

# Data Preparation and Analysis

## Final Exam

### Multiple Choice Questions

#### 1) For a smooth fit

#### PART-II

indications of population and signal such

1.  $\lambda$  is the tuning parameter which controls the effective degree of freedom of smoothing spline. As  $\lambda$  varies between 0 to infinity, the value of degree of freedom varies from  $n$  to 2.

To fit a smooth curve to a dataset, we need to find a function  $g(x)$  such that,

$$\text{RSS} = \sum_{i=1}^n (y_i - g(x_i))^2 \text{ is minimum}$$

One way to find such a smooth function is to minimize,

$$S(\lambda) = \sum_{i=1}^n (y_i - g(x_i))^2 + \lambda \int g'(t)^2 dt.$$

where  $\lambda$  is a nonnegative tuning parameter. The function that minimizes this is called as smoothing spline. The first part is a loss function and the second term is called penalty.

$\lambda \rightarrow$  larger the value of  $\lambda$ , smoother the  $g$  as well.

$\rightarrow$  when  $\lambda = 0$ , the given model will be very flexible and interpretable the training data.

$\rightarrow$  when  $\lambda \rightarrow \infty$ , the model corresponds to least squares linear regression.

$\lambda$  controls the effective degree of freedom (df $\lambda$ )

$\rightarrow$  As  $\lambda$  increases from 0 to  $\infty$ , the effective degree of freedom (df $\lambda$ ) decreases from  $n$  to 2.

where,  $n$  is number of parameters in smoothing spline and hence  $n$  is normal degree of freedom.

$\lambda$  can be chosen by cross-validation

The way RSS is calculated is

$$RSS(\lambda) = \sum_{i=1}^n (y_i - g_n^{(i)}(x_i))^2 = \sum_{i=1}^n \left[ \frac{y_i - \hat{g}_n(x_i)}{1 - s_n} \right]^2$$

The matrix  $S_n$  can be completed as

$$\hat{g}_n = S_n y$$

The effective degree of freedom for the smoothing spline is given as:

$$df_n = \sum_{i=1}^n (S_n)_{ii}$$

The maximum value that  $df_n$  can take is  $n$ , i.e., the number of parameters in smoothing spline.

The minimum value that  $df_n$  can take is 2 and it represents a linear model.

## PART-II

2. Slack variables that allows individual observations to be on the wrong side of the margin or the hyperplane.

- when  $\xi_i = 0$ ,  $i^{th}$  observation is on correct side of margin.

- when  $\xi_i > 0$ ,  $i^{th}$  observation is on wrong side of margin and we can say it has violated the margin.

- when  $\xi_i \geq 1$ ,  $i^{th}$  observation is on wrong side of hyperplane as budget  $c$  increase we become tolerant of violation to the violation to margin and so margin will be widen.

- when  $c$  is larger, margin allows more violation to it, so we can have many support vectors.

-  $c$  also amounts to fitting the less data and obtaining a classifier that is potentially more biased but lower variance.

## PART-II

and hence better knowledge of each  
and hence above all with more information

3. As per the question decision tree  $n=100$ ,

split -  $C_1, C_2$  which is prepared

100 observation in a node split evenly

between the 2 classes  $C_1$  and  $C_2$

So  $C_1$  has 50 observations and  $C_2$  has 50  
observation.

Gini coefficient =  $1 - \sum [p(j)]^2$  where  $p(j)$   
is the probability of getting class j

observation thus entropy value =  $1 - [(50/100)^2 + (50/100)^2]$

$$= 1 - (1/4 + 1/4) = 1/2 = 0.5$$

Entropy value =  $-\sum p(j) * \log_2(p(j))$

where  $p(j)$  is probability of getting class  
j observation thus = entropy value

$$= -(50/100) * \log_2(50/100) +$$

$$0 = [(50/100) * \log_2(50/100)] + [(50/100) * \log_2(50/100)]$$

$$= 0.5 + 0.5 = 1$$

Now if optimal split occurs into 2 leaf nodes then both the nodes will be pure i.e., both will have only observations belonging to single class only.

