

Student Name: Shrilakshmi S K

Roll Number: 211012

Date: September 15, 2023

The given optimization function for class c is

$$L_c(\mathbf{w}_c, \mathbf{M}_c) = \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} (\mathbf{x}_n - \mathbf{w}_c)^\top \mathbf{M}_c (\mathbf{x}_n - \mathbf{w}_c) - \log |\mathbf{M}_c|$$

Since \mathbf{M}_c is positive definite matrix, Hessian of $(\mathbf{x}_n - \mathbf{w}_c)^\top \mathbf{M}_c (\mathbf{x}_n - \mathbf{w}_c)$ will be positive definite, which is a sufficient condition for strict convexity. Hence, we will use first order optimality on the objective function to obtain closed form optimal values of \mathbf{w}_c and \mathbf{M}_c .

1. Finding $\hat{\mathbf{w}}_c$

Taking the gradient of $L_c(\mathbf{w}_c, \mathbf{M}_c)$ with respect to \mathbf{w}_c and equating it to zero will give us $\hat{\mathbf{w}}_c$.

$$\frac{\partial L_c(\mathbf{w}_c, \mathbf{M}_c)}{\partial \mathbf{w}_c} = \frac{\partial}{\partial \mathbf{w}_c} \left(\sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} (\mathbf{x}_n - \mathbf{w}_c)^\top \mathbf{M}_c (\mathbf{x}_n - \mathbf{w}_c) - \log |\mathbf{M}_c| \right)$$

Using $\frac{\partial}{\partial \mathbf{s}} (\mathbf{x} - \mathbf{s})^\top \mathbf{A} (\mathbf{x} - \mathbf{s}) = -(\mathbf{A} + \mathbf{A}^\top)(\mathbf{x} - \mathbf{s})$,

$$\begin{aligned} \frac{\partial L_c(\mathbf{w}_c, \mathbf{M}_c)}{\partial \mathbf{w}_c} &= \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} \frac{\partial}{\partial \mathbf{w}_c} \left((\mathbf{x}_n - \mathbf{w}_c)^\top \mathbf{M}_c (\mathbf{x}_n - \mathbf{w}_c) \right) \\ \frac{\partial L_c(\mathbf{w}_c, \mathbf{M}_c)}{\partial \mathbf{w}_c} &= - \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} \left((\mathbf{M}_c + \mathbf{M}_c^\top)(\mathbf{x}_n - \mathbf{w}_c) \right) \end{aligned}$$

Now let us equate $\frac{\partial L_c(\mathbf{w}_c, \mathbf{M}_c)}{\partial \mathbf{w}_c}$ to 0

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

(\mathbf{M}_c is positive definite $\iff \mathbf{M}_c^\top$ is positive definite) $\implies \mathbf{M}_c + \mathbf{M}_c^\top$ is positive definite.
 Since $\mathbf{M}_c + \mathbf{M}_c^\top$ is positive definite, the solution of the above equation will be:

$$\begin{aligned} \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} (\mathbf{x}_n - \mathbf{w}_c) &= 0 \\ \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} \mathbf{w}_c &= \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} \mathbf{x}_n \\ \frac{N_c}{N_c} \mathbf{w}_c &= \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} \mathbf{x}_n \\ \hat{\mathbf{w}}_c &= \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} \mathbf{x}_n \end{aligned}$$

2. Finding $\hat{\mathbf{M}}_c$

Taking the gradient of $L_c(\mathbf{w}_c, \mathbf{M}_c)$ with respect to \mathbf{M}_c and equating it to zero will give us $\hat{\mathbf{M}}_c$.

$$\begin{aligned}\frac{\partial L_c(\mathbf{w}_c, \mathbf{M}_c)}{\partial \mathbf{M}_c} &= \frac{\partial}{\partial \mathbf{M}_c} \left(\sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} (\mathbf{x}_n - \mathbf{w}_c)^\top \mathbf{M}_c (\mathbf{x}_n - \mathbf{w}_c) - \log |\mathbf{M}_c| \right) \\ \frac{\partial L_c(\mathbf{w}_c, \mathbf{M}_c)}{\partial \mathbf{M}_c} &= \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} \frac{\partial}{\partial \mathbf{M}_c} \left((\mathbf{x}_n - \mathbf{w}_c)^\top \mathbf{M}_c (\mathbf{x}_n - \mathbf{w}_c) \right) - \frac{\partial}{\partial \mathbf{M}_c} \log |\mathbf{M}_c|\end{aligned}$$

