

Assignment-1

- (A) 1. True 5. True 9. True
2. True 6. False 10. True
3. False 7. False
4. True 8. True
- (B) 3: Model Loss function Regularizer
- | | | |
|------------|------------------------------------|----------------------------------|
| 1. SVM - | Hinge Loss $\max(1 - y\hat{y}, 0)$ | ℓ_2 Norm ($\ \omega\ ^2$) |
| 2. LASSO - | Squared Loss $(y - \hat{y})^2$ | ℓ_1 Norm ($\ \omega\ $) |
| 3. RIDGE - | Squared Loss $(y - \hat{y})^2$ | ℓ_1 Norm ($\ \omega\ $) |
- 3
- Loss functions that are differentiable, can be optimized using gradient descent. Examples include mean squared loss and hinge loss (due to existence of sub-gradient).
 - Loss functions that are at least twice differentiable, can be optimized using Newton's method. Squared loss can be optimized using Newton's method.
- (C)
- One of the major reason why underfitting occurs is when the model is too simple so it can't capture the complex ~~nature~~ of features in data. This can happen due to using a strong regularization or using overly simple model like ~~the~~ linear regression which may not be able to capture non-linear features.
 - If both training and testing errors are ~~high~~ high, then it implies that either the data is too complex for the model to extract and learn some features or, the ~~data~~ in the data, the important features are missing. ~~is not well fitted~~

3. Bagging reduces variance by training multiple models on different samples of data (obtained by random sampling with ~~repetition~~ repetition) and then, averaging their predictions, which may help in cancelling out individual model fluctuations and reduce noise in output.
4. Boosting reduces bias by sequentially training models such that each new model focuses more on wrong predictions made by previous models. In this way, by repeatedly correcting these wrong predictions, it finally helps build a strong model with weights that can now capture the features of our complex data. Bias decrease as model is now not simple and properly capture the features in the data.

Variance generally increases as boosting sequentially leads to overfitting.

(E)

1. Let at a leaf node, all samples reaching the leaf have target values $\{y_1, y_2, \dots, y_n\}$.

Let the predicted value by model be λ .

$$\text{So, Squared Loss} \Rightarrow L(\lambda) = \sum_{i=1}^n (y_i - \lambda)^2$$

So, optimal prediction (λ) will be such that the squared loss is minimum.

$$\text{So, } \frac{dL(\lambda)}{d\lambda} = -\sum_{i=1}^n 2(y_i - \lambda) = 0$$

$$\lambda = \bar{y} = \frac{\sum y_i}{n} \rightarrow \text{mean}$$

2. We know, $\text{Gini} = 1 - \sum_{i=1}^n p_i^2$ where, p_i = probability of class i .
- * $(\text{Gini})_{\min} \rightarrow$ when all the samples belong to one target class (So one of the class has $p=1$ and others 0)

For $n=3 \rightarrow 3$ classes,

$$(Gini)_{min} = 1 - 1^2 = 0$$

* $(Gini)_{max} \rightarrow$ When ~~the~~ the samples are equally likely to occur in each class (so all class has some probability)

For $n=3 \rightarrow 3$ classes,

$$(Gini)_{max} = 1 - 3 \times \left(\frac{1}{3}\right)^2 = \frac{2}{3}$$

3.

Decision trees work in such a way that they select the ~~the~~ split that gives maximum immediate reduction in impurity without looking on the long-term impact of that split. And since, once a split is made, the decision becomes irreversible. So a better local split may affect the global output. Therefore, because of this greedy way the decision trees work, they are called myopic learners.

4. ~~and values of other splits~~

Pre-pruning :- This method avoids overfitting by restricting tree growth using constraints like maximum depth (no. of levels in a tree), minimum samples per leaf (so that tree don't create leaves that have ~~the~~ very less data points, noise) and minimum impurity decrease.

Post-pruning :- This method reduces overfitting by first letting a tree ~~the~~ grow completely and then removing branches that do not significantly improve performance, to improve generalization. (removes noise).

(F)

1. Yes, in random forests, each tree is only trained on a subset of that data so many data points are there that are not yet seen by that tree. So, we can have an algorithm to sample the data ^(by sampling) in such a way that we make 2 datasets for each tree, that have quite less data points common between them and hence these datasets can be used as training and testing.

2. ~~both bagging and boosting are used in decision trees~~

Bagging

Boosting

- Trains many models independently
- Uses randomly sampled datasets for each model.
- Mainly tries to reduce variance
- Trains model sequentially
- Improves weights of data points (in each iteration)
- Main role is to reduce bias