for lead node with class C<sub>1</sub>,

$$\text{thus gini coefficient} = 1 - (50/50)^2 - (0/50)^2 \\ = 1 - 1 - 0 \\ = 0 \quad (\text{Ans})$$

$$\text{Entropy value} = -[(50/50) \log_2(50/50) + (0/50) * \log_2(0/50)] \\ = -[\log_2(1) + 0 \log_2(0)] = 0. \quad (\text{Ans})$$

for leaf node with class C<sub>2</sub>,

thus gini coefficient =

$$1 - (0/50)^2 - (50/50)^2 \\ = 1 - 0 - 1 = 0 \quad (\text{Ans})$$

Entropy value

$$= -[(0/50) \log_2(0/50) + (50/50) * \log_2(50/50)] \\ = -[0 \log_2 0 + \log_2 1] = 0$$

## PART-II

4. Hierarchical clustering process with points  $P_1(4, 5)$  and  $P_2(6, 13)$ ,  $P_3(1, 1)$

Euclidean distance  $f(x, y), (a, b)] =$

$$\sqrt{(x-a)^2 + (y-b)^2}$$

$$\text{Distance } (P_1, P_2) = \sqrt{(4-6)^2 + (5-13)^2} = 8.246$$

$$(P_1, P_3) = \sqrt{(4-1)^2 + (5-1)^2} = 5$$

$$(P_2, P_3) = \sqrt{(6-1)^2 + (13-1)^2} = 13$$

we perform dissimilarity matrix

	P <sub>1</sub>	P <sub>2</sub>	P <sub>3</sub>
P <sub>1</sub>	0	8.246	5
P <sub>2</sub>	8.246	0	13
P <sub>3</sub>	5	13	0

So cluster 1 is formed between P<sub>1</sub> & P<sub>3</sub>. Distance

matrix for complete linkage:

$$1.5 \begin{bmatrix} 1.5 & 2 \\ 0 & 13 \end{bmatrix}$$

Stand by 2nd value 13 and 0 values as step of 1

$$\max [\text{dist}(P_1, P_3), P_2]$$

$$\text{dist}(P_1, P_2) [P_2, P_1]$$

$$\max [13, 8.24] \Rightarrow 13.5$$

Single linkage:  $\min[\text{dis}(P_1, P_2, P_3)] = \min[(P_1, P_2), (P_3, P_2)]$

$$\text{centroid for cluster} = \left( \frac{1+4}{2}, \frac{5+1}{2} \right)$$

$$(P_1, P_3) = (2.5, 3)$$

$$= \min(13, 8.24)$$

$$= 8.24$$

## PART-II

~~Explain forward and backward search methods~~

### 1.1 Maximal-Margin Classifier :

Maximum margin hyperplane is the farthest from the training observation which we compute the distance of each training observation from given hyperplane. The smallest distance is margin. Based on this we can classify the data, on which side it lies. The distance between a line and closest data points of the trained data is known as maximum-marginal classifier.

The maximal margin classifier a test observation  $x^*$  based on sign of

$$\text{Test sign } f(x^*) = \beta_0 + \beta_1 x_1^* + \dots + \beta_p x_p^*$$

where  $\beta_0, \beta_1, \dots, \beta_p$  are coefficient of maximal margin hyperplane. The maximal margin hyperplane represents mid-line of gap

between the two classes.

There are relatively few training observations which are constituted of  $p$ -dimensional vectors, are those that would cause the maximal hyperplane to move if they were moved in some dimension. The maximal margin hyperplane is the solution to the optimization problem of choosing  $\beta_0, \dots, \beta_p$  to maximize  $M$  such that:

$$\sum_{j=1}^p \beta_j^2 = 1$$

$$y_i(\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}) \geq M$$

- ensuring all margin at both having

Support vector classifier.

The margin in support vector machine is computed as the perpendicular distance from line to only closest points of trained data.

Support vector classifier separates points that are relevant in defining the line and in construction of classifier.

and to minimize the error function

and to minimize the error function

SVC is generalization of maximum margin classifier to non separable case is known as SVC. A hyperplane is line that splits the input variable space. The role of separating hyperplane is to accurately classify all the learning examples for splitting the line.

Optimal hyperplane is used in optimization estimation of enhance the generalized ability of support vector machine as shown

(i) Non linear machine learning produces classifiers that have a lower bias but higher variance.

