

CS771A- Homework Assignment 1

Jaivardhan Kapoor - 150300

August 29, 2016

Question 1

Disance from means

Given training data $\{x_n, y_n\}_{n=1}^N$ for a classification problem with $y_n \in \{-1, +1\}$, the goal is to find $f(\mathbf{x})$ such that $y = \text{sign}[f(\mathbf{x})]$ and $f(\mathbf{x})$ can be written as $\sum_{n=1}^N \langle \mathbf{x}_n, \mathbf{x} \rangle + b$ where \mathbf{x} is the test example feature vector.

Let N_+ be the class with $y = +1$ and N_- be the class with $y = -1$. Then,

$$\boldsymbol{\mu}_+ = \frac{\sum_{y_n \in N_+} \mathbf{x}_n}{|N_+|} \text{ and } \boldsymbol{\mu}_- = \frac{\sum_{y_n \in N_-} \mathbf{x}_n}{|N_-|}$$

Predict +1 if $(\mathbf{x} - \boldsymbol{\mu}_+)^T(\mathbf{x} - \boldsymbol{\mu}_+) < (\mathbf{x} - \boldsymbol{\mu}_-)^T(\mathbf{x} - \boldsymbol{\mu}_-)$, and vice versa.

That is,

$$y = \text{sign}[(\mathbf{x} - \boldsymbol{\mu}_-)^T(\mathbf{x} - \boldsymbol{\mu}_-) - (\mathbf{x} - \boldsymbol{\mu}_+)^T(\mathbf{x} - \boldsymbol{\mu}_+)] = \text{sign}[f(\mathbf{x})]$$

Substituting the value of $\boldsymbol{\mu}_+$ and $\boldsymbol{\mu}_-$ in $f(\mathbf{x})$ and further simplifying, we get

$$\begin{aligned} f(\mathbf{x}) &= \frac{1}{N_+^2 N_-^2} \left(N_+^2 (N_- \mathbf{x} - \sum_{n \in N_-} \mathbf{x}_n)^T (N_- \mathbf{x} - \sum_{n \in N_-} \mathbf{x}_n) - N_-^2 (N_+ \mathbf{x} - \sum_{n \in N_+} \mathbf{x}_n)^T (N_+ \mathbf{x} - \sum_{n \in N_+} \mathbf{x}_n) \right) \\ &= \frac{1}{N_+^2 N_-^2} \left[\left(N_+^2 \sum_{i,j \in N_-} \mathbf{x}_i^T \mathbf{x}_j - N_-^2 \sum_{i,j \in N_+} \mathbf{x}_i^T \mathbf{x}_j \right) + \left(N_+^2 \mathbf{x}^T \left(\sum_{j \in N_-} \mathbf{x}_j \right) - N_-^2 \mathbf{x}^T \left(\sum_{j \in N_+} \mathbf{x}_j \right) \right) \right] \\ &= \frac{1}{N_+^2 N_-^2} \left(d(\text{const.}) + \mathbf{x}^T \left(\sum_{j \in N} \beta_j \mathbf{x}_j \right) \right) \end{aligned}$$

$$\text{where } \beta_j = \begin{cases} N_+^2 & j \in N_- \\ N_-^2 & j \in N_+ \end{cases}$$

Rewriting this, we get:

$$f(\mathbf{x}) = \sum_{n=1}^N \alpha_n \langle \mathbf{x}_n, \mathbf{x} \rangle + b$$

$$\text{where } b = \frac{\left(N_+^2 \sum_{i,j \in N_-} \mathbf{x}_i^T \mathbf{x}_j - N_-^2 \sum_{i,j \in N_+} \mathbf{x}_i^T \mathbf{x}_j \right)}{N_+^2 N_-^2} \text{ and } \alpha_n = \begin{cases} \frac{1}{N_-^2} & j \in N_- \\ \frac{1}{N_+^2} & j \in N_+ \end{cases}$$

