

Name:- Kalyan Kumar Paladugula

NetID:- Kpalad4

UIN:- 679025059.

Final Exam

i) No. Generally speaking, Regression problems are easier than Classification problems because Regression is continuous optimization and classification is discrete optimization. Continuous optimization is easier than the discrete. The constraint of the classification is that the o/p of the classification should not take an arbitrary value.

Ex:-  $y = w^T x$  - Regression

$y = \text{Sign}(w^T x)$  - classification.

The Sign function makes the <sup>associated</sup> cost function a non-convex. So, we could end up in the local minima.

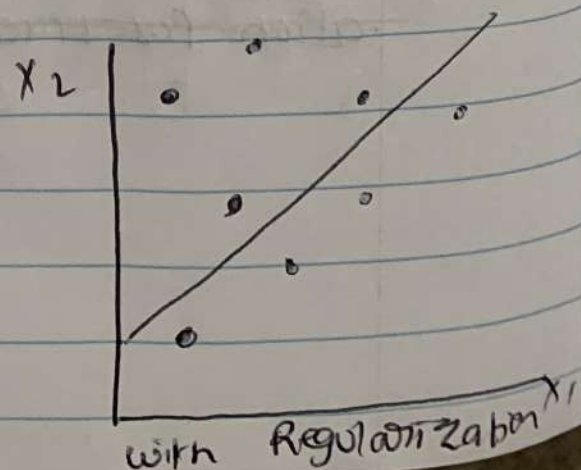
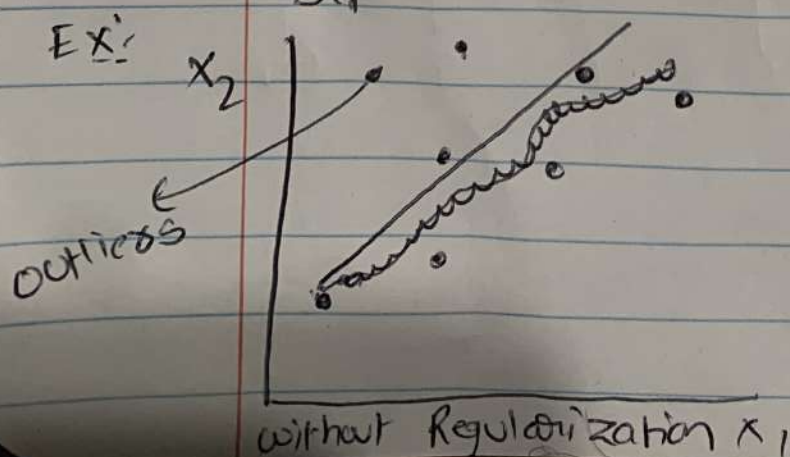
~~Even when we need to apply Kernel transfor-~~  
~~mation for linear~~

b) False.

With higher  $c$ , the margin would become small (the error / hinge loss ~~it~~ would be less).  
So, the No of support vectors would be lower than the original SVM.

c) False.

Regularization is done to increase the accuracy / decrease the error in the testing dataset i.e. to make model more generalizable which might increase the error in training set.

