

CSE317A

Machine Learning (SP18)

Written Homework 1.

Chunqun Li, Jianui Xing

Problem 1.

$$\begin{aligned} (a) \because L(\vec{w}) &= \sum_{i=1}^n (\vec{w}^T \vec{x}_i - y_i)^2 + \lambda \|\vec{w}\|_2^2 \\ &= \sum_{i=1}^n [(\vec{w}^T \vec{x}_i)^2 - 2y_i \vec{w}^T \vec{x}_i + y_i^2] + \lambda \vec{w}^T \vec{w} \\ &= \sum_{i=1}^n [\vec{w}^T \vec{x}_i \vec{x}_i^T \vec{w} - 2y_i \vec{w}^T \vec{x}_i + y_i^2] + \lambda \vec{w}^T \vec{w} \end{aligned}$$

$$\begin{aligned} \therefore \frac{\partial}{\partial \vec{w}} L(\vec{w}) &= \sum_{i=1}^n (2 \vec{x}_i \vec{x}_i^T \vec{w} - 2y_i \vec{x}_i) + 2\lambda \vec{w} \\ &= \sum_{i=1}^n 2 \vec{x}_i (\vec{x}_i^T \vec{w} - y_i) + 2\lambda \vec{w} \end{aligned}$$

(b) Note $\nabla \vec{w}$:

$$\nabla \vec{w} = \frac{\partial}{\partial \vec{w}} \|\vec{w}\|$$

$$\nabla w_d = \frac{\partial}{\partial \vec{w}} \|w_d\| = \frac{\partial}{\partial w_d} |w_d| = \begin{cases} 1, & w_d > 0 \\ -1, & w_d < 0 \\ \text{not defined}, & w_d = 0 \end{cases}, d=1, \dots, d$$

Since in practice the chance that $w_d = 0$ is very rare, we can define $\nabla \vec{w}$:

$$\nabla w_d = \begin{cases} 1, & w_d \geq 0 \\ -1, & w_d \leq 0 \end{cases}, d=1, 2, \dots, d$$

Therefore,

$$\frac{\partial}{\partial \vec{w}} L(\vec{w}) = \sum_{i=1}^n 2 \vec{x}_i (\vec{x}_i^T \vec{w} - y_i) + \lambda \nabla \vec{w},$$

$$\text{where } \nabla \vec{w} = \nabla w_d = \begin{cases} 1, & w_d \geq 0 \\ -1, & w_d < 0 \end{cases}, d=1, \dots, d$$

$$\begin{aligned} (c) \frac{\partial}{\partial \vec{w}} L(\vec{w}) &= \sum_{i=1}^n \frac{1}{1 + e^{-y_i \vec{w}^T \vec{x}_i}} \cdot e^{-y_i \vec{w}^T \vec{x}_i} \cdot (-y_i \vec{x}_i) \\ &= - \sum_{i=1}^n \frac{e^{-y_i \vec{w}^T \vec{x}_i}}{1 + e^{-y_i \vec{w}^T \vec{x}_i}} y_i \vec{x}_i \\ &= - \sum_{i=1}^n \frac{y_i \vec{x}_i}{1 + e^{y_i \vec{w}^T \vec{x}_i}} \end{aligned}$$

$$\begin{aligned} (d) \because L(\vec{w}) &= C \sum_{i=1}^n \max\{1 - y_i \vec{w}^T \vec{x}_i, 0\} + \lambda \|\vec{w}\|_2^2 \\ &= C \sum_{i=1}^n \begin{cases} 1 - y_i \vec{w}^T \vec{x}_i, & y_i \vec{w}^T \vec{x}_i < 1 \\ 0, & y_i \vec{w}^T \vec{x}_i \geq 1 \end{cases} + \lambda \vec{w}^T \vec{w} \end{aligned}$$

$$\therefore \frac{\partial}{\partial \vec{w}} L(\vec{w}) = C \sum_{i=1}^n \begin{cases} -y_i \vec{x}_i, & y_i \vec{w}^T \vec{x}_i < 1 \\ 0, & y_i \vec{w}^T \vec{x}_i \geq 1 \end{cases} + 2\lambda \vec{w}$$

Problem 2.

