

SUBHODIP SAHA  
5485317  
HW3

$$f(w) = \frac{1}{n} \sum_{i=1}^n \{-y_i w^T x_i + \log(1 + \exp(w^T x_i))\} + \frac{\lambda}{2} \|w\|_2^2$$

Let's define  $\sigma_w(x_i) = \frac{1}{1 + e^{-w^T x_i}}$

Let's at first prove,  
 $y_i w^T x_i + \log(1 + \exp(w^T x_i)) = -y_i \log \sigma_w(x_i) - (1 - y_i) \log(1 - \sigma_w(x_i))$

then it would be easier to ~~offer~~ calculate gradient.

$$\begin{aligned} & y_i \log \sigma_w(x_i) + (1 - y_i) \log (1 - \sigma_w(x_i)) \\ &= y_i \log \sigma_w(x_i) + \log \sigma_w(x_i) - y_i \log (1 - \sigma_w(x_i)) + \log (1 - \sigma_w(x_i)) \\ &= y_i [\log \sigma_w(x_i) - \log (1 - \sigma_w(x_i))] + \log (1 - \sigma_w(x_i)) \\ &= y_i \frac{\sigma_w(x_i)}{1 - \sigma_w(x_i)} + \log (1 - \sigma_w(x_i)) \end{aligned}$$

$$\sigma_w(x_i) = \frac{1}{1 + e^{-w^T x_i}} = \frac{e^{w^T x_i}}{e^{w^T x_i} + 1}$$

$$1 - \sigma_w(x_i) = \frac{1}{e^{w^T x_i} + 1}$$

$$\frac{\sigma_w(x_i)}{1 - \sigma_w(x_i)} = e^{w^T x_i}$$

$$y_i \frac{\sigma_w(x_i)}{1 - \sigma_w(x_i)} + \log (1 - \sigma_w(x_i)) = y_i w^T x_i - \log (1 + e^{w^T x_i})$$

So, we have proved.

$$\left[ \begin{aligned} & \sum_{i=1}^n -y_i w^T x_i + \log(1 + \exp(w^T x_i)) + \frac{\lambda}{2} \|w\|_2^2 \\ &= -\sum_{i=1}^n y_i \log h_w(x_i) + (1-y_i) \log(1-h_w(x_i)) + \frac{\lambda}{2} \|w\|_2^2 \end{aligned} \right]$$

Now we will differentiate this expression.

$$\begin{aligned} \frac{\partial f(w)}{\partial w_j} &= -\frac{1}{n} \sum_{i=1}^n y_i \frac{\partial \log h_w(x_i)}{\partial w_j} + (1-y_i) \frac{\partial \log(1-h_w(x_i))}{\partial w_j} \\ &= -\frac{1}{n} \sum_{i=1}^n y_i \frac{\frac{\partial h_w(x_i)}{\partial w_j}}{h_w(x_i)} + (1-y_i) \frac{\frac{\partial (1-h_w(x_i))}{\partial w_j}}{1-h_w(x_i)} \\ h_w(x_i) &= \frac{1}{1+e^{-w^T x_i}} \quad \frac{\partial h_w(x_i)}{\partial w_j} = + \frac{e^{-w^T x_i} x_j^i}{(1+e^{-w^T x_i})^2} \\ &= \frac{h_w(x_i)(1-h_w(x_i)) x_j^i}{h_w(x_i)(1-h_w(x_i))} \end{aligned}$$

$$\begin{aligned} &= -\frac{1}{n} \sum_{i=1}^n y_i (1-h_w(x_i)) x_j^i - (1-y_i) h_w(x_i) x_j^i + \lambda w_j \\ &= -\frac{1}{n} \sum_{i=1}^n x_j^i (y_i - y_i h_w(x_i) - h_w(x_i) + y_i h_w(x_i)) + \lambda w_j \\ &= \frac{1}{n} \sum_{i=1}^n (h_w(x_i) - y_i) x_j^i + \lambda w_j \end{aligned}$$

$$\left[ \frac{\partial f(w)}{\partial w_j} = \frac{1}{n} \sum_{i=1}^n (h_w(x_i) - y_i) x_j^i + \lambda w_j \right]$$

Expression for gradient.



