

Solutions :

Name : Achyut Kulkarni

For all values here I have taken $W = \theta \Rightarrow L(w) \Rightarrow L(\theta)$ and so on ...

1) Logistic Regression :

a) Lets take $W = \theta \Rightarrow L(w) \Rightarrow L(\theta)$

w.k.t for Binary logistic regression :

$$P(y=1|x; \theta) = h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}} \Rightarrow P(y=0|x; \theta) = 1 - h_{\theta}(x) = 1 - \frac{1}{1+e^{-\theta^T x}}$$

Since there are only two classes we can define this as :

$$\begin{aligned} L(\theta) &= -\log\left(\prod_{i=1}^n p(y = y_i | x_i)\right) = -\log\left(\prod_{i=1}^n (h_{\theta}(x)^{y_i} * (1 - h_{\theta}(x))^{1-y_i})\right) \\ \Rightarrow & -y_{mi} \sum_{i=1}^n \log(h_{\theta}(x)) - (1 - y_i) \sum_{i=1}^n \log((1 - h_{\theta}(x))) \end{aligned}$$

b)

From the final equation in above answer ., we have

$$-y_i \sum_{i=1}^n \log(h_{\theta}(x)) - (1 - y_i) \sum_{i=1}^n \log((1 - h_{\theta}(x))) \text{-----(1)}$$

And we also know that

$$P(y=1|x; \theta) = h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}} \Rightarrow P(y=0|x; \theta) = 1 - h_{\theta}(x) = 1 - \frac{1}{1+e^{-\theta^T x}}$$

Substituting in 1, we get

However w.k.t Gradient descent is

$$\theta_j : \theta_j - \alpha \frac{\partial}{\partial \theta} J(\theta)$$

$$\text{For } n \geq 1, \text{ we have } J(\theta) = \left(-y_i \sum_{i=1}^n \log\left(\frac{1}{1+e^{-\theta^T x}}\right) - (1 - y_i) \sum_{i=1}^n \log\left(1 - \frac{1}{1+e^{-\theta^T x}}\right) \right)$$

Hence;

$$\frac{\partial}{\partial \theta} J(\theta) = \left(-y_i \sum_{i=1}^n \log(\sigma(\theta^T x_n)) - (1-y_i) \sum_{i=1}^n \log((1-\sigma(\theta^T x_n))) \right)$$

$$\frac{\partial}{\partial \theta} J(\theta) = \left(-\sum_{i=1}^n y_n (1 - \sigma(\theta^T x_n) x_n) - (1-y_n)(\sigma(\theta^T x_n) x_n) \right)$$

$$\frac{\partial}{\partial \theta} J(\theta) = \sum_{i=1}^n ((\sigma(\theta^T x_n) - y_n) x_n)$$

Hence,

$$\theta_j : \theta_j - \alpha \frac{\partial}{\partial \theta} J(\theta)$$

$$\theta_j : \theta_j - \alpha \sum_{i=1}^n ((\sigma(\theta^T x_n) - y_n) x_n)$$

Intuitively, a (strictly) convex function has a “bowl shape”, and hence has a unique global minimum θ^* corresponding to the bottom of the bowl. Hence its second derivative must be positive everywhere, $d^2 f(\theta) > 0$. A twice-continuously differentiable, multivariate function f is $d^2 \theta$ convex if its Hessian is positive definite for all θ .

Now let's differentiate $\sum_{i=1}^n ((\sigma(\theta^T x_n) - y_n) x_n)$ second time and see if it's > 0

Hence ;'

$$\frac{\partial}{\partial \theta} \sum_{i=1}^n ((\sigma(\theta^T x_n) - y_n) x_n) = \sigma(\theta^T x_n)(1 - \sigma(\theta^T x_n))x_n^2 \quad \dots\dots\dots 1$$

Since , the sigmoid function lies between zero and 1 $\Rightarrow \sigma(\theta^T x_n) > 0$

$\Rightarrow (1 - \sigma(\theta^T x_n)) > 0$ too hence

The whole equation is greater than zero hence, the

Hessian H is > 0 .

