1) SHORT ANSWERS

a) "All of the following answers are my own work"

S. Anusha

b) Regularization finds a model where it does not fit the training model so well so as to reduce overfitting. i.e., it makes the model less complex so as to decrease the variance in the unseen data

L1 Regularization → Lasso → $\sim\sim\sim + \lambda \sum_i |w_i|$

L2 Regularization → Ridge → $\sim\sim\sim + \lambda \sum_i (w_i)^2$

Both the models give some penalty to the model - In L1, it adds absolute value to the coefficient as penalty & in L2, it adds squared value to the coefficient as penalty.
To be more precise, lasso can only either totally remove or add feature coefficients depending on their correlation with the model. But ridge can have weights for each of the feature's coefficients. (high or low)

c) Decision trees are quite unstable & each model can give different results for the same dataset. i.e, they have high independence b/w the models & their errors.
The ensemble work well with unstable models as they combine each weak learners together & predict based on voting or giving weights to each of the model based on the accuracy of each model.

∴ Decision trees are unstable & can give different results everytime (i.e., this can lead to overfitting), ensemble methods like bagging & boosting can be used to decrease overfitting

thereby decreasing variance & boosting can reduce bias without specifally increasing variance.

d) Usually the activation fn is added to the Neural Network to achieve some kind of non-linearity in the function.

In NN, during feedforward, we multiply (dot product) input vector with random weights vector. & so on, till the output node. Generally, this dot product is linear operation & gives similar o/ps. If we don't provide activation fn like RELU & sigmoid (which are basically based on confidence values) functions, our NN would just give some linear output & will not able to label or categorize the outputs.
So, to have some non-linearity in the outputs, we use activation fns so that we can get desired output labels.

e) Yes, we can apply kernel transformation to NN. It helps in creating much more complex models that a required. For example, in image classification based on pixel values, having very high degree for kernel can make the computation more complex & requires high computation time. It depends on the CPU availability & depends on the problem that we are trying to solve.

In simple digit classification, with less no. of colors & their intensities, we need not require kernel transformation because from the rough distingutation b/w pixels, we will be able to predict the label. But for some complex images that are hard to read & identify, kernel transf., may be able to distinguish pixel by pixel (basically in the low confidence regions).

f) SVM would take longest to retrain.

① In Nearest Neighbours, we don't actually have any training. So, retraining would not be needed or happen.

② In Neural N/ws, if we encounter a new point, we can just train that single point by passing it to the neural network that we've already built & the previous points that we've trained can still be retained. So, the retraining is not at all complex or time consuming here.

③ In SVM, if we encounter a new point, we need to change our seperator & margins depending on where the point lies. Suppose we have hard margin SVM, & the new point falls within the margin, we need to shift our boundary line & recalculate all the values like margin, support vectors etc. So, this is the longest to retrain.
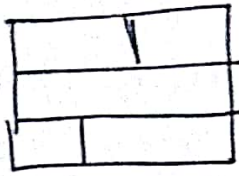
g) Bagging works on the concept of creating equal sized bags by drawing samples from the training data & creating different learners for each bootstrap that is created.

For eg, if we just have a single model on the whole trai -ning data, our model might overfit the data & sometimes we might even fit our model for the outliers that are present which will obviously lead to overfitting & high variance.

In bagging, all bags don't contain all sample points i.e., the model doesn't memorize all training points which makes each model have some amount of bias to the points that are not included. Also, outliers can easily be identified and the results can exclude them.

h) In decision trees, the boundaries are continuous & is rectangular form. & then each block can further be subdivided using continuous boundary.
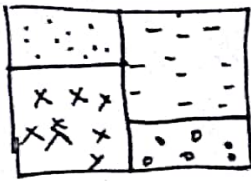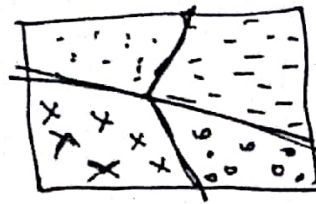
For eg:- boundaries created by decision trees are · is,

 } rectilinear decision boundaries.

