

# Assignment 4 { Theory }

Q5

Ans

(a) Convex Sets

→ A Subset  $C$  of  $R^n$  is called convex if

$$\alpha x + (1-\alpha)y \in C$$

$$x, y \in C \text{ & } 0 \leq \alpha \leq 1$$

→ operations that preserve convexity  
↳ Intersection, scalar multiplication  
vector sum, closure, interior  
Linear transformation

A convex set is a collection of points in which the line  $AB$  connecting any 2 points  $A, B$  in the set lies completely within the set

## Conven functions

Let  $C$  be a conven subset of  $\mathbb{R}^n$ . A function  $F: C \rightarrow \mathbb{R}$  is called conven if

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

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

Definition: A conven function is a continuous function whose values at the midpoint of every interval in its domain does not exceed the arithmetic mean of its values at the ends of interval

A function  $f(x)$  is said to be strictly conven if for every  $(A, B)$  in the domain

the line segment obtained by joining these strictly lies above the curve except the 2 end points which would be  $(A, f(A))$ ,  $(B, f(B))$  that are common

### Conveniency of Lasso & Ridge regression

For lasso, the lemma has been proved that, For any  $x, y \neq \lambda \geq 0$  and lasso solutions  $B(1) + B(2)$  must satisfy  $B_{i(1)} \cdot B_{i(2)} \geq 0$  for  $i = 1, \dots, p$ . In other words, any 2 lasso solutions must have the 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 equal which means  $\infty$  no. of points could be common.

Ridge regression Minimization problem is

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

Alternatively, it can be said that  
 $\hat{\beta} \in \arg \min F(\beta) + \lambda \|\beta\|$  --- (2)

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

By (1) + (2) we can say that ridge regression is strictly convex.

Q5(b)

Ans  $k=3$  would be better than  $k=2$  because we don't want to be in a situation where 1 neighbour suggests one class & another neighbour suggests another class for  $k=3$ , there will always be

majority because atleast 2 out of the 3 neighbours will belong to the same class, and there will be majority.

Q5(c) { Plz see next Page }

Q5(C)

Assuming training set  $D$  consists of  $N$  points  $(x_i, y_i)$  sampled  $i.i.d$

Let the classifier trained on  $D$  be  $H_D$

$y = H_D(x)$  is a distribution with mean ' $\mu$ ' & variance  $\sigma^2$  { given }

→ Loss over single  $n$

$$\text{Total Loss} = E_{nD} [L(H_D(x), y)]$$

Average loss over all ' $n$ '

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

Squared Loss.

$$E_{nD} [(H_D(x) - y)^2] = E_{nD} [(H_D(x) - E_D[H_D(x)]) + E_D[H_D(x)] - y]^2$$

$$\Rightarrow E_{x,D} [(H_D(x) - E_D[H_D(x)])^2 + (E_D[H_D(x)] - y)^2]$$

$$2(H_D(n) - E_D[H_D(n)]) (E_D[H_D(n)] - \gamma)$$

$\Rightarrow E_{n,D}[(H_D(n) - E_D[H_D(n)])^2] + E_n[(E_D[H_D(n)] - \gamma)^2]$

Variance                          bias

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

The bias term, given some data point  $(n, y)$ , depends on  $E_D[H_D(n)]$  only.

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

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

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

$$= \mu$$

So the bias term does not change with ensembling.

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

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

$$\left\{ \begin{array}{l} \text{Var}(Kn) \\ = K^2 \text{Var}(n) \end{array} \right\} = \frac{1}{K^2} \text{Var}\left(\sum_{i=1}^K H_{D,i}(n)\right)$$

$$\left\{ \begin{array}{l} \text{Var}(n_1 + n_2 + \dots) \\ = \text{Var}(n_1) + \text{Var}(n_2) \end{array} \right\} = \frac{1}{K^2} \left( \sum_{i=1}^K \text{Var}(H_{D,i}(n)) \right)$$

$$\text{when } n_i \text{ iid} \left\{ \begin{array}{l} \\ \end{array} \right. = \frac{1}{K^2} \times KL\delta^2$$

$$= \frac{L}{K} \delta^2$$

Now, as the bias term does not change, we need to reduce the variance term to reduce

the loss.

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

$\downarrow$

Variance of  
single Model trained  
on Entire Data

Variance of ensemble model.

$$L < K$$

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