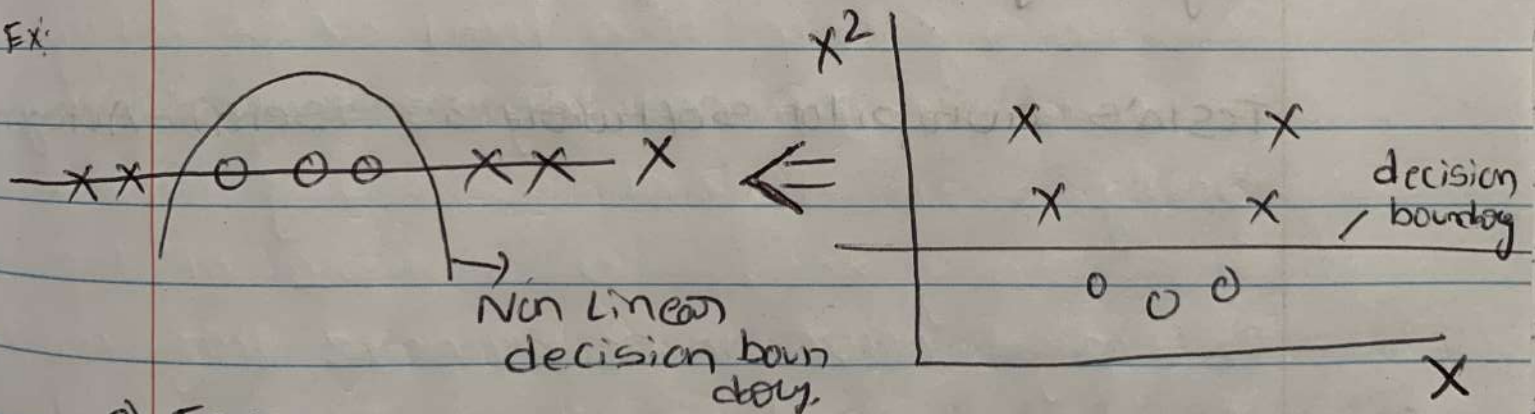




d) True:

Kernel transformation produces higher dimensional version of the data in which the data could become linearly separable and SVM produces a linear decision boundary in the transformed space. If we project the data along with decision boundary, back to the original dimensions, the discriminator would be nonlinear.

Ex:



e) False.

KNN is a lazy learner - no training.

So, the training time has no correlation with the 'K' value.

2) a) Self-driving car programmes by tech giants like Google, Tesla et.c

Ex:

In 2018, a Tesla car crashed by hitting a highway barrier and Tesla auto-pilot system found probably at fault in this crash. In this case, the Tesla software failed to detect the highway barrier.

Tesla's autopilot software is Semi-Autonomous



## b) Leave-one-out-cross-validation: (LOOCV)

In this, we leave one data point for validation/testing and use the rest for training.

LOOCV can be used for the KNN. as KNN is a lazy learner - No training, LOOCV would be very efficient for KNN.

## c) ~~Very less/zero~~

- c) i) ~~Zero~~ <sup>Very low</sup> training error and high testing error
- ii) high Variance, Low bias.
- iii) When training loss continues to decrease but Validation loss started to increase boundaries