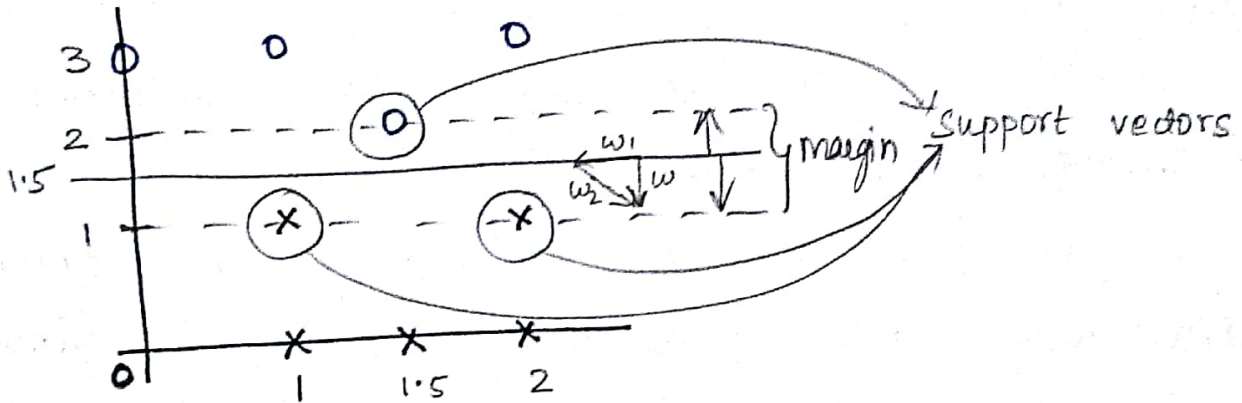So, & would look how the decision boundaries are.

Eg:-



Decision Trees        Neural Networks.

# 2) SUPPORT VECTOR MACHINES.

## a)



## b)

$$\beta_0 + \beta_1 x_1 + \beta_2 x_2 = y.$$

$$\boxed{w_1 = 0}$$

$\Rightarrow$ when $x_2 = 1.5 \Rightarrow w_0 + w_2(1.5) = 0 \Rightarrow \boxed{w_0 = -1.5 w_2}$

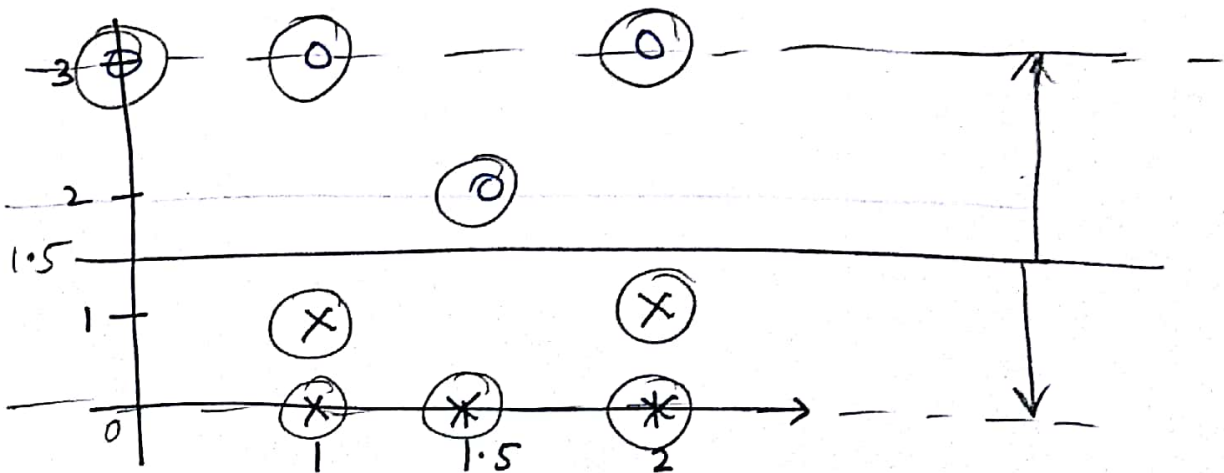$\Rightarrow$ when $x_2 = 1 \Rightarrow w_0 + w_2(1) = 1 \Rightarrow -1.5 w_2 + w_2 = 1$

$$\Rightarrow -0.5 w_2 = 1 \Rightarrow \boxed{w_2 = -2}$$

Now $w_0 = -1.5(w_2) \Rightarrow w_0 = -1.5(-2) \Rightarrow w_0 = -\dfrac{3}{2} \times -2$

