

Q1) Distance from mean.

In this given a test example we say the label is $y = \text{sign}(f(x))$. $\rightarrow x$ feature vector of test
 $\rightarrow y$ predicted label.
 where

$$f(x) = \|M_- - x\|^2 - \|M_+ - x\|^2$$

$$= \underbrace{M_-^2 - M_+^2}_b + 2M_+^T x - 2M_-^T x.$$

$$= b + 2[M_+ - M_-]^T x$$

$$= b + 2 \left[\frac{\sum_{n \in N_+} x_n}{|N_+|} - \frac{\sum_{n \in N_-} x_n}{|N_-|} \right]^T x.$$

$$= \sum_{n=1}^N \alpha_n \langle x_n, x \rangle + b$$

where $\alpha_n = \begin{cases} \frac{2y_n}{|N_+|} & \text{if } y_n > 0 \\ \frac{2y_n}{|N_-|} & \text{if } y_n < 0 \end{cases}$

$$b = M_-^2 - M_+^2$$

hence we get the form

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

Q2) Gaussian Distribution

$$N(x) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2} (x-M)^T \Sigma^{-1} (x-M)\right]$$

taking the ratio

$$\frac{N_+(x)}{N_-(x)}$$

Now we want a function f such that label of a testing feature vector x is given by

$$y = \text{sign}(f(x))$$

Define

$$f(x) = \ln \frac{N_+(x)}{N_-(x)}$$

$$= \ln \left[\frac{\frac{1}{(2\pi)^{D_+} \sqrt{|\Sigma|}} \exp\left[-\frac{1}{2} (x-M_+)^T \Sigma^{-1} (x-M_+)\right]}{\frac{1}{(2\pi)^{D_-} \sqrt{|\Sigma|}} \exp\left[-\frac{1}{2} (x-M_-)^T \Sigma^{-1} (x-M_-)\right]} \right]$$

$$= \frac{1}{2} \ln \frac{(2\pi)^{D_+}}{(2\pi)^{D_-}} - \frac{1}{2} (x-M_+)^T \Sigma^{-1} (x-M_+) + \frac{1}{2} (x-M_-)^T \Sigma^{-1} (x-M_-)$$

$$= \frac{D_+ - D_-}{2} \ln 2\pi + \frac{1}{2} \left[(x^T \Sigma^{-1} - M_-^T \Sigma^{-1}) (x-M_-) - (x^T \Sigma^{-1} - M_+^T \Sigma^{-1}) (x-M_+) \right]$$

(\therefore using $x^T \Sigma^{-1} M_- = M_-^T \Sigma^{-1} x$)

$$= \underbrace{\frac{D_+ - D_-}{2} \ln 2\pi + \frac{1}{2} [M_-^T \Sigma^{-1} M_- - M_+^T \Sigma^{-1} M_+]}_b + \underbrace{[M_+^T \Sigma^{-1} - M_-^T \Sigma^{-1}] x}_{w^T x}$$

$$= b + w^T x$$

$$w = (M_+^T \Sigma^{-1} - M_-^T \Sigma^{-1})^T$$

$$b = \frac{D_+ - D_-}{2} \ln 2\pi + \frac{1}{2} [M_-^T \Sigma^{-1} M_- - M_+^T \Sigma^{-1} M_+]$$

$\begin{cases} D_+ = D_- \\ \text{since } D \text{ is gaussian} \\ \text{distribution is the} \\ \text{dimension.} \end{cases}$