Support vector machine algorithm has low bias and high variance. The bias can be increased by increased the parameters that affects the numbers of violations of margin that are allowed in data to be trained.

$$(1) \text{Def} \left( \text{Margin} = \sum_{j=1}^p \beta_j^2 \right) = 1 \quad (\beta_1, \beta_2, \dots, \beta_p)$$

$$= Y_i(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_ip) - M(1 - \varepsilon_i) \leq 0 \\ -(d - x) + (c - \varepsilon_i) \geq 0 \\ \sum_{i=1}^n \varepsilon_i \leq C$$

- The maximal marginal classifier depends on smaller set of points called support present in data.
- In SVC, effects of points on right side of plane is small, but effects of points on wrong side of plane is more. Hence they are called support vectors.
- SVC is robust in nature, because SVC works with inseparable class data and hyperplane depends upon support vectors.
- In SVC, when  $c$  is large, it will have low variance and high bias; when  $c$  is small, it will have biasing & high variance.

## PART-I

$$\text{Given } f(x) = (x-2)^2 + (x-3)^2, \text{ i.e. } f(x) = (x-2)(x-3)$$

1. 3-cut points in our data, we fit piece wise polynomial for regression -

→ The first constraint is that the fitted curve must be continuous at the knots.

→ Second constraint first derivative of the piece wise polynomials has to be continuous.

→ Third constraint second derivative of the piece wise polynomials has to be continuous.

→ To get a cubic spline with k-knots uses  $K+4$  degree of freedom

$$\Rightarrow 3 \text{ knots} \quad (P_1, P_2, P_3)$$

$$\text{Degrees of freedom} = 4 + 3 = 7.$$

$$\text{param} = (P_1, P_2, P_3)$$

$$P_1, P_2, P_3 =$$

$$(2, 1, 4) \text{ and } (3, 2, 1)$$

$$(3, 2, 1) \rightarrow (3, 1)$$

## PART - I

2.  $d = 250$  features

$N = 5,000,000$  observations.

see home,  $d = 250$  features

length of covariance matrix =  $250 \times 250$ .

number of eigen vectors =  $250 \times 250$

number of eigen values =  $250 \times 1$

Now, we take 10% of eigen vector and

values,

$\therefore$  dimensions of projected data,

sample matrix =  $25 \times 5,000,000$ .

$$(25 \times 5,000,000) =$$

## PART - I

4. \* Lasso performs shrinkage as ridge, but additional Lasso shrinks all up Predictors in final model exactly to zero.

\* Unlike ridge where it wont set the predictors to zero especially when ( $\lambda = 0$ )

\* Lasso has  $\beta_j$  coefficients as  $|\beta_j|$  in terms of penalty. Overall lasso performs variable selection. As a result, it is easier to interpret model generated in lasso compared to ridge. Hence, lasso yields sparse model, that is model that involve only a subset of variables.

5. When performing regression, the bagging we constructed  $B$  regressions trees using  $B$  bootstrapped training sets and finally we averaged resulting predictions to compute out prediction.

predictions for new sample observations.  
 But considering classification for a given observation, we record the class predicted by each of B trees and maximum, overall prediction is most commonly occurring class among B predictors of individual

6.  $P = 225$  predictors  
 number of splits offered by each tree is  $m = \sqrt{P} = 15$ .

In order to find no. of splits for given tree in stand or forest, will not consider strong predictors.

$$\therefore \text{On average} = \frac{(P-m)}{P}$$

$$= \frac{225-15}{225} = 0.9333\ldots$$

$\therefore$  On average, 0.93 of splits will not even consider by random forest. we can think of this process as decorrelating the trees, thereby making the average of resulting tree

Resulting trees less variable and  
hence more reliable.

f.  $d=5$  columns, variance = 100

eigen values  $\{35, 25, 20, 15, 5\}$

$\therefore$  here  $\lambda_1 = 35, \lambda_2 = 25, \lambda_3 = 20, \lambda_4 = 15$

$$\lambda_5 = 5.$$

also  $\sum_{i=1}^5 \lambda_i = \lambda_1 + \lambda_2 + \lambda_3 + \lambda_4 + \lambda_5 = 100$

so.  $\frac{\lambda_1}{\sum \lambda_i}$  explain the variance of 1<sup>st</sup> P.C

$\frac{\lambda_1 + \lambda_2}{\sum \lambda_i}$  explain variance of 1<sup>st</sup> & 2<sup>nd</sup> P.C