Therefore we can say that the function is convex and thus converges at only global minimum.

c) Let's take $w = \theta$

Hence for K classes :

$$P(y=k|x; \theta) = h_{\theta}(x) = \frac{e^{\theta^T x}}{1 + \sum_{i=1}^{K-1} e^{\theta^T x}} \text{ for } K = 1, 2, 3, \dots, K-1$$

$$\Rightarrow P(y=K|x; \theta) = 1 - h_{\theta}(x) = 1 - \frac{1}{1 + \sum_{i=1}^{K-1} e^{\theta^T x}} \Rightarrow \frac{1}{1 + \sum_{i=1}^{K-1} e^{\theta^T x}} \text{ for } k = K$$

Negative Log Likelihood of the above equation is $L(\theta_1, \theta_2, \theta_3, \dots, \theta_k)$

$$L(\theta_1, \theta_2, \theta_3, \dots, \theta_k) = -\log\left(\prod_{i=1}^n p(y=y_i|x_i)\right) = -\sum_{i=1}^n \sum_{k=1}^{k-1} \left(\log\left(\frac{e^{\theta_{kX}^T}}{1+\sum_{t=1}^{k-1} e^{\theta_t^T X}}\right)\right) - \sum_{i=1}^n \log\left(\frac{1}{1+\sum_{t=1}^{k-1} e^{\theta_t^T X}}\right)$$

Also; we can do ;

$$L(\theta_1, \theta_2, \theta_3, \dots, \theta_k) = -\log\left(\prod_{i=1}^n p(y=y_i|x_i)\right) = -\left(\log\left(\prod_{i=1}^n \prod_{k=1}^{k-1} \frac{e^{\theta_t^T X}}{1+\sum_{t=1}^{k-1} e^{\theta_t^T X}}\right) \left(\frac{1}{1+\sum_{t=1}^{k-1} e^{\theta_t^T X}}\right)\right) \text{-----}$$

$$\frac{1}{1+\sum_{t=1}^{k-1} e^{\theta_t^T X}} = \frac{e^0}{1+\sum_{t=1}^{k-1} e^{\theta_t^T X}} = \frac{e^{\theta_{kX}^T}}{1+\sum_{t=1}^{k-1} e^{\theta_t^T X}} \text{ as we know } \theta_{kX}^T = 0$$

$$\begin{aligned} & -\sum_{i=1}^n \sum_{k=1}^{k-1} \left(\log\left(\frac{e^{\theta_{kX}^T}}{1+\sum_{t=1}^{k-1} e^{\theta_t^T X}}\right)\right) + \log\left(\frac{e^{\theta_{kX}^T}}{1+\sum_{t=1}^{k-1} e^{\theta_t^T X}}\right) = -\sum_{i=1}^n \sum_{k=1}^k \left(\log\left(\frac{e^{\theta_{kX}^T}}{1+\sum_{t=1}^{k-1} e^{\theta_t^T X}}\right)^{y_{ik}}\right); \text{ where } y \\ & = -\sum_{i=1}^n \sum_{k=1}^k \left(\log\left(\frac{e^{\theta_{kX}^T}}{1+\sum_{t=1}^{k-1} e^{\theta_t^T X}}\right)^{y_{ik}}\right) \\ & = -\sum_{i=1}^n \sum_{k=1}^k \theta_{kX}^T y_{ik} - y_{ik} \log\left(1 + \sum_{t=1}^{k-1} e^{\theta_t^T X}\right) \end{aligned}$$

This reduces to :

$$-\sum_{i=1}^n \sum_{k=1}^k \theta_{kX}^T y_{ik} - y_{ik} \log\left(1 + \sum_{t=1}^{k-1} e^{\theta_t^T X}\right)$$

d)

Now w.k.t for gradient descent

$$\theta_j : \theta_j - \alpha \frac{\partial}{\partial \theta} J(\theta)$$

And

$$J(\theta) = -\sum_{i=1}^n \sum_{k=i}^k \theta_{kX}^T y_{ik} - y_{ik} \log\left(1 + \sum_{t=1}^{k-1} e^{\theta_t^T X}\right) - \sum_{i=1}^n \sum_{k \neq i}^{k-1} \theta_{kX}^T y_{ik} - y_{ik} \log\left(1 + \sum_{t=1}^{k-1} e^{\theta_t^T X}\right)$$

Hence

$$\frac{\partial}{\partial \theta_i} J(\theta) = -\sum_{i=1}^n \sum_{k=i}^k x_{ij} y_{ik} - y_{ik} \left(\frac{\sum_{t=1}^k e^{\theta_t^T X} x_{ij}}{1+\sum_{t=1}^k e^{\theta_t^T X}}\right) + \sum_{k \neq i}^k -y_{ik} \left(\frac{\sum_{t=1}^k e^{\theta_t^T X} x_{ij}}{1+\sum_{t=1}^k e^{\theta_t^T X}}\right)$$

$$\frac{\partial}{\partial \theta_i} J(\theta) = -\sum_n \begin{matrix} \square \\ \square \end{matrix} y_{ik} x_i - \sum_k \frac{x_{ij} y_{ik} e^{\theta_k^T X_i}}{1+\sum_{t=1}^{k-1} \exp(w_t^T x_n)} \begin{matrix} \square \\ \square \end{matrix}$$

