# Machine Learning Assignment 2: Naive Bayes & ANNs

30.04.2022
—

Done By

Maitreya Manohar (2021A7PS2663H)

Ayush Bhauwala (2021A7PS0180H)

Arunachala Amuda Murugan (2021A7PS0205H)

# Part 1: Naive Bayes

## Model 1: Naive Bayes

Notes

1. Conversion of continuous columns to categorical using **pandas.cut**, allowed us to speed up the training process as well as increase the accuracy from 80% to 83.5%.
2. Data preprocessing: imputed categorical columns using its mode and numerical columns with its mean.
3. Smoothing techniques allowed in the model: Laplace smoothing, Ignoring features not present in distribution.
4. We create a Distribution object for every class (here 2, but extendable to n classes as well), where we store its prior probability and the likelihood of each feature in a distribution in a dictionary, whose keys are in the format **(column, value)**.

Accuracies

| split 1 | 0.8335194490973385 |
| split 2 | 0.8289596128792108 |
| split 3 | 0.8324958123953099 |
| split 4 | 0.8322166387493021 |
| split 5 | 0.8334263912153359 |
| split 6 | 0.8351944909733855 |
| split 7 | 0.8273776288851665 |
| split 8 | 0.8331472175693281 |
| split 9 | 0.8351944909733855 |
| split 10 | 0.8335194490973385 |

| Mean Accuracy | 0.8325051181835101 |
| Standard deviation | 0.0025106281502815394 |

| Average precision | 0.8253973136540113 |
|---|---|
| Average recall | 0.8029594125987934 |

## Model 2: KNN

Accuracies

| split 1 | 0.8323096966313047 |
|---|---|
| split 2 | 0.828680439233203 |
| split 3 | 0.8306346547552578 |
| split 4 | 0.8259817606551275 |
| split 5 | 0.824213660897078 |
| split 6 | 0.8310999441652708 |
| split 7 | 0.8325888702773125 |
| split 8 | 0.8327749860413177 |
| split 9 | 0.830820770519263 |
| split 10 | 0.8306346547552578 |

| Mean Accuracy | 0.8299739437930391 |
|---|---|
| Standard deviation | 0.0028633513601374907 |
| Average precision | 0.8606935915399496 |
| Average recall | 0.9257837842067571 |

## Model 3 : Logistic Regression

Accuracies

| split 1 | 0.8528754885538805 |
| split 2 | 0.844965568583659 |
| split 3 | 0.8517587939698492 |
| split 4 | 0.8484087102177554 |
| split 5 | 0.8467336683417085 |
| split 6 | 0.8463614368136981 |
| split 7 | 0.85324772008189 |
| split 8 | 0.8503629257398102 |
| split 9 | 0.8481295365717476 |
| split 10 | 0.8494323469197841 |

| Mean Accuracy | 0.8492276195793783 |
| Standard deviation | 0.0028221762132697053 |
| Average precision | 0.8787840738371813 |
| Average recall | 0.9257837842067571 |

| | index | accuracy | precision | recall | f1 score | knn accuracy | knn recall | knn f1 score | knn precision | logreg accuracy | logreg recall | logreg f1 score | logreg precision |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Accuracy 1 | 0.833519 | 0.849771 | 0.922736 | 0.884752 | 0.832310 | 0.923707 | 0.893448 | 0.865109 | 0.852875 | 0.931165 | 0.905966 | 0.882094 |
| 1 | Accuracy 2 | 0.828960 | 0.844551 | 0.922321 | 0.881725 | 0.828680 | 0.920850 | 0.891187 | 0.863376 | 0.844966 | 0.925369 | 0.900939 | 0.877766 |
| 2 | Accuracy 3 | 0.832496 | 0.845331 | 0.929412 | 0.885380 | 0.830635 | 0.925163 | 0.892320 | 0.861730 | 0.851759 | 0.929088 | 0.904833 | 0.881812 |
| 3 | Accuracy 4 | 0.832217 | 0.796911 | 0.617779 | 0.696004 | 0.825982 | 0.926305 | 0.889205 | 0.854962 | 0.848409 | 0.931366 | 0.902566 | 0.875493 |
| 4 | Accuracy 5 | 0.833426 | 0.800391 | 0.615038 | 0.695578 | 0.824214 | 0.926364 | 0.888456 | 0.853528 | 0.846734 | 0.928703 | 0.901560 | 0.875958 |
| 5 | Accuracy 6 | 0.835194 | 0.849188 | 0.928438 | 0.887046 | 0.831100 | 0.928343 | 0.892737 | 0.859761 | 0.846361 | 0.926991 | 0.901344 | 0.877079 |
| 6 | Accuracy 7 | 0.827378 | 0.843742 | 0.922480 | 0.881356 | 0.832589 | 0.922429 | 0.893557 | 0.866437 | 0.853248 | 0.927681 | 0.905935 | 0.885185 |
| 7 | Accuracy 8 | 0.833147 | 0.783888 | 0.621049 | 0.693032 | 0.832775 | 0.927064 | 0.894288 | 0.863750 | 0.850363 | 0.931089 | 0.904717 | 0.879797 |
| 8 | Accuracy 9 | 0.835194 | 0.794942 | 0.621539 | 0.697627 | 0.830821 | 0.930215 | 0.892642 | 0.857986 | 0.848130 | 0.933785 | 0.902892 | 0.873978 |
| 9 | Accuracy 10 | 0.833519 | 0.845257 | 0.928803 | 0.885063 | 0.830635 | 0.927398 | 0.892587 | 0.860296 | 0.849432 | 0.929973 | 0.903599 | 0.878679 |

