

Student Name: Prashant Kumar

Roll Number: 231110036

Date: September 15, 2023

The Given optimization problem is-

$$(\hat{w}_c, \hat{M}_c) = \operatorname{argmin}_{w_c, M_c} \sum_{x_n: y_n=c} \frac{1}{N_c} (x_n - w_c)^T M_c (x_n - w_c) - \log |M_c| \quad (1)$$

To find the optimal value of  $(\hat{w}_c)$  and  $(\hat{M}_c)$  (i.e. value on which loss will be minimum) I have to partially differentiate the loss function with respect to  $(w_c)$  and  $(M_c)$ , and equating the results with zero will give me the optimal values for  $(\hat{w}_c)$  and  $(\hat{M}_c)$  correspondingly.

First partially differentiating given equation w.r.t  $w_c$  and equating it to zero-

$$\frac{\partial}{\partial w_c} \left( \operatorname{argmin}_{w_c, M_c} \sum_{x_n: y_n=c} \frac{1}{N_c} (x_n - w_c)^T M_c (x_n - w_c) - \log |M_c| \right) = 0 \quad (2)$$

partially differentiating  $\log |M_c|$  w.r.t  $w_c$  will results in zero as  $\log |M_c|$  is a constant w.r.t  $w_c$ ,

$$\frac{-2 * M_c}{N_c} \sum_{x_n: y_n=c} (x_n - w_c) = 0 \quad (3)$$

Simplifying the expression,

$$\sum_{x_n: y_n=c} x_n - \sum_{x_n: y_n=c} w_c = 0 \quad (4)$$

Since there are  $N_c$  numbers of input from class C so this expression becomes-

$$\sum_{x_n: y_n=c} x_n - N_c * w_c = 0 \quad (5)$$

Further Simplifying,

$$w_c = \frac{1}{N_c} \sum_{x_n: y_n=c} x_n \quad (6)$$

Again partially differentiating given equation w.r.t  $M_c$  and equating it to zero-

$$\frac{\partial}{\partial M_c} \left( \sum_{x_n: y_n=c} \frac{1}{N_c} (x_n - w_c)^T M_c (x_n - w_c) - \log |M_c| \right) = 0 \quad (7)$$

Since  $\frac{\partial}{\partial M_c} \log |M_c| = (M_c^{-1})^T$  but positive definite matrix is symmetric so  $(M_c^{-1})^T = M_c^{-1}$ , hence after partially differentiating -

$$\frac{1}{N_c} \sum_{x_n: y_n=c} (x_n - w_c)^T (x_n - w_c) - M_c^{-1} = 0 \quad (8)$$

$$\frac{1}{N_c} \sum_{x_n: y_n=c} (x_n - w_c)^T (x_n - w_c) = M_c^{-1} \quad (9)$$

Pre-multiply by  $M_c$  on both sides and simplify to calculate for  $M_c$ ,

$$M_c * \sum_{x_n: y_n=c} (x_n - w_c)^T (x_n - w_c) = N_c * I_{d*d} \quad (10)$$

$$M_c * \sum_{x_n: y_n=c} \|x_n - w_c\|_2^2 = N_c * I_{d*d} \quad (11)$$

$$M_c = \frac{N_c * I_{d*d}}{\sum_{x_n: y_n=c} \|x_n - w_c\|_2^2} \quad (12)$$

**the second part of the question** when  $M_c$  is reduced to identity matrix then the given optimization expression becomes-

$$(\hat{w}_c, \hat{M}_c) = \operatorname{argmin}_{w_c, M_c} \sum_{x_n: y_n=c} \frac{1}{N_c} (x_n - w_c)^T I (x_n - w_c) - \log |I| \quad (13)$$

simplifying,

$$(\hat{w}_c) = \operatorname{argmin}_{w_c} \sum_{x_n: y_n=c} \frac{1}{N_c} ((x_n - w_c)^T (x_n - w_c)) \quad (14)$$

$$(\hat{w}_c) = \operatorname{argmin}_{w_c} \sum_{x_n: y_n=c} \frac{1}{N_c} \|x_n - w_c\|_2^2 \quad (15)$$