$$(a) \therefore P_r(y=1|x) = \frac{e^{\vec{w}^T \vec{x}}}{1 + e^{\vec{w}^T \vec{x}}} = \text{Sigm}(\vec{w}^T \vec{x})$$

$$P_r(y=0|x) = \frac{1}{1 + e^{\vec{w}^T \vec{x}}} = 1 - \text{Sigm}(\vec{w}^T \vec{x})$$

$$\therefore P_r(y|x) = [\text{Sigm}(\vec{w}^T \vec{x})]^y \cdot [1 - \text{Sigm}(\vec{w}^T \vec{x})]^{1-y}$$

$$\therefore L(\vec{w}) = -\log[P_r(Y|X)]$$

$$= -\log\left[\prod_{i=1}^n P_r(y_i|x_i)\right]$$

$$= -\sum_{i=1}^n \log[P_r(y_i|x_i)]$$

$$= -\sum_{i=1}^n \left[\log[\text{Sigm}(\vec{w}^T \vec{x}_i)]^{y_i} + \log[1 - \text{Sigm}(\vec{w}^T \vec{x}_i)]^{1-y_i} \right]$$

$$= -\sum_{i=1}^n \left[y_i \log[\text{Sigm}(\vec{w}^T \vec{x}_i)] + (1-y_i) \log[1 - \text{Sigm}(\vec{w}^T \vec{x}_i)] \right]$$

$$(b) \frac{\partial}{\partial \vec{w}} L = -\sum_{i=1}^n \left[y_i \cdot \frac{1}{\text{Sigm}(\vec{w}^T \vec{x}_i)} \cdot \text{Sigm}(\vec{w}^T \vec{x}_i) [1 - \text{Sigm}(\vec{w}^T \vec{x}_i)] \cdot \vec{x}_i \right. \\ \left. - (1-y_i) \frac{1}{1 - \text{Sigm}(\vec{w}^T \vec{x}_i)} \text{Sigm}(\vec{w}^T \vec{x}_i) [1 - \text{Sigm}(\vec{w}^T \vec{x}_i)] \cdot \vec{x}_i \right]$$

$$= -\sum_{i=1}^n \left[y_i \vec{x}_i [1 - \text{Sigm}(\vec{w}^T \vec{x}_i)] - (1-y_i) \vec{x}_i \text{Sigm}(\vec{w}^T \vec{x}_i) \right]$$

$$= -\sum_{i=1}^n \left[y_i \vec{x}_i - \vec{x}_i \text{Sigm}(\vec{w}^T \vec{x}_i) \right]$$

$$= -\sum_{i=1}^n \vec{x}_i [y_i - \text{Sigm}(\vec{w}^T \vec{x}_i)]$$

(c) By definition,

$$H = \frac{\partial}{\partial \vec{w}} L(\vec{w})$$

$$= \frac{\partial}{\partial \vec{w}} \left(\frac{\partial}{\partial \vec{w}} L(\vec{w}) \right)$$

$$= \frac{\partial}{\partial \vec{w}} (-1) \sum_{i=1}^n \vec{x}_i [y_i - \text{Sigm}(\vec{w}^T \vec{x}_i)]$$

$$= \sum_{i=1}^n \vec{x}_i \frac{\partial}{\partial \vec{w}} \text{Sigm}(\vec{w}^T \vec{x}_i)$$

$$= \sum_{i=1}^n \vec{x}_i \cdot \text{Sigm}(\vec{w}^T \vec{x}_i) [1 - \text{Sigm}(\vec{w}^T \vec{x}_i)] \vec{x}_i^T$$

$$= \sum_{i=1}^n \vec{x}_i^T \vec{x}_i \text{Sigm}(\vec{w}^T \vec{x}_i) [1 - \text{Sigm}(\vec{w}^T \vec{x}_i)]$$

$$X^T W X = \sum_{i=1}^n \vec{x}_i^T W_{ii} \vec{x}_i$$

$$= \sum_{i=1}^n \vec{x}_i^T \vec{x}_i \text{Sigm}(\vec{w}^T \vec{x}_i) [1 - \text{Sigm}(\vec{w}^T \vec{x}_i)]$$

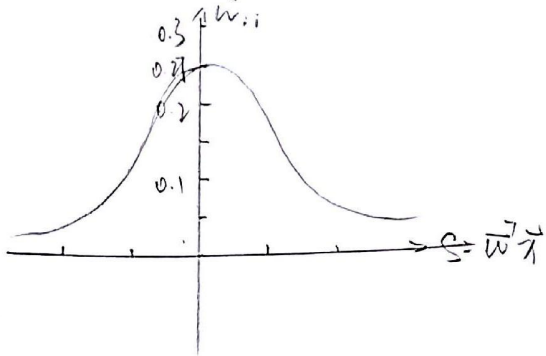
$$= H$$

Problem 2 (continue)

(c) (cont.)

