

DDA2020 Assignment 1

Name: 费祥 Student ID: 120090414

T1.

$$\begin{aligned} KL(P_{emp} || q) &= \sum_{x \in X} P_{emp}(x) \log \frac{P_{emp}(x)}{q(x|\theta)} \\ &= \sum_{x \in X} P_{emp}(x) \log P_{emp}(x) - \sum_{x \in X} P_{emp}(x) \log q(x|\theta) \end{aligned}$$

$$\begin{aligned} \operatorname{argmin}_\theta KL(P_{emp} || q) &= \operatorname{argmin}_\theta - \sum_{x \in X} P_{emp}(x) \log q(x|\theta) \\ &= \operatorname{argmax}_\theta \sum_{x \in X} P_{emp}(x) \log q(x|\theta) \end{aligned}$$

so we just need to show that:

$$\operatorname{argmax}_\theta \sum_{x \in X} P_{emp}(x) \log q(x|\theta) = \operatorname{argmax}_\theta \sum_{x \in X} \log q(x|\theta).$$

Suppose n_x is the number of x samples, N is the total number, by the definition of empirical distribution, we have:

$$P_{emp}(x) = \frac{n_x}{N}$$

$$\begin{aligned} \therefore \operatorname{argmax}_\theta \sum_{x \in X} P_{emp}(x) \log q(x|\theta) &= \operatorname{argmax}_\theta \sum_{x \in X} \frac{n_x}{N} \log q(x|\theta) \\ &= \operatorname{argmax}_\theta \sum_{x \in X} \log q(x|\theta)^{n_x} = \operatorname{argmax}_\theta \prod_{x \in X} q(x|\theta)^{n_x} \end{aligned}$$

$$= \operatorname{argmax}_\theta \prod_{x \in X} q(x|\theta) = \operatorname{argmax}_\theta \sum_{x \in X} \log q(x|\theta)$$

$\therefore Q.E.D.$

T₂,

$$\begin{aligned} J(w, w_0) &= (y^T - w^T X^T - w_0 I^T)(y - Xw - w_0 I) + \lambda w^T w \\ &= y^T y - y^T Xw - y^T w_0 I - w^T X^T y + w^T X^T Xw + w^T X^T w_0 I - w_0 I^T y \\ &\quad + w_0 I^T Xw + w_0^2 I^T I + \lambda w^T w \end{aligned}$$

$$\frac{\partial J(w, w_0)}{\partial w_0} = -y^T I + w^T X^T I - I^T y + I^T Xw + 2w_0 I^T I = 0$$

since $\bar{X} = 0 \Rightarrow X^T I = 0, I^T X = 0$

we also have: $y^T I = I^T y = n\bar{y}, |I^T| = n$
 $\therefore -2n\bar{y} + 2nw_0 = 0 \Rightarrow w_0 = \bar{y}$

$$\frac{\partial J(w, w_0)}{\partial w} = 2\lambda Iw - 2X^T y + 2X^T Xw + w_0 I^T X + w_0 X^T I = 0$$

since $\bar{X} = 0 \Rightarrow X^T I = 0, I^T X = 0$

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

since $X^T X + \lambda I$ is invertible

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

\therefore Q.E.D.

$$\begin{aligned}
T_3 &= \sum_{i=1}^N \log P(y_i | x_i, w) - \lambda \sum_{c=1}^C \|w_c\|_2^2 \\
&= \sum_{i=1}^N \log \frac{\exp(w_{y_{i0}} + w_{y_i}^T x_i)}{\sum_{k=1}^C \exp(w_{k0} + w_k^T x_i)} - \lambda \sum_{c=1}^C \|w_c\|_2^2 \\
&= \sum_{i=1}^N \left(w_{y_{i0}} + w_{y_i}^T x_i - \log \sum_{k=1}^C \exp(w_{k0} + w_k^T x_i) \right) - \lambda \sum_{c=1}^C \|w_c\|_2^2
\end{aligned}$$

