# COL 774: Machine Learning
# Assignment 3

**BRAJ RAJ NAGAR**

**2022AIB2682**

aib222682@scai.iitd.ac.in

1.  **Decision Trees, Random Forests and Gradient Boosted Trees:**

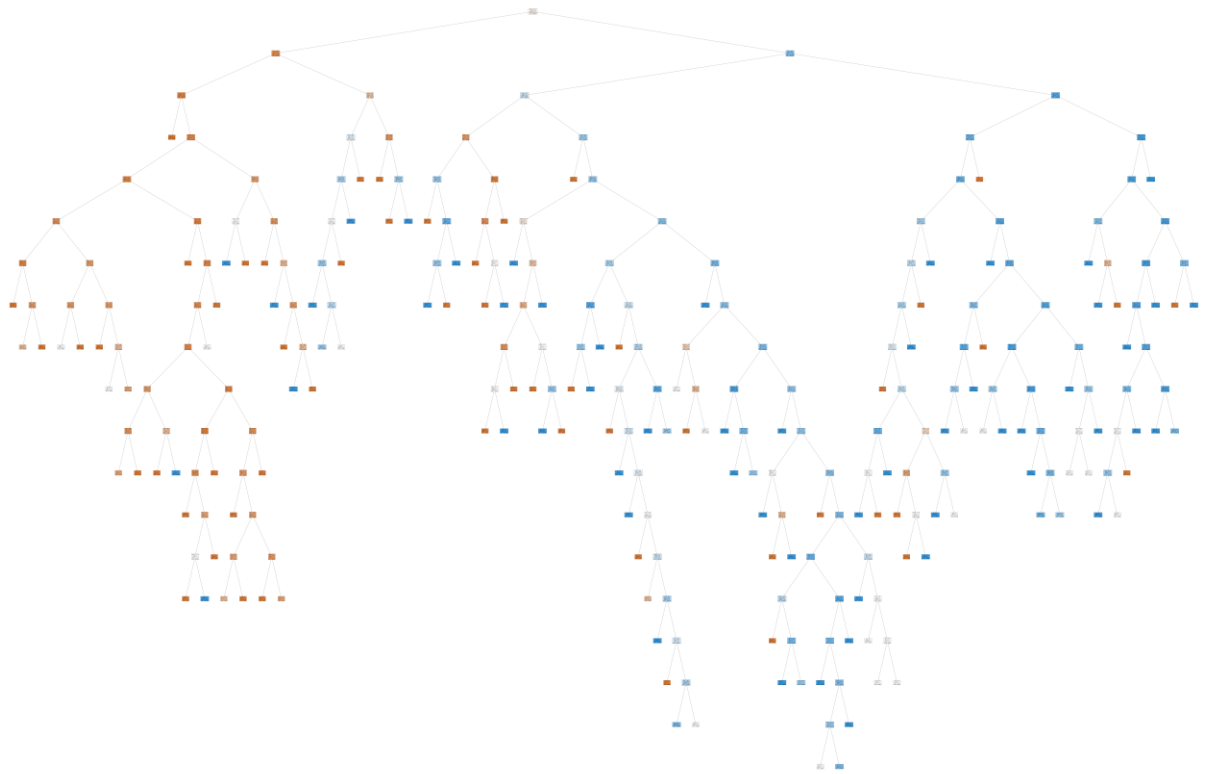    **Dataset 1:**

    a.  **Decision Tree Construction and Visualization:**
        Missing values is ignored in examples.
        Training accuracy: 92.52747252747253%
        Test accuracy: 68.77470355731226 %
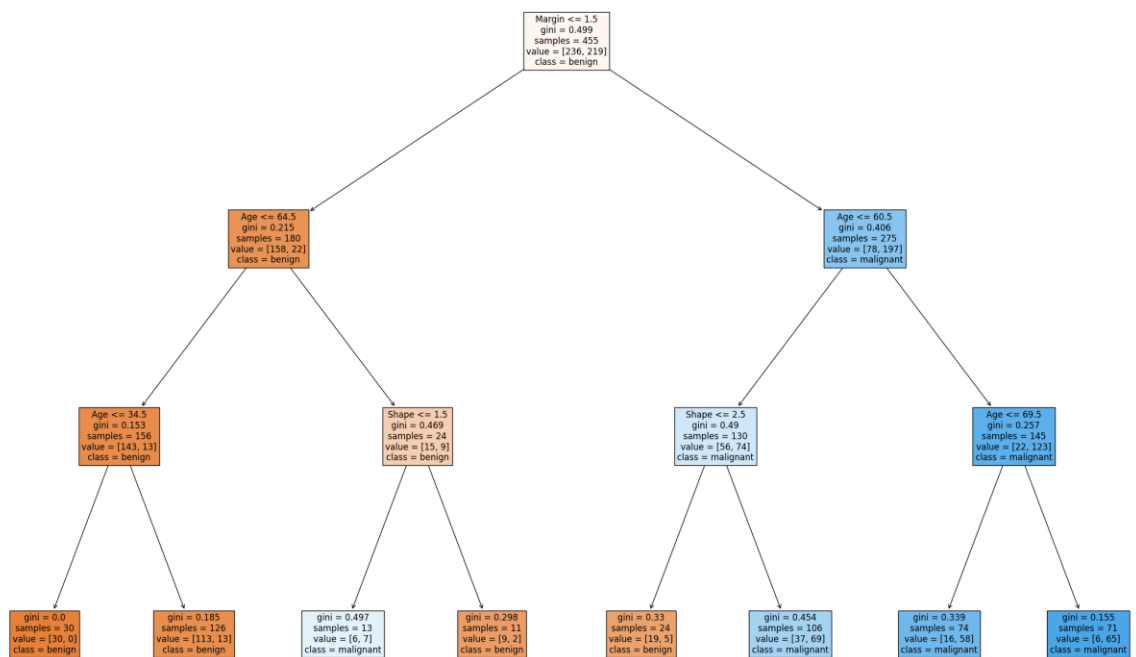        Validation accuracy: 76.03305785123967%