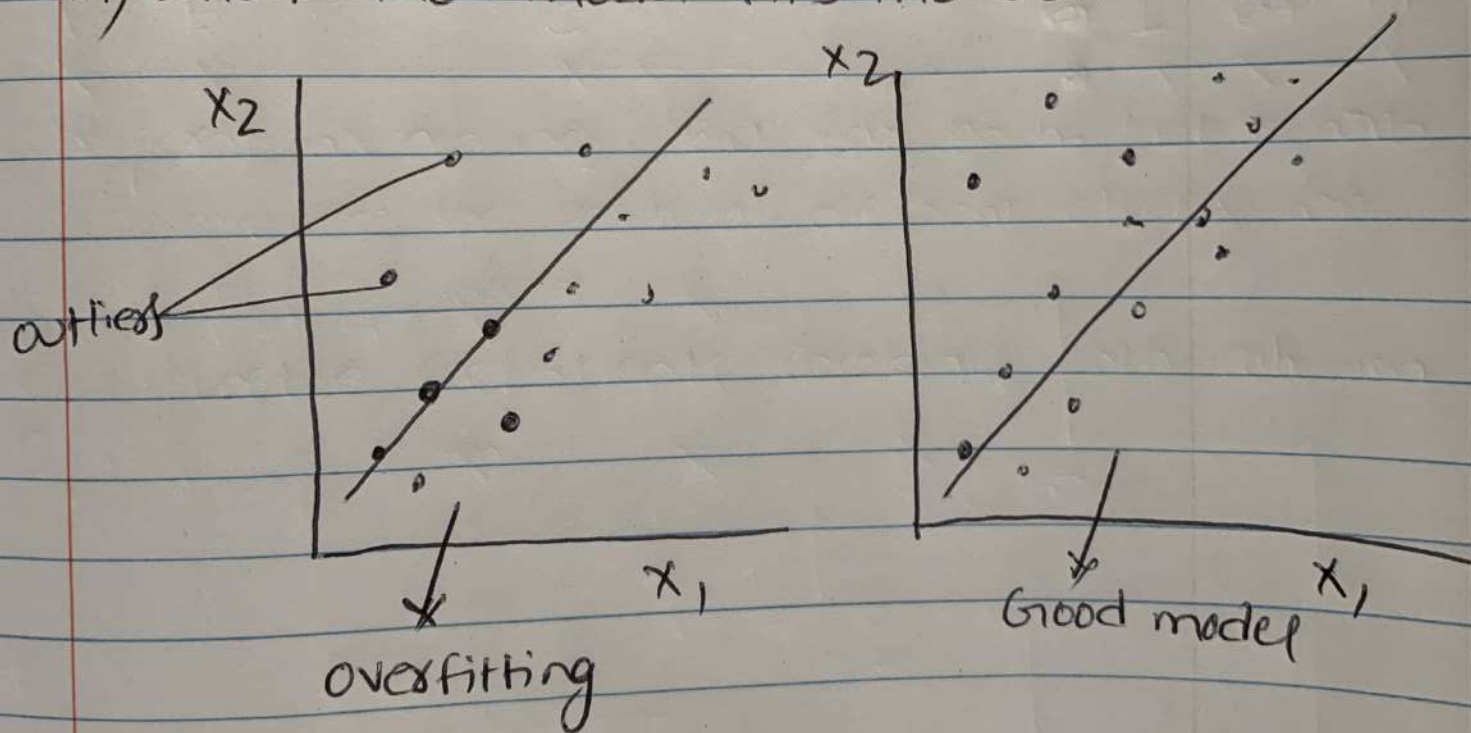
a) c) i) Very less/zero training error and high testing error

- Low generalizability.

ii) In the error vs complexity graph, if the training error continues to decrease but the testing/validation error starts increasing

- overfitting region

iii) When the model fits the outliers also



iv) Very complex decision boundaries

v) If the model has high variance and low bias, then we can say it is overfitting.



#### d) Vanishing Gradient

when the  $|\text{gradient update}| < 1$ , during the backpropagation, since we multiply the sensitivity of the connected nodes, the gradient updates would become very small for the initial layers of the network. This is called "Vanishing Gradient".

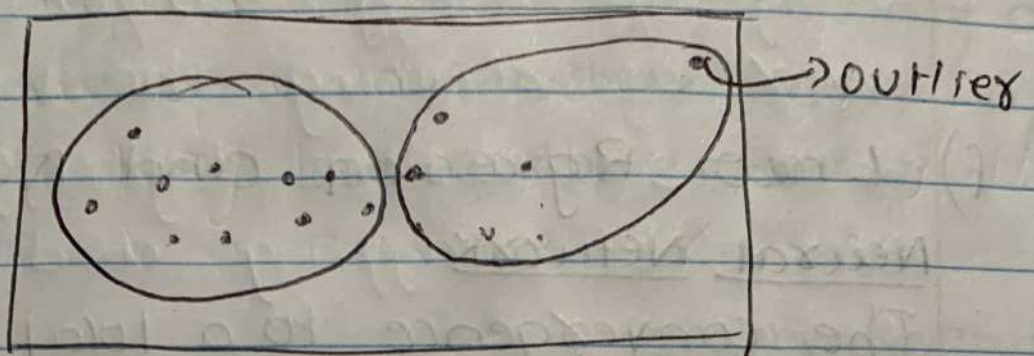
→ This causes initial layers to learn very slowly.

This might occur when we have large number of hidden layers/ nodes. i.e. Deep Neural Networks or when  $|\text{gradient update}| \approx 0$ . or when due to the Sigmoid Activation function.