Using $\frac{\partial}{\partial \mathbf{A}} (\mathbf{x}^\top \mathbf{A} \mathbf{x}) = \mathbf{x} \mathbf{x}^\top$

$$\begin{aligned}\frac{\partial L_c(\mathbf{w}_c, \mathbf{M}_c)}{\partial \mathbf{M}_c} &= \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} (\mathbf{x}_n - \mathbf{w}_c) (\mathbf{x}_n - \mathbf{w}_c)^\top - \frac{\partial}{\partial \mathbf{M}_c} \log |\mathbf{M}_c| \\ \frac{\partial L_c(\mathbf{w}_c, \mathbf{M}_c)}{\partial \mathbf{M}_c} &= \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} (\mathbf{x}_n - \mathbf{w}_c) (\mathbf{x}_n - \mathbf{w}_c)^\top - \frac{1}{|\mathbf{M}_c|} \frac{\partial}{\partial \mathbf{M}_c} |\mathbf{M}_c|\end{aligned}$$

Using $\frac{\partial}{\partial \mathbf{A}} |\mathbf{A}| = |\mathbf{A}| (\mathbf{A}^{-1})^\top$

$$\begin{aligned}\frac{\partial L_c(\mathbf{w}_c, \mathbf{M}_c)}{\partial \mathbf{M}_c} &= \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} (\mathbf{x}_n - \mathbf{w}_c) (\mathbf{x}_n - \mathbf{w}_c)^\top - \frac{1}{|\mathbf{M}_c|} |\mathbf{M}_c| (\mathbf{M}_c^{-1})^\top \\ \frac{\partial L_c(\mathbf{w}_c, \mathbf{M}_c)}{\partial \mathbf{M}_c} &= \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} (\mathbf{x}_n - \mathbf{w}_c) (\mathbf{x}_n - \mathbf{w}_c)^\top - (\mathbf{M}_c^{-1})^\top\end{aligned}$$

Now let us equate $\frac{\partial L_c(\mathbf{w}_c, \mathbf{M}_c)}{\partial \mathbf{M}_c}$ to 0

$$\begin{aligned}\sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} (\mathbf{x}_n - \mathbf{w}_c) (\mathbf{x}_n - \mathbf{w}_c)^\top - (\mathbf{M}_c^{-1})^\top &= 0 \\ (\mathbf{M}_c^{-1})^\top &= \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} (\mathbf{x}_n - \mathbf{w}_c) (\mathbf{x}_n - \mathbf{w}_c)^\top \\ \mathbf{M}_c^{-1} &= \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} (\mathbf{x}_n - \mathbf{w}_c) (\mathbf{x}_n - \mathbf{w}_c)^\top \\ \mathbf{M}_c^{-1} &= \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} (\mathbf{x}_n - \mathbf{w}_c) (\mathbf{x}_n - \mathbf{w}_c)^\top \\ \hat{\mathbf{M}}_c &= \left(\sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} (\mathbf{x}_n - \hat{\mathbf{w}}_c) (\mathbf{x}_n - \hat{\mathbf{w}}_c)^\top \right)^{-1}\end{aligned}$$

Substituting $\hat{\mathbf{w}}_c = \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} \mathbf{x}_n$, we get

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

3. If \mathbf{M}_c is an identity matrix

$$\mathbf{M}_c = \mathbf{I}$$

$$L_c(\mathbf{w}_c, \mathbf{M}_c) = \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} (\mathbf{x}_n - \mathbf{w}_c)^\top \mathbf{M}_c (\mathbf{x}_n - \mathbf{w}_c) - \log |\mathbf{M}_c|$$

$$L_c(\mathbf{w}_c, \mathbf{M}_c = \mathbf{I}) = \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} (\mathbf{x}_n - \mathbf{w}_c)^\top \mathbf{I} (\mathbf{x}_n - \mathbf{w}_c) - \log |\mathbf{I}|$$

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

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

$$L_c(\mathbf{w}_c, \mathbf{M}_c = \mathbf{I}) = \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} (\|\mathbf{x}_n - \mathbf{w}_c\|)^2$$