$$\Rightarrow \boxed{w_0 = +3}$$

$\therefore \quad \beta_0 = +3, \quad \beta_1 = 0, \quad \beta_2 = -2$

## c)

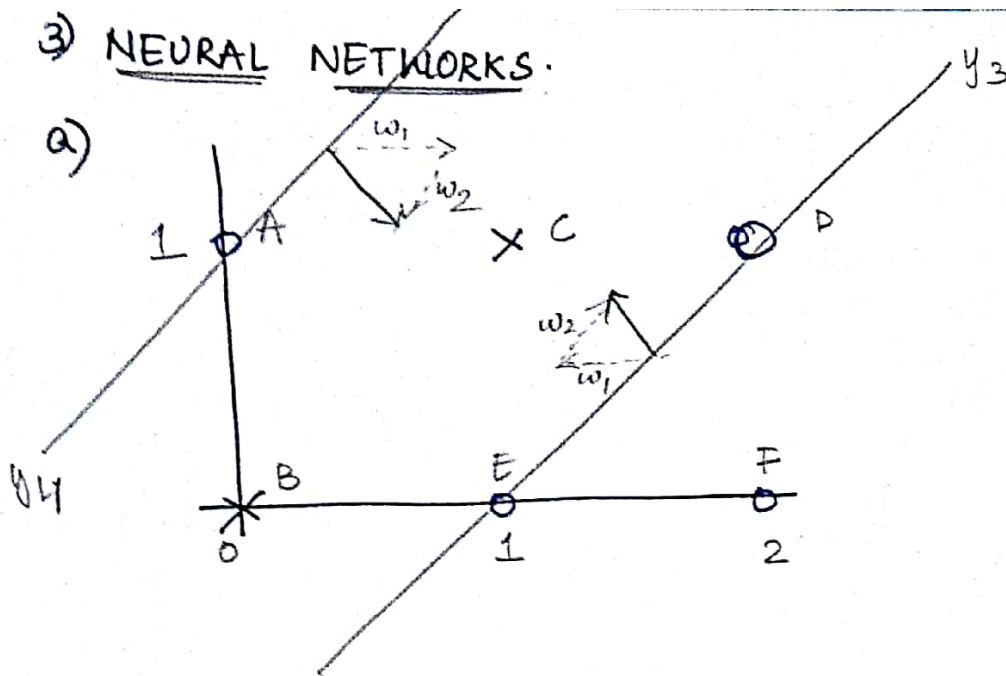d) The model is almost same where the seperator is still at 1.5

But the soft margin increases, which includes high no. of misclassification.

Because for high values of c, we generally tend towards hard margin & as 'c' value becomes lower, the margin increases allowing some misclassifications & as 'c' keeps decreasing the no. of misclassifications increase.
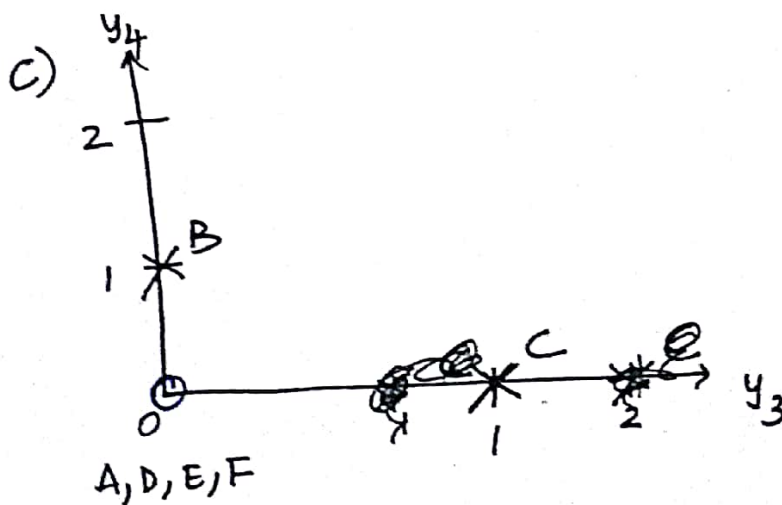
③ NEURAL NETWORKS.

a)



b) $y_3 = x_1 + x_2 + bias$

$y_3 = \underset{w_1}{-1} + \underset{w_2}{(+1)} + (+1)$

∴ Weights are $-1, +1, +1$

$y_4 = x_1 + x_2 + bias$

$y_4 = \underset{w_1}{+1} + \underset{w_2}{(-1)} + (+1)$

∴ Weights are $-1, +1, +1$

c)



A, D, E, F

⎱ Transformed $y_3$ & $y_4$
⎰      Space.