# Part 2: ANNs

## Dataset: MNIST

Used the inbuilt MNIST (80-20 split) train and test datasets provided in **keras.datasets**

## Models Used

Using the various combinations from the given parameter grid, we have trained, tested and documented the results of 15 models on 10 splits each.

The output layer is a 10 neuron layer, with a softmax activation function.

## Hyperparameters

| Number of hidden layers | Number of neurons | Activation functions |
|---|---|---|
| 2 | 100 | relu |
| 3 | 150 | tanh |
| | | sigmoid |

## Model details with accuracies

| Model no. | No. of Inner Layers | Neurons | Activation function | Accuracy |
|---|---|---|---|---|
| 1 | 2 | 150, 150 | sigmoid, sigmoid | 95.77% |
| 2 | 2 | 100, 100 | sigmoid, sigmoid | 95.21% |
| 3 | 2 | 100, 150 | sigmoid, sigmoid | 95.52% |
| 4 | 2 | 150, 100 | sigmoid, sigmoid | 95.68% |
| 5 | 2 | 150, 150 | relu, relu | 97.48% |
| 6 | 2 | 100, 100 | relu, relu | 97.85% |
| 7 | 2 | 100, 150 | relu, relu | 97.99% |
| 8 | 2 | 150, 100 | relu, relu | 97.96% |
| 9 | 2 | 150, 150 | tanh, tanh | 94.06% |
| 10 | 2 | 100, 100 | tanh, tanh | 94.16% |
| 11 | 2 | 100, 150 | tanh, tanh | 94.48% |
| 12 | 2 | 150, 100 | tanh, tanh | 94.56% |
| 13 | 3 | 150, 150, 150 | tanh, tanh, tanh | 94.11% |
| 14 | 3 | 150, 150, 150 | relu, relu, relu | 98.25% |
| 15 | 3 | 150, 150, 150 | sigmoid, sigmoid, sigmoid | 95.92% |

## Analysis

On analysing the accuracies and confusion matrices of all 15 models, we find that the models with the best performance are the ones that use the ReLU activation function.

The best performing model is model 14 which has 3 inner layers and uses the ReLU activation function.