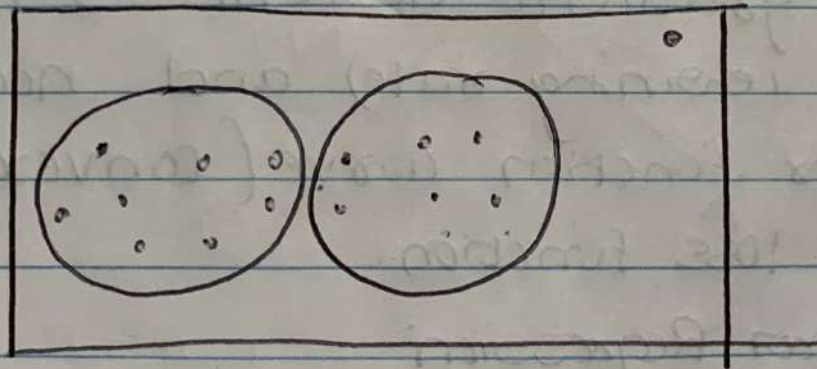


e) Kmeans is vulnerable to the Random start points.

i) If we choose the outliers as starting points, Kmeans takes a long time to converge or might produce incorrect clusters



Predicted clusters.



desired clusters.

ii) Some random initializations cause Kmeans to stuck in a local minima rather than Global minima



## f) Linear Regression and Neural Networks.

### Neural Network

The convergence to a local or global optima is contingent on the learning rate (as gradient descent is heavily dependent on learning rate) and nature of the error function curve (convex or non-convex) i.e. loss function.

### Linear Regression:-

In case of gradient descent for Linear regression, the cost function curve has only one minimum that is global minimum. Hence, here gradient descent always converge to Global minimum. Keeping learning rate to neither too small or too large.

g) By separating the data into training and testing, we are avoiding "Data snooping" which causes good accuracy in the Seen data but poor accuracy in the unseen data.

By building the model on the training data, we can check its generalizability by applying it on the testing data.

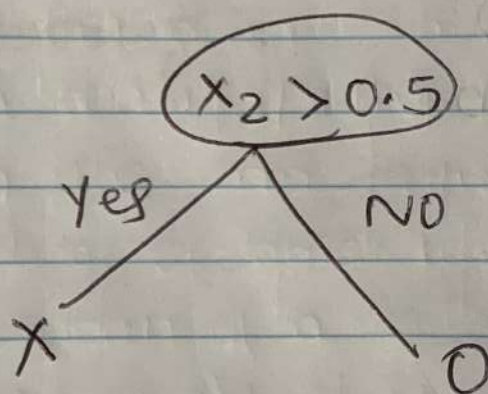
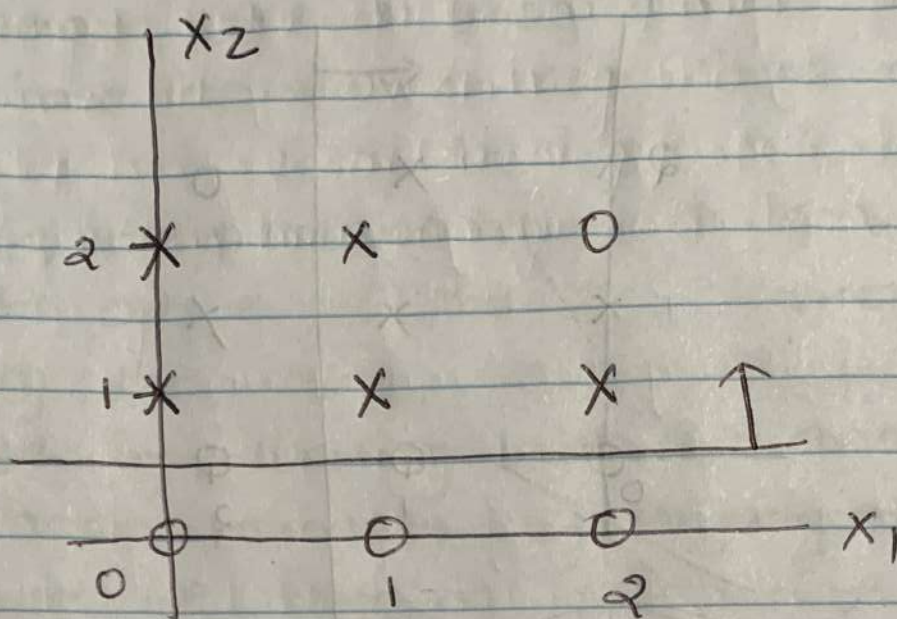
We can also find the best hyperparameters by applying the model on validation/testing data.



