

CSE 517A Homework 1

$$1. (a) \quad XX^T = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \end{pmatrix} \begin{pmatrix} x_{11} & x_{21} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{1n} & x_{2n} \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^n x_{ii}^2 & \sum_{i=1}^n x_{1i}x_{2i} \\ \sum_{i=1}^n x_{1i}x_{2i} & \sum_{i=1}^n x_{2i}^2 \end{pmatrix}$$

$$w^T X X^T w = (w_1, w_2) \begin{pmatrix} \sum_{i=1}^n x_{ii}^2 & \sum_{i=1}^n x_{1i}x_{2i} \\ \sum_{i=1}^n x_{1i}x_{2i} & \sum_{i=1}^n x_{2i}^2 \end{pmatrix} \begin{pmatrix} w_1 \\ w_2 \end{pmatrix}$$

$$= \left(w_1, \sum_{i=1}^n x_{ii}^2 + w_2 \sum_{i=1}^n x_{1i}x_{2i} \quad w_1 \sum_{i=1}^n x_{1i}x_{2i} + w_2 \sum_{i=1}^n x_{2i}^2 \right) \begin{pmatrix} w_1 \\ w_2 \end{pmatrix}$$

$$= w_1^2 \sum_{i=1}^n x_{ii}^2 + w_1 w_2 \sum_{i=1}^n x_{1i}x_{2i} + w_1 w_2 \sum_{i=1}^n x_{1i}x_{2i} + w_2^2 \sum_{i=1}^n x_{2i}^2$$

$$= w_1^2 \left(\sum_{i=1}^n x_{ii}^2 \right) + 2w_1 w_2 \left(\sum_{i=1}^n x_{1i}x_{2i} \right) + w_2^2 \left(\sum_{i=1}^n x_{2i}^2 \right)$$

$$(b) \quad (X^T w)^T = \left(\begin{pmatrix} x_{11} & x_{21} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{1n} & x_{2n} \end{pmatrix} \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} \right)^T = \begin{pmatrix} w_1 x_{11} + w_2 x_{21} \\ \vdots \\ w_1 x_{1n} + w_2 x_{2n} \end{pmatrix}^T = (w_1 x_{11} + w_2 x_{21}, \dots, w_1 x_{1n} + w_2 x_{2n})$$

$$w^T X = (w_1, w_2) \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \end{pmatrix} = (w_1 x_{11} + w_2 x_{21}, \dots, w_1 x_{1n} + w_2 x_{2n})$$

Thus, $(X^T w)^T = w^T X$

$$(c) \quad \sum_{i=1}^n x_{ii} y_i = \sum_{i=1}^n (x_{ii} y_i) = \left(\sum_{i=1}^n x_{ii} y_i \right)$$

$$Xy = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \left(\sum_{i=1}^n x_{ii} y_i \right)$$

Thus, $\sum_{i=1}^n x_{ii} y_i = Xy$

$$(d) \quad \sum_{i=1}^n x_{ii} x_{ii}^T w = \sum_{i=1}^n (x_{ii}) (x_{ii} x_{ii}) \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} = \sum_{i=1}^n \begin{pmatrix} x_{ii}^2 & x_{ii} x_{ii} \\ x_{ii} x_{ii} & x_{ii}^2 \end{pmatrix} \begin{pmatrix} w_1 \\ w_2 \end{pmatrix}$$

$$= \sum_{i=1}^n \begin{pmatrix} w_1 x_{ii}^2 + w_2 x_{ii} x_{ii} \\ w_1 x_{ii} x_{ii} + w_2 x_{ii}^2 \end{pmatrix} = \begin{pmatrix} w_1 \sum_{i=1}^n x_{ii}^2 + w_2 \sum_{i=1}^n x_{ii} x_{ii} \\ w_1 \sum_{i=1}^n x_{ii} x_{ii} + w_2 \sum_{i=1}^n x_{ii}^2 \end{pmatrix}$$

$$XX^T w = \begin{pmatrix} \sum_{i=1}^n x_{ii}^2 & \sum_{i=1}^n x_{ii} x_{ii} \\ \sum_{i=1}^n x_{ii} x_{ii} & \sum_{i=1}^n x_{ii}^2 \end{pmatrix} \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} = \begin{pmatrix} w_1 \sum_{i=1}^n x_{ii}^2 + w_2 \sum_{i=1}^n x_{ii} x_{ii} \\ w_1 \sum_{i=1}^n x_{ii} x_{ii} + w_2 \sum_{i=1}^n x_{ii}^2 \end{pmatrix}$$