Since, we want to reduce dimension with 80%. So we need to find pairs where

sum is 80

$$\text{So, } \lambda_1 + \lambda_2 + \lambda_3 = 35 + 25 + 20 = 80$$

$$\lambda_1 + \lambda_2 + \lambda_4 + \lambda_5 = 35 + 25 + 15 + 5 = 80.$$

So we have pairs

$$[\lambda_1 + \lambda_2 + \lambda_3, \lambda_4, \lambda_5] \text{ and } [\lambda_1 + \lambda_2 + \lambda_4 + \lambda_5, \lambda_3]$$

$$\begin{bmatrix} 20, 15, 5 \end{bmatrix}_{1 \times 3} \text{ and } \begin{bmatrix} 20, 20 \end{bmatrix}_{1 \times 2}$$

Here 2 dimensions  
are reduced

here 3 dimensions  
and reduced

Now,

$$\sigma_1 = \left\{ \frac{1}{3} \left( 80 - \frac{100}{3} \right)^2 + \left( 15 - \frac{100}{3} \right)^2 + \left( 5 - \frac{100}{3} \right)^2 \right\}^{1/2}$$

$$\sigma_1 = \left[ \frac{1}{3} \left\{ \left( \frac{140}{3} \right)^2 + \left( \frac{55}{3} \right)^2 + \left( \frac{85}{3} \right)^2 \right\} \right]^{1/2}$$

$$2.0 \sigma_1 = \frac{1}{3\sqrt{3}} \left[ (140)^2 + (55)^2 + (85)^2 \right]^{1/2}$$

$$0.8 \sigma_1 = \frac{1}{3\sqrt{3}} [19600 + 3025 + 7225]^{1/2}$$

$$\text{Now } \sigma_1 = \frac{1}{3\sqrt{3}} [29850]^{1/2} = \underline{\underline{172.771}}$$

similarly finding 3rd and 4th dimension

$$\sigma_1 = \frac{57.520}{\sqrt{3}} = 33.24$$

Similarly  $\sigma_2 = \sqrt{60+20+20} = \sqrt{100} = 10$

$$\sigma_2 = \left\{ \frac{1}{3} \left[ \left( 80 - \frac{100}{2} \right)^2 + \left( 20 - \frac{100}{2} \right)^2 \right] \right\}^{1/2}$$

$$= \left[ \frac{1}{3} \left\{ \left( \frac{60}{2} \right)^2 + \left( \frac{60}{2} \right)^2 \right\} \right]^{1/2}$$

$$\sigma_2 = \frac{1}{2\sqrt{3}} [(60)^2 + (60)^2]^{1/2}$$

$$\sigma_2 = \frac{1}{2\sqrt{3}} [2(60)^2]^{1/2}$$

$$\text{and } \sigma_2 = \frac{\sqrt{2} \cdot 60}{2 \cdot \sqrt{3}}$$

$$\sigma_2 = \frac{30\sqrt{2}}{\sqrt{3}} = 24.434$$

So Analysis able to reduce 2 dim in  $[80, 15, 5]$

reduce 3 dim in  $[80, 20]$

Standard deviation of  $[80, 15, 5]$   $\sigma_1 = 33.24$

Standard deviation of  $[80, 20]$   $\sigma_2 = 24.494$

8. If  $K=2$ , what are the clusters?

$$C_1 = \{1, 2, 3, 4\}$$

$$C_2 = \{-9, -8, -7, -6\}$$

Centroids of clusters are

$$\begin{aligned} & \left\{ \frac{1+2+3+4}{4} \right\}, \left\{ \frac{-9-8-7-6}{4} \right\} \\ & = (2.5, -7.5) \end{aligned}$$

we have to check

$$2 \leq \sum_{i \in C_k} \sum_{j=1}^k (x_{ij} - \bar{x}_{nj})^2$$
 reaches

optimum values.

- Result obtained will depend on the initial (random) cluster assignment of each observation.

## Lucky 7

1. 2662
2. Kazhdan - Lusztig polynomial.
3. Knot Theory and Representation theory
4. Capsule neural network
5. Under 18 years of age
6. IBM
7. PYTHIA
8. IKEA Chair

### Part - I

3. First layer of CNN - 100 units

gray scale images -  $30 \times 30$  pixel

size of weight matrix  $w_1$ ,

(Image size)

$$N \times N = 30 \times 30.$$

F filter size

$$F \times F = 10 \times 10 \text{ (100 units)}$$

$$i/p = 30 \times 30 = 900.$$

$$\text{weight} = [N - F + 1] \times [N - F + 1]$$

$$= (30 - 10 + 1) \times (30 - 10 + 1)$$

$$= 21 \times 21$$

Rows columns .

~~Output~~