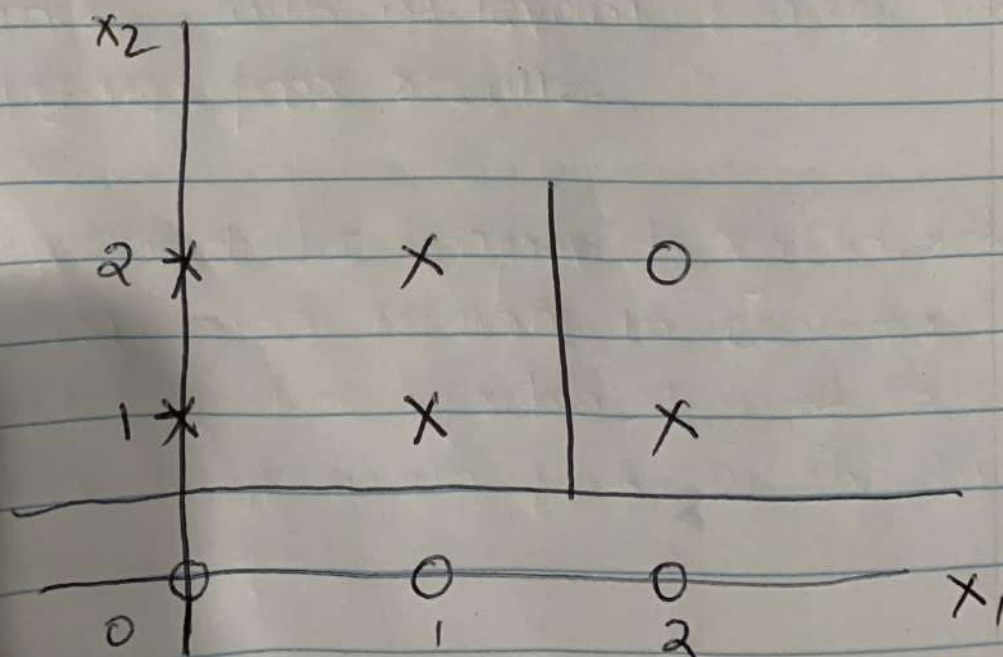
h) Ensemble methods use unstable learners as base learners to make them independent.

Since, there would be a lot of variance in the outputs of the base learners, by taking the voting, we would reduce the error by reducing the variance of the base learners by aggregation/ensembling in case of bagging and bias in case of boosting.