    b.  **Decision Tree Grid Search**
        Optimal parameter found in grid search: {'max_depth': 3, 'min_samples_leaf': 3, 'min_samples_split': 4}
        Training accuracy:  81.31868131868132%
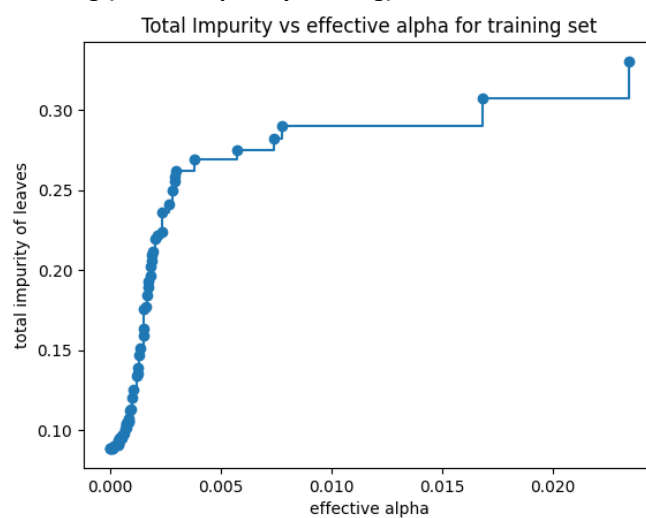        Test accuracy:  75.49407114624506%
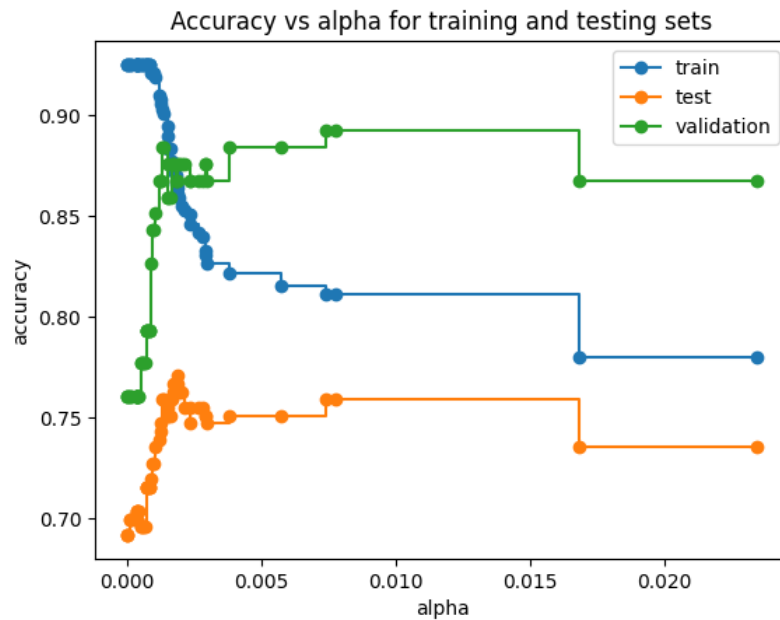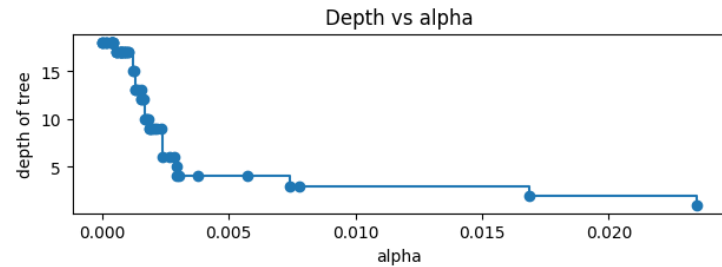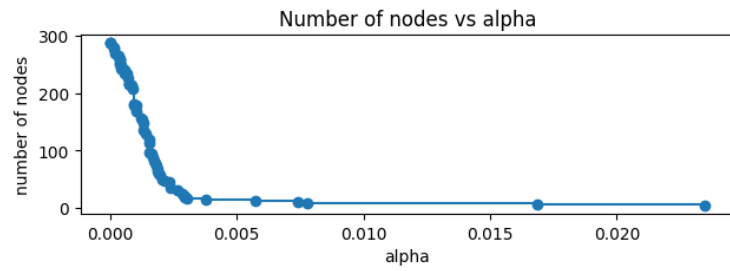        Validation accuracy: 87.60330578512396%

Observation:
With optimal parameters test and validation accuracy improved. But train data accuracy is decreased which indicate that in previous part model was overfitting.

c. **Decision Tree Post Pruning (Cost Complexity Pruning):**

Number of nodes vs alpha


Depth vs alpha


Accuracy vs alpha for training and testing sets

Observations:
Effective alpha increases total impurities of leaves increases.
Effective alpha increases no. of nodes decreases.
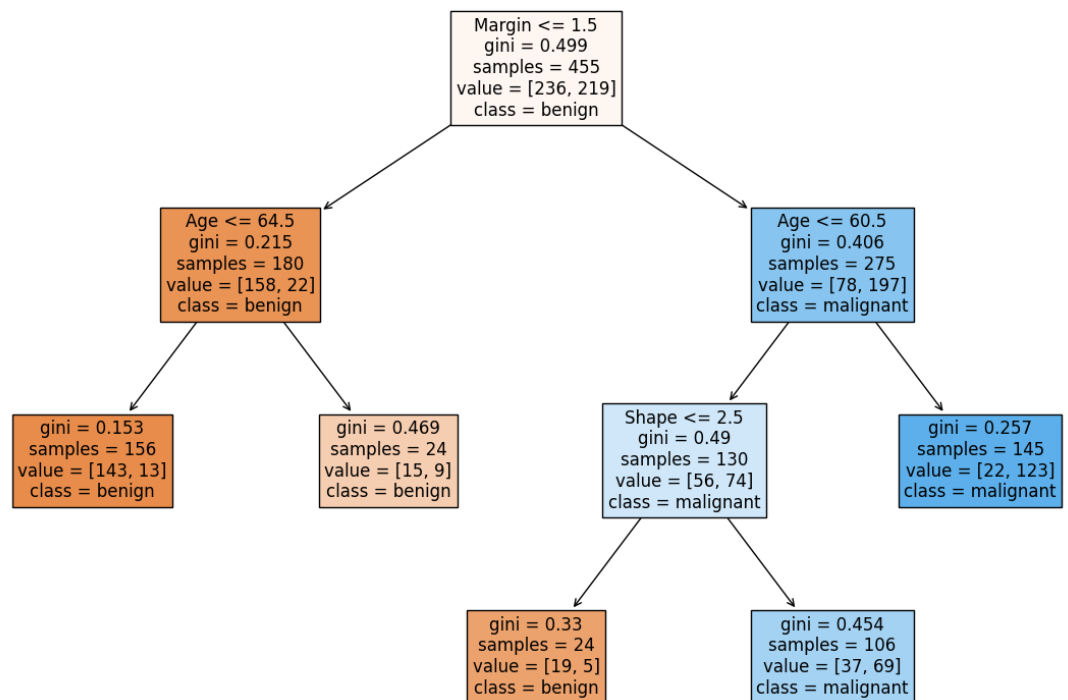Effective alpha increases depth of the tree decreases.
Accuracy is maximum only at optimal value of effective alpha.

Optimal value of alpha: 0.00738705738705739

Train accuracy with respect to the best tree: 81.0989010989011%
Test accuracy with respect to the best tree: 75.8893280632411%
Validation accuracy with respect to the best tree: 89.25619834710744%

Margin <= 1.5
gini = 0.499
samples = 455
value = [236, 219]
class = benign

Age <= 64.5
gini = 0.215
samples = 180
value = [158, 22]
class = benign

Age <= 60.5
gini = 0.406
samples = 275
value = [78, 197]
class = malignant

gini = 0.153
samples = 156
value = [143, 13]
class = benign

gini = 0.469
samples = 24
value = [15, 9]
class = benign

Shape <= 2.5
gini = 0.49
samples = 130
value = [56, 74]
class = malignant

gini = 0.257
samples = 145
value = [22, 123]
class = malignant

gini = 0.33
samples = 24
value = [19, 5]
class = benign

gini = 0.454
samples = 106
value = [37, 69]
class = malignant

Comparison with part (a) and part (b):
Tree in this has lesser depth compare to part (a) and (b) due to cost complexity pruning. And accuracy is also improved because cost complexity pruning avoid overfitting of data.

d. **Random Forests:**
Out-of-bag accuracy (as explained in the class) to tune to the optimal values for these parameters.