$$\Delta w_j = -\eta \frac{\partial f}{\partial w_j}$$

$\eta \rightarrow$  step size

1

$$= +\eta \left[ \frac{1}{n} \sum_{i=1}^n (y_i - h_w(x_i)) x_{ij} - \lambda w_j \right]$$

$$w_j^{\text{new}} = w_j + \eta \left[ \frac{1}{n} \sum_{i=1}^n (y_i - h_w(x_i)) x_{ij} - \lambda w_j \right]$$

$\rightarrow$  we should update the ~~value~~  $w_j$  according to that.

b Let's derive 2nd derivative of  $f(w)$ .

$$\frac{\partial f(w)}{\partial w_j} = \frac{1}{n} \sum_{i=1}^n (h_w(x_i) - y_i) x_{ij} + \lambda w_j$$

$$\frac{\partial^2 f(w)}{\partial w_j \partial w_k} = \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial w_k} (h_w(x_i) - y_i) x_{ij} + \lambda \delta_{jk}$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{\partial h_w(x_i)}{\partial w_k} x_{ij} + \lambda \delta_{jk}$$

we know,  $\frac{\partial h_w(x_i)}{\partial w_k} = h_w(x_i)(1 - h_w(x_i)) x_{ik}$

$$= \frac{1}{n} \sum_{i=1}^n x_{ij} x_{ik}$$

$$\frac{\partial^2 f(w)}{\partial w_j \partial w_k} = \frac{1}{n} \sum_{i=1}^n h_w(x_i)(1 - h_w(x_i)) x_{ik} x_{ij} + \lambda \delta_{jk}$$

$$\frac{\partial^2 f}{\partial w_j \partial w_k} = \frac{1}{n} \sum_{i=1}^n h_w(x_i) (1 - h_w(x_i)) x_i^j x_i^k + \lambda \delta_{jk}$$

→ The 2nd derivative of objective func.

For strong convexity we need to show  
 $\nabla^2 f(w) \succeq \alpha \mathbb{I}$  for some  $\alpha > 0$ .

Let's take  $x$  is 1dim, we can generalize to multidimension.

then,

$$\frac{\partial^2 f}{\partial w^2} = \frac{1}{n} \sum_i h_w(x_i) (1 - h_w(x_i)) x_i^2 + \lambda$$

we know,  $\mathbb{I} \succeq h_w(x_i) \succeq 0$  so,  $\mathbb{I} \succeq 1 - h_w(x_i) \succeq 0$   
 $x_i^2 \succeq 0$   $\lambda > 0$

$$\frac{\partial^2 f}{\partial w^2} = \frac{1}{n} \sum_i \underbrace{h_w(x_i)}_{\succeq 0} \underbrace{(1 - h_w(x_i))}_{\succeq 0} \underbrace{x_i^2}_{\succeq 0} + \underbrace{\lambda}_{> 0}$$

So,  $\frac{\partial^2 f}{\partial w^2} > 0$   ~~$\nabla^2 f(w) \succeq \alpha$~~

Can be written as  $\nabla^2 f(w) \succeq \alpha$

If it is true for 1-dim, it should be true for multidim.

So, the func is strongly convex

$$\boxed{\nabla^2 f(w) - \alpha \mathbb{I} \succeq 0}$$



<sup>b</sup> I am proving the strong convexity of LR using another method, (for any dim - general case)

$$f(w) = \sum_{i=1}^n \underbrace{y_i \log h_w(x_i)}_{\text{Need to prove it is convex func.}} + \underbrace{(1-y_i) \log(1-h_w(x_i))}_{\text{Need to prove it is convex func.}} + \underbrace{\frac{\lambda}{2} \|w\|_2^2}_{\text{Need to prove convex func.}}$$

$$= \sum_{i=1}^n -y_i \log h_w(x_i) - (1-y_i) \log(1-h_w(x_i)) + \frac{\lambda}{2} \|w\|_2^2$$