Substituting optimal value of \mathbf{w}_c ,

$$L_c(\mathbf{w}_c = \hat{\mathbf{w}}_c, \mathbf{M}_c = \mathbf{I}) = \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} \left(\|\mathbf{x}_n - \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} \mathbf{x}_n\| \right)^2$$

$$L_c(\mathbf{w}_c = \hat{\mathbf{w}}_c, \mathbf{M}_c = \mathbf{I}) = \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} \|\mathbf{x}_n - \boldsymbol{\mu}_c\|^2$$

where $\boldsymbol{\mu}_c = \mathbf{w}_c = \sum_{\mathbf{x}_n: y_n=c} \frac{1}{N_c} \mathbf{x}_n$. Here $\boldsymbol{\mu}_c$ is centroid of class c inputs

In this special case, optimizing this loss function is equivalent to solving **LwP** as both measure the Euclidean distances from a centroid of the class. $\hat{\mathbf{w}}_c$ will not change in this case.

Student Name: Shrilakshmi S K

Roll Number: 211012

Date: September 15, 2023

In the given setting, one-nearest-neighbour algorithm will be **consistent**.

One-nearest-neighbour classifier has an error rate of at most twice of Bayes optimal error rate.

$$E_{Bayes} \leq E_{1NN} \leq 2E_{Bayes}(1 - E_{Bayes}) \leq 2E_{Bayes}$$

- Cover and Hart (1967)

As number of training examples approach infinity and since they are labelled correctly in the given setting, both E_{Bayes} and $2E_{Bayes}$ approach 0. This implies that E_{1NN} approaches 0, hence it is consistent in this setting where noise is absent. Note that it is not necessarily consistent when the noise is present.

Since there are infinite training examples, for every test input, we can always find a training example labelled correctly in close proximity of test input. Therefore, we can classify the test inputs with a minimal error rate. This error rate tends to 0 as the number of training examples approach infinity because test input will tend to coincide with its nearest training input.

Student Name: Shrilakshmi S K

Roll Number: 211012

Date: September 15, 2023

In Decision Tree Classifier, we decided the features for nodes based on the Information Gain. Information Gain was defined as difference in entropy before split and after the split of that node. The entropy captured the homogeneity/diversity of the labels. If the split was such that more homogenous labels were grouped together, the purity was high and the entropy was low. If the split was such that more diverse labels were grouped together, the purity was low and the entropy was high.

Extending the analogy of entropy to regression, we want a metric that can measure the homogeneity/diversity of a set of real valued labels. One such metric could be **variance**. If the variance of the labels grouped together is low, then their homogeneity/purity is high. If the variance of the labels grouped together is high, then their diversity is high and purity is low.

Let the variance of labels of set S be $Var(S)$. Suppose we assign a feature F to a node and it splits S into smaller disjoint sets S_1, S_2, \dots, S_k . Here k is the total number of distinct values taken by feature F in training data inputs and S_i is the set of training data inputs whose feature F takes i^{th} distinct value.

Similar to Information Gain in classification, we can define Information Gain in regression as

$$IG_{Regression}(F) = Var(S) - \sum_{i=1}^k \frac{|S_i|}{|S|} Var(S_i)$$

The feature F which gives highest $IG_{Regression}(F)$ will be chosen as a feature to split the data for the node at that level.

Student Name: Shrilakshmi S K

Roll Number: 211012

Date: September 15, 2023

Given a closed form solution $\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ for an unregularized linear regression, the predicted output $y_* = f(\mathbf{x}_*)$ for test input \mathbf{x}_* can be obtained by:

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

$$f(\mathbf{x}_*) = \langle \hat{\mathbf{w}}, \mathbf{x}_* \rangle = \langle \mathbf{x}_*, \hat{\mathbf{w}} \rangle$$

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

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

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

$$f(\mathbf{x}_*) = \mathbf{w}_* \mathbf{y}$$

Note that \mathbf{w}_* is an $1 \times N$ matrix as the dimensions of \mathbf{x}_*^\top is $1 \times D$, $(\mathbf{X}^\top \mathbf{X})^{-1}$ is $D \times D$ and \mathbf{X}^\top is $D \times N$. Dimension of \mathbf{y} is $N \times 1$. Let $\mathbf{y} = [y_1, y_2 \dots y_N]^\top$ and $\mathbf{w}_* = [w_1, w_2 \dots w_N]$. Then we can write