Optimal set of parameters: {'max_features': 1, 'min_samples_split': 7, 'n_estimators': 51, 'oob_score': True}

Train accuracy:  87.25274725274725%
Out-of-bag accuracy:  78.46153846153847%
Test accuracy:  76.6798418972332%
Validation accuracy:  84.29752066115702%

e. **Missing Data Imputation:**
Simple imputer from sklearn library used to replace the missing value with median of that data group.
Results of Part (a) to (d) after imputation:
Part (a)
Train accuracy:  91.80633147113594%
Test accuracy:  72.91666666666666%
Validation accuracy:  73.33333333333333%

Part (b)
Train accuracy: 81.56424581005587%
Test accuracy: 80.90277777777778%
Validation accuracy: 87.40740740740741%

Part (c)
Train accuracy with respect to the best tree: 8212290502793296%
Test accuracy with respect to the best tree: 79.86111111111112%
Validation accuracy with respect to the best tree: 86.66666666666667%

Part (D)
Train accuracy: 82.30912476722533%
Out-of-bag accuracy: 79.70231914157149%
Test accuracy: 80.55555555555556%
Validation accuracy: 85.92592592592593%

Observation:
Test set accuracy is improved in every part after imputation. Simply because now we have trained our model with more information by imputing. Although in some cases validation accuracy is less compare to other parts maybe because imputation function (median) may not be the best choice for validation data set.

f. **Gradient Boosted Trees:**
XGBoost (Extreme Gradient Boosting) is used with the goal to construct a model with less bias, and better predictive performance.

Optimal set of parameters: {'max_depth': 4, 'n_estimators': 20, 'subsample': 0.1}

Train accuracy: 79.51582867783985%
Test accuracy: 75.69444444444444%
Validation accuracy: 85.92592592592593%

Observations:
XGBoost provides better accuracy compare part (a) and even missing data imputation part (a) because XGBoost uses the approach of weak learner which gives less bias and more variance which is very helpful to avoid overfitting.

## Dataset 2:

a. **Decision Tree Construction:**
Train accuracy: 99.94863074361428%
Test accuracy: 57.49730312837109%
Validation accuracy: 57.75899481287069%
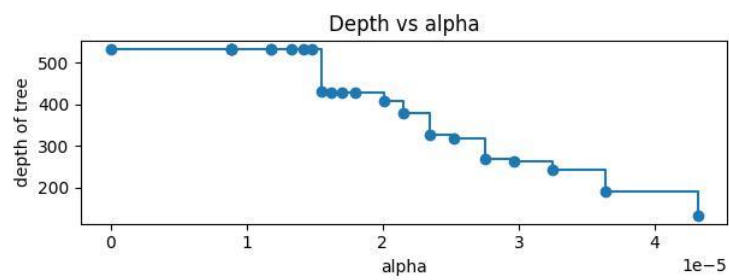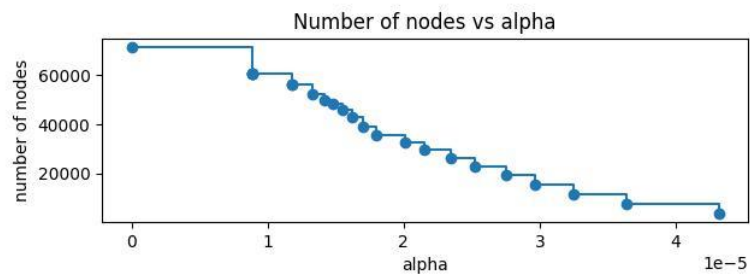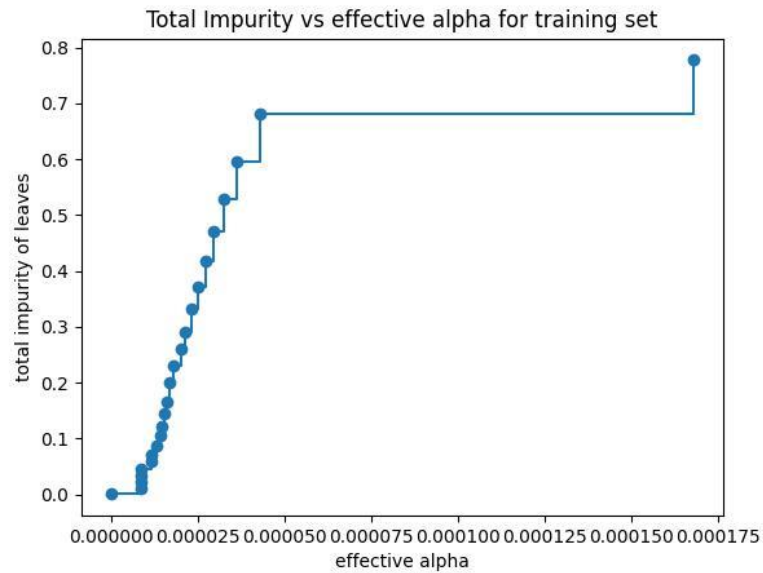
b. **Decision Tree Grid Search:**

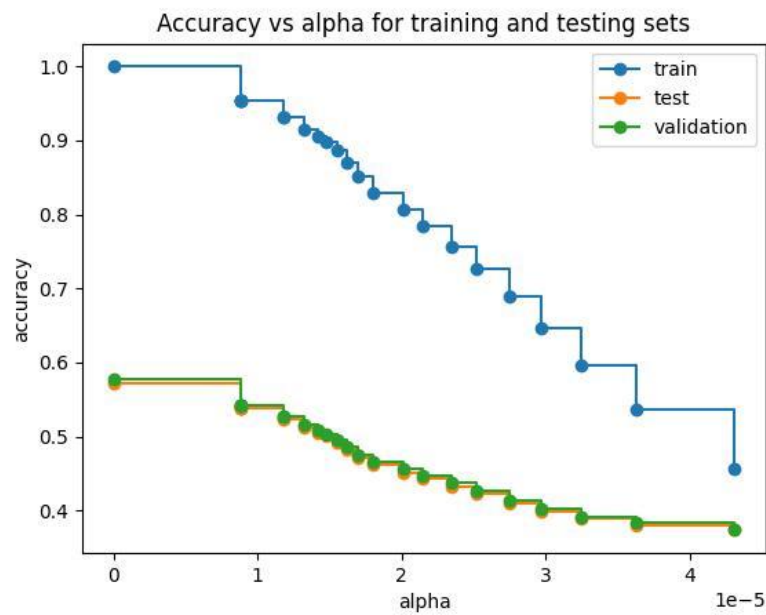Optimal parameter found in grid search: {'max_depth': 50, 'min_samples_leaf': 1, 'min_samples_split': 2}