Let's take,

$$h(z) = \frac{1}{1 + \exp(-z)}$$

$$f_1(z) = -\log(h(z))$$

$$f_2(z) = -\log(1-h(z))$$

$$\frac{d}{dz} f_1(z) = -\log(1/(1+\exp(-z))) = -1 + 1/(1+\exp(-z)) = -1 + h(z)$$

$$\frac{d^2}{dz^2} f_1(z) = \frac{d}{dz} h(z) = h(z)(1-h(z)) \geq 0 \rightarrow f_1(z) \text{ is convex.}$$

$$\frac{d}{dz} f_2(z) = \log(1+\exp(-z)) = f_1(z) + z$$

$$\frac{d^2}{dz^2} f_2(z) = \frac{d^2}{dz^2} f_1(z) = h(z)(1-h(z)) \geq 0 \rightarrow f_2(z) \text{ is also convex.}$$

$$\frac{d^2}{dz^2} \frac{\lambda}{2} \|w\|_2^2 = \frac{\lambda}{2} \rightarrow \text{is also convex.}$$

Now, we need to prove  $f(w)$  is strongly convex.

Let's take,  $g(y) = f(Ay+b)$

$$\nabla_y g(y) = A^T \nabla_x f(Ay+b) \in \mathbb{R}^n$$

$$\nabla_y^2 g(y) = A^T \nabla_y^2 f(Ay+b) A \in \mathbb{R}^{n \times n}$$

$$A \in \mathbb{R}^{m \times n}$$

$$b \in \mathbb{R}^m$$

QED

Since  ~~$f_1 \neq f_2$~~   $f_1 \neq f_2$  (sigmoid) is convex func,  
 having  $\nabla_x^2 f(x) \geq 0$   

$$z^T \nabla_y^2 g(y) z = z^T A^T \nabla_y^2 f(Ay+b) Az$$

$$= (Az)^T \nabla_x^2 f(Ay+b) (Az) \geq 0.$$

So,  $\nabla_y^2 g(y)$  is also positive semi definite matrix.

~~If  $\nabla_w^2 g(w)$  is +ve semi definite, i.e.  $\nabla_y^2 g(y) \geq \alpha$   
 $f(w)$  is strongly convex.~~

If  $\nabla_w^2 g(w)$  is +ve semi-definite, i.e.  $\nabla_w^2 g(w) \geq \alpha$   
 then  $f(w)$  is strongly convex.



The 2nd derivative of objective fun'

$$\frac{\partial^2 f}{\partial w_j \partial w_k} = \frac{1}{n} \sum_{i=1}^n h_w(x_i) (1 - h_w(x_i)) x_i^j x_i^k + \lambda \delta_{jk}$$

To establish smoothness we need to show,  
 $\nabla^2 f(w) \leq \beta \mathbb{I}$  for some  $\beta$   
 $\beta < \infty$

As,  $0 < h_w(x_i) \leq 1$

Then,  $h_w(x_i) (1 - h_w(x_i)) \leq 0.25$  (for  $h_w(x_i) = 0.5$ )

Let's prove for 1 dim.

$$\frac{\partial^2 f}{\partial w^2} = \frac{1}{n} \sum_i h_w(x_i) (1 - h_w(x_i)) x_i^2 + \lambda$$

$$\max_w h_w(x_i) (1 - h_w(x_i)) = 0.25$$

$x_i^2 < \infty$  (as  $x_i$  is finite)

$\lambda < \infty$  ( $\lambda$  is also finite)

$$\frac{\partial^2 f}{\partial w^2} = \frac{1}{n} \sum_i \underbrace{h_w(x_i) (1 - h_w(x_i))}_{\max = 0.25} \underbrace{x_i^2}_{< \infty} + \underbrace{\lambda}_{< \infty}$$

~~So,  $\frac{\partial^2 f}{\partial w^2}$  is bounded & we can always~~

So,  $\frac{\partial^2 f}{\partial w^2}$  is bounded & we can always  
 write  $\frac{\partial^2 f}{\partial w^2} \leq \beta \mathbb{I}$ . As true for 1 dim, it  
 must be true for multidim.

So the fun' is smooth.

$\beta \mathbb{I} - \nabla^2 f(w) \geq 0$

d For smooth, strongly convex func.