We know that $P(Y = k | X = x) = \frac{\exp(w_k^T x)}{1 + \sum_1^{K-1} \exp(w_i^T x)}$,

$$\frac{\partial L(w_1, w_2, \dots, w_k)}{\partial w_i} = - \sum_n x_n \left(\left[y_{nk} - \sum_k P(Y = k | X = x) y_{nk} \right] \right)$$

2) Linear/Gaussian Discriminant analysis :

a)

W.k.t :

Lets denote $x_n = x$ and $y_n = y$ for the ease

$$P(x_n, y_n) = P(y_n) * P(x_n | y_n)$$

$$\Rightarrow P(x_n, y_n) = \left\{ p_1 \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp\left(-\frac{(x_n - \mu_1)^2}{2\sigma_1^2}\right) \text{ if } y_n = 1 \right\} \text{ and}$$

$$P(x_n, y_n) = \left\{ p_2 \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{(x_n - \mu_2)^2}{2\sigma_2^2}\right) \text{ if } y_n = 2 \right\}$$

Where $p_1 + p_2 = 1$

Hence $p_2 = 1 - p_1$

Log likelihood is :

$$L(D) = \sum_1^2 \log(P(x_n, y_n))$$

$$\Rightarrow \sum_{n: y_n=1} \log\left(\left\{ p_1 \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp\left(-\frac{(x_n - \mu_1)^2}{2\sigma_1^2}\right) \right\} \right) + \sum_{n: y_n=2} \log\left(\left\{ p_2 \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{(x_n - \mu_2)^2}{2\sigma_2^2}\right) \right\} \right)$$

Now to calculate MLE's of parameters:

i) p_1^*

$$p_1^* = \frac{\partial}{\partial p_1} \left(\sum_{n: y_n=1} \log\left(\left\{ p_1 \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp\left(-\frac{(x_n - \mu_1)^2}{2\sigma_1^2}\right) \right\} \right) + \sum_{n: y_n=2} \log\left(\left\{ p_2 \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{(x_n - \mu_2)^2}{2\sigma_2^2}\right) \right\} \right) \right) = 0$$

$$p_1^* = \sum_{n: y_n=1} \left(\frac{1}{p_1}\right) + \sum_{n: y_n=1} \left(\frac{1}{1-p_1}\right) \Rightarrow 0$$

Lets denote

$$\sum_{n: y_n=1} 1 = n_1 \quad \text{and} \quad \sum_{n: y_n=2} 1 = n_2 \dots\dots\text{hence}$$

$$\frac{n_1}{p_1} - \frac{n_2}{1-p_1} = 0 \quad \Rightarrow \quad p_1 = \frac{n_1}{n_1+n_2}$$

$$p_1^* = \frac{n_1}{n}$$

ii) p_2^*

$$p_2^* = \frac{\partial}{\partial p_1} \left(\sum_{n: y_n=1} \log \left(\left\{ p_1 \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp \left(-\frac{(x_n - \mu_1)^2}{2\sigma_1^2} \right) \right\} \right) + \sum_{n: y_n=2} \log \left(\left\{ p_2 \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp \left(-\frac{(x_n - \mu_2)^2}{2\sigma_2^2} \right) \right\} \right) \right) = 0$$

$$p_2^* = \sum_{n: y_n=1} \left(\frac{1}{p_1} \right) + \sum_{n: y_n=1} \left(\frac{1}{p_2} \right) \Rightarrow 0$$

Lets denote

$$\sum_{n: y_n=1} 1 = n_1 \quad \text{and} \quad \sum_{n: y_n=2} 1 = n_2 \dots\dots\text{hence}$$

$$\frac{n_1}{1-p_1} - \frac{n_2}{p_2} = 0 \quad \Rightarrow \quad p_2 = \frac{n_2}{n_2+n_1}$$

$$p_2^* = \frac{n_2}{n}$$

iii) μ_1^*

$$\mu_1^* = \frac{\partial}{\partial \mu_1} \left(\sum_{n: y_n=1} \log \left(\left\{ p_1 \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp \left(-\frac{(x_n - \mu_1)^2}{2\sigma_1^2} \right) \right\} \right) + \sum_{n: y_n=2} \log \left(\left\{ p_2 \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp \left(-\frac{(x_n - \mu_2)^2}{2\sigma_2^2} \right) \right\} \right) \right) = 0$$