Train accuracy: 74.84943493817976%
Test accuracy: 48.0098947290109%
Validation accuracy: 48.221703279671%

c. **Decision Tree Post Pruning:**

Accuracy vs alpha for training and testing sets

Optimal value of alpha: 0

Train accuracy with respect to the best tree: 99.94863074361428%
Test accuracy with respect to the best tree: 57.25923446043968%
Validation accuracy with respect to the best tree: 57.74866188596582%

d. **Random Forests:**
   max_features=0.4, min_samples_split=10, n_estimators=100
   Train accuracy:  0.9964395791263684
   Test accuracy:  0.6385448052672693
   Validation accuracy:  0.6413441071317861

e. **XGBoost:**
   At default parameters
   Train accuracy:  83.25096538774932 %
   Test accuracy:  55.241230517427375 %
   Validation accuracy:  55.252226745748004 %

f. **LightGBM:**
   At default parameters
   Train accuracy:  53.55333545895774%
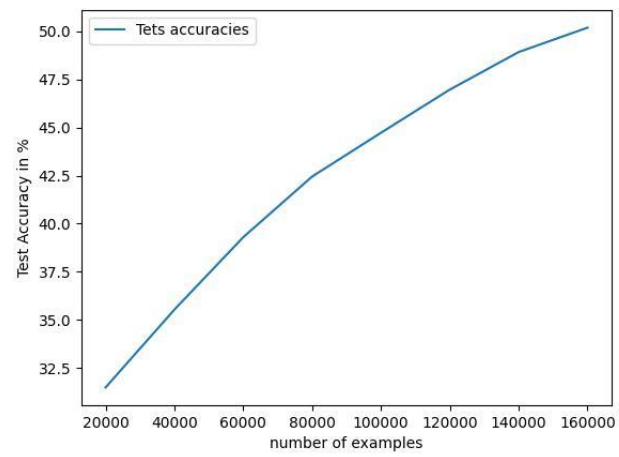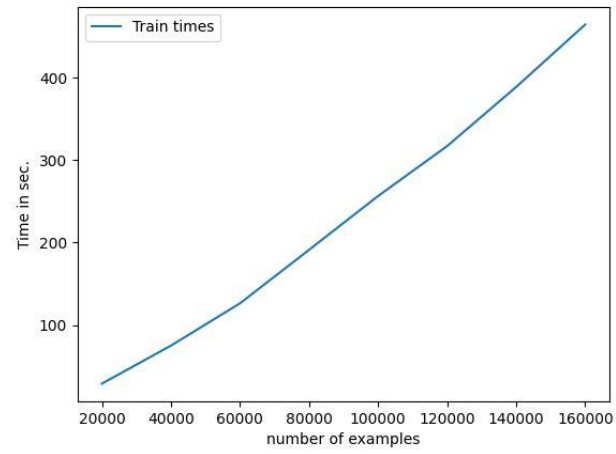   Test accuracy:  43.98132648885913%
   Validation accuracy:  44.206327884436547%
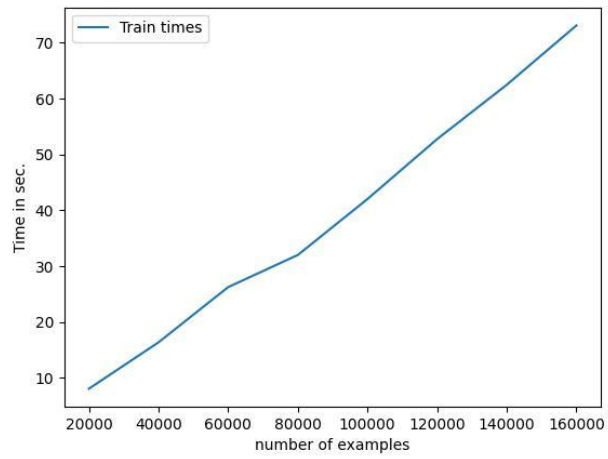
g. **Training with Varying amount of data:**
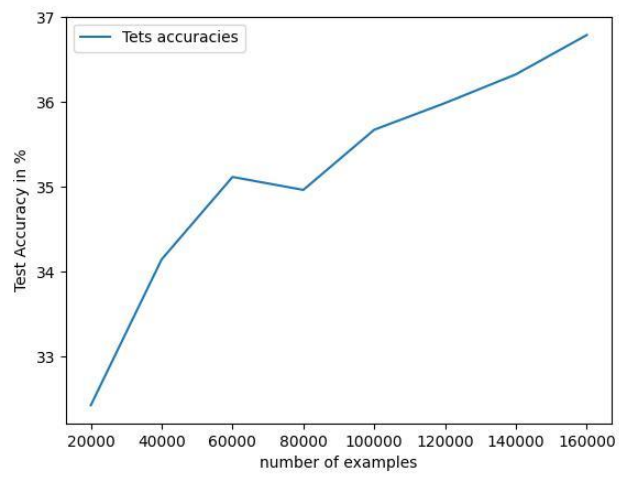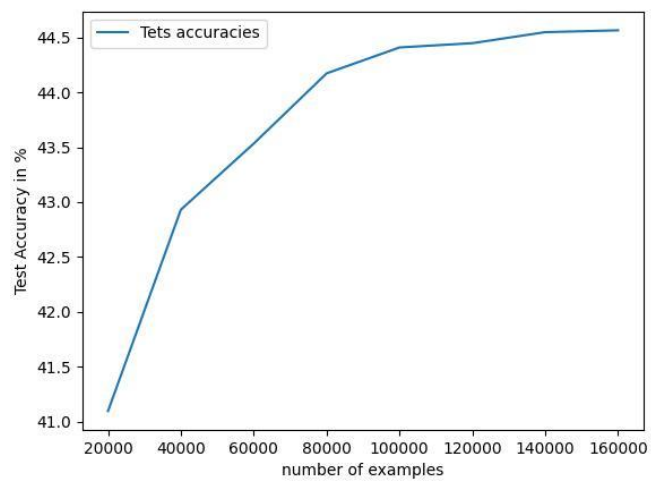   n ∈ {20k, 40k, 60k, 80k, 100k, 120k, 140k, 160k}
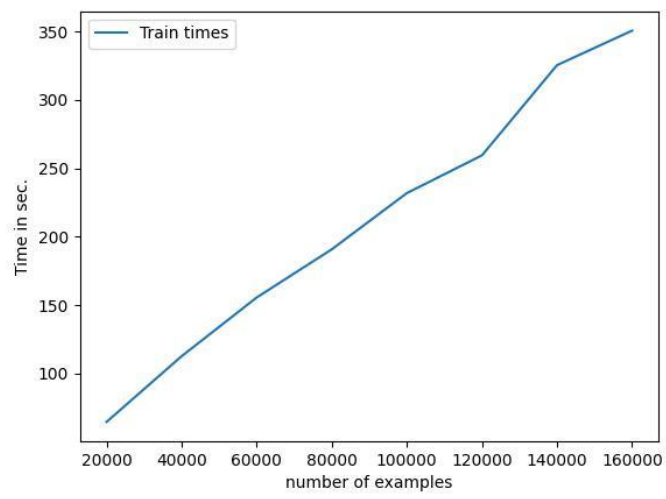
Part (a):





