

$$f(w) = \frac{1}{n} \sum_{i=1}^n \{-\eta_i w^T x_i + \log(1 + \exp(w^T x_i))\} + \frac{\lambda}{2} \|w\|^2, \text{ assume } \|x_i\| \leq R$$

$$\nabla_w f(w) = \frac{1}{n} \sum_{i=1}^n (-\eta_i x_i + \frac{x_i \exp(w^T x_i)}{1 + \exp(w^T x_i)}) + \lambda \|w\|$$

$$= \frac{1}{n} \sum_{i=1}^n (\sigma(w^T x_i) - \eta_i) x_i + \lambda \|w\|, \text{ where } \sigma(w^T x_i) = \frac{\exp(w^T x_i)}{1 + \exp(w^T x_i)}$$

Steps: i) initialize  $w_0$  values as random values

(j) = 0, 1, ..., d) where d is the number of features,  $w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$

ii) compute  $\nabla_w f(w)$  as above

iii) compute  $y_{k+1} = w_k - \eta \nabla_w f(w_k)$  until it converges

iv) if  $\|y_{k+1}\| \leq R$ ,  $w_{k+1} = y_{k+1}$ , if  $\|y_{k+1}\| > R$ ,  $w_{k+1} = \frac{R}{\|y_{k+1}\|} y_{k+1}$

b) Yes, for every  $w_1 > w_2$ , to prove  $f(w_1) \geq f(w_2) + (w_1 - w_2)^T \nabla f(w_2) + \frac{\beta}{2} \|w_1 - w_2\|^2$

$$\text{we prove } \frac{1}{n} \sum_{i=1}^n \{-\eta_i w_1^T x_i + \log(1 + \exp(w_1^T x_i))\} + \frac{\lambda}{2} \|w_1\|^2 \geq \frac{1}{n} \sum_{i=1}^n \{-\eta_i w_2^T x_i + \log(1 + \exp(w_2^T x_i))\} + \frac{\lambda}{2} \|w_2\|^2$$

$$+ (w_1 - w_2)^T \cdot \frac{1}{n} \sum_{i=1}^n (\sigma(w_2^T x_i) - \eta_i) x_i + \lambda \|w_1\| + \frac{\lambda}{2} \|w_1 - w_2\|^2$$

To prove  $f(w)$  is strongly convex, we can prove  $f(w) - \frac{\lambda}{2} \|w\|^2$  is convex

$$\nabla^2 (f(w) - \frac{\lambda}{2} \|w\|^2) = \frac{1}{n} \sum_{i=1}^n \left\{ \frac{\exp(w^T x_i)}{(1 + \exp(w^T x_i))^2} \right\} x_i x_i^T \geq 0$$

since  $\exp(w^T x_i) > 0$  and  $x_i \geq 0$  ( $i=0, 1, \dots, n$ )

so the objective function  $f(w)$  is strongly convex

We prove each part of objective function is  $\beta$ -smooth

(c) Yes, There exists  $w_1 > w_2$  where  $f(w_1) \leq f(w_2) + (w_1 - w_2)^T \nabla f(w_2) + \frac{\beta}{2} \|w_1 - w_2\|^2$

if  $\beta > \lambda$  then it is true since  $\log(1 + \exp(x))$  is  $\beta$ -smooth since

$$\log\left(\frac{1 + \exp(x)}{1 + \exp(x)}\right) + (x - x_1)^T \frac{\exp(x_1)}{1 + \exp(x_1)} + \frac{\beta}{2} \|x - x_1\|^2 \geq 0 \text{ if } \beta \geq (x - x_1)^T \frac{\exp(x_1)}{1 + \exp(x_1)} > 0 \text{ and}$$

$\log\left(\frac{1 + \exp(x)}{1 + \exp(x_1)}\right)$  is growing less than quadratic function, and  $-\eta w^T x$  is  $\beta$ -smooth since

$$-\eta(w_1 - w_2)^T x + (w_1 - w_2)^T (-\eta x) + \frac{\beta}{2} \|w_1 - w_2\|^2 \geq 0 \text{ if } \beta \geq \eta x \text{ and obviously } \frac{\lambda}{2} (\|w_1\|^2 - \|w_2\|^2) + (w_1 - w_2)^T \lambda \|w_1\| + \frac{\beta}{2} \|w_1 - w_2\|^2 \geq 0$$

if  $\beta > \lambda$  since  $\|w_1 - w_2\| \geq \|w_2\| - \|w_1\|$  based on triangle inequality.

d) for  $\beta$ -smooth and  $\alpha$ -strongly convex function,

$$\text{strong-convexity, } f(w_T) - f(w^*) \leq (w_T - w^*)^T \nabla f(w^*) + \frac{\beta}{2} \|w_T - w^*\|^2 = \frac{\beta}{2} \|w_T - w^*\|^2 \text{ since } \nabla f(w^*) = 0$$

$$w_{T+1} = w_T - \eta_T \nabla f(w_T), \|w_{T+1} - w^*\|^2 = \|w_T - w^*\|^2 + \eta_T^2 \|\nabla f(w_T)\|^2 - 2\eta_T \nabla f(w_T)^T (w_T - w^*)$$

$$\leq \|w_T - w^*\|^2 + \eta_T^2 \|\nabla f(w_T)\|^2 - 2\eta_T \left( \frac{\alpha \beta}{\alpha + \beta} \|w_T - w^*\|^2 + \frac{1}{\alpha + \beta} \|\nabla f(w_T)\|^2 \right)$$