Hence, when  $M_c$  is reduced to the identity matrix then the problem becomes the minimization of the sum of  $l_2$  norm of each input (for a class C) w.r.t  $w_c$

---

One-nearest-neighbor algorithm predicts a class label based on the class label of the nearest training input from the test data point. Since it has access to infinite amounts of training data, in this case, the algorithm will always find the nearest training point having the same label as the test data so the algorithm will approach zero error which is the condition for consistency. Hence, **the one-nearest-neighbor algorithm will be consistent in this setting.**

---

In Decision tree construction for classification, we split on a feature that gives the highest information gain (high entropy reduction) on splitting. (i.e. a notion for deciding the purity of children nodes, moreover purity refers to the homogeneity of labels present in that node) similarly, for regression using a decision tree, we can **use standard deviation or variance reduction** for ensuring the purity of children nodes having real-valued labels.

**how standard deviation or variance quantify the homogeneity/diversity of the set of real-valued labels-**

standard deviation (i.e. the square root of the variance) is the measure of the dispersion of values, standard deviation will be zero for a set of completely homogeneous real-valued labels and will be high for completely diverse labels.

$$\text{standard deviation}(\sigma) = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2}$$

$$\text{variance} \sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

**Idea for building the decision tree using standard deviation/variance reduction**

On each step to decide the feature for splitting:

- (i) First I calculate the standard deviation for output labels on that node then split data using every feature that is present (remaining at the corresponding step) and again calculate the standard deviation for output labels.
- (ii) subtracting each resulting standard deviation from the standard deviation before the split will give me a standard deviation reduction. (kind of similar to the highest entropy reduction in classification)
- (iii) Now I will Split on the feature that is giving the highest standard deviation reduction.

Since the decision tree is a recursive algorithm, to implement a regression model using a decision tree and also to ensure that the model is not over-fitting I will stop splitting further when the standard deviation of output labels is below a predefined threshold (also stop when there is no more feature or very few data points in a node)

Student Name: Prashant Kumar

Roll Number: 231110036

Date: September 15, 2023

For the unregularized linear regression model the solution for  $\hat{w}$  is-

$$\hat{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (16)$$

and Prediction for an input  $x_*$  at test time is written as-

$$y_* = \hat{w}^\top x_* = x_*^\top \hat{w} \quad (17)$$

Substituting the value of  $\hat{w}$  in  $y_* = x_n^\top \hat{w}$ ,

$$y_* = x_*^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (18)$$

This expression can also be written as-

$$y_* = \sum_{n=1}^N (x_*^\top (x_n^\top x_n)^{-1} x_n^\top) y_n \quad (19)$$

OR

$$y_* = \sum_{n=1}^N \mathbf{W} y_n \quad (20)$$

The overall equation can be written as-

$$y_* = \sum_{n=1}^N w_n y_n \quad (21)$$

Here each  $w_n = x_*^\top (x_n^\top x_n)^{-1} x_n^\top$  or  $\mathbf{W} = x_*^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  and by the expression for the  $w_n$ , i can say that  $w_n$  is dependent on  $n^{th}$  training input and test input  $x_*$  (also note that to make a prediction  $y_*$  all the  $(w_1, w_2, \dots, w_N)$  vectors and all the  $(y_1, y_2, \dots, y_N)$  vectors are needed)

**The difference in this  $w_n$  and the weights in the weighted K nearest neighbors** is that  $w_n$  is calculated by  $n^{th}$  training input (and all the weights i.e.  $(w_1, w_2, \dots, w_N)$  needed for prediction) but in weighted KNN weight depends on only K nearest input and weights are assigned by the inverse of their distance from the test point. Also in KNN weights are normalized by dividing the sum and that is not the case in calculating  $w_n$ .