Part (b):

Part (f):

## 2. Neural Networks:

**a.** Neural Network algorithm has been implemented from scratch for Multi-class Classification. In this outputs are converted into one hot encoded form. We have used Mini-Batch Stochastic Gradient Descent to train the network. A fully connected architecture is assumed while implemented the algorithm.
The implementation is generic for following parameters:
(a) Mini-Batch Size
(b) Number of features/attributes
(c) Hidden Layer Architecture : List of numbers denoting the number of perceptrons in the corresponding hidden layer.
(d) Number of target classes

**b.** To train the network following parameter are fixed.
Number of Hidden Layer = 1
$\eta = 0.1$
Hidden Layer Units = [5, 10, 15, 20, 25]
Mini Batch Sizer = 100

Convergence Criteria: (current average cost − previous average cost) < 5e-5 and no. of epochs > 50

Train accuracies....
Train accuracy for hidden layer of 5 units is: 79.77166666666666 %
Train accuracy for hidden layer of 10 units is: 87.37833333333333 %
Train accuracy for hidden layer of 15 units is: 88.54833333333333 %
Train accuracy for hidden layer of 20 units is: 88.81666666666666 %
Train accuracy for hidden layer of 25 units is: 89.61666666666666 %

Test accuracies....
Test accuracy for hidden layer of 5 units is: 78.31 %
Test accuracy for hidden layer of 10 units is: 84.92 %
Test accuracy for hidden layer of 15 units is: 86.15 %
Test accuracy for hidden layer of 20 units is: 86.33 %
Test accuracy for hidden layer of 25 units is: 87.15 %
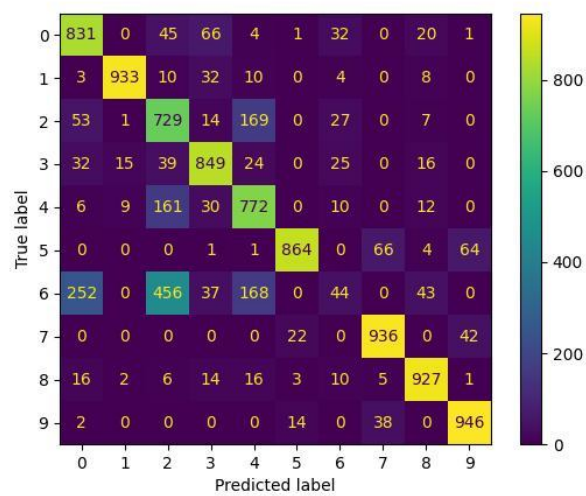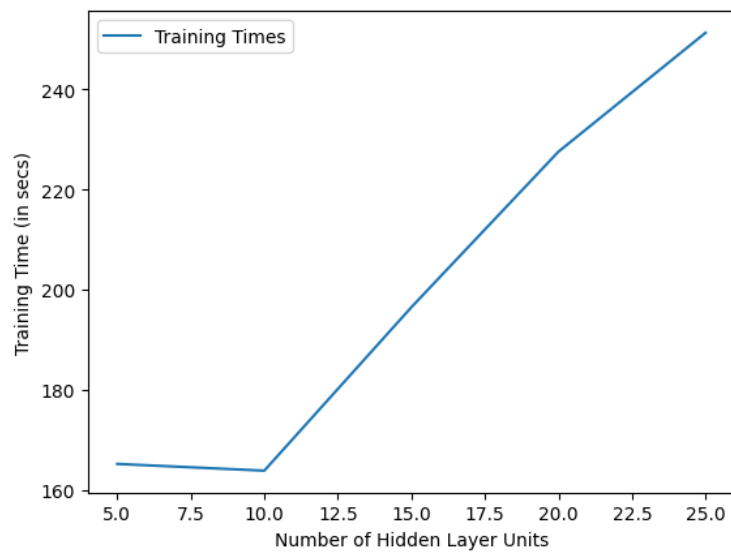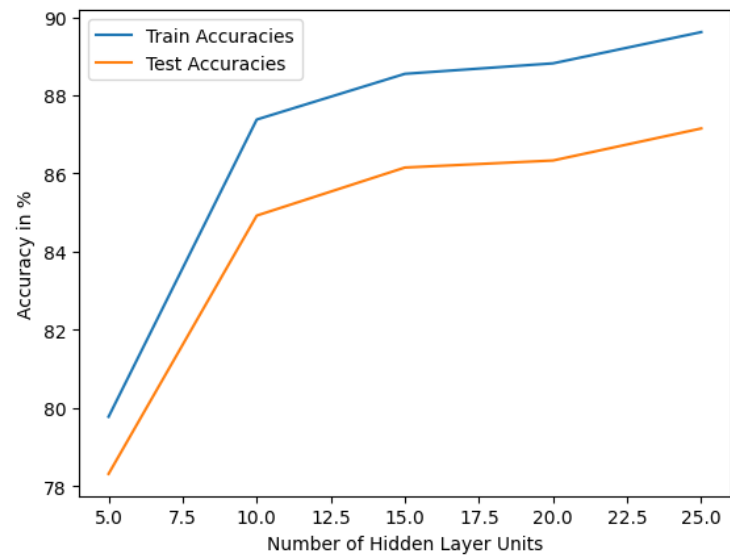
Training Time....
Training time for hidden layer of 5 units is: 165.18345522880554 sec
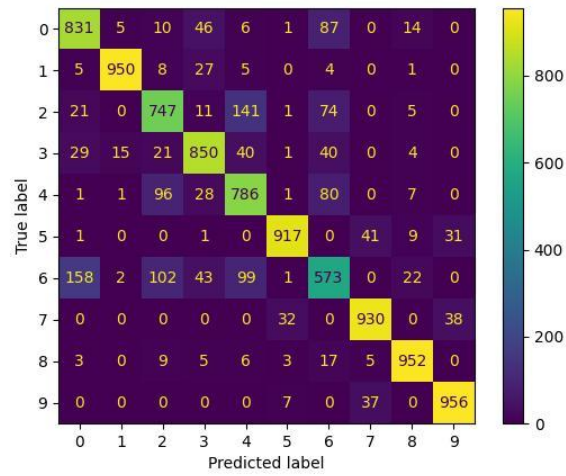Training time for hidden layer of 10 units is: 163.8339102268219 sec
Training time for hidden layer of 15 units is: 196.49265480041504 sec
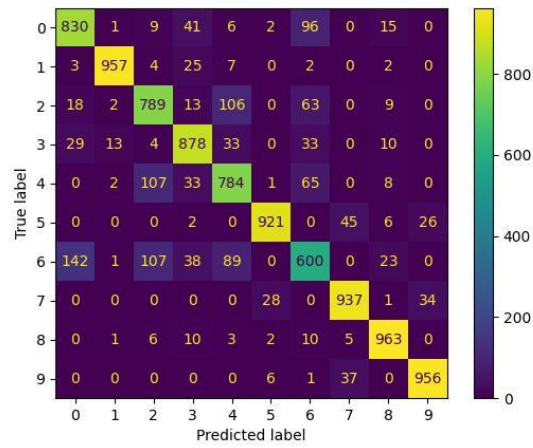Training time for hidden layer of 20 units is: 227.5731236934662 sec
Training time for hidden layer of 25 units is: 251.32813835144043 sec
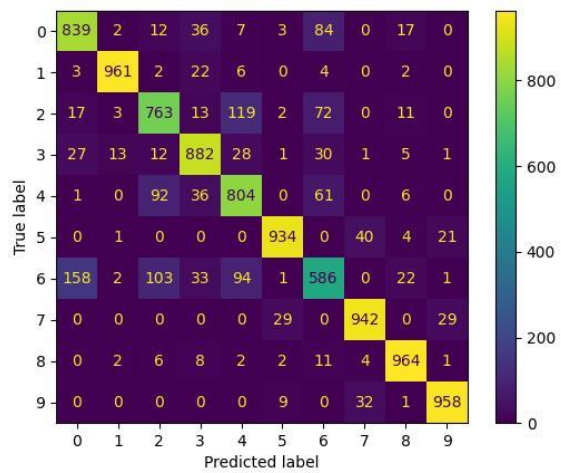
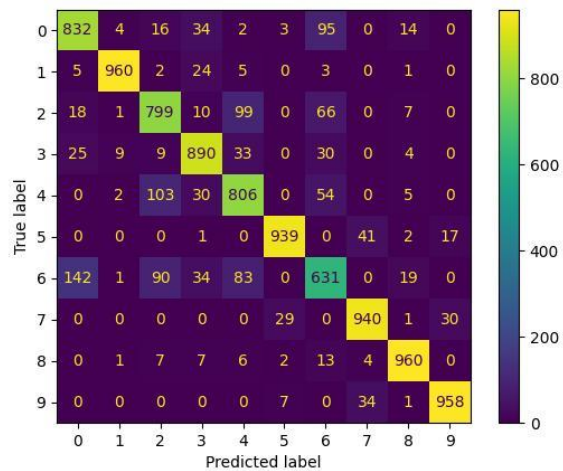Confusion Matrix for hidden layer of 5 units