$$= \left(1 - 2\eta_T \frac{\alpha \beta}{\alpha + \beta}\right) \|w_T - w^*\|^2 + \eta_T \left(\eta_T - \frac{2}{\alpha + \beta}\right) \|\nabla f(w_T)\|^2$$

suppose  $\eta_T \leq \frac{2}{\alpha + \beta}$

$$\|w_{T+1} - w^*\|^2 \leq \left(1 - 2\eta_T \frac{\alpha \beta}{\alpha + \beta}\right) \|w_T - w^*\|^2$$

$$\leq \|w_0 - w^*\|^2 \prod_{t=1}^T \left(1 - 2\eta_t \frac{\alpha \beta}{\alpha + \beta}\right) \text{ and } 1 - x \leq e^{-x}$$

$$\text{so, } f(w_T) - f(w^*) \leq \frac{\beta}{2} \|w_T - w^*\|^2 \leq \frac{\beta}{2} \prod_{t=1}^T \left(1 - \frac{2\eta_t \alpha \beta}{\alpha + \beta}\right) \|w_0 - w^*\|^2 \leq \frac{\beta}{2} \|w_0 - w^*\|^2 e^{-\sum_{t=1}^T \eta_t}$$

initialize  $\Phi^0$  with mean vector of zero vector and covariance matrix of  $I$

2. (a) E-step: assign class labels for every data point by using Bayes' rule to compute the likelihood that it belongs to certain class  $G_i$ , given the data point and the old set of parameters  $\Phi^l$  and set the label  $G_i$  with the largest likelihood  $P(G_i|X^0, \Phi^l) = h_i^0$  to that data point.

M-step: recompute the prior probability  $P(G_i)$ , mean vector and the covariance matrix with the posterior  $h_i^0$ , the data point  $X^0$  and the total number of class  $N$ , update the new set of parameters to be  $\Phi^{l+1}$ , maximum likelihood is used to get new parameters.

keep doing E-step & M-step until  $h_i^0$  converges

$$b) \bar{\lambda}_i^{(L)} = \frac{\sum_{j=1}^N h_j^0}{N}, \quad M_i^{(L)} = \frac{\sum_{j=1}^N h_j^0 X^0}{\sum_{j=1}^N h_j^0}$$

$$\Sigma_i^{(L)} = \frac{\sum_{j=1}^N h_j^0 (X^0 - M_i^{(L)}) (X^0 - M_i^{(L)})^T}{\sum_{j=1}^N h_j^0}$$

where  $L$  is the step number

$$c) h_i^0 \equiv P(\lambda_i | X^0, M_i^L, \Sigma_i^L) = \frac{P(X^0 | \pi_i, M_i^L, \Sigma_i^L) P(\lambda_i)}{\sum_{j=1}^K P(X^0 | \pi_j, M_j^L, \Sigma_j^L) P(\lambda_j)}, \text{ where } L \text{ is step}$$

$$\text{and } P(X^0 | \lambda_i, M_i^L, \Sigma_i^L) = \lambda_i \cdot \frac{1}{(2\pi)^{D/2} |\Sigma_i^L|^{D/2}} \exp\left(-\frac{1}{2} (X^0 - M_i^L)^T \Sigma_i^{L-1} (X^0 - M_i^L)\right)$$

### Q3

#### Summary

I use the gradient descent method with iteration number of 850 and step size of 0.00001 based on my experiment and some suggestions on the textbook.

#### MyLogisticReg2 with Boston50

K=0	K=1	K=2	K=3	K=4	Mean	Std
0.207920792079	0.188118811881	0.128712871287	0.247524752475	0.147058823529	0.18386721025	0.0425342027945

#### MyLogisticReg2 with Boston75

K=0	K=1	K=2	K=3	K=4	Mean	Std
0.207920792079	0.178217821782	0.227722772277	0.138613861386	0.107843137255	0.172063676956	0.0439652339392

#### LogisticRegression with Boston50

K=0	K=1	K=2	K=3	K=4	Mean	Std
0.128712871287	0.108910891089	0.0891089108911	0.277227722772	0.117647058824	0.144321490973	0.0677075481368

#### LogisticRegression with Boston75

K=0	K=1	K=2	K=3	K=4	Mean	Std
0.0891089108911	0.128712871287	0.138613861386	0.108910891089	0.0490196078431	0.102873228499	0.0318471419286

### Q4

#### Summary

The feature matrix gets preprocessed by a standardized method in a way that for each feature column vector, every element gets subtracted by the mean of the column and divided by the standard deviation of the column so that every feature vector has mean of zero and standard deviation of one.

The gradient descent method with iteration number of 300 and step size of 0.00001 based on my experiment and some suggestions on the textbook. And I consider the lambda value of one to penalize big weight in each step.

#### MyLogisticRegGen with Digits

K=0	K=1	K=2	K=3	K=4	Mean	Std
0.075208913649	0.122562674095	0.0611111111111	0.050139275766	0.0388888888889	0.0695821727019	0.0290876487138

#### LogisticRegression with Digits

K=0	K=1	K=2	K=3	K=4	Mean	Std
0.0640668523677	0.116991643454	0.0611111111111	0.0389972144847	0.0166666666667	0.0595666976168	0.0334178913221