

Convex Set:

Q5) a) A set C is convex if the line segment b/w any 2 pts in C lies in C ,
i.e. $\forall x_1, x_2 \in C, \forall \theta \in [0, 1]$

$$\theta x_1 + (1-\theta)x_2 \in C$$

Convex Function:

Let C be a convex function of \mathbb{R}^n .

A function $F: C \rightarrow \mathbb{R}$ is called convex, if

$$F(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y), \forall x, y \in C$$

If F is a convex function, then all its level sets $\{x \in C \mid F(x) \leq a\}$, where 'a' is a scalar, are convex.

Defⁿ: A convex function is a continuous function whose values at the midpoint of every interval in its domain doesn't exceed the arithmetic mean of its values at the ends of interval.

A function $f(x)$ is said to be strictly convex if for every (A, B) in the domain, the line segment obtained by joining these strictly lies above the curve except the two end points which would be $(A, f(A))$,

$(B, f(B))$ that are common.

Commonality of Lasso & Ridge Regression.

For LASSO the lemma has been proved that, for any y, x and $\lambda \geq 0$ and lasso solutions $B(1)$ and $B(2)$ must satisfy $B_i^{(1)}$ and $B_i^{(2)} \geq 0$ for $i = 1 \dots p$. In other words, any two LASSO solutions must have same signs over the common support.

Hence, LASSO penalty is not strict because if A, B have the same sign then the line segment & curve are eq. which means infinite no of points could be common.

In ridge regression, the minimisation problem is

$$J(w) = \|y - Xw\|_2^2 + \lambda \|w\|_2^2 \quad \text{--- (1)}$$

Another defn for convex loss function is that $\hat{\beta} \in \argmin_{\beta} F(X\beta) + \lambda \|\beta\|_1$ --- (2)

where the loss function $F: \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable & strictly convex.

By joining (1) & (2) we can say that Ridge regression is strictly convex.

Q5(b)

We prefer odd values for K because if we take an even value of K then a situation can arise when there is no tie breaker.

If we are dealing with problem of binary classification & 50% of neighbours are of its class & rest 50% of class 2, then it would be impossible to decide which class current point belongs.

But, if we have a odd K then these kind of situations can be avoided.

So $K=3$ would be better than $K=2$.

Q5(c)

Assuming training set D consists of N pts (x_i, y_i) sampled IID,

let the classifier trained on D be H_D

Given that:

$y = H_D(x)$ is a distribution with mean ' μ ' & variance σ^2

$$\text{Total loss} = E_{x,D} [\underbrace{L(H_D(x), y)}_{\substack{\text{loss over a single } x \\ \text{Avg loss over all 'x'}}}]$$

$\text{let } L(h(x), y) = \frac{1}{2} (h(x) - y)^2$

 Squared loss

$$\begin{aligned} E_{x,D} [(H_D(x) - y)^2] &= E_{x,D} [(H_D(x) - E_D[H_D(x)] + E_D[H_D(x)] - y)^2] \\ &= E_{x,D} [(H_D(x) - E_D[H_D(x)])^2 + (E_D[H_D(x)] - y)^2 + \\ &\quad 2(H_D(x) - E_D[H_D(x)])(E_D[H_D(x)] - y)] \end{aligned}$$

$$\Rightarrow \underbrace{E_{x,D} [H_D(x) - E_D[H_D(x)]]^2}_{\text{Variance}} + \underbrace{E_x [(E_D[H_D(x)] - y)^2]}_{\text{bias}}$$

Suppose K models are trained on K subsets of $D \{D_i\}_{i=1}^K$

The bias term, given some data point (x, y) depends on $E_n[H_D(x)]$ only

$$H_D(x) = \frac{1}{K} \sum_{i=1}^K H_{D_i}(x)$$

$$E_D(H_D(x)) = \frac{1}{K} \sum_{i=1}^K E_{D_i}(H_{D_i}(x))$$

$$= \frac{1}{K} \times K \mu$$

$$= \mu$$

Thus, ensembling doesn't change the bias term.

$$\text{Now, } H_D(x) = \frac{1}{K} \sum_{i=1}^K H_{D_i}(x)$$

$$\Rightarrow \text{Var}(H_D(x)) = \text{Var}\left(\frac{1}{K} \sum_{i=1}^K H_{D_i}(x)\right)$$

$$\left\{ \text{Var}(Kx) = K^2 \text{Var}(x) \right\} = \frac{1}{K^2} \text{Var}\left(\sum_{i=1}^K H_{D_i}(x)\right)$$

$$= \frac{1}{K^2} \left(\sum_{i=1}^K \text{Var}(H_{D_i}(x)) \right)$$

$$= \frac{1}{K^2} \times K L \sigma^2 = \frac{1}{K} \sigma^2$$

We now reduce the variance term to reduce the loss as the bias term doesn't change.

thus $\frac{L}{K} < \sigma^2$

$\frac{L}{K}$ variance of ensemble models

σ^2 variance of single model trained on Entire data

$$L < K$$

thus, there must be more than L learners for the ensemble to perform better than the single bigger model.