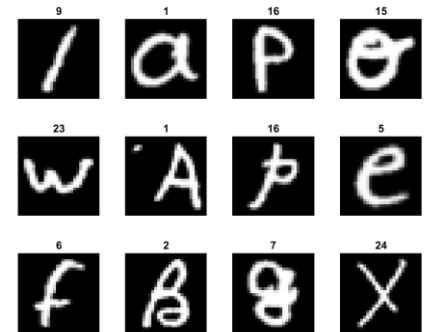# Optical Character Recognition from the handwritten EMNIST character set

## Introduction
Optical Character Recognition is the process of converting handwritten text into machine-encoded text through mediums such as documents, written papers, or images. Developing automated methods for this using machine learning and deep learning of large datasets is important for producing AI models that can appropriately recognise and convert this text. In this work, we aim to compare varying classification models for their efficiency and accuracy in creating an AI model for Optical Character Recognition. This has many uses such as data entry in transcribing written records into a digital form.

## Data and Preparation
The EMNIST dataset is a collection of 26,000 images. The format of this data is 26,000 28x28 pixel images with each pixel represented by a numerical value. Each image is also provided with a corresponding label as to the correct digital character that the handwritten image refers to. Data is provided as a 1x784 vector and variable labels stored as a column vector. We convert the images into a double type so that it can be scaled accurately unlike an integer. We also convert the labels into a categorical array for further processing. The data is split into training and testing subsets to avoid overfitting. A random permutation of the 26,000 images is created and used to make a 50/50 split of training and testing subsets.
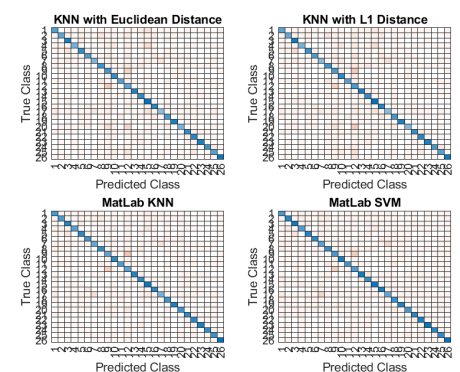


## Methodology
In this study we compare various classification models and distance measures to classify the test data. We implement the K-nearest neighbour classification model. This model measures the distance from our test point to the nearest points. It then looks at the k-nearest of these and sets the value of the test point to the most common labelled point nearby. We use the k-value of 28 as a compromise for efficiency and accuracy. In this method we compare the impact of using Euclidean distance against using L1 or Manhattan distance. We then also compare two models from MatLab using pre-defined functions for K-nearest neighbour and SVM for Multiclass. SVM for Multiclass works by splitting data into linearly separated groups.

## Results

| Classification Model | Time of Processing (s) | Correct Predicts(/13000) | Accuracy (%) |
|---|---|---|---|
| KNN Custom Euclidean | 424 | 9496 | 73.05 |
| KNN Custom Manhattan | 423 | 9202 | 70.78 |
| KNN MatLab | 21.5 | 10096 | 77.66 |
| SVM MatLab | 41.3 | 9518 | 73.22 |

The Custom KNN model is far slower to process than the built-in functions which is to be expected as pre-built functions are optimised in machine code already without the need for compilation or interpretation. The distance measures have very similar time of processing. Manhattan is quicker as it doesn't need calculating of square numbers however this is a large dataset, and this advantage doesn't apply. The MatLab KNN is ~20 times quicker and has a greater accuracy making it more efficient and better. SVM has a slower processing time as it classifies the entire data which takes a lot of time on the dataset. This is also not rewarded as the accuracy is less. SVM is often used for text-processing and is efficient on small datasets, so this observation is expected on a large dataset.



## Conclusion
In conclusion, our study shows that the K-nearest neighbour classification model is both the most efficient and accurate algorithm for correctly classifying these handwritten images. KNN is very good for implementing on multiclass models however it can be slow on large datasets and struggle with imbalanced data. An accuracy of 77.66% is not good enough for making this OCR model viable in commercial use and so further testing and training of other models is important for increasing these numbers. Further evaluation