Scanned by CamScanner

d)



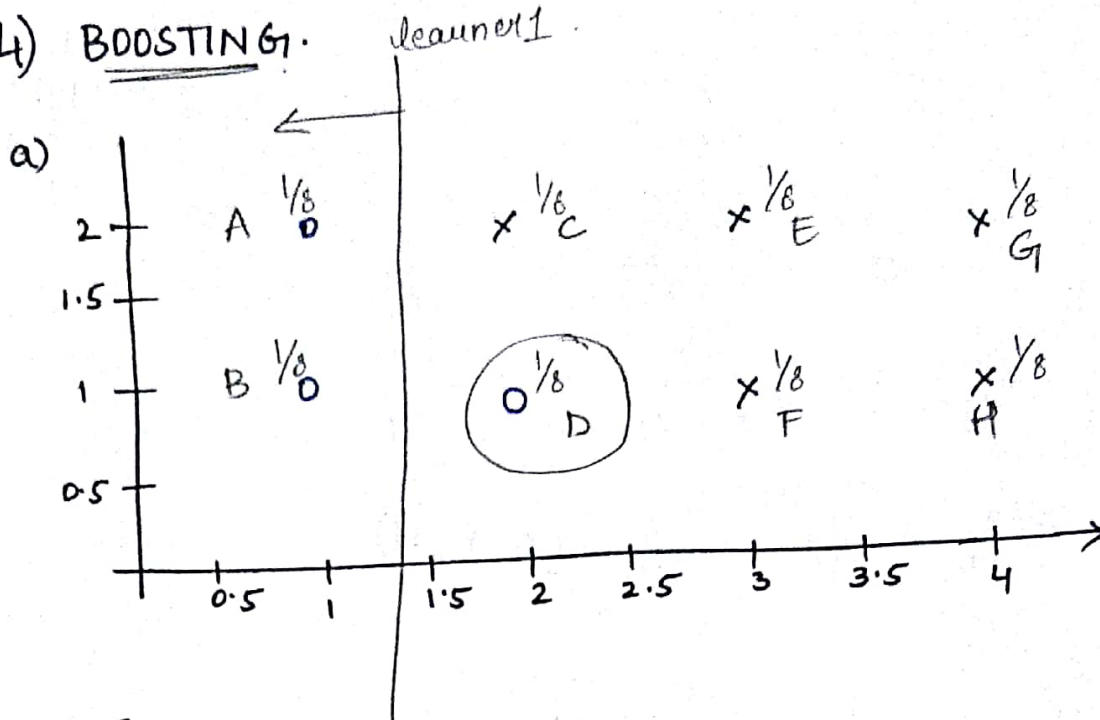e)  $y_5 = x_1 + x_2 + bias$

$y_5 = (+1) + (+1) + (-1)$
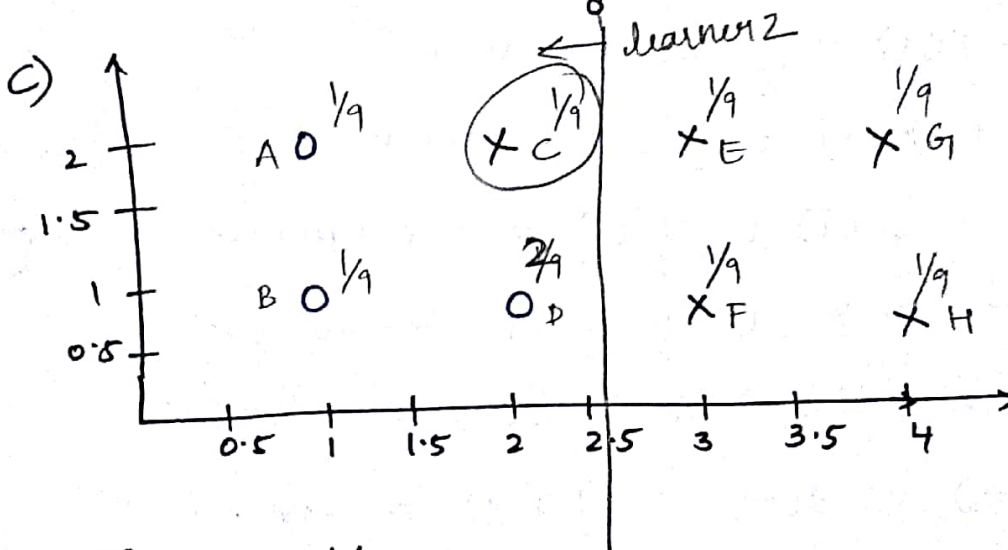
$\therefore$ Weights are $+1, +1, -1$

