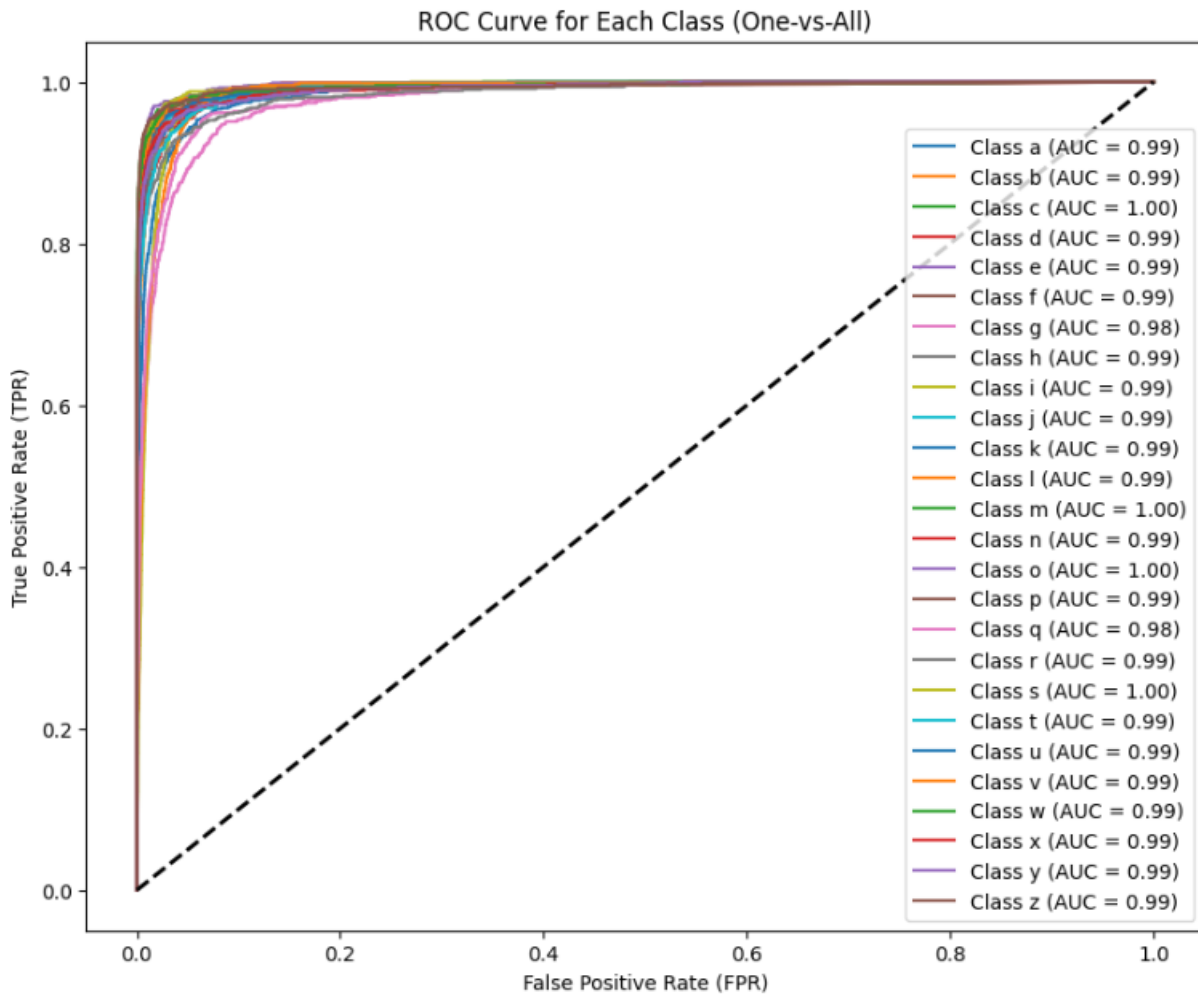
```

# param_grid = {
#   'criterion': ['gini', 'entropy'],
#   'splitter': ['best', 'random'],
#   'max_depth': [None, 10, 20, 30],
#   'min_samples_split': [2, 5, 10],
#   'min_samples_leaf': [1, 2, 5],
#   'max_features': [None, 'sqrt', 'log 2']
# }
# kf = KFold(n_splits=5, shuffle=True, random_state=12)
# classifier = DecisionTreeClassifier(random_state=12)
# grid_search = GridSearchCV(classifier, param_grid, cv=kf, scoring='accuracy')
  
```

Decision Tree Merged Dataset ROC Curve:



Random Forest Merged Dataset ROC Curve:



• Analysis, Discussion and Future work:

- By analyzing the results of both algorithms, we got insights about the efficiency of both algorithms then we found that performance of results of Random Forest algorithm is better than Decision tree.
- Random Forest gives overall better performance but needs more computational power. It consists of a collection of decision trees to improve its performance
- Future work: Adding functionality to generalize on any handwritten English characters

٦- Development Platform:

Tools: Jupyter Notebook, VScode

Programming Languages: Python.

Python Libraries:

- Pandas => Dealing with csv files.
- NumPy => Dealing with arrays.
- sklearn => Used Models.
- joblib => saving models.

- Github Repo:

[SeifMohamed99/AI_ML_DecisionTree_RF: Decison Tree and Random Forest models done on MINST English letters Dataset \(github.com\)](https://github.com/SeifMohamed99/AI_ML_DecisionTree_RF_Decision_Tree_and_Random_Forest_models_done_on_MINST_English_letters_Dataset)