Thus, $\sum_{i=1}^n x_{ii} x_{ii}^T w = XX^T w$

$$(e) \quad \mathcal{L}(w) = \sum_{i=1}^n (w^T x_i - y_i)^2 = \sum_{i=1}^n ((w^T x_i)^2 - 2w^T x_i y_i + y_i^2)$$

$$= \sum_{i=1}^n w^T x_i x_i^T w - \sum_{i=1}^n 2 w^T x_i y_i + \sum_{i=1}^n y_i^2$$

$$= w^T \left(\sum_{i=1}^n x_i x_i^T \right) w - 2 w^T \left(\sum_{i=1}^n x_i y_i \right) + \sum_{i=1}^n y_i^2$$

$$= w^T X X^T w - 2 w^T X y + y^T y$$

$$= (X^T w - y)^T (X^T w - y)$$

$$2. (a) L(w) = \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|_2^2$$

$$= (X^T w - y)^T (X^T w - y) + \lambda w^T w$$

$$\frac{\partial L(w)}{\partial w} = \frac{\partial}{\partial w} (w^T X X^T w - 2w^T X y + y^T y + \lambda w^T w)$$

$$= 2X X^T w - 2X y + 2\lambda w = 2X(X^T w - y) + 2\lambda w$$

$$w_{t+1} = w_t - C \frac{\partial L(w_t)}{\partial w_t} = w_t - C (2X(X^T w_t - y) + 2\lambda w_t)$$

$$(b) L(w) = \sum_{i=1}^n (w^T x_i - y_i)^2 + \lambda \|w\|_1$$

$$= w^T X X^T w - 2w^T X y + y^T y + \lambda \|w\|_1$$

$$\frac{\partial \|w\|_1}{\partial w} = \text{sign}(w) = \begin{cases} +1, & w > 0 \\ -1, & w < 0 \\ \text{any value in } [-1, +1], & w = 0 \end{cases}$$

$$\frac{\partial L(w)}{\partial w} = 2X X^T w - 2X y + \text{sign}(w) = 2X(X^T w - y) + \text{sign}(w)$$

$$w_{t+1} = w_t - C (2X(X^T w_t - y) + \text{sign}(w_t))$$

$$(c) L(w) = \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$$

$$\frac{\partial}{\partial w} \log(1 + \exp(-y_i w^T x_i)) = \frac{1}{1 + \exp(-y_i w^T x_i)} \cdot \exp(-y_i w^T x_i) \cdot (-y_i x_i) = \frac{-y_i x_i}{\exp(y_i w^T x_i) + 1}$$

$$\frac{\partial L(w)}{\partial w} = \sum_{i=1}^n \frac{\partial}{\partial w} \log(1 + \exp(-y_i w^T x_i))$$

$$= \sum_{i=1}^n \frac{-y_i x_i}{\exp(y_i w^T x_i) + 1}$$

$$w_{t+1} = w_t - C \sum_{i=1}^n \frac{-y_i x_i}{\exp(y_i w^T x_i) + 1}$$

$$(d) L(w) = C \sum_{i=1}^n \max\{1 - y_i w^T x_i, 0\} + \|w\|_2^2$$

$$1) \text{ when } 1 - y_i w^T x_i > 0, \max\{1 - y_i w^T x_i, 0\} = 1 - y_i w^T x_i, \frac{\partial}{\partial w} (1 - y_i w^T x_i) = -y_i x_i$$

$$2) \text{ when } 1 - y_i w^T x_i < 0, \max\{1 - y_i w^T x_i, 0\} = 0, \frac{\partial}{\partial w} (0) = 0$$

$$\text{Thus, } \frac{\partial}{\partial w} \sum_{i=1}^n \max\{1 - y_i w^T x_i, 0\} = \sum_{i=1}^n \mathbb{1}(1 - y_i w^T x_i > 0) \cdot (-y_i x_i)$$

$$\frac{\partial L(w)}{\partial w} = \frac{\partial}{\partial w} \left(C \sum_{i=1}^n \max\{1 - y_i w^T x_i, 0\} + \|w\|_2^2 \right)$$

$$= C \sum_{i=1}^n \frac{\partial}{\partial w} \max\{1 - y_i w^T x_i, 0\} + 2w$$

$$= C \sum_{i=1}^n \mathbb{1}(1 - y_i w^T x_i > 0) (-y_i x_i) + 2w$$

$$w_{t+1} = w_t - C \left[C \sum_{i=1}^n \mathbb{1}(1 - y_i w_t^T x_i > 0) (-y_i x_i) + 2w_t \right]$$

