

Student Name: Akash Shivaji Varude

Roll Number: 231110006

Date: September 15, 2023

We can find optimal values of \mathbf{w}_c and \mathbf{M}_c , by finding the first-order derivative of our objective function and equating it with zero.

Our cost function is

$$(\hat{\mathbf{w}}_c, \hat{\mathbf{M}}_c) = \arg \min_{\mathbf{w}_c, \mathbf{M}_c} \sum_{(\mathbf{x}_n: y_n=c)} \frac{1}{N_c} (\mathbf{x}_n - \mathbf{w}_c)^T \mathbf{M}_c (\mathbf{x}_n - \mathbf{w}_c) - \log |\mathbf{M}_c| \quad (1)$$

Partially differentiating above function with respect to \mathbf{w}_c we get,

$$\frac{\partial L}{\partial \mathbf{w}_c} = -\frac{2}{N_c} \sum_{(\mathbf{x}_n: y_n=c)} \mathbf{M}_c (\mathbf{x}_n - \mathbf{w}_c)$$

Equating it with zero to get optimal value of \mathbf{w}_c

$$-\frac{2\mathbf{M}_c}{N_c} \sum_{(\mathbf{x}_n: y_n=c)} (\mathbf{x}_n - \mathbf{w}_c) = 0$$

$$\sum_{(\mathbf{x}_n: y_n=c)} (\mathbf{x}_n - \mathbf{w}_c) = 0$$

$$\sum_{(\mathbf{x}_n: y_n=c)} \mathbf{x}_n - \sum_{(\mathbf{x}_n: y_n=c)} \mathbf{w}_c = 0$$

$$\sum_{(\mathbf{x}_n: y_n=c)} \mathbf{x}_n = \sum_{(\mathbf{x}_n: y_n=c)} \mathbf{w}_c$$

$$\sum_{(\mathbf{x}_n: y_n=c)} \mathbf{x}_n = \mathbf{w}_c N_c$$

$$\mathbf{w}_c = \frac{1}{N_c} \sum_{(\mathbf{x}_n: y_n=c)} \mathbf{x}_n \quad (2)$$

Partially differentiating equation (1) with respect to \mathbf{M}_c we get,

$$\frac{\partial L}{\partial \mathbf{M}_c} = \frac{1}{N_c} \sum_{(\mathbf{x}_n: y_n=c)} (\mathbf{x}_n - \mathbf{w}_c)^T (\mathbf{x}_n - \mathbf{w}_c) - \mathbf{M}_c^{-T}$$

Equating it with zero to get optimal value of \mathbf{M}_c

$$\frac{1}{N_c} \left(\sum_{(\mathbf{x}_n: y_n=c)} (\mathbf{x}_n - \mathbf{w}_c)^T (\mathbf{x}_n - \mathbf{w}_c) \right) - \mathbf{M}_c^{-T} = 0$$

$$\frac{1}{N_c} \sum_{(\mathbf{x}_n: y_n=c)} (\mathbf{x}_n - \mathbf{w}_c)^T (\mathbf{x}_n - \mathbf{w}_c) = \mathbf{M}_c^{-T}$$

$$\mathbf{M}_c = N_c \sum_{(\mathbf{x}_n: y_n=c)} [(\mathbf{x}_n - \mathbf{w}_c) \cdot (\mathbf{x}_n - \mathbf{w}_c)^T]^{-1} \quad (3)$$

Special Case: (When \mathbf{M}_c is an Identity matrix)

In this case our cost function becomes

$$L_{\hat{\mathbf{w}}_c} = \frac{1}{N_c} \sum_{(\mathbf{x}_n: y_n=c)} \|\mathbf{x}_n - \mathbf{w}_c\|^2 \quad (4)$$

This means we are minimizing squared distance between data points and class specific \mathbf{w}_c . Hence optimal value of \mathbf{w}_c remains same i.e.

$$\mathbf{w}_c = \frac{1}{N_c} \sum_{(\mathbf{x}_n: y_n=c)} \mathbf{x}_n$$

and as \mathbf{M}_c is Identity matrix, it is fixed and it doesn't require optimization

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

Student Name: Akash Shivaji Varude

Roll Number: 231110006

Date: September 15, 2023

QUESTION

2

Yes, the one-nearest-neighbor (1-NN) algorithm will be consistent in the noise-free setting where all training inputs are correctly labeled.