Question 2

Classes as Gaussians

Given 2 Gaussian distributions as $\mathcal{N}_+(\mathbf{x}|\mu_+, \Sigma)$ and $\mathcal{N}_-(\mathbf{x}|\mu_-, \Sigma)$ and test example \mathbf{x} , probability that \mathbf{x} belongs to class K ($K \in \{+1, -1\}$) is $\mathcal{N}_K(\mathbf{x}|\mu_K, \Sigma)$. // Since these are disjoint, $\mathcal{N}_+ + \mathcal{N}_- = 1$. For \mathbf{x} to belong in class +1, $\frac{\mathcal{N}_+}{\mathcal{N}_-} > 1$, and vice-versa.

Taking log on both sides, $\text{sign}\left[\log\left(\frac{\mathcal{N}_+}{\mathcal{N}_-}\right)\right]$ determines y . Writing both distributions in exponential form,

we get

$$\begin{aligned} f(\mathbf{x}) &= \log\left(\frac{\mathcal{N}_+}{\mathcal{N}_-}\right) = \log\left(\frac{e^{\frac{-1}{2}(\mathbf{x}-\boldsymbol{\mu}_+)^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}_+)}}{e^{\frac{-1}{2}(\mathbf{x}-\boldsymbol{\mu}_-)^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}_-)}}\right) \\ &= \frac{1}{2}\left((\mathbf{x}-\boldsymbol{\mu}_-)^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}_-) - (\mathbf{x}-\boldsymbol{\mu}_+)^T \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu}_+)\right) \end{aligned}$$

Note that if the covariance matrix is of the form $|\sigma|I$, this reduces to the Distance from Means problem. Expanding terms and removing constant factors from the equation, we get

$$\begin{aligned} f(\mathbf{x}) &= 2\mathbf{x}^T \Sigma^{-1}(\boldsymbol{\mu}_+ - \boldsymbol{\mu}_-) + (\boldsymbol{\mu}_-^T \Sigma^{-1} \boldsymbol{\mu}_- - \boldsymbol{\mu}_+^T \Sigma^{-1} \boldsymbol{\mu}_+) \\ &= \mathbf{x}^T \left[2 \Sigma^{-1}(\boldsymbol{\mu}_+ - \boldsymbol{\mu}_-) \right] + (\boldsymbol{\mu}_-^T \Sigma^{-1} \boldsymbol{\mu}_- - \boldsymbol{\mu}_+^T \Sigma^{-1} \boldsymbol{\mu}_+) \\ &= \mathbf{w}^T \mathbf{x} + b \end{aligned}$$

where $\mathbf{w} = \left[2 \Sigma^{-1}(\boldsymbol{\mu}_+ - \boldsymbol{\mu}_-) \right]$ and $b = (\boldsymbol{\mu}_-^T \Sigma^{-1} \boldsymbol{\mu}_- - \boldsymbol{\mu}_+^T \Sigma^{-1} \boldsymbol{\mu}_+)$

Question 3

Importance-Weighted linear regression

Let $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$ and correspondingly $Y = [y_1, y_2, \dots, y_N]^T$. Also $\langle \mathbf{x}_n, y_n \rangle$ carries penalty c_n . Let C be a $N \times N$ diagonal matrix with $C_{ii} = c_i$. Define empirical loss function

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

where \mathbf{w} is the parametric vector.