This reduces to ;

$$\mu_1^* = \sum_{n_1} \log \left(\left\{ \exp \left(-\frac{(x_n - \mu_1)^2}{2\sigma_1^2} \right) \right\} \right) + 0 = 0$$

$$\frac{\partial}{\partial \mu_1} \left(-\frac{1}{2\sigma_1^2} \sum_{n_1} ((x_n - \mu_1)^2) \right) + 0 = 0$$

$$\sum_{n1} \left((x_n - \mu_1) \right) = 0$$

$$\Rightarrow \mu_1^* = \frac{\sum x_n}{n1}$$

$$\text{iv) } \mu_2^* = \frac{\partial}{\partial \mu_2} \left(\sum_{n: y_n=1} \log \left(\left\{ p1 \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp\left(-\frac{(x_n - \mu_1)^2}{2\sigma_1^2}\right) \right\} \right) + \sum_{n: y_n=2} \log \left(\left\{ p2 \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{(x_n - \mu_2)^2}{2\sigma_2^2}\right) \right\} \right) \right) = 0$$

This reduces to ;

$$\mu_2^* = \sum_{n2} \log \left(\left\{ \exp \left(-\frac{(x_n - \mu_2)^2}{2\sigma_2^2} \right) \right\} \right) + 0 = 0$$

$$\frac{\partial}{\partial \mu_2} \left(-\frac{1}{2\sigma_2^2} \sum_{n2} ((x_n - \mu_2)^2) \right) + 0 = 0$$

$$\sum_{n2} \left((x_n - \mu_2) \right) = 0$$

$$\Rightarrow \mu_2^* = \frac{\sum x_n}{n2}$$

$$\text{v) } \sigma_1^*$$

$$\sigma_1^* = \frac{\partial}{\partial \sigma_1^*} \left(\sum_{n: y_n=1} \log \left(\left\{ p1 \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp\left(-\frac{(x_n - \mu_1)^2}{2\sigma_1^2}\right) \right\} \right) + \sum_{n: y_n=2} \log \left(\left\{ p2 \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left(-\frac{(x_n - \mu_2)^2}{2\sigma_2^2}\right) \right\} \right) \right) = 0$$

This reduces to ;

$$\sigma_1^* = \frac{\partial}{\partial \sigma_1^*} \left(\sum_{n: y_n=1} \log \left(\left\{ p1 \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp\left(-\frac{(x_n - \mu_1)^2}{2\sigma_1^2}\right) \right\} \right) \right) + 0 = 0$$

$$\sigma_1^* = \frac{\partial}{\partial \sigma_1^*} \left(\sum_{n: y_n=1} \left\{ \log \left\{ \frac{1}{\sqrt{2\pi\sigma_1^2}} \right\} + \log \left\{ \exp\left(-\frac{(x_n - \mu_1)^2}{2\sigma_1^2}\right) \right\} \right\} \right) = 0$$

$$\sigma_1^* = \frac{\partial}{\partial \sigma_1^*} \sum_{n: y_n=1} \left\{ \frac{-1}{2} \log \sigma_1^2 \right\} = \frac{(x_n - \mu_1)^2}{2\sigma_1^2}$$

$$\sigma_1^* = \frac{\partial}{\partial \sigma_1^*} \sum_{n: y_n=1} \{ \log \sigma_1^2 = \frac{(x_n - \mu_1)^2}{\sigma_1^2} = 0 \}$$

$$\sigma_1^* = 2 \frac{n1}{\sigma_1} = \frac{-2(x_n - \mu_1)^2}{\sigma_1^3}$$

$$\sigma_1^* = \frac{\sum (x_n - \mu_1)}{n1 \sqrt{n1}}$$

vi)

$$\sigma_2^* = \frac{\partial}{\partial \sigma_2^*} (\sum_{n: y_n=1} \log(\left\{ p1 \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp(-\frac{(x_n - \mu_1)^2}{2\sigma_1^2}) \right\}) + \sum_{n: y_n=2} \log(\left\{ p2 \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp(-\frac{(x_n - \mu_2)^2}{2\sigma_2^2}) \right\}) = 0$$