1NN assigns the label based on nearest neighbour in training data. In a noise-free environment, where all training inputs are accurately labeled, the nearest neighbor will share the same class label as the input being classified. Hence, as the training dataset size increases infinitely, 1-NN will consistently make zero errors, converging to the Bayes optimal error rate of zero.

Student Name: Akash Shivaji Varude

Roll Number: 231110006

Date: September 15, 2023

If the decision trees are used for regression problem, 'Reduction in Variance' would be great criteria for choosing a feature to split nodes. So variance can be measure of homogeneity. If a node is entirely homogeneous, then the variance is zero.

At first we will calculate variance of target variable in each subset after the split. We will choose that feature which results in the greatest reduction in variance (i.e., the difference between the variance before and after the split).

Formula for Reduction in Variance given by:

$$RV = V(D) - \left(\frac{|D|}{|D_L|} \cdot V(D_L) + \frac{|D|}{|D_R|} \cdot V(D_R) \right)$$

where,

$V(D)$: Variance of the target variable in the current group of data

$|D_L|$: Number of data points that end up in the left child node after the split

$|D_R|$: Number of data points that end up in the right child node after the split

$V(D_L)$: Variance of the target variable in the left child node after the split

$V(D_R)$: Variance of the target variable in the right child node after the split

So, we find the split that will maximize this reduction in variance as it will result in more homogeneous groups. Hence tree will be better at predicting values within each group.

Student Name: Akash Shivaji Varude

Roll Number: 231110006

Date: September 15, 2023

In unregularized linear regression, the prediction at a test input \mathbf{x}_* is given by, $f(\mathbf{x}_*) = \hat{\mathbf{w}}^\top \mathbf{x}_*$ where $\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$.

Substituting this expression of $\hat{\mathbf{w}}$ into the prediction formula we get,

$$f(\mathbf{x}_*) = \left((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \right)^\top \mathbf{x}_*$$

$$f(\mathbf{x}_*) = (\mathbf{X}^\top \mathbf{y})^\top ((\mathbf{X}^\top \mathbf{X})^{-1})^\top \mathbf{x}_*$$

$$f(\mathbf{x}_*) = (\mathbf{X}^\top \mathbf{y})^\top ((\mathbf{X}^\top \mathbf{X})^\top)^{-1} \mathbf{x}_*$$

as $\mathbf{X}^\top \mathbf{X}$ is symmetric we can write

$$f(\mathbf{x}_*) = (\mathbf{X}^\top \mathbf{y})^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_*$$

$$f(\mathbf{x}_*) = \mathbf{y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_*$$

$$f(\mathbf{x}_*) = \mathbf{y}^\top \hat{\mathbf{w}}_{\text{new}}$$

where $\hat{\mathbf{w}}_{\text{new}} = \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_*$ and it is $N \times 1$ vector.

as we can see, \mathbf{y}^\top is $1 \times N$ vector and $\hat{\mathbf{w}}_{\text{new}}$ is $N \times 1$ vector, the expression $\mathbf{y}^\top \hat{\mathbf{w}}_{\text{new}}$ represents their inner product. Hence we can write it as

$$f(\mathbf{x}^*) = \sum_{n=1}^N w_n y_n$$

Discussion about weights:

In our problem, each weight w_n represents importance of each feature for modeling linear relationship between features and the training responses. So, value w_n is higher for more important feature than others. While in KNN, each weight w_n emphasize the importance of nearby training examples, giving more importance to nearby training example. hence value of w_n for training example nearer to \mathbf{x}_* will be more than farther training example.