could also be done to recognise areas of weakness. For example, looking at the confusion charts we see the models have issues with distinguishing between the letter I and L.
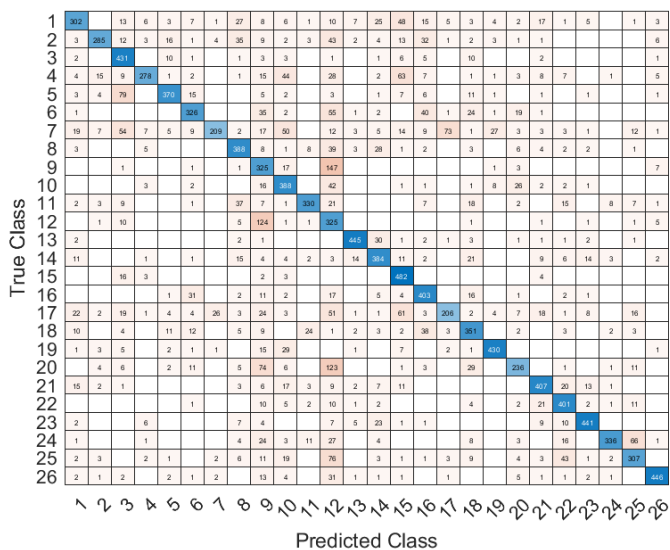
KNN with Euclidean Distance:

Test 1: 435s, 9496/13000, 73.05%    Test 2: 413s, 9472/13000, 72.86%
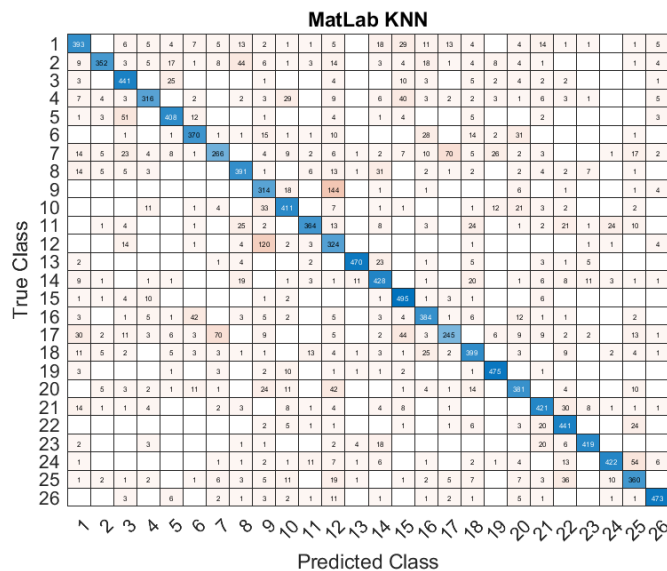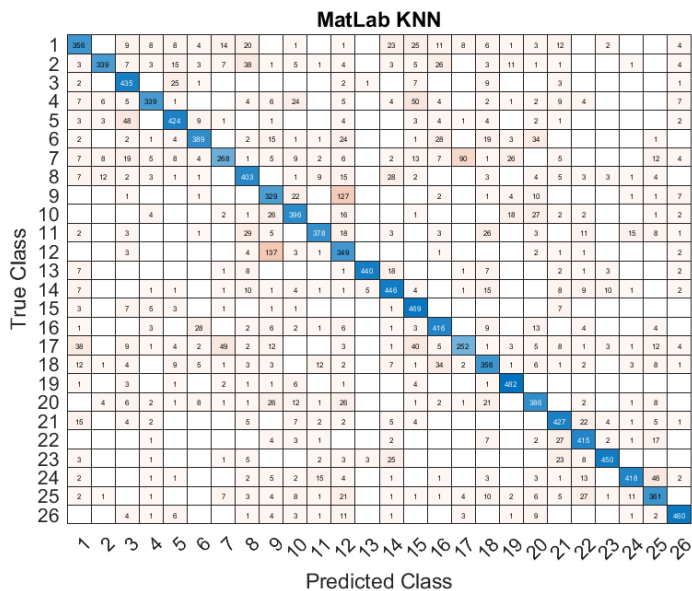


KNN with Manhattan Distance:

Test 1: 420s, 9202/13000, 70.78%     Test 2: 426s, 9232/13000, 71.02%

Test 1: 21.5s, 10096/13000, 77.66%     Test 2: 21.05s, 10163/13000, 78.18%



**MatLab KNN**



**MatLab KNN**

KNN with SVM:

Test 1: 41.3s, 9518/13000, 73.22%     Test 2: 38.9s, 9543/13000, 73.41%

Full Results:

**KNN with Euclidean Distance**

**KNN with L1 Distance**

**MatLab KNN**

**MatLab SVM**