This reduces to ;

$$\sigma_2^* = \frac{\partial}{\partial \sigma_2^*} (\sum_{n: y_n=2} \log(\left\{ p2 \frac{1}{\sqrt{2\pi\sigma_2^2}} \exp(-\frac{(x_n - \mu_2)^2}{2\sigma_2^2}) \right\}) + 0 = 0$$

$$\sigma_2^* = \frac{\partial}{\partial \sigma_2^*} (\sum_{n: y_n=2} \left\{ \log \left\{ \frac{1}{\sqrt{2\pi\sigma_2^2}} \right\} + \log \left\{ \exp(-\frac{(x_n - \mu_2)^2}{2\sigma_2^2}) \right\} \right\} = 0)$$

$$\sigma_2^* = \frac{\partial}{\partial \sigma_2^*} \sum_{n: y_n=2} \{ \frac{-1}{2} \log \sigma_2^2 = \frac{(x_n - \mu_2)^2}{2\sigma_2^2} \}$$

$$\sigma_2^* = \frac{\partial}{\partial \sigma_2^*} \sum_{n: y_n=2} \{ \log \sigma_2^2 = \frac{(x_n - \mu_2)^2}{\sigma_2^2} = 0 \}$$

$$\sigma_2^* = 2 \frac{n2}{\sigma_2} = \frac{-2(x_n - \mu_2)^2}{\sigma_2^3}$$

$$\sigma_2^* = \frac{\sum (x_n - \mu_2)}{n2 \sqrt{n2}}$$

2)

b)

W.k.t Naive Bayes

$p(x|y = c_1)$ and $p(x|y = c_2)$ follows a multivariate gaussian distributions

$N(\mu_1, \Sigma)$ and $N(\mu_2, \Sigma)$ respectively.

To prove:

$$p(y = 1|x) = \frac{1}{1 + \exp(-\theta^T x)} \text{ for some } \theta$$

$$\begin{aligned} P(Y = 1|X) &= \frac{P(X|Y=1)P(Y=1)}{P(X|Y=0)P(Y=0) + P(X|Y=1)P(Y=1)} \\ &= \left(\frac{1}{1 + \frac{P(X|Y=0)P(Y=0)}{P(X|Y=1)P(Y=1)}} \right) \\ &= \frac{1}{\frac{p_2}{\sqrt{(2\pi)^2 |\Sigma|}} \exp \left[-\frac{1}{2} (x - \mu_2)^T (\Sigma)^{-1} (x - \mu_2) \right] + \frac{p_1}{\sqrt{(2\pi)^2 |\Sigma|}} \exp \left[-\frac{1}{2} (x - \mu_1)^T (\Sigma)^{-1} (x - \mu_1) \right]} \end{aligned}$$

By substituting $p_2 = 1 - p_1$ in the above equation,

$$P(Y = 1|X) = \frac{1}{\frac{(1-p_1)}{\sqrt{(2\pi)^2 |\Sigma|}} \exp \left[-\frac{1}{2} (x - \mu_2)^T (\Sigma)^{-1} (x - \mu_2) \right] + \frac{p_1}{\sqrt{(2\pi)^2 |\Sigma|}} \exp \left[-\frac{1}{2} (x - \mu_1)^T (\Sigma)^{-1} (x - \mu_1) \right]}$$

$$P(Y = 1|X) = \frac{1}{\frac{(1-p_1) \exp \left[-\frac{1}{2} (x - \mu_2)^T (\Sigma)^{-1} (x - \mu_2) \right]}{\sqrt{(2\pi)^2 |\Sigma|}} + \frac{p_1 \exp \left[-\frac{1}{2} (x - \mu_1)^T (\Sigma)^{-1} (x - \mu_1) \right]}{\sqrt{(2\pi)^2 |\Sigma|}}}$$

Since $x = e^{\log(x)}$ we get,

$$P(Y = 1|X) = \frac{1}{1 + \exp \left(\log \left(\frac{(1-p_1) \exp \left(-\frac{1}{2} (x - \mu_2)^T (\Sigma)^{-1} (x - \mu_2) \right) + \frac{1}{2} (x - \mu_1)^T (\Sigma)^{-1} (x - \mu_1)}{(p_1)} \right) \right)}$$