Confusion Matrix for hidden layer of 10 units



Confusion Matrix for hidden layer of 15 units



Confusion Matrix for hidden layer of 20 units

Confusion Matrix for hidden layer of 25 units

Conclusion:
For large hidden layer units, time taken to train the model is increasing substantially as we can see it from the graphs above.
Accuracy increases as the no. of hidden layer units increases.

c. **Using Adaptive Learning Rate**
Learning rate ($\eta$) would get updated after every epoch
$$\eta_e = \eta_0 / \sqrt{e}$$
Where $\eta_0 = 0.1$ and e is the current epoch number.

Convergence criteria: (current average cost – previous average cost) < 5e-5 and no. of epochs > 50

Train accuracies....
Train accuracy for hidden layer of 5 units is: 72.75%
Train accuracy for hidden layer of 10 units is: 83.42333333333333%
Train accuracy for hidden layer of 15 units is: 84.11%
Train accuracy for hidden layer of 20 units is: 84.37333333333333%
Train accuracy for hidden layer of 25 units is: 84.30666666666666%

Test accuracies....
Test accuracy for hidden layer of 5 units is: 72.16%
Test accuracy for hidden layer of 10 units is: 82.03%
Test accuracy for hidden layer of 15 units is: 82.71%
Test accuracy for hidden layer of 20 units is: 83.08%
Test accuracy for hidden layer of 25 units is: 83.07%
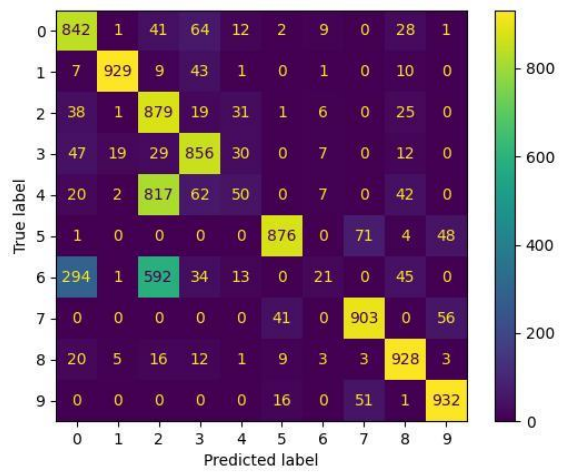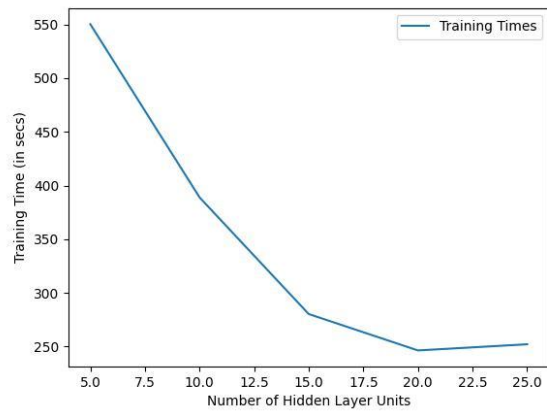
Training Time....
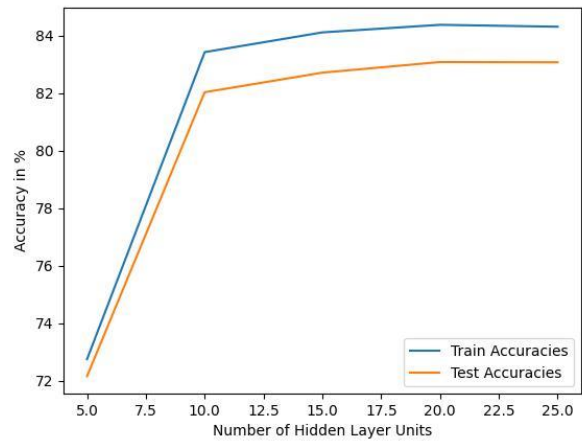Training time for hidden layer of 5 units is: 550.1844565868378 sec
Training time for hidden layer of 10 units is: 388.9296541213989 sec
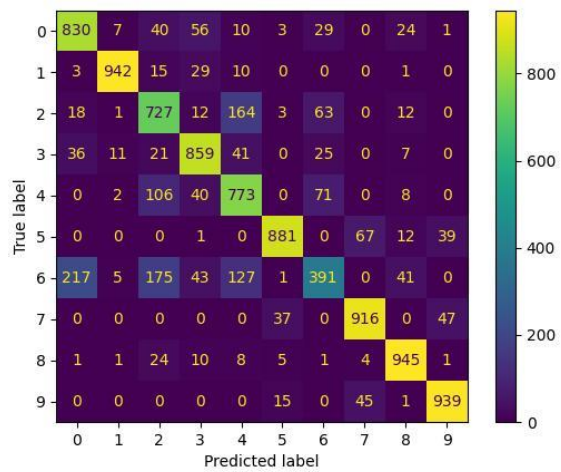Training time for hidden layer of 15 units is: 280.2832624912262 sec
Training time for hidden layer of 20 units is: 246.40840601921082 sec
Training time for hidden layer of 25 units is: 252.13452792167664 sec
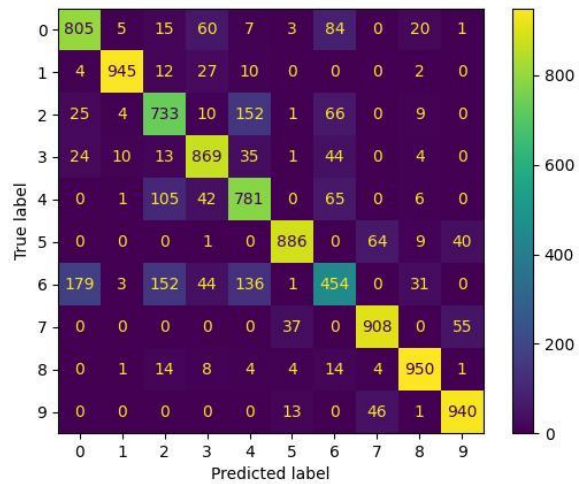
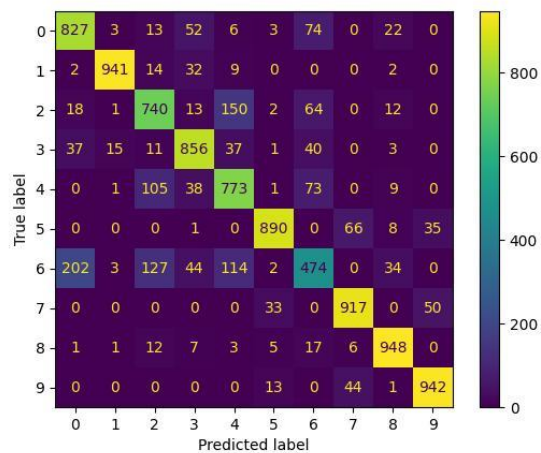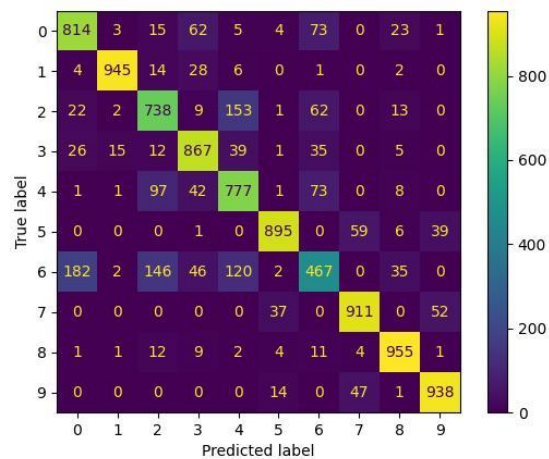Confusion Matrix for hidden layer of 5 units