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

w_n in Linear Regression

$$w_n = \langle \mathbf{x}_*^\top, ((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)_n \rangle$$

where $((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top)_n$ is n^{th} column of $(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$

w_n measures the cosine similarity between the test input and n^{th} example in training data input with the weighted matrix $(\mathbf{X}^\top \mathbf{X})^{-1}$

w_n in KNN

$$w_n = \frac{1}{d(\mathbf{x}_*, \mathbf{x}_n)} = \frac{1}{\sqrt{(\mathbf{x}_* - \mathbf{x}_n)^\top (\mathbf{x}_* - \mathbf{x}_n)}}$$

w_n in KNN too measures the similarity between the test input and n^{th} example in training data input, but this similarity is different than that of Linear Regression. Here it measures similarity based on Euclidean Distance.

Student Name: Shrilakshmi S K

Roll Number: 211012

Date: September 15, 2023

The new loss function is

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

where $\tilde{\mathbf{x}}_n = \mathbf{x}_n \circ \mathbf{m}_n$

Calculating Expectation of \mathbf{m}_{nd} using the property of Bernoulli variable,

$$E(m_{nd}) = p$$

Calculating Expectation of \tilde{x}_{nd} ,

$$\tilde{x}_{nd} = x_{nd} \cdot m_{nd}$$

$$E(\tilde{x}_{nd}) = E(x_{nd} \cdot m_{nd}) = x_{nd} \cdot E(m_{nd}) = px_{nd}$$

Calculating Expectation of loss function,

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

By linearity of Expectation,

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

$$E(L(\mathbf{w})) = \sum_{n=1}^N E(y_n^2 - 2y_n \mathbf{w}^\top \tilde{\mathbf{x}}_n + (\mathbf{w}^\top \tilde{\mathbf{x}}_n)^2)$$

$$E(L(\mathbf{w})) = \sum_{n=1}^N (E(y_n^2) - 2E(y_n \mathbf{w}^\top \tilde{\mathbf{x}}_n) + E((\mathbf{w}^\top \tilde{\mathbf{x}}_n)^2))$$

$$E(L(\mathbf{w})) = \sum_{n=1}^N (y_n^2 - 2y_n E(\mathbf{w}^\top \tilde{\mathbf{x}}_n) + E((\mathbf{w}^\top \tilde{\mathbf{x}}_n)^2))$$

Calculating Expectation of $\mathbf{w}^\top \tilde{\mathbf{x}}_n$,

$$E(\mathbf{w}^\top \tilde{\mathbf{x}}_n) = E\left(\sum_{d=1}^D w_d \tilde{x}_{nd}\right) = \sum_{d=1}^D E(w_d \tilde{x}_{nd}) = \sum_{d=1}^D w_d E(\tilde{x}_{nd}) = \sum_{d=1}^D p w_d x_{nd} = p \sum_{d=1}^D w_d x_{nd} = p \mathbf{w}^\top \mathbf{x}_n$$

Calculating Expectation of $(\mathbf{w}^\top \tilde{\mathbf{x}}_n)^2$,

$$\text{Var}(z) = E(z^2) - (E(z))^2$$

$$E(z^2) = Var(z) + (E(z))^2$$

$$E((\mathbf{w}^\top \tilde{\mathbf{x}}_n)^2) = Var(\mathbf{w}^\top \tilde{\mathbf{x}}_n) + (E(\mathbf{w}^\top \tilde{\mathbf{x}}_n))^2$$

Finding $Var(\mathbf{w}^\top \tilde{\mathbf{x}}_n)$,

$$Var(\mathbf{w}^\top \tilde{\mathbf{x}}_n) = Var\left(\sum_{d=1}^D w_d \tilde{x}_{nd}\right)$$

Dimensions are independent (covariance is equal to 0). So linearity of variance holds.

$$Var(\mathbf{w}^\top \tilde{\mathbf{x}}_n) = \sum_{d=1}^D Var(w_d \tilde{x}_{nd}) = \sum_{d=1}^D Var(w_d x_{nd} m_{nd}) = \sum_{d=1}^D w_d^2 x_{nd}^2 Var(m_{nd})$$