$\alpha$  smooth func ~~strongly convex func~~ satisfy,

$$f(w) \geq f(v) + (w-v)^T \nabla f(v) + \frac{\alpha}{2} \|w-v\|_2^2$$

$\beta$  strongly convex func satisfy,

$$f(w) \leq f(v) + (w-v)^T \nabla f(v) + \frac{\beta}{2} \|w-v\|_2^2$$

Initial point  $x_0$ .

$$x_{t+1} = x_t - \eta \nabla f(x_t)$$

$\leadsto$  step size.

For smooth, strongly convex func

$$\boxed{\eta = \frac{2}{\alpha + \beta}}$$

$$\boxed{f(x_t) - f(x^*) \leq \frac{\beta}{2} \exp\left(-\frac{4T}{\frac{\beta}{\alpha} + 1}\right) \|x_0 - x^*\|^2}$$

This is the bound  
on difference.

fixed step size  $\eta$ , <sup>optimization</sup> is substantially faster, exponential.



2 EM algorithm for learning gaussian mixture model: EM is a type of clustering algorithm similar to k-means. Rather than having hard assignment into clusters like k-means we have soft assignment. As a result each gaussian distribution has some responsibility for generating ~~part~~ particular data point.

EM algo in high level.

$$\ln p(x|\theta) = \ln \left\{ \sum_z p(x, z|\theta) \right\}$$

Our goal is to maximize MLE of  $x$  given parameters  $\theta$ . ( $x$  is observed,  $z$  is hidden)

E step: Estimate posterior distribution of responsibilities of each gaussian  $p(z_i|x_i)$  depending on weight ( $\pi$ ), mean ( $\mu$ ) & covariance ( $\Sigma$ ).  
i.e. estimate  $p(z_i|x_i) = f(\pi, \mu, \Sigma)$

M step: Use  $p(z_i|x_i)$  to maximize likelihood w.r.t the parameters  $\theta$   
i.e. maximize  $\ln p(x|\mu, \Sigma, \pi)$

Repeat EM step until converge.

b M step. (calculate mean, covariance, prior)

~~$\mu_h^{new} = \frac{1}{N}$~~  We define,  $n_h = \sum_{i=1}^n p(G_h | x_i)$

$$\mu_h^{new} = \frac{1}{n_h} \sum_{i=1}^n p(G_h | x_i) x_i$$

$$\Sigma_h^{new} = \frac{1}{n_h} \sum_{i=1}^n p(G_h | x_i) (x_i - \mu_h^{new})(x_i - \mu_h^{new})^T$$

$$\pi_h^{new} = \frac{n_h}{n} \quad \text{for all component, } h=1, \dots, k$$

where,  $n_h = \sum_{i=1}^n p(G_h | x_i)$

(posterior prob.)

Gaussian

with  $\mu_h$ ,  
mean covar  $\Sigma_h$ .

c E step.

$$\pi_h \mathcal{N}(x_i | \mu_h, \Sigma_h)$$

$$p(G_h | x_i) = \frac{\pi_h \mathcal{N}(x_i | \mu_h, \Sigma_h)}{\sum_{j=1}^k \pi_j \mathcal{N}(x_i | \mu_j, \Sigma_j)}$$

$$= \frac{\pi_h |\Sigma_h|^{-1/2} \exp \left[ -\frac{1}{2} (x_i - \mu_h)^T \Sigma_h^{-1} (x_i - \mu_h) \right]}{\sum_{j=1}^k \pi_j |\Sigma_j|^{-1/2} \exp \left[ -\frac{1}{2} (x_i - \mu_j)^T \Sigma_j^{-1} (x_i - \mu_j) \right]}$$



## Problem 3

1. Error rates for myLogisticReg2 with Boston 50

F1	F2	F3	F4	F5	Mean	SD
0.17	0.22	0.24	0.22	0.16	0.20	0.03

2. Error rates for myLogisticReg2 with Boston 75

F1	F2	F3	F4	F5	Mean	SD
0.18	0.27	0.24	0.22	0.24	0.23	0.02

3. Error rates for LogisticRegression with Boston 50

F1	F2	F3	F4	F5	Mean	SD
0.12	0.20	0.25	0.25	0.17	0.20	0.04

4. Error rates for LogisticRegression with Boston 75

F1	F2	F3	F4	F5	Mean	SD
0.11	0.11	0.09	0.08	0.09	0.14	0.01