3. (a) In logistic regression with $y_i \in \{0, 1\}$, we model

$$P(y_i=1 | x_i) = \text{sigm}(w^T x_i) \quad P(y_i=0 | x_i) = 1 - \text{sigm}(w^T x_i)$$

$$\text{Then, we have } P(y_i | x_i) = (\text{sigm}(w^T x_i))^{y_i} (1 - \text{sigm}(w^T x_i))^{1-y_i}$$

Given a dataset $\{(x_i, y_i)\}_{i=1}^n$, then we have

$$\begin{aligned} \log P(\{y_i\} | \{x_i\}) &= \sum_{i=1}^n \log [(\text{sigm}(w^T x_i))^{y_i} (1 - \text{sigm}(w^T x_i))^{1-y_i}] \\ &= \sum_{i=1}^n (y_i \log(\text{sigm}(w^T x_i)) + (1-y_i) \log(1 - \text{sigm}(w^T x_i))) \end{aligned}$$

$$\text{Thus, } L(w) = - \sum_{i=1}^n (y_i \log(\text{sigm}(w^T x_i)) + (1-y_i) \log(1 - \text{sigm}(w^T x_i)))$$

$$\begin{aligned} (b) \quad \frac{\partial L(w)}{\partial w} &= \frac{\partial}{\partial w} \left[- \sum_{i=1}^n (y_i \log(\text{sigm}(w^T x_i)) + (1-y_i) \log(1 - \text{sigm}(w^T x_i))) \right] \\ &= - \sum_{i=1}^n \frac{\partial}{\partial w} (y_i \log(\text{sigm}(w^T x_i)) + (1-y_i) \log(1 - \text{sigm}(w^T x_i))) \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial w} (y_i \log(\text{sigm}(w^T x_i))) &= y_i \frac{1}{\text{sigm}(w^T x_i)} \cdot \frac{\partial \text{sigm}(w^T x_i)}{\partial w} \\ &= y_i \cdot \frac{1}{\text{sigm}(w^T x_i)} \cdot \text{sigm}(w^T x_i)(1 - \text{sigm}(w^T x_i)) \cdot x_i \\ &= y_i (1 - \text{sigm}(w^T x_i)) \cdot x_i \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial w} (1-y_i) \log(1 - \text{sigm}(w^T x_i)) &= -(1-y_i) \frac{1}{1 - \text{sigm}(w^T x_i)} \cdot \frac{\partial \text{sigm}(w^T x_i)}{\partial w} \\ &= -(1-y_i) \frac{1}{1 - \text{sigm}(w^T x_i)} \cdot \text{sigm}(w^T x_i)(1 - \text{sigm}(w^T x_i)) \cdot x_i \\ &= -(1-y_i) \text{sigm}(w^T x_i) \cdot x_i \end{aligned}$$

$$\text{Thus, } \frac{\partial}{\partial w} L(w) = - \sum_{i=1}^n y_i (1 - \text{sigm}(w^T x_i)) \cdot x_i - (1-y_i) \cdot \text{sigm}(w^T x_i) \cdot x_i$$

$$\begin{aligned} &= - \sum_{i=1}^n y_i x_i - y_i \text{sigm}(w^T x_i) x_i + y_i \text{sigm}(w^T x_i) x_i - \text{sigm}(w^T x_i) x_i \\ &= - \sum_{i=1}^n (y_i - \text{sigm}(w^T x_i)) x_i \end{aligned}$$

$$\begin{aligned} (c) \quad H &= \frac{\partial^2}{\partial w^2} L(w) = \frac{\partial}{\partial w} \sum_{i=1}^n (\text{sigm}(w^T x_i) - y_i) x_i \\ &= \sum_{i=1}^n x_i \text{sigm}(w^T x_i) (1 - \text{sigm}(w^T x_i)) x_i^T \\ &= \sum_{i=1}^n x_i W_{ii} x_i^T = X W X^T \end{aligned}$$