Variance of Bernoulli variable is $p(1-p)$

$$Var(\mathbf{w}^\top \tilde{\mathbf{x}}_n) = p(1-p) \sum_{d=1}^D w_d^2 x_{nd}^2$$

Substituting,

$$E((\mathbf{w}^\top \tilde{\mathbf{x}}_n)^2) = Var(\mathbf{w}^\top \tilde{\mathbf{x}}_n) + (E(\mathbf{w}^\top \tilde{\mathbf{x}}_n))^2$$

$$E((\mathbf{w}^\top \tilde{\mathbf{x}}_n)^2) = p(1-p) \sum_{d=1}^D w_d^2 x_{nd}^2 + (p \mathbf{w}^\top \mathbf{x}_n)^2$$

Substituting in $E(L(\mathbf{w}))$,

$$E(L(\mathbf{w})) = \sum_{n=1}^N (y_n^2 - 2y_n E(\mathbf{w}^\top \tilde{\mathbf{x}}_n) + E((\mathbf{w}^\top \tilde{\mathbf{x}}_n)^2))$$

$$E(L(\mathbf{w})) = \sum_{n=1}^N (y_n^2 - 2y_n p \mathbf{w}^\top \mathbf{x}_n + p(1-p) \sum_{d=1}^D w_d^2 x_{nd}^2 + (p \mathbf{w}^\top \mathbf{x}_n)^2)$$

$$E(L(\mathbf{w})) = \sum_{n=1}^N \left((y_n - p \mathbf{w}^\top \mathbf{x}_n)^2 + p(1-p) \sum_{d=1}^D w_d^2 x_{nd}^2 \right)$$

$$E(L(\mathbf{w})) = \sum_{n=1}^N (y_n - p \mathbf{w}^\top \mathbf{x}_n)^2 + p(1-p) \sum_{n=1}^N \sum_{d=1}^D w_d^2 x_{nd}^2$$

$$E(L(\mathbf{w})) = L'(\mathbf{w}) = L'_{unreg}(\mathbf{w}) + L'_{reg}(\mathbf{w})$$

Here, we can see that expectation value of the given loss function is in the form of a regularized loss function. So minimizing this given loss function would be equivalent to minimizing the regularized loss function.

The new unregularized loss function is

$$L'_{unreg}(\mathbf{w}) = \sum_{n=1}^N (y_n - p \mathbf{w}^\top \mathbf{x}_n)^2$$

And the regularizer function is

$$L'_{reg}(\mathbf{w}) = p(1-p) \sum_{n=1}^N \sum_{d=1}^D w_d^2 x_{nd}^2$$

$$L'_{reg}(\mathbf{w}) = \mathbf{w}^\top (p(1-p) \sum_{n=1}^N \text{diag}(\mathbf{x}_n \mathbf{x}_n^\top)) \mathbf{w}$$

Let $\lambda = p(1-p) \sum_{n=1}^N \text{diag}(\mathbf{x}_n \mathbf{x}_n^\top)$

$$L'_{reg}(\mathbf{w}) = \lambda \|\mathbf{w}\|^2$$

$$L'(\mathbf{w}) = \sum_{n=1}^N (y_n - p \mathbf{w}^\top \mathbf{x}_n)^2 + \lambda \|\mathbf{w}\|^2$$

where $\lambda = p(1-p) \sum_{n=1}^N \text{diag}(\mathbf{x}_n \mathbf{x}_n^\top)$

Let $b \sim \text{Bernoulli}(p)$

$$E(b) = 1 \cdot p + 0 \cdot (1-p) = p$$

$$E(b^2) = 1^2 \cdot p + 0^2 \cdot (1-p) = p$$

$$\text{Var}(b) = E(b^2) - (E(b))^2 = p - p^2 = p(1-p)$$

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

QUESTION

6

Student Name: Shrilakshmi S K

Roll Number: 211012

Date: September 15, 2023

Method 1

Classification Accuracy on Test Set = **46.89320388349515 %**

Method 2

Classification Accuracy on Test Set for Different Lambdas

λ	Accuracy
0.01	58.090614886731395 %
0.1	59.54692556634305 %
1	67.39482200647248 %
10	73.28478964401295 %
20	71.68284789644012 %
50	65.08090614886731 %
100	56.47249190938511 %

Maximum Accuracy (73.28478964401295 %) is given by $\lambda = 10$