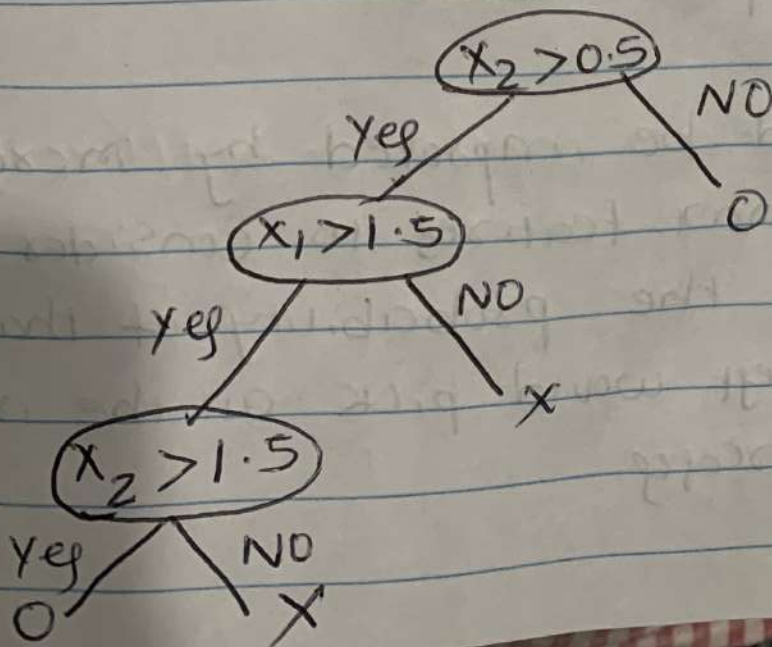
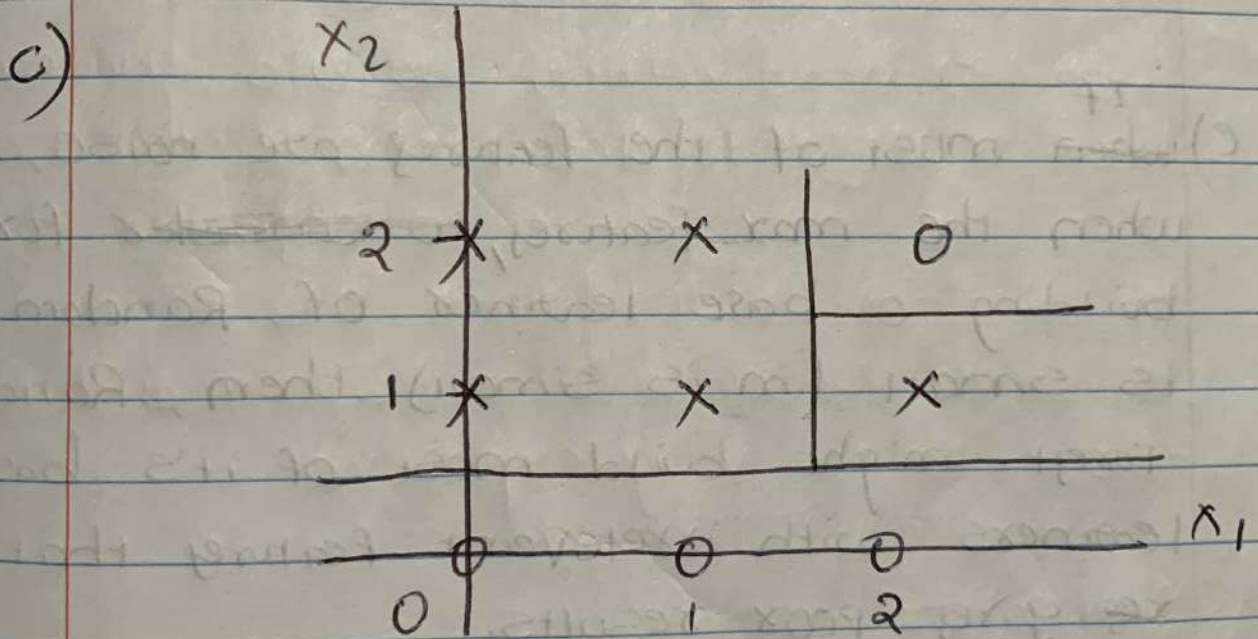
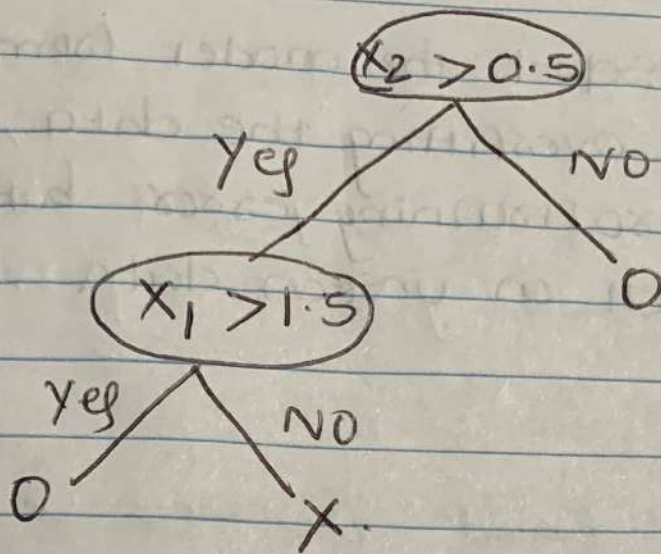
3) a)