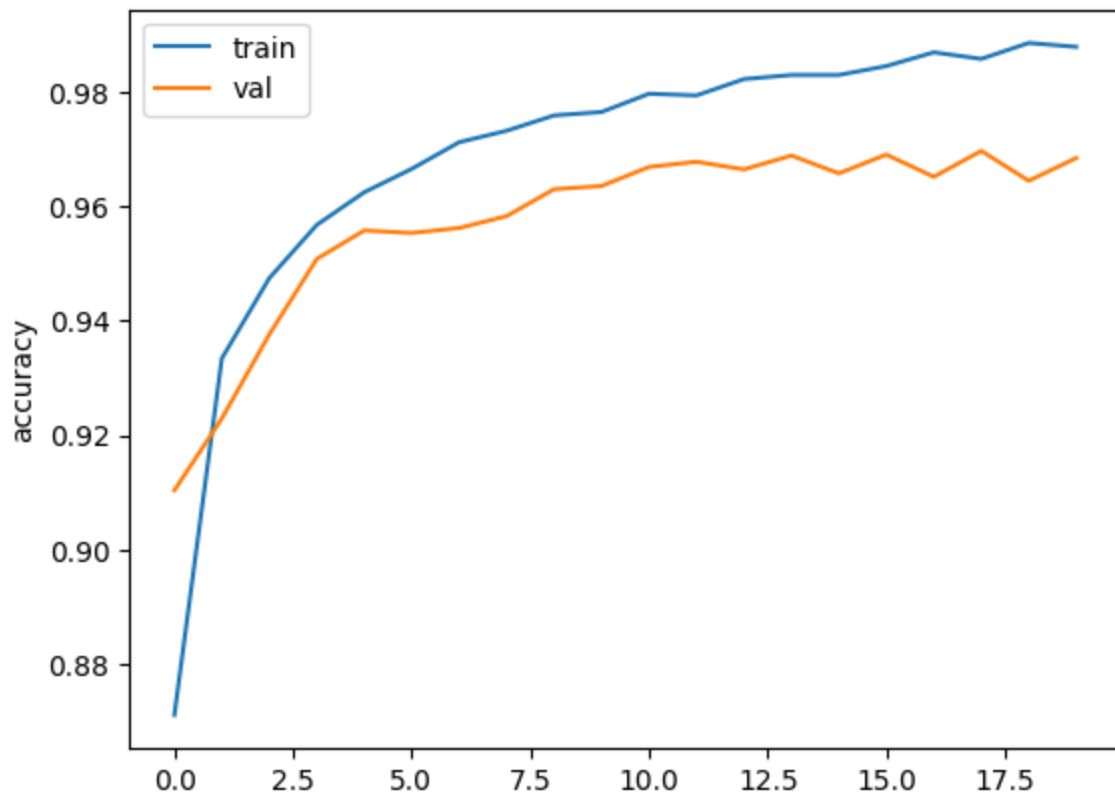We find that after ReLU activation function, sigmoid function performs the best followed by tanh function.