Q3) linear regression

We change then L_{emp} to incorporate the weight c_n

$$L_{emp} = \frac{1}{2} \sum_{n=1}^N (y_n - w^T x_n)^2 c_n$$

$$\Rightarrow \frac{dL_{emp}}{dw} = \sum_{n=1}^N (y_n - w^T x_n) \cdot c_n x_n$$

Equating $\frac{dL_{emp}}{dw}$ to zero

$$\sum_{n=1}^N y_n c_n x_n = \left(\sum_{n=1}^N w^T x_n \right) \cdot c_n x_n$$

$$\Rightarrow \sum_{n=1}^N y_n c_n x_n = \left(\sum_{n=1}^N x_n^T w \right) c_n x_n$$

$$\Rightarrow \sum_{n=1}^N y_n c_n x_n = \left[\sum_{n=1}^N c_n x_n x_n^T \right] w$$

$$w = \left[\sum_{n=1}^N c_n x_n x_n^T \right]^{-1} \left[\sum_{n=1}^N y_n c_n x_n \right]$$

$$= [X^T X]^{-1} [X^T Y]$$

$$X = [\sqrt{c_1} x_1, \sqrt{c_2} x_2, \dots, \sqrt{c_N} x_N]^T$$

$$Y = [\sqrt{c_1} y_1, \sqrt{c_2} y_2, \dots, \sqrt{c_N} y_N]^T$$

Q. Noise as regulariser.

$$L(w) = \sum_{n=1}^N (y_n - w^T x_n)^2$$

Now replacing x_n with $x_n + \epsilon_n$

$$\tilde{L}(w) = \sum_{n=1}^N (y_n - w^T (x_n + \epsilon_n))^2$$

Also

$$E(\epsilon_n) = 0 \quad E(\epsilon_n \epsilon_m) = \delta_{nm} \sigma^2 I$$

$$\begin{aligned} E(\tilde{L}(w)) &= E \left[\sum_{n=1}^N (y_n - w^T (x_n + \epsilon_n))^2 \right] \\ &= E \left[\sum_{n=1}^N \left[y_n^2 + (w^T (x_n + \epsilon_n))^2 - 2 y_n w^T (x_n + \epsilon_n) \right] \right] \\ &= E \left[\sum_{n=1}^N \left[y_n^2 + (w^T x_n)^2 + (w^T \epsilon_n)^2 - 2 (w^T x_n) (w^T \epsilon_n) - 2 y_n w^T (x_n + \epsilon_n) \right] \right] \\ &= \sum_{n=1}^N y_n^2 + (w^T x_n)^2 + \cancel{2 w^T x_n (w^T E(\epsilon_n))} \\ &\quad + w^T E(\epsilon_n \epsilon_n^T) w - \cancel{2 y_n w^T x_n} - \cancel{2 y_n w^T E(\epsilon_n)} \\ &= \sum_{n=1}^N y_n^2 + (w^T x_n)^2 + w^T E(\epsilon_n \epsilon_n^T) w - 2 y_n w^T x_n \\ &= \sum_{n=1}^N \left[(y_n - w^T x_n)^2 + \sigma^2 w^T w \right] \\ &= \sum_{n=1}^N (y_n - w^T x_n)^2 + N \sigma^2 w^T w \end{aligned}$$

$\Rightarrow E(\tilde{L}(w))$ is the same as the regularized objective function for linear regression.

Q5) Decision tree for Regression

Let the problem has n attributes x_1, x_2, \dots, x_n
Now for the decision tree we need to split the dataset based on some attribute x_i at some value k . So to find x_i and k we use the following method.

Define

$$D = \sum_{\substack{j \\ x_{ij} > k}} (M_1 - y_j)^2 + \sum_{\substack{j \\ x_{ij} \leq k}} (M_2 - y_j)^2$$

Note my nomenclature
 x_{ij} is the value of x_i feature for j th example.

M_1 = average over all y such that the corresponding feature x_i for that y is greater than k .

M_2 = average over all y such that the corresponding feature x_i for that y is less than k .

→ We are finding a standard deviation like measure.

→ We find D for different values of i and k and take that pair of (i, k) for which D is minimum.

→ keep doing process recursively.

→ finding condition - keep doing this until you get perfect branches or you can keep doing it for fixed number of iterations.

Q 6).

Initially the accuracy increases rapidly by increasing so more training data but after one more the accuracy does not change so much.

This is mostly because initially on add so more training data the mean of each class change a lot but after some time the mean moves by very small distance and so we don't have much change in the accuracy.