Our goal is to optimize \mathbf{w} such that \mathcal{L} is minimum. writing the above equation in vector form, we get

$$\mathcal{L} = (\mathbf{X}\mathbf{w} - Y)^T C (\mathbf{X}\mathbf{w} - Y)$$

to minimize \mathcal{L} , we use matrix calculus. Recall that $\nabla_{\mathbf{x}} \mathbf{x}^T A \mathbf{x} = 2A\mathbf{x}$, and $\nabla_{\mathbf{x}} \mathbf{b}^T \mathbf{x} = \mathbf{b}$. Applying these results on the expanded Loss function \mathcal{L} , we get

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 2\mathbf{X}^T C \mathbf{X} \mathbf{w} - 2\mathbf{X}^T C Y = 0$$

Since \mathcal{L} is of quadratic form, $\mathcal{L} = 0$ gives us the global optimum (minimum).

$$\mathbf{X}^T C \mathbf{X} \mathbf{w} = \mathbf{X}^T C Y \mathbf{w} = (\mathbf{X}^T C \mathbf{X})^{-1} \mathbf{X}^T C Y$$

Similarly, adding l_2 norm results in

$$\mathcal{L} = (\mathbf{X}\mathbf{w} - Y)^T C (\mathbf{X}\mathbf{w} - Y) + \lambda \mathbf{w}^T \mathbf{w}$$

which on minimisation gives

$$\mathbf{w} = (\mathbf{X}^T C \mathbf{X} + \lambda I_D)^{-1} \mathbf{X}^T C Y$$

where D is the dimensions of the feature vector. Also note that we have assumed that the bias units and bias parameters have already been added to the feature vector \mathbf{X} and the parameter vector \mathbf{w} .

Question 4

Noise as regularizer

For a standard non-regularised linear regression model, we define the loss function as

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

and minimise \mathcal{L} with respect to \mathbf{w} . Note that this is identical to maximising

$$\prod_{n=1}^N \mathcal{N}(\mathbf{w}^T \mathbf{x} | \mathbf{w})$$

. This is the maximum likelihood function and translates to

$$\prod_{n=1}^N \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(\mathbf{w}^T \mathbf{x}_n - y_n)^2}{2\sigma_0^2}} = e^{-\frac{\left(-(X\mathbf{w} - Y)^T(X\mathbf{w} - Y)\right)}{2\sigma_0^2}}$$

Also, after adding prior $\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{\mathbf{w}^T \mathbf{w}}{2\sigma^2}}$ the MAP function becomes

$$e^{-\frac{\left(-(X\mathbf{w} - Y)^T(X\mathbf{w} - Y)\right)}{\sigma_0^2} + \frac{\mathbf{w}^T \mathbf{w}}{2\sigma^2}}$$

Taking logarithm of the MAP function and maximizing this leads to the empirical loss function

$$\mathcal{L} = (X\mathbf{w} - Y)^T C (X\mathbf{w} - Y) + \lambda \mathbf{w}^T \mathbf{w}$$

where $\lambda = \frac{\sigma_0^2}{\sigma^2}$

As we have already seen in the previous question, this is the optimisation function of regularised Linear Regression and thus the prior in this case acts as the regulariser, befitting our intuition.

Question 5

Decision trees for regression

We basically have to make our selection purer as we go down the tree. Suppose X is the feature vector matrix with rows as feature vector (which may consist of real valued features as well as labeled features). If all of the features are labeled features, we split the data according to the feature which on splitting reduces the variance of the data the most. That is, the entropy, or randomness of the real valued prediction is interpreted by its variance. So the feature \mathcal{F} which minimises the variance of the resulting split subclasses of the data the most is selected for that node.

If the feature is a real valued one, we calculate a boundary value for the values of that feature, for which the variance is minimized for the resulting subclasses of X . In this way we tackle splitting real valued features as well as discrete valued ones, and complete the tree (doing pruning of branches either during the process or after).