All models though had different performances, were statistically significant, as the losses were reducing with increase in epochs. (as seen below)
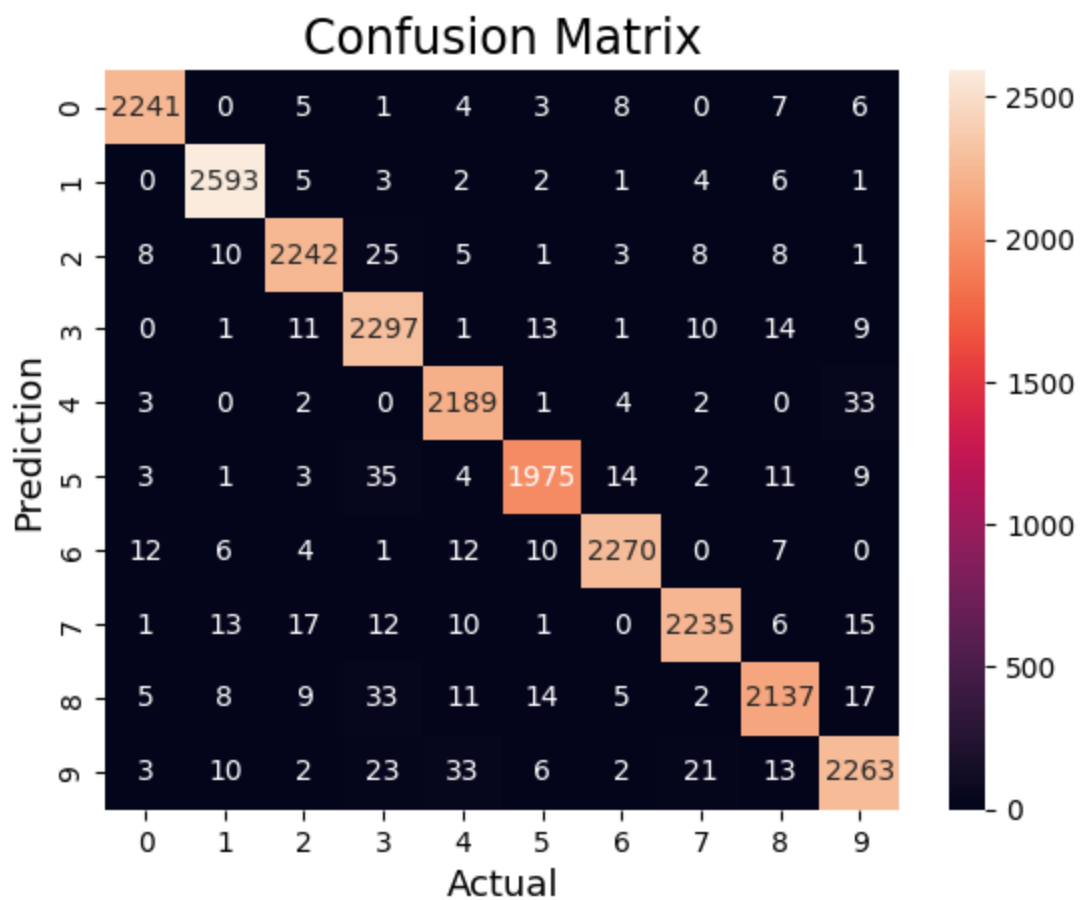
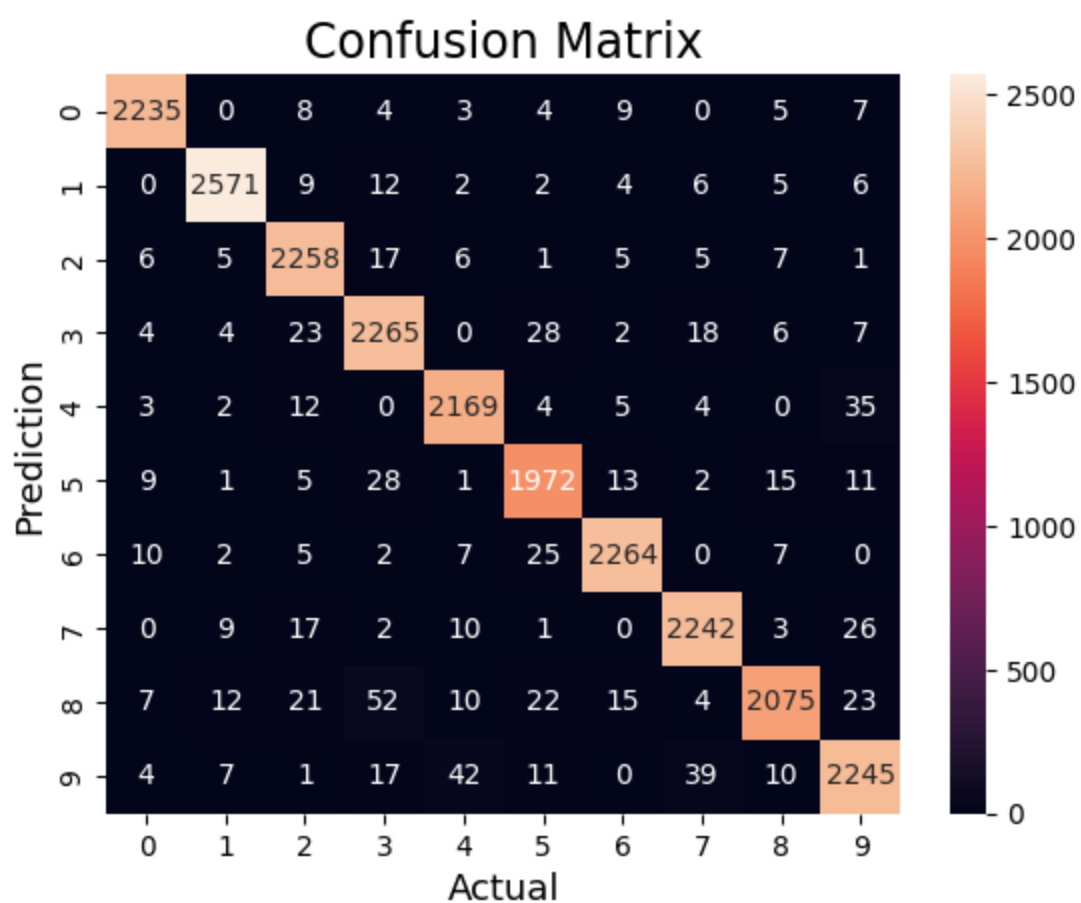# Train and Test plots