$$\begin{aligned}
P(Y = 1|X) &= \frac{1}{1+\exp\left(\log\left(\frac{(1-p_1)}{(p_1)}\right)+\left[\frac{-1}{2}(x-\mu_2)^T(\Sigma)^{-1}(x-\mu_2)+\frac{-1}{2}(x-\mu_1)^T(\Sigma)^{-1}(x-\mu_1)\right]\right)} \\
P(Y = 1|X) &= \frac{1}{1+\exp\left(\log\left(\frac{(1-p_1)}{(p_1)}\right)+\left[\left(\frac{-1}{2}\right)(-x^T(\Sigma)^{-1}\mu_2-\mu_2^T(\Sigma)^{-1}x+\mu_2^T(\Sigma)^{-1}x+\left(\frac{1}{2}\right)(-x^T(\Sigma)^{-1}\mu_1-\mu_1^T(\Sigma)^{-1}x+\mu_1^T(\Sigma)^{-1}\mu_1\right)\right]\right)} \\
P(Y = 1|X) &= \frac{1}{1+\exp\left(\log\left(\frac{(1-p_1)}{(p_1)}\right)+\left[\frac{\mu_1^T(\Sigma)^{-1}\mu_1-\mu_2^T(\Sigma)^{-1}\mu_2}{2}\right]-[\mu_1^T-\mu_2^T](\Sigma)^{-1}x\right)} \\
P(Y = 1|X) &= \frac{1}{1+\exp\left(-\left[-\log\left(\frac{(1-p_1)}{(p_1)}\right)-\left[\frac{\mu_1^T(\Sigma)^{-1}\mu_1-\mu_2^T(\Sigma)^{-1}\mu_2}{2}\right]+[\mu_1^T-\mu_2^T](\Sigma)^{-1}x\right]\right)}
\end{aligned}$$

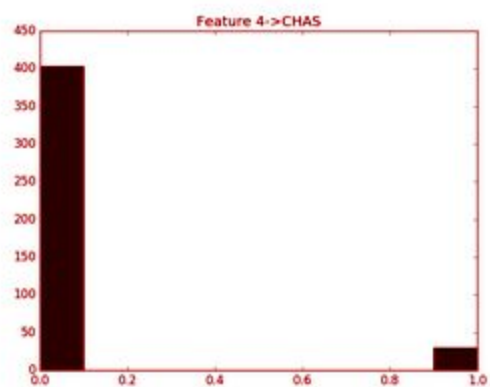
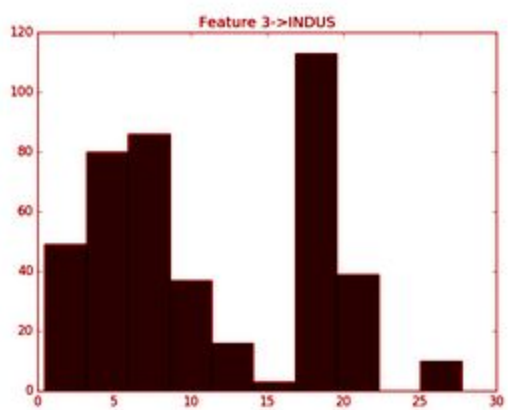
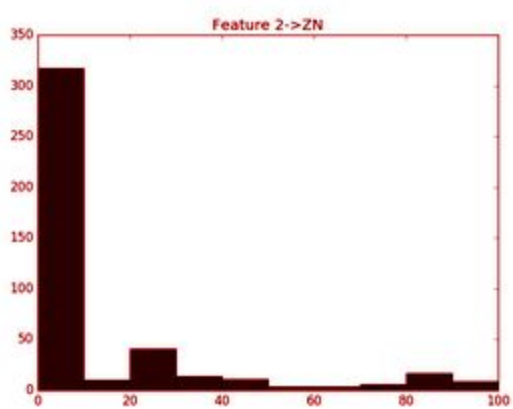
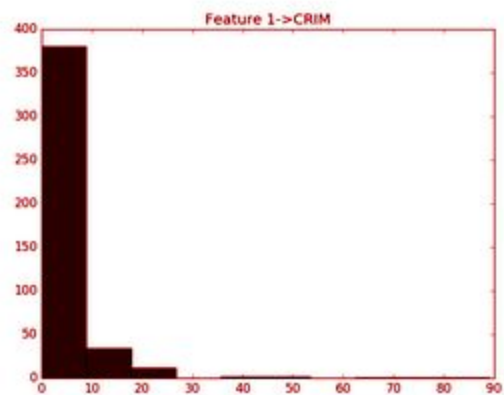
The above equation is of the required form