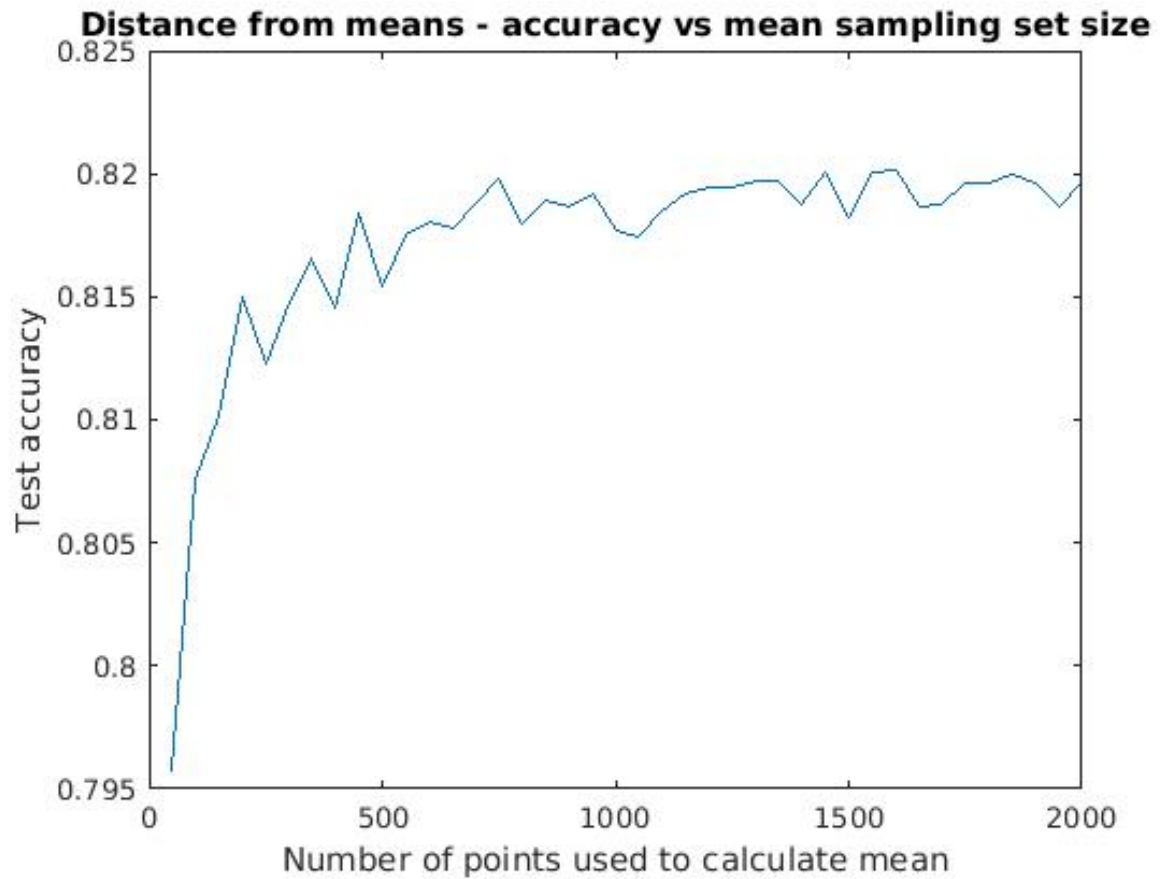
This is just a crude algorithm, and the accuracy of it I can barely predict. However, This is one way of "purifying" the features.

Question 6

Programming assignment- Distance from means classifier on MNIST dataset

This question requires a programming element. The program is written in MATLAB. The goal is to use the *Distance From Means* algorithm, as discussed in class to classify handwritten digits as a part of the MNIST dataset provided. The code is available on this [Github repository](#).

The algorithm is a primitive one that uses the concept of representing one class with the mean of its constituents, and predicts new datapoints' classes by computing euclidean distances from the means as a measure of the similarity from the classes. The size of subset of each class used to compute the mean is varied and the corresponding test accuracy is plotted against it, as shown in the following figure.



We notice that as the mean sampling size increases, the performance, on average, increases (barring some abrupt fluctuations due to random sampling, which have been smoothed out by averaging accuracy over multiple iterations). This may be explained on the basis that the increased number of data points better approximate the underlying central tendency of the class. As a result, better results are obtained when the sampling set is larger.