Since $W_{ii} = \text{sigm}(w^T x_i)(1 - \text{sigm}(w^T x_i))$, when $\text{sigm}(w^T x_i) = \frac{1}{2}$, i.e. $w^T x_i = 0$, W_{ii} achieves the largest value $\frac{1}{4}$

when $\text{sigm}(w^T x_i) = 0$ or 1 i.e. $w^T x_i \rightarrow -\infty$ or $+\infty$, W_{ii} achieves the smallest value 0

$$(d) \quad V \in \mathbb{R}^{d \times d}, \quad V^T H V = V^T X^T W X V = (X_V)^T W (X_V)$$

since W is diagonal matrix and all diagonal entries $W_{ii} = \text{sign}(w^T x_i) / (1 - \text{sign}(w^T x_i))$

which are non-negative, it follows that $(X_V)^T W (X_V) \geq 0$

Thus, $V^T H V \geq 0$ which means H is positive semi-definite

$$(e) \quad \text{Newton's method } w_{\text{new}} = w - H^{-1} \frac{\partial L}{\partial w}$$

$$= w + (X^T W X)^{-1} \sum_{i=1}^n (y_i - \text{sign}(w^T x_i)) x_i$$

since $\varepsilon_i = x_i^T w + \frac{1}{W_{ii}} (y_i - \text{sign}(w^T x_i))$, then we have $(y_i - \text{sign}(w^T x_i)) = W_{ii} (\varepsilon_i - x_i^T w)$

$$\text{Thus, } w_{\text{new}} = w + (X^T W X)^{-1} \sum_{i=1}^n W_{ii} (\varepsilon_i - x_i^T w) x_i$$

$$= w + (X^T W X)^{-1} \left[\sum_{i=1}^n x_i^T W_{ii} \varepsilon_i - \sum_{i=1}^n x_i^T W_{ii} x_i w \right]$$

$$= w + (X^T W X)^{-1} (X W \varepsilon - X^T W X w)$$

$$= w + (X^T W X)^{-1} X W \varepsilon - w = (X^T W X)^{-1} X W \varepsilon$$

$$4(a) \quad \mathcal{L}(w) = \sum_{i=1}^n p_i (w^T x_i - y_i)^2 + \lambda w^T w$$

$$= (X_w^T - y)^T P (X_w^T - y) + \lambda w^T w$$

$$(b) \quad \frac{\partial \mathcal{L}(w)}{\partial w} = \frac{\partial}{\partial w} [(X_w^T - y)^T P (X_w^T - y) + \lambda w^T w]$$

$$= \frac{\partial}{\partial w} ((w^T X - y^T) P (X_w^T - y)) + 2\lambda w$$

$$= \frac{\partial}{\partial w} [w^T X P X^T w - w^T X P y - y^T P X^T w + y^T P y] + 2\lambda w$$

$$= \frac{\partial}{\partial w} [w^T X P X^T w - 2w^T X P y + y^T P y] + 2\lambda w$$

$$= X P X^T w + (X P X^T)^T w - 2 X P y + 2\lambda w$$

$$= (2 X P X^T + 2\lambda I) w - 2 X P y = 0$$

$$\Rightarrow w^* = (X P X^T + \lambda I)^{-1} X P y$$

$$(c) \quad \text{Let } \lambda = 0, \quad w^* = (X P X^T)^{-1} X P y$$

$$\text{In problem 3(e)}, \quad w_{\text{new}} = (X W X^T)^{-1} X W z$$

Newton's Method can be shown to solve each iteration by fitting a new weighted least squares problem, where the weights get updated at every iteration. Because one is reweighting the least squares objective at each step, that procedure is known as iteratively reweighted least squares.

5. (a) Using $y \in \{-1, +1\}$ is a good choice for implementation project 1 because it simplifies the mathematics, particularly in gradient calculations. And this label choice directly reflects the directionality of classification, enabling the model to more clearly learn the boundary

(b) For regression tasks, the 'true' regression loss is often considered to be the Mean Squared Error ($MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$). This is because it directly measures the average squared difference between the predicted and true values. MSE is widely used due to its mathematical properties, interpretability and its strong connection to statistical models. Other loss functions, like Absolute-loss and Huber-loss may serve as surrogates to handle specific scenarios.