## Loss curves
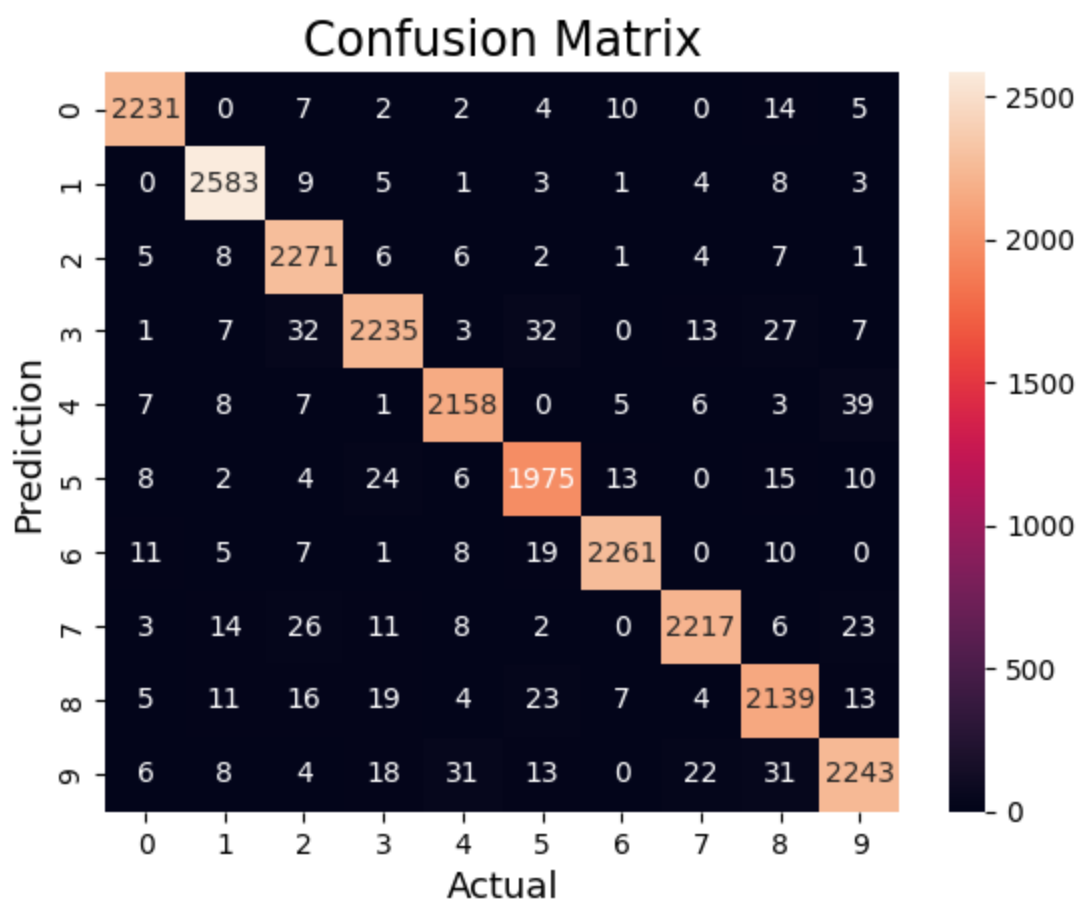
Accuracy curves

# Confusion Matrices
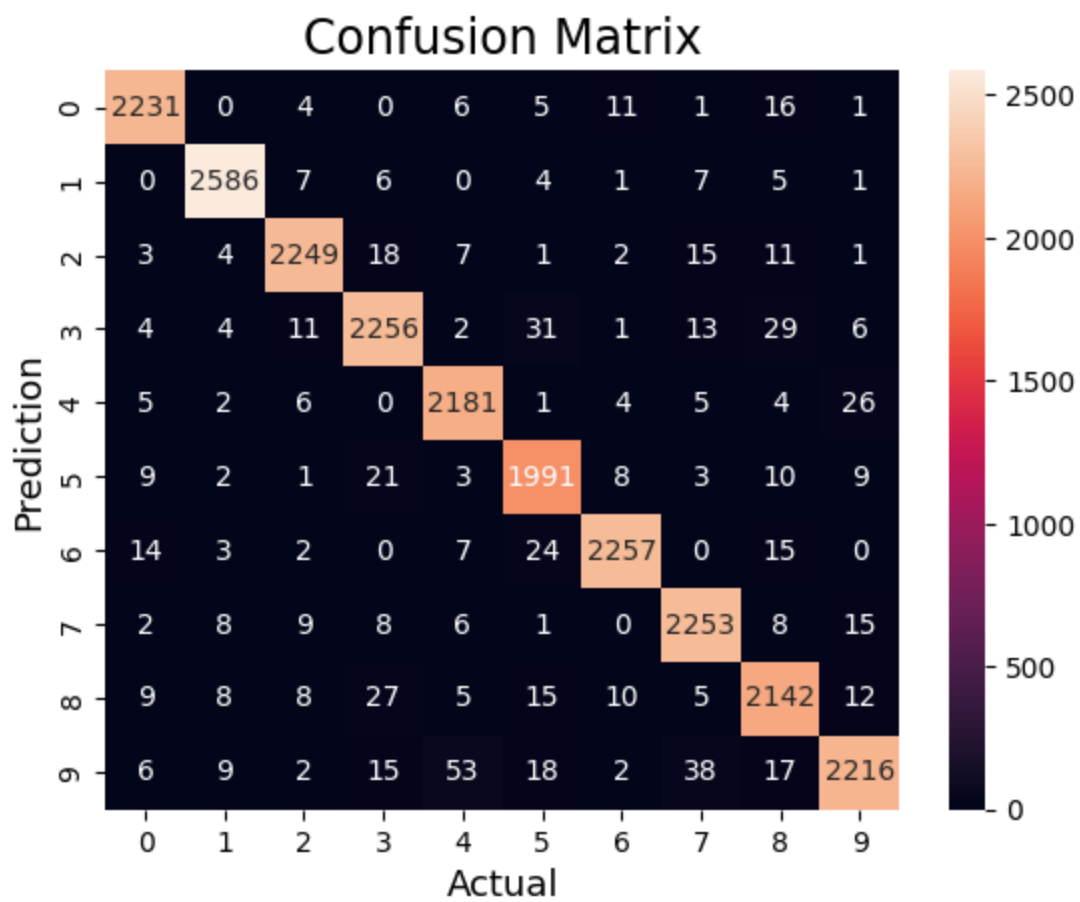
## Model 1



Confusion Matrix

Model 2



Confusion Matrix

Model 3



Confusion Matrix

Model 4



Confusion Matrix

Model 5



Confusion Matrix

Model 6

## Confusion Matrix

| Prediction \ Actual | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2181 | 0 | 1 | 1 | 1 | 2 | 3 | 0 | 86 | 0 |
| 1 | 0 | 2595 | 2 | 0 | 0 | 3 | 0 | 3 | 13 | 1 |
| 2 | 0 | 0 | 2276 | 2 | 1 | 0 | 0 | 1 | 31 | 0 |
| 3 | 0 | 0 | 2 | 2282 | 0 | 8 | 0 | 1 | 61 | 3 |
| 4 | 0 | 0 | 1 | 0 | 2215 | 0 | 0 | 1 | 12 | 5 |
| 5 | 0 | 0 | 0 | 2 | 0 | 2004 | 1 | 0 | 49 | 1 |
| 6 | 0 | 0 | 3 | 0 | 4 | 30 | 2231 | 0 | 54 | 0 |
| 7 | 1 | 0 | 5 | 1 | 0 | 0 | 0 | 2280 | 17 | 6 |
| 8 | 0 | 2 | 1 | 5 | 0 | 3 | 0 | 1 | 2228 | 1 |
| 9 | 1 | 1 | 0 | 1 | 13 | 1 | 0 | 2 | 72 | 2285 |