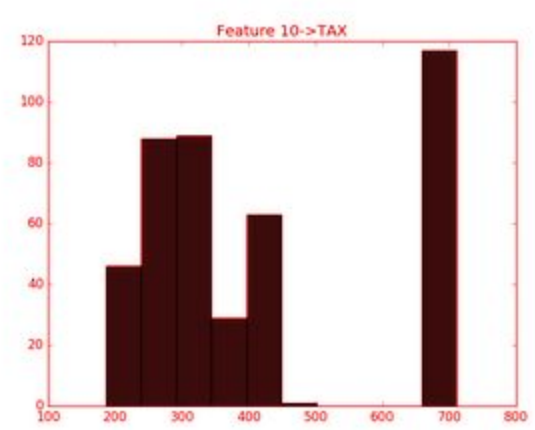
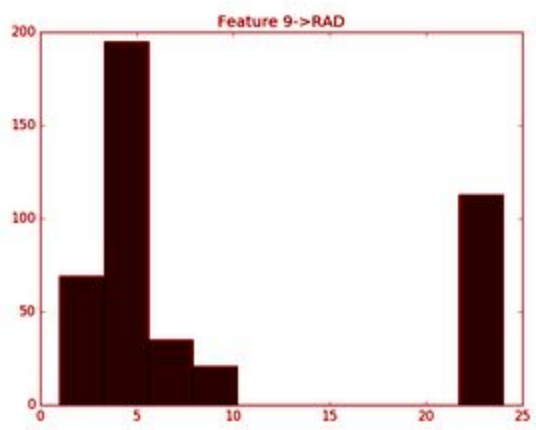
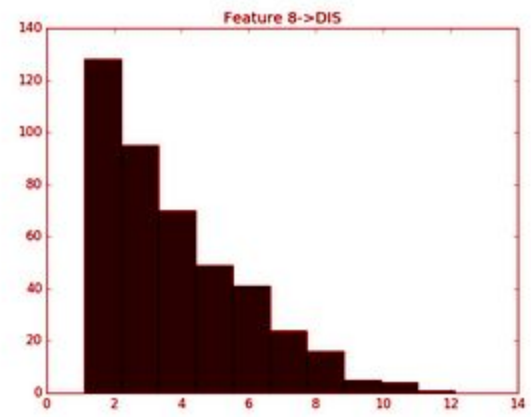
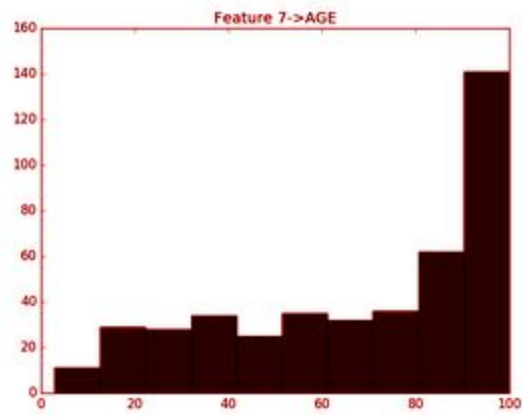
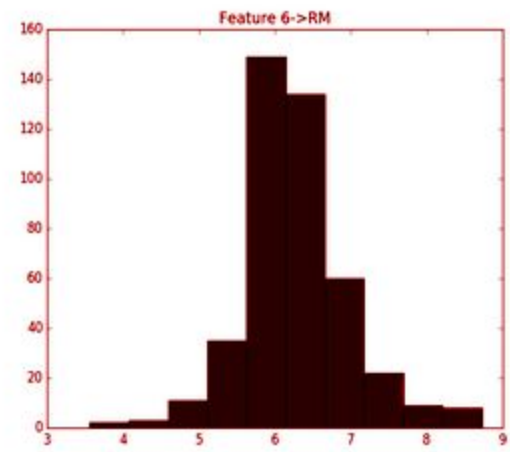
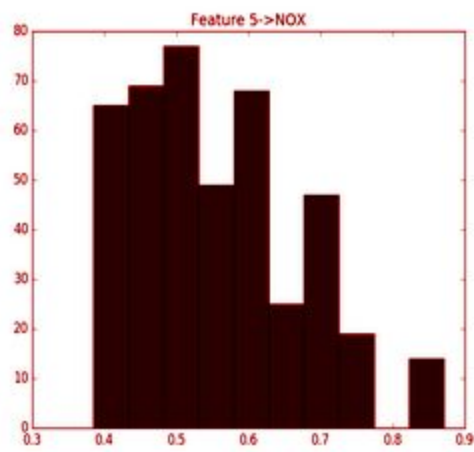
$$p(y = 1|x) = \frac{1}{1+\exp(-\theta^T x)}$$

