

Assignment - 1

- (A)
- ① True. Training set assignment
 - ② True. \hat{y}_i is either 0 or 1
 - ③ False. making "gradient descent more accurate" is very vague. purpose of boosting is to focus more on previously misclassified pts.
 - ④ True. bottom constraint
 - ⑤ True. features are both positive
 - ⑥ False. Simplifying model increases bias
 - ⑦ False. More memory \rightarrow more overfitting
 - ⑧ True
 - ⑨ True
 - ⑩ True

<u>Model</u>	<u>Loss Function</u>	<u>Regularizer</u>
SVM	Hinge loss	L_1 - L_2
LASSO	Squared loss	L_1
RIDGE	Squared loss	L_2

- ② Insufficient DATA

③ (a) Loss functions can only be optimised by gradient descent if they are differentiable.

eg: squared loss. → works

eg: hinge loss → gradient descent does NOT work

(b) Newton's method works ^{but} for loss functions that are at least twice differentiable.

eg. it works here

(c) ① One major reason for underfitting is model being too simple (v. high bias). Such a model might not be able to properly learn data.

② If model error is high during both testing and training, problem might be in data. Data might be too noisy, or might have important features missing, features which would have helped model learn.

③ Bagging trains multiple models on different bootstrap samples of dataset and averages their predictions thus reducing noise and variance.

④ Boosting builds a strong model by combining sequentially training models on failures of previous models, giving importance to

Bias decreases due to boosting since model goes from simple to more expressive.

Variance may increase, since later models focus heavily on specific data points, which can be problematic if data is too noisy.

(D)

(E) (1) We know Loss (where λ is prediction),

$$\text{Loss } L(\lambda) = \sum_{i=1}^n \ell(\lambda - y_i)^2$$

To minimise loss, $\frac{dL}{d\lambda} = 0$

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

$$\Rightarrow \lambda = \frac{\sum y_i}{n} \rightarrow \underline{\text{mean}} \quad \left\{ \begin{array}{l} \text{Hence} \\ \text{proved} \end{array} \right\}$$

$$(b) G_i = 1 - \sum_{k=1}^K p_k^2 \quad \text{and} \quad 1 - \sum_{i=1}^n p_i^2$$

$n=3$ with 3 classes

$G_{\min} \rightarrow$ when ~~all~~ all samples belong to one class

$\rightarrow P_K = 1, \text{ rest } 0$

$$\therefore G_{\min} = 1 - 1^2 = 0$$

$G_{\max} \rightarrow$ When samples equally dist.

$P_i = 1/3 \text{ for all } i$

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

(3) Myopia means short-sightedness.

Decision trees only select best split at current node, they do not care look ahead multiple levels.

Basically it is a greedy type algo which looks for immediate impurity reduction. Hence it's a "myopic" learner.

(4) Pre / post pruning are two methods which can avoid overfitting.

Pre - Pruning: Stop tree growth when set condition is triggered (e.g. max depth). This prevents overfitting by memorization of training data.

Post - Pruning: First fully let a tree grow, then remove unnecessary branches which don't improve its performance on a validation dataset.

This removes noise or noisy branches.

(F) (1) Using same data for training and testing is not completely optimal, but it can work well. Since in random forest, each tree is trained on multiple trees are trained, and each tree is trained on a random subset of data, hence a data point can't be seen by all trees (in most cases). Hence data can be reused, as predictions of trees are averaged. Although using a separate dataset would be better.

(2) The main difference b/w bagging and boosting is in bagging, multiple models are trained independently on bootstrap samples, and their prediction are averaged out, while in boosting, models are trained sequentially, with each new model focusing on earlier mistakes.

Bagging \rightarrow reduces variance

Boosting \rightarrow reduces bias