Model 7

## Confusion Matrix

| Prediction \ Actual | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2241 | 0 | 2 | 0 | 0 | 1 | 1 | 0 | 29 | 1 |
| 1 | 0 | 2607 | 0 | 0 | 1 | 0 | 0 | 2 | 6 | 1 |
| 2 | 0 | 0 | 2276 | 3 | 1 | 1 | 0 | 2 | 28 | 0 |
| 3 | 0 | 2 | 5 | 2294 | 0 | 8 | 0 | 1 | 40 | 7 |
| 4 | 0 | 0 | 1 | 0 | 2197 | 0 | 0 | 1 | 15 | 20 |
| 5 | 0 | 0 | 0 | 0 | 0 | 2032 | 1 | 0 | 21 | 3 |
| 6 | 3 | 2 | 0 | 0 | 1 | 13 | 2274 | 0 | 29 | 0 |
| 7 | 0 | 0 | 5 | 2 | 0 | 1 | 0 | 2273 | 7 | 22 |
| 8 | 0 | 2 | 1 | 2 | 0 | 1 | 0 | 0 | 2228 | 7 |
| 9 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 4 | 17 | 2352 |

Model 8

## Confusion Matrix

Model 9



Confusion Matrix

Model 10

## Confusion Matrix

Model 11


Confusion Matrix

Model 12



Confusion Matrix

Model 13


Confusion Matrix

Model 14

## Confusion Matrix

Model 15