Student Name: Prashant Kumar

Roll Number: 231110036

Date: September 15, 2023

Given the new loss function(after masking inputs) is-

$$L(w) = \sum_{n=1}^N (y_n - \mathbf{w}^T \tilde{\mathbf{x}}_n)^2 \quad (22)$$

When  $\mathbf{x}$  is dropped s.t. any input dimension is retained with probability  $p$ , then the expected value of  $L(w)$ -

$$L(w) = E_{R \sim \text{Bernoulli}(N, p)}[L(w)] \quad (23)$$

Replacing  $\tilde{x}_n$  by  $x_n \odot m_n$  in  $L(w)$  and minimizing the objective,

$$\mathbf{L}(\mathbf{w}) = \text{argmin}_w \{E_{R \sim \text{Bernoulli}(n, p)}[\sum_{n=1}^N (y_n - \mathbf{w}^T (\mathbf{x}_n \odot \mathbf{m}_n))^2]\} \quad (24)$$

This can be written as-

$$L(w) = \|\mathbf{y} - (\mathbf{M} * \mathbf{X})\mathbf{w}\|^2 \quad (25)$$

We minimize this w.r.t.  $\mathbf{w}$ , so new objective function will be-

$$\mathbf{L}(\mathbf{w}) = \text{argmin}_w \{E_{R \sim \text{Bernoulli}(n, p)}[\|\mathbf{y} - (\mathbf{M} * \mathbf{X})\mathbf{w}\|^2]\} \quad (26)$$

$$\mathbf{L}(\mathbf{w}) = \text{argmin}_w \{E_{R \sim \text{Bernoulli}(n, p)}[(\mathbf{y} - (\mathbf{M} * \mathbf{X})\mathbf{w})^T (\mathbf{y} - (\mathbf{M} * \mathbf{X})\mathbf{w})]\} \quad (27)$$

Let  $\mathbf{b} = \mathbf{y} - (\mathbf{M} * \mathbf{X})\mathbf{w}$  and  $\mu = E[R \sim \text{Bernoulli}(n, p)]$ . Then, we have  $[\mathbf{b}] = \mathbf{y} - p\mathbf{X}\mathbf{w}$ .

$$\mathbf{L}(\mathbf{w}) = \arg \min_w (E[R \sim \text{Bernoulli}(n, p)] \|\mathbf{b}^T \mathbf{b}\|) \quad (28)$$

$$\mathbf{L}(\mathbf{w}) = \arg \min_w (\mu^T \mu + \text{TRACE}((\mathbf{b} - \mu)(\mathbf{b} - \mu)^T)) \quad (29)$$

$$\mathbf{L}(\mathbf{w}) = \arg \min_w ((\mathbf{y} - p\mathbf{X}\mathbf{w})^T (\mathbf{y} - p\mathbf{X}\mathbf{w}) + p(1-p)\text{TRACE}((\mathbf{X}\mathbf{w})(\mathbf{X}\mathbf{w})^T)) \quad (30)$$

$$\mathbf{L}(\mathbf{w}) = \arg \min_w (\|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + p(1-p)\text{TRACE}(\mathbf{X}\mathbf{w}\mathbf{w}^T \mathbf{X}^T)) \quad (31)$$

$$\mathbf{L}(\mathbf{w}) = \arg \min_w \left( \|\mathbf{y} - p\mathbf{X}\mathbf{w}\|^2 + p(1-p) \|\sqrt{\text{diag}(\mathbf{X}^T \mathbf{X} \mathbf{w})}\| \right) \quad (32)$$

Now when we compare this with ridge regression objective function

$$\mathbf{L}_{\text{ridge}}(\mathbf{w}) = \arg \min_w ((y - Xw)^T (y - Xw) + \lambda w^T w) \quad (33)$$

$$= \arg \min_w (\|y - \mathbf{X}\mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2) \quad (34)$$

By comparing this with the ridge regression objective we can say that  $\lambda \|\mathbf{w}\|^2$  is acting as a regularizer term.

**Introduction to ML (CS771), Autumn 2023**  
**Indian Institute of Technology Kanpur**  
**Homework Assignment Number 1**

*Student Name:* Prashant Kumar

*Roll Number:* 231110036

*Date:* September 15, 2023

**QUESTION**

**6**

---

A '**231110036.zip**' file is submitted with code and readme file.