Note $\text{sigm}(\cdot)$ as $\theta(\cdot)$, $\vec{w}^T \vec{x}_i$ as s ,

$$\begin{aligned} \text{then: } W_{ii} &= \theta(s) [1 - \theta(s)] \\ &= \frac{e^s}{1+e^s} \cdot \frac{1}{1+e^s} \\ &= \frac{e^s}{(1+e^s)^2} \end{aligned}$$



From the plot we can see that W_{ii} gets larger when $|\vec{w}^T \vec{x}_i| \rightarrow 0$, and gets

smaller when $|\vec{w}^T \vec{x}_i| \rightarrow \infty$

$$\text{Since } P(y_i=1 | x_i) = \frac{e^{\vec{w}^T \vec{x}_i}}{1 + e^{\vec{w}^T \vec{x}_i}}$$

$$P(y_i=0 | \vec{x}_i) = \frac{1}{1 + e^{\vec{w}^T \vec{x}_i}},$$

$$\text{When } \vec{w}^T \vec{x}_i = 0, P(y_i=1 | \vec{x}_i) = P(y_i=0 | \vec{x}_i) = \frac{1}{2}$$

which means \vec{x}_i is "ambiguous" or "hard to classify". In other words,

W_{ii} is larger for examples harder to be classified and smaller for examples easier to be classified.

(e) By definition of Newton's method:

$$\vec{w}_{\text{new}} = \vec{w} - \frac{\frac{\partial}{\partial \vec{w}} L(\vec{w})}{H}$$

$$= \vec{w} + \frac{\sum_{i=1}^n \vec{x}_i [y_i - \text{sigm}(\vec{w}^T \vec{x}_i)]}{H}$$

$$= \frac{\sum_{i=1}^n \vec{x}_i [y_i - \text{sigm}(\vec{w}^T \vec{x}_i)] + \vec{X}^T \vec{W} \vec{X} \cdot \vec{w}}{\vec{X}^T \vec{W} \vec{X}}$$

$$= \frac{\sum_{i=1}^n \vec{x}_i [y_i - \text{sigm}(\vec{w}^T \vec{x}_i)] + \sum_{i=1}^n \vec{x}_i \cdot \vec{x}_i^T W_{ii} \cdot \vec{w}}{\vec{X}^T \vec{W} \vec{X}}$$

$$= \frac{\sum_{i=1}^n \vec{x}_i [y_i - \text{sigm}(\vec{w}^T \vec{x}_i) + W_{ii} \vec{x}_i^T \vec{w}]}{\vec{X}^T \vec{W} \vec{X}}$$

$$= \frac{\sum_{i=1}^n \vec{x}_i W_{ii} \left[\frac{1}{W_{ii}} (y_i - \text{sigm}(\vec{w}^T \vec{x}_i)) + \vec{x}_i^T \vec{w} \right]}{\vec{X}^T \vec{W} \vec{X}}$$

$$= \frac{\sum_{i=1}^n \vec{x}_i W_{ii} \vec{z}}{\vec{X}^T \vec{W} \vec{X}}$$

$$= \frac{\vec{X} \vec{W} \vec{z}}{\vec{X}^T \vec{W} \vec{X}}$$

$$= (\vec{X}^T \vec{W} \vec{X})^{-1} \vec{X} \vec{W} \vec{z}$$

Problem 3

$$(a) L(\vec{w}) = (\vec{w}^T X - Y) P (\vec{w}^T X - Y)^T + \lambda \vec{w}^T \vec{w}$$

b) Note $L_1(\vec{w})$:

$$\begin{aligned} L_1(\vec{w}) &= (\vec{w}^T X - Y) P (\vec{w}^T X - Y)^T \\ &= (\vec{w}^T X P - Y P) (X^T \vec{w} - Y^T) \\ &= \vec{w}^T X X^T \vec{w} - \vec{w}^T X P Y^T - Y P X^T \vec{w} + Y P Y^T \\ &= \vec{w}^T X X^T \vec{w} - 2 \vec{w}^T X P Y^T + Y P Y^T \end{aligned}$$

$$\frac{\partial}{\partial \vec{w}} L_1(\vec{w}) = 2 X P X^T \vec{w} - 2 X P Y^T$$

$$\begin{aligned} \therefore \frac{\partial}{\partial \vec{w}} L(\vec{w}) &= \frac{\partial}{\partial \vec{w}} L_1(\vec{w}) + 2 \lambda \vec{w} \\ &= 2 X P X^T \vec{w} - 2 X P Y^T + 2 \lambda \vec{w} \\ &= 2 (X P X^T + \lambda I) \vec{w} - 2 X P Y^T = 0 \\ \Rightarrow \vec{w} &= (X P X^T + \lambda I)^{-1} X P Y^T \end{aligned}$$

(c) Comparing the equation above and that in Problem 2(c), we can see that using Newton's method to solve logistic regression is equal to solve a weighted ridge regression at each iteration, where the weight p_i is given by w_{i1} and examples that are harder to be classified will get larger weights.