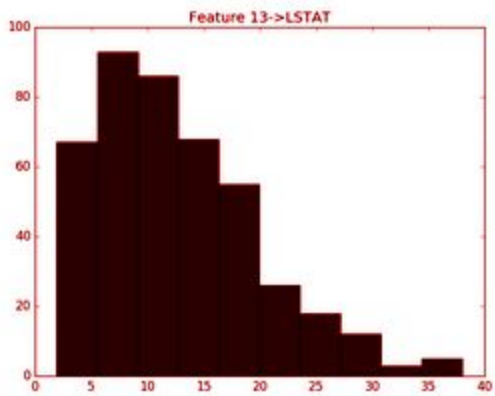
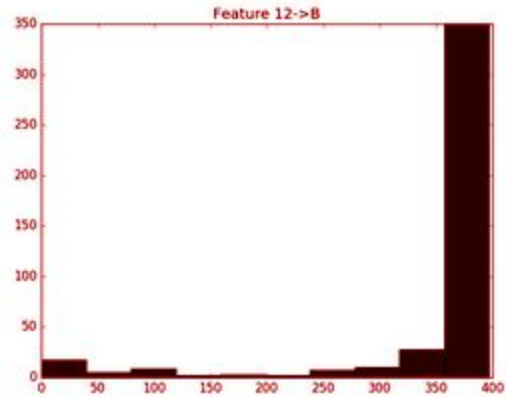
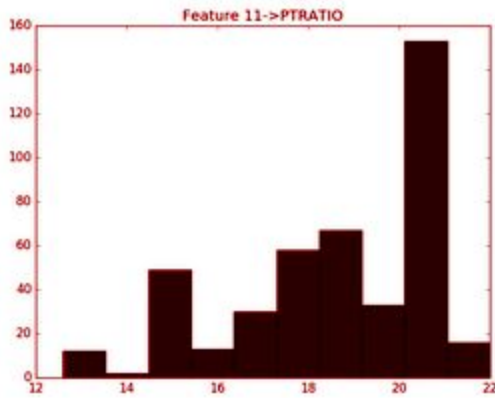
Where,

$$\theta^T x = \theta^0 + w^T x = -\log\left(\frac{(1-p_1)}{(p_1)}\right) - \frac{\mu_1^T(\Sigma)^{-1}\mu_1 - \mu_2^T(\Sigma)^{-1}\mu_2}{2} + [\mu_1^T - \mu_2^T](\Sigma)^{-1}x$$

3) Histograms







Pearson Correlation : We can see from the data that the Pearson Correlation is highest for the feature 'LSTAT' which is -0.073

Programming :

1) Linear Regression and Ridge Regression :

MSE for Training Data : 20.950144508

MSE for testing Data : 28.4179164975

2) Ridge Regression :

MSE for training Data : 20.9502783985

MSE for testing data : 28.4253883068

2) Ridge Regression with Cross Validation - 10 fold :

lambda = 0.0001 Mean CV MSE : 23.5822099346

lambda = 0.001 Mean CV MSE : 23.5821939596
lambda = 0.01 Mean CV MSE : 23.5820343881
lambda = 0.1 Mean CV MSE : 23.5804563838
lambda = 1 Mean CV MSE : 23.5663434685
lambda = 10 Mean CV MSE : 23.5276152939

Here we can see that as Lambda Value increases the MSE decreases. However after a certain point in the curve even though when the lambda increases, the MSE will not decrease. For example in my output for lambda value at 100, the MSE is

lambda = 100 Mean CV MSE : 24.7837340983

Hence, we can conclude that the Lambda lies in the interval [1,10]

Performance :

We can see from our output that the performance of the Linear regression is better than the Ridge Regression in calculating MSE's.

Ridge Regression reduces overfitting and helps in Regularisation when the design matrix has features which are linear combination of each other making the matrix Singular or Non-invertible, however the data set we are provided with does neither has such feature nor does it reduce overfitting hence Ridge Regression does not improve the MSE of Linear Regression.

Feature Selection :

a)

The top 4 features selected are :

[INDUS RM PTRATIO LSTAT]

Train-MSE : 26.4066042155

Test-MSE : 31.4962025449

b)

The top 4 features selected are :

CHAS RM PTRATIO LSTAT

train-MSE : 25.1060222464

test-MSE : 34.6009447175

Selecting one feature at a time improved the MSE

Brute Force :

The top 4 features selected are :

CHAS RM PTRATIO LSTAT

train-MSE : 25.1060222464

test-MSE : 34.6009447175

Brute Force always selects the top 4 features as it exhausts all the combinations and gets the best results.

However Brute force performs solwest as the algorithm is itself brute force.

Feature Expansion :

Train MSE : 5.05978429711

Test MSE : 14.555304

We can see that as we expand the features the number of features, we end up increasing the points in the data space and thus we find better ways to fit the data with the new features which might have not been possible without the expansion of the features.

Thus we can say feature expansion is one to improve the accuracy of our Linear Regression