b)







d) I might not select the model from c because it is overfitting the data, this model has zero training error but might not perform well on unseen data

e) <sup>If</sup> when most of the features are noise, when the max-features, ~~to consider~~ for building a base learner of Random Forest, is small ( $m$  is small) then, Random-Forest might build most of its base learners with irrelevant features that ~~gives~~ gives poor results.

This could be improved by increasing the  $m$  (max-features to consider for base learner) as the probability of that the Random Forest would pick all the irrelevant features decreases



# 4) unsupervised Learning

a)

A	B	C	D	E	F	G
X	X	X	X	X	X	X
-2	-1.5	0	0.25	3	4	6

$$AD = 0.25$$

$$AG = 8$$

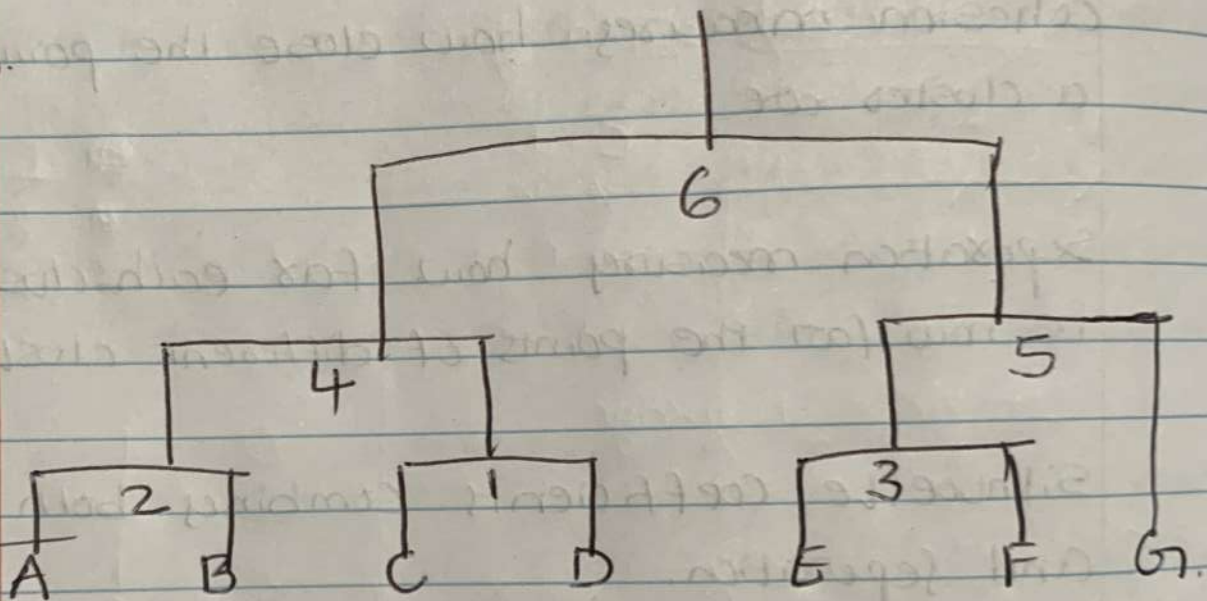
$$AF = 6$$

~~$$AG = 8$$~~

$$CF = 4$$

$$CG = 6$$

$$EG = 3$$



$$AF = 6$$

$$AG = 8$$

$$EG = 3$$

#### 4) b) Semi-supervised learning

Semi-supervised learning falls between unsupervised learning and supervised learning.