Confusion Matrix for hidden layer of 10 units



Confusion Matrix for hidden layer of 15 units



Confusion Matrix for hidden layer of 20 units

Confusion Matrix for hidden layer of 25 units

Observations:
Network trained using adaptive learning rate takes more time than what it had taken in non-adaptive learning.
Accuracies are slightly less than what we had earlier for each hidden layer units.
Adaptive learning rate has taken more time than non-adaptive learning rate as learning rate keeps on decreasing with each epoch and that is leading to its slow convergence.
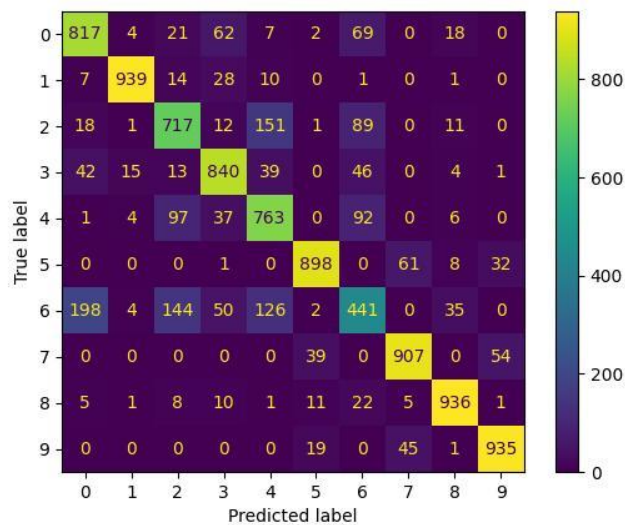
d. **ReLU and sigmoid activation unit comparison:**
A network is implemented with 2 hidden layers with 100 units each.

**Results with sigmoid Activation Units**

Train accuracy for hidden layer of [100, 100] units is: 83.45833333333333%
Test accuracy for hidden layer of [100, 100] units is: 81.93%
Training time for hidden layer of [100, 100] units is: 872.1873600482941 sec
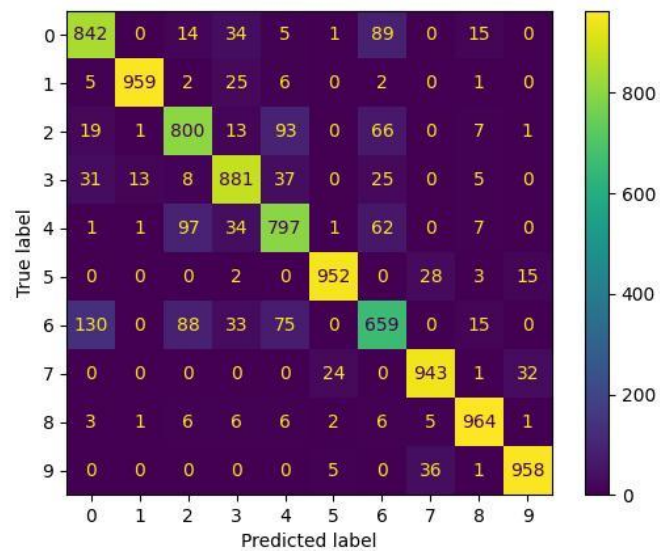
**Result with ReLu Activation Units**

Train accuracy for hidden layer of [100, 100] units is: 90.275%
Test accuracy for hidden layer of [100, 100] units is: 87.55%
Training time for hidden layer of [100, 100] units is: 1426.159077167511 sec



Observations:
ReLu activation function for hidden layers gives better results.
Multiple hidden layers gives better accuracies compare to single hidden layer in part (b).

e. **Experiment by increasing the number of hidden layers:**

Hidden layers in network from 2 to 5. Keeping the number of hidden units per layer to be 50.

