

## SECTION A - True/False

1. Bagging reduces the chance of overfitting making the model more adaptable to unseen data. (True)
2. Averaging predictions reduces fluctuation in data (True)
3. In Boosting, the weights of data points change at every step to make the next gradient descent one, accurate. (True)
4. Sampling without randomness can introduce Sampling Bias. (True)
5. If a model is overfit, training error can be 0% while testing error can be 100%. (True)
6. Regularization simplifies the model, which decreases bias. (~~True~~) (False)
7. Increasing the depth of a decision tree always prevents overfitting. (~~True~~) (False)
8. Every classifier makes assumptions about the data. (True)
9. Random forests are more accurate than single decision tree because they combine bagged trees. (True)
10. Irreducible error is the lower bound on error due to inherent noise in the data. (True)

## SECTION-15 - ERM and EVM

1. fill the following table

Model	Loss function	Regularizer
SVM	Hinge Loss	L2
LASSO	Squared Loss	
RIDGE	Squared Loss	L2

2. Identify the Loss-function A-E

A:

B:

C:

D:

E:

3. Short Answer

(a) which loss function can be optimized using gradient descent? why?

→ Squared Loss function can be optimized using gradient descent because it is differentiable function (smooth with ~~loss function~~)

(b) which loss functions can be optimized using Newton's method

→ Squared Loss and Logistic Loss function can be optimized using Newton's method because these functions are twice differentiable (smooth).

## SECTION-C Bias and Variance

1. Explain one major reason why underfitting occurs?

→ Underfitting occurs when our model is too linear or simple ~~for complex~~ to capture the complex patterns in the data, it means our model has not learned anything.

2. If both training and test errors remains high, what does this imply ~~to~~ about the data?

→ It means our model fails the ~~exact~~ training set and test set ~~and~~ our model has not learned anything. So if our model is not in underfitting then there is a chance that no pattern in the data to be found.

3. Explain how bagging reduces variance.

→ Bagging creates many copies of the same model. Each copy is trained on a random subset of the data. So average of all their results is usually very accurate and drastically reduces variance.

4. Explain the effect of boosting on bias and variance.

- Boosting Primarily reduces Bias by Sequentially training models to correct the errors of previous ones.  
It ~~also~~ reduces variance.  
It can also reduce Variance.

#### SECTION D - KNN & Curse of Dimensionality

1. what step can be taken to reduce KNN computation time?

- we can ~~speed up~~ reduce KNN computation time by reducing the number of features (using dimension reduction or shrinking the dataset).

Instead of checking every single point, using ~~Space~~ ~~kd~~-Trees or Ball trees allows the KNN algorithm to efficiently skip distant neighbours.

2. (a) Does Squared Euclidean distance change predictions? Explain.

No, it doesn't change predictions because the nearest neighbours remain the same regardless of whether we compare raw distances or squared distances.

b) Does it affect the previous dimensionality ~~conclusion~~ conclusion  
→ No, it doesn't affect the previous dimensionality conclusion - because high dimensional data remains sparse and distances still lose meaning, so we have still problem of high dimensionality.

3. Why does KNN perform poorly in high dimensions with few data points.

→ In high dimensional space available data points are scattered far apart. This breaks distance based algorithm ~~KNN~~ because 'nearest' neighbours are no longer actually close in very high dimensions. So every point roughly the same distance away from every other point.

4. How do bias and variance change with  $K$ ?

→ with a small  $K$ , we get low bias but high variance, meaning the model captures local details perfectly but is very sensitive to noise. And as we increase  $K$  bias increases and variance decreases.

5. When is KNN preferred over linear SVM?

→ KNN is preferred when preferred when decision boundary is highly non-linear or irregular, as a linear SVM is restricted to drawing straight lines to separate classes if our data has complex clusters rather than a clear linear split. KNN is much more effective.

## SECTION - Decision Trees

Q1. Derive that the optimal prediction at a leaf (with squared loss) is the mean.

→ actual target values  $\Rightarrow y_1, y_2, y_3, \dots, y_n$   
 $n =$  number of points in the leaf  
 Single Value  $c$

$$J(c) = \sum_{i=1}^n (y_i - c)^2$$

To Diff. w.r.t. to  $c$

$$\frac{dJ}{dc} = \frac{d}{dc} \sum_{i=1}^n (y_i - c)^2$$

$$= -2 \sum_{i=1}^n (y_i - c)$$

Set the derivative to zero.

$$-2 \sum_{i=1}^n (y_i - c) = 0$$

$$\sum_{i=1}^n (y_i - c) = 0$$

$$n_c = \sum_{i=1}^m y_i$$

$$C = \frac{1}{m} \sum_{i=1}^m y_i$$

So the value of  $C$  that minimizes the Squared loss is exactly the arithmetic mean of the target values in that lead

2. what are the max/min Gini impurity values for 3 classes?

→ minimum Gini impurity = 0

Node  $i$  is completely pure; meaning all data points belong to a single

Class So if class 1 has prob.  $\frac{1}{3}$  and others

$$\text{Gini} = 1 - (\left(\frac{1}{3}\right)^2 + \left(\frac{1}{3}\right)^2 + \left(\frac{1}{3}\right)^2) = 0.$$

maximum Gini impurity =  $\frac{2}{3}$

here data is ~~evenly~~ evenly split among all 3 classes. Each class has a prob. of  $\frac{1}{3}$ .

$$\text{Gini} = 1 - \left(\left(\frac{1}{3}\right)^2 + \left(\frac{1}{3}\right)^2 + \left(\frac{1}{3}\right)^2\right)$$

$$= \frac{2}{3}$$

3. why are decision trees myopic learners?

→ decision trees are considered myopic learners. because they use a greedy algorithm. At each step they choose the split that maximize the immediate gain based on the current state without caring about future values.

4. Explain two methods to avoid overfitting in decision trees

→ (1) pre-pruning ⇒ stop the tree from growing before it become too complex by setting limits like maximum depth or minimum sampling.

(2) Post - Pruning ⇒ Allow the tree to grow fully to fit the training data then trim back branches that do not improve accuracy

on a validation set.

## Section F:- Boosting and Bagging.

1. Can Random forests ~~use~~ use the same data for training and testing?

Justify.

→ No, we can't use the same data for training and testing. If we test a random forest on the exact same data it was trained on, the model will likely achieve 100% accuracy. This is because the model has memorized the patterns in the data set. So ~~we~~ we can't ~~general~~ generalize our model ~~to~~ to new unseen data.

2. Explain the key difference between bagging and boosting.

→ bagging:- bagging creates many copies of the same model. Each copy is trained on a random subset of the data.

boosting:- it trains model sequentially. The second model focuses heavily on the data points the first model got wrong.

So key difference is ~~how they~~ working principle.