4) BOOSTING.

a)



learner1

Error $= \frac{1}{8}$ (point D)

b) $\alpha = \frac{1-E}{E} = \frac{1-\frac{1}{8}}{\frac{1}{8}} = \frac{7}{1} = 7.$   $\boxed{\alpha = 7}$

(Double weight of D)

c)



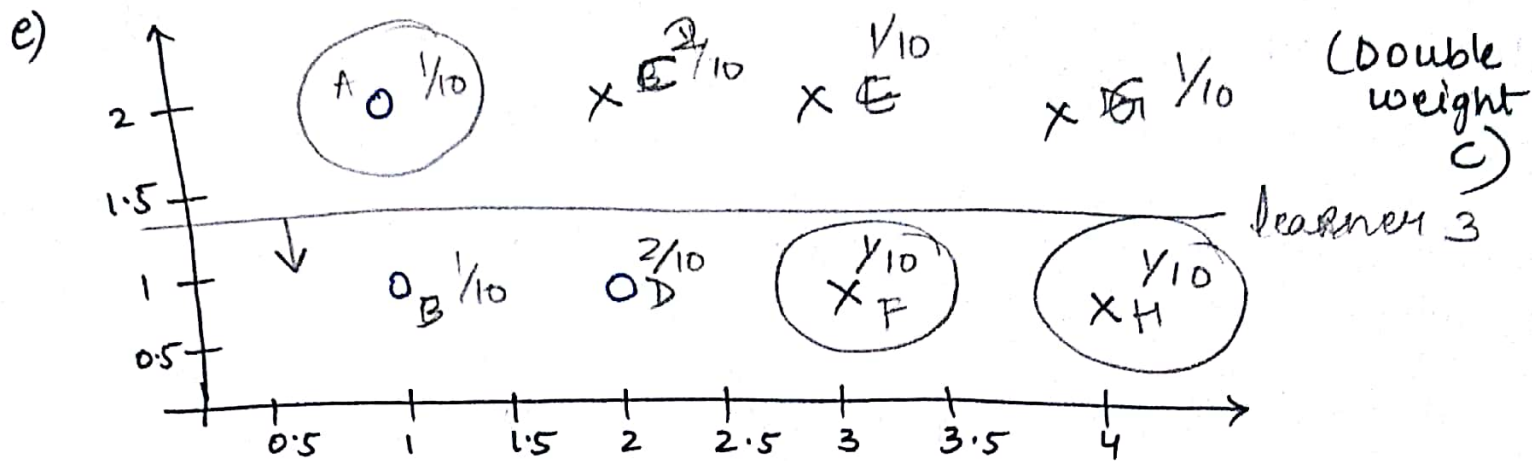learner 2

Error $= \frac{1}{9}$

d) $\alpha = \frac{1-\frac{1}{9}}{\frac{1}{9}} = \frac{8}{1} = 8.$

$\boxed{\alpha = 8}$

e)



$$\text{Error} = \frac{1}{10} + \frac{1}{10} + \frac{1}{10} = \frac{3}{10} \quad (A, F, H)$$

f)
$$\alpha = \frac{1 - \frac{3}{10}}{\frac{3}{10}} = \frac{7}{3} \implies \boxed{\alpha = \frac{7}{3}}$$

g) Classification for $\overset{(+ve)}{A} \implies 7(1) + 8(1) + \frac{1}{3}(-1) = +ve$ ✓

$(+ve) B \implies 7(1) + 8(1) + \frac{1}{3}(1) = +ve$ ✓

$(-ve) C \implies 7(-1) + 8(1) + \frac{1}{3}(-1) = -ve$ ✓

$(+ve) D \implies 7(-1) + 8(1) + \frac{1}{3}(1) = +ve\ value$ ✓

$(-ve) E \implies 7(-1) + 8(1) + \frac{1}{3}(1) = -ve$ ✓

$(-ve) F \implies 7(-1) + 8(-1) + \frac{1}{3}(+1) = -ve$ ✓

$(-ve) G \implies 7(-1) + 8(-1) + \frac{1}{3}(-1) = -ve.$ ✓

$(-ve) H \implies 7(-1) + 8(-1) + \frac{1}{3}(+1) = -ve$ ✓

∴ All the points are classified correctly after 3 iterations.

∴ Training Error = 0 //.