**Accuracies with ReLu function for Hidden Layers:**

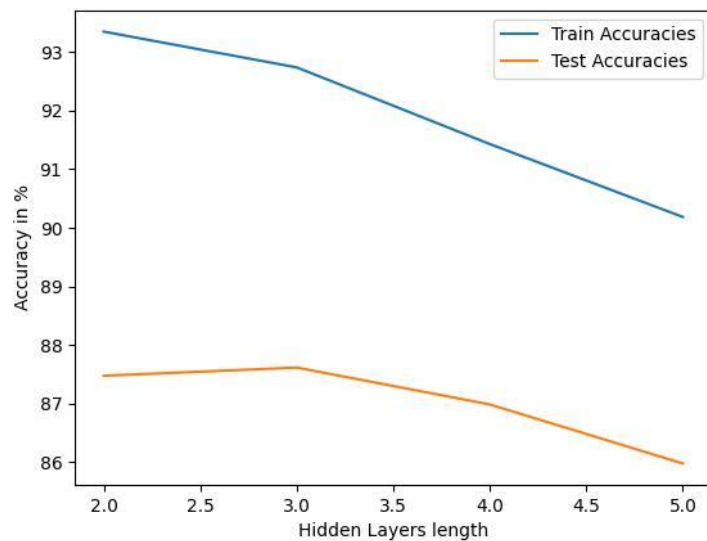Train accuracy for hidden layer of [50, 50] units is: 93.35333333333334 %
Train accuracy for hidden layer of [50, 50, 50] units is: 92.74166666666666 %
Train accuracy for hidden layer of [50, 50, 50, 50] units is: 91.43166666666667 %
Train accuracy for hidden layer of [50, 50, 50, 50, 50] units is: 90.185 %

Test accuracy for hidden layer of [50, 50] units is: 87.47 %
Test accuracy for hidden layer of [50, 50, 50] units is: 87.61 %
Test accuracy for hidden layer of [50, 50, 50, 50] units is: 86.98 %
Test accuracy for hidden layer of [50, 50, 50, 50, 50] units is: 85.97 %

Training time for hidden layer of [50, 50] units is: 489.35870456695557 sec
Training time for hidden layer of [50, 50, 50] units is: 412.4283878803253 sec
Training time for hidden layer of [50, 50, 50, 50] units is: 334.13640904426575 sec
Training time for hidden layer of [50, 50, 50, 50, 50] units is: 356.448077917099 sec
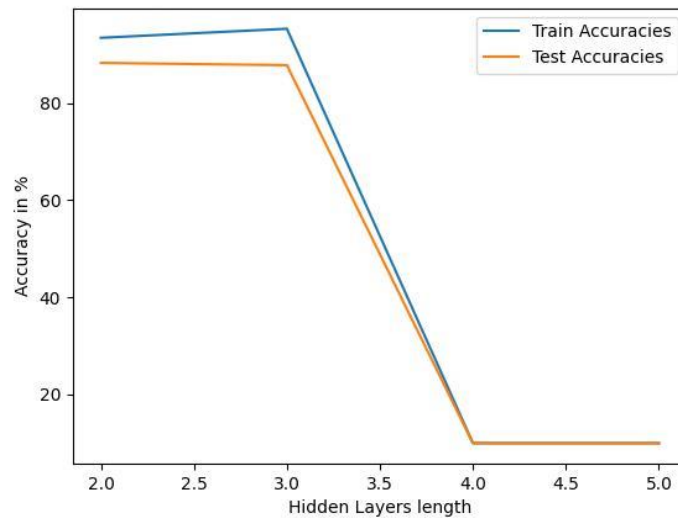


**Accuracies with Sigmoid function for Hidden Layers:**

Train accuracy for hidden layer of [50, 50] units is: 93.4 %
Train accuracy for hidden layer of [50, 50, 50] units is: 95.26 %
Train accuracy for hidden layer of [50, 50, 50, 50] units is: 10.0 %
Train accuracy for hidden layer of [50, 50, 50, 50, 50] units is: 10.0 %

Test accuracy for hidden layer of [50, 50] units is: 88.23 %
Test accuracy for hidden layer of [50, 50, 50] units is: 87.77000000000001 %
Test accuracy for hidden layer of [50, 50, 50, 50] units is: 10.0 %
Test accuracy for hidden layer of [50, 50, 50, 50, 50] units is: 10.0 %

Training time for hidden layer of [50, 50] units is: 663.4633610248566 sec
Training time for hidden layer of [50, 50, 50] units is: 986.2991437911987 sec
Training time for hidden layer of [50, 50, 50, 50] units is: 74.34880018234253 sec
Training time for hidden layer of [50, 50, 50, 50, 50] units is: 83.22205805778503 sec

Observations:

ReLu activation function more suitable and gives better results for multiple hidden layer neural network.
Sigmoid function in multiple hidden layers network suffers from vanishing gradient problem.
The gradient for the sigmoid function will saturate and when using the chain rule, it will shrink. By contrast, the derivative for ReLU is always 1 or 0.
Because of this as the no. of hidden layer increases we can see the drastic reduction in accuracies of model.

**f.   Cross-Entropy loss:**

$$\mathbf{h} = \mathbf{w}^T \mathbf{X}$$

$$\text{Logistic regression: } \mathbf{z} = \sigma(\mathbf{h}) = \frac{1}{1 + e^{-\mathbf{h}}}$$

$$\text{Cross-entropy loss: } J(\mathbf{w}) = -(\mathbf{y}log(\mathbf{z}) + (1 - \mathbf{y})log(1 - \mathbf{z}))$$

$$\text{Use chain rule: } \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \frac{\partial J(\mathbf{w})}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{w}}$$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{z}} = -\left(\frac{\mathbf{y}}{\mathbf{z}} - \frac{1 - \mathbf{y}}{1 - \mathbf{z}}\right) = \frac{\mathbf{z} - \mathbf{y}}{\mathbf{z}(1 - \mathbf{z})}$$

$$\frac{\partial \mathbf{z}}{\partial \mathbf{h}} = \mathbf{z}(1 - \mathbf{z})$$

$$\frac{\partial \mathbf{h}}{\partial \mathbf{w}} = \mathbf{X}$$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \mathbf{X}^T(\mathbf{z} - \mathbf{y})$$

$$\text{Gradient descent: } \mathbf{w} = \mathbf{w} - \alpha \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}}$$

Train accuracy for hidden layer of [50, 50, 50] units is: 92.43666666666667 %

Test accuracy for hidden layer of [50, 50, 50] units is: 87.83 %

Training time for hidden layer of [50, 50, 50] units is: 334.8329563140869 sec

Total no. of epochs for hidden layer of [50, 50, 50] units is: 85


g. **MLPClassifier:**

Parameters for MLPClassifier:
hidden_layer_sizes=(50,50,50)
activation="relu"
solver='sgd'
max_iter=85
learning_rate_init=0.1

Train accuracy:  93.49833333333333 %
Test accuracy:  87.42 %
Training time:  215.6363410949707 sec

Observation:
Accuracies in this part and part (f) are almost similar as we kept all parameters same. Which gives validation to my neural network model that I've build from scratch.


**Libraries used for this assignment:**

pandas

numpy

matplotlib

sklearn

time

sys

scipy

xgboost

lightgbm