In SSL, we first use the unlabeled dataset. we first use unsupervised learning on the unlabeled data and cluster the similar data. Then we use the existing labeled data to label the unlabeled data.

It is most useful when the acquisition of unlabeled data is cheap and labeling the data is very expensive.

#### Applications

- i) Alexa
- ii) Movie Recommendation - Netflix, Prime
- iii) Book Recommendation



- c) We can judge/assess the quality of an unsupervised learning approach by calculating the
- i) cohesion - Intra-cluster distance
  - ii) Separation - Inter-cluster

Cohesion measures how close the points in a cluster are

Separation measures how far each cluster is, i.e. how far the points of different clusters are

Silhouette coefficients combine both cohesion and separation.

### Q) 5) a):

The hyperparameters for my final neural network are:

**No of Hidden Layers = 1, No of hidden nodes = 10, Learning Rate = 0.1, no of iterations/epochs = 1000, solver = 'adam', activation function = 'relu'**

To make the algorithm simple, I used rectangular weight matrix i.e. the number of nodes in all hidden layers is same

I used the following grid for the Grid Search method to select the best hyperparameters:

1. hidden = [1,2,5]
2. nodes = [2,5,10]
3. lr=[0.1,1]
4. af=['tanh', 'relu']
5. solver=['sgd','adam']

I tried every combination of the above (72 combinations) and calculated the cross-validation error. Then, I picked the hyperparameters which gave the least cross validation error.

I applied the best model on the test dataset and got **77** percent accuracy

**While doing the search, I got a few convergence warnings which say that the no of iterations = 1000 is not enough for the model to converge. This may be due to the 0.1 learning rate**

### Q) 5) b):

The concentration bound obtained at 95% confidence interval using Hoeffding bound is **0.096**

### Q) 5) c):

I am getting different results (accuracy) every time I run my neural network because of

- 1) Random initialization of weights
- 2) Randomness in the training data (I set shuffle = True in the MLPClassifier)
- 3) Due to the adaptive Learning Rate (I set Learning Rate = 'adaptive' in the MLP Classifier)
- 4) Multiple Local Minima (in case of Non-Convex cost functions)

### Q) 5) d):

We can say whether a model is overfitting or not by checking the Training error.

I saved the training error of all the models I built using the cross\_validate function of the Sklearn.model\_selection. To get the training error, I set the return\_train\_score parameter of cross\_validate to True.

For the model with the least validation error (the model that I picked as the best) has training error of **0.17**.

### Q) 5) c)

```
print('The training error obtained with my best model = ' +str(t_error))
```

```
The training error obtained with my best model = 0.17041666666666666
```



Since, the training error is not zero, I can say that my model is not overfitting.

I checked the testing accuracy of the model that has low training error:

```
b_hl2, b_hn2, b_lr2, b_af2, b_s2, error, t_error2 = minimum7(d, e, f, g, h, c, b)
print(t_error2)
b_hid2 = [b_hn2 for n in range(b_hl2)]
classifier2 = MLPClassifier(hidden_layer_sizes= b_hid2, activation = b_af2, learning_rate_init = b_lr2, learning_rate = 'adaptive',
classifier2.fit(X_train, y_train.values.ravel())
y_pred2 = classifier2.predict(X_test)
nn_accuracy2 = accuracy_score(y_test, y_pred2)
print('The testing accuracy obtained with the model that has lower training error = ' + str(nn_accuracy2))
```

```
0.01597222222222221
```

```
The testing accuracy obtained with the model that has lower training error = 0.755
```

I got lesser testing accuracy of 0.755 with the model that has low training error of 0.016.

**Note: I used random\_state = 100 in the MLP classifier in the code to check the testing accuracies of my best model and the model that has low training error of 0.016**