New loss function is defined as:

$$L_{\text{new}} = \sum_{n=1}^N (y_n - \mathbf{w}^\top \tilde{\mathbf{x}}_n)^2 \quad (1)$$

Expectation of this function with respect to the random mask vectors \mathbf{m}_n is:

$$E[L_{\text{new}}] = E \left[\sum_{n=1}^N (y_n - \mathbf{w}^\top \tilde{\mathbf{x}}_n)^2 \right]$$

$$E[L_{\text{new}}] = \sum_{n=1}^N \left[E(y_n - \mathbf{w}^\top \tilde{\mathbf{x}}_n)^2 \right]$$

$$E[L_{\text{new}}] = \sum_{n=1}^N \left[E(y_n - \mathbf{w}^\top (\mathbf{x}_n \circ \mathbf{m}_n))^2 \right]$$

$$E[L_{\text{new}}] = \sum_{n=1}^N \left[E \left(y_n^2 - 2y_n(\mathbf{w}^\top (\mathbf{x}_n \circ \mathbf{m}_n)) + (\mathbf{w}^\top (\mathbf{x}_n \circ \mathbf{m}_n))^2 \right) \right]$$

$$E[L_{\text{new}}] = \sum_{n=1}^N \left[y_n^2 - 2y_n \mathbf{w}^\top E[(\mathbf{x}_n \circ \mathbf{m}_n)] + (\mathbf{w}^\top)^2 E[(\mathbf{x}_n \circ \mathbf{m}_n)^2] \right] \quad (2)$$

Now, Each element of $\mathbf{x}_n \circ \mathbf{m}_n$ is either 0 with probability $(1 - p)$ or 1 with probability p . Hence every element of \mathbf{m}_n also follows a Bernoulli distribution with probability p . Hence

$$E[(\mathbf{x}_n \circ \mathbf{m}_n)] = \mathbf{x}_n p \quad (3)$$

Also, for $E[(\mathbf{x}_n \circ \mathbf{m}_n)^2]$, elements of $\mathbf{x}_n \circ \mathbf{m}_n$ are either 0 or \mathbf{x}_n . Hence,

$$E[(\mathbf{x}_n \circ \mathbf{m}_n)^2] = E[(\mathbf{x}_n \circ \mathbf{m}_n)(\mathbf{x}_n \circ \mathbf{m}_n)]$$

$$E[(\mathbf{x}_n \circ \mathbf{m}_n)^2] = E[(\mathbf{x}_n^2 \circ \mathbf{m}_n^2)] \quad (4)$$

Hence substituting (3), (4) in equation (2) we get,

$$\begin{aligned}
E[L_{\text{new}}] &= \sum_{n=1}^N \left[y_n^2 - 2y_n \mathbf{w}^\top \mathbf{x}_n p + (\mathbf{w}^\top)^2 E[(\mathbf{x}_n^2 \circ \mathbf{m}_n^2)] \right] \\
E[L_{\text{new}}] &= \sum_{n=1}^N \left[y_n^2 - 2py_n \mathbf{w}^\top \mathbf{x}_n + (\mathbf{w}^\top)^2 E[(\mathbf{x}_n^2 \circ \mathbf{m}_n^2)] \right] \tag{5}
\end{aligned}$$

So for converting it in regularized form, we can see, y_n^2 represents square of target values. So it is constant with respect to model weights w . $2py_n \mathbf{w}^\top \mathbf{x}_n$ this term represents the contribution of the data to expected loss. Because $\mathbf{w}^\top \mathbf{x}_n$ is predicted output and y_n is actual output. So, it's the mean squared error (MSE) term. While the term, $(\mathbf{w}^\top)^2 E[(\mathbf{x}_n^2 \circ \mathbf{m}_n^2)]$ penalizes large values of the model's weights \mathbf{w} using L2 regularization.

Hence we can write equation (5) in regularized form as,

$$E[L_{\text{new}}] = \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \lambda \sum_{d=1}^D \mathbf{w}_d^2 \tag{6}$$

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

Student Name: Akash Shivaji Varude

Roll Number: 231110006

Date: September 15, 2023

QUESTION

6

Method 1(Convex):

Test Accuracy = 46.89320388349515

Method 2(Regress):

for $\lambda=[0.01, 0.1, 1, 10, 20, 50, 100]$

Corresponding Test Accuracies are

For $\lambda=0.01$ Test accuracy = 58.090614886731395 %

For $\lambda=0.1$ Test accuracy = 59.54692556634305 %

For $\lambda=1$ Test accuracy = 67.39482200647248 %

For $\lambda=10$ Test accuracy = 73.28478964401295 %

For $\lambda=20$ Test accuracy = 71.68284789644012 %

For $\lambda=50$ Test accuracy = 65.08090614886731 %

For $\lambda=100$ Test accuracy = 56.47249190938511 %