for any y_i , we let $w_{y_{i0}} = 0$, and $w_{k0} = 0$

$$J(w) = - \sum_{i=1}^N \left(w_{y_i}^T x_i - \log \sum_{k=1}^C \exp(w_k^T x_i) \right) + \lambda \sum_{c=1}^C \|w_c\|_2^2$$

$$\begin{aligned}
\text{we have } w_{y_i}^T x_i &= \sum_{c=1}^C I(y_i=c) w_c^T x_i \\
\frac{\partial \sum_{c=1}^C I(y_i=c) w_c^T x_i}{\partial w_c} &= I(y_i=c) x_i
\end{aligned}$$

$$\begin{aligned}
\frac{\partial \log \sum_{k=1}^C \exp(w_k^T x_i)}{\partial w_c} &= \frac{1}{\sum_{k=1}^C \exp(w_k^T x_i)} \cdot \frac{\partial \sum_{k=1}^C \exp(w_k^T x_i)}{\partial (w_c^T x_i)} \cdot \frac{\partial (w_c^T x_i)}{\partial w_c} \\
&= \frac{\exp(w_c^T x_i)}{\sum_{k=1}^C \exp(w_k^T x_i)} \cdot x_i
\end{aligned}$$

$$\frac{\partial \sum_{c=1}^C \|w_c\|_2^2}{\partial w_c} = \frac{\partial w_c^T w_c}{\partial w_c} = 2w_c$$

$$\frac{\partial J(w)}{\partial w_c} = - \sum_{i=1}^N I(y_i=c) x_i + \sum_{i=1}^N \frac{\exp(w_c^T x_i)}{\sum_{k=1}^C \exp(w_k^T x_i)} x_i + 2\lambda w_c$$

$$\sum_{c=1}^C \frac{\partial J(w)}{\partial w_c} = 0, \quad \sum_{c=1}^C \sum_{i=1}^N I(y_i=c) x_i = \sum_{i=1}^N x_i$$

$$\text{and in fact: } \frac{\exp(w_c^T x_i)}{\sum_{k=1}^C \exp(w_k^T x_i)} = P(y_i=c | x_i, w)$$

$$\begin{aligned}
\therefore \sum_{c=1}^C \sum_{i=1}^N P(y_i=c | x_i, w) \cdot x_i &= \sum_{i=1}^N x_i \quad \therefore \sum_{c=1}^C \frac{\partial J(w)}{\partial w_c} = 2\lambda \sum_{c=1}^C w_c = 0 \\
\therefore \sum_{c=1}^C \hat{w}_c &= 0 \Rightarrow \sum_{c=1}^C \hat{w}_{cj} = 0 \text{ for } j=1:D. \quad Q.E.D.
\end{aligned}$$

T4.
a, F.

$$\begin{aligned}
 J(w) &= -\frac{1}{|D|} \sum_{i \in D} \log \delta(y_i x_i^T w) + \lambda \|w\|_2^2 \\
 &= -\frac{1}{|D|} \sum_{i \in D} \log \frac{1}{1 + e^{-y_i x_i^T w}} + \lambda w^T w \\
 &= \frac{1}{|D|} \sum_{i \in D} \log(1 + e^{-y_i x_i^T w}) + \lambda w^T w \\
 \frac{\partial (w^T w)}{\partial w} &= 2w, \quad \frac{\partial^2 (w^T w)}{\partial w^2} = 2I \\
 \frac{\partial \log(1 + e^{-y_i x_i^T w})}{\partial w} &= \frac{\partial \log(1 + e^{-y_i x_i^T w})}{\partial (1 + e^{-y_i x_i^T w})} \cdot \frac{\partial (1 + e^{-y_i x_i^T w})}{\partial w} \\
 &= \frac{1}{1 + e^{-y_i x_i^T w}} \cdot e^{-y_i x_i^T w} \cdot (-y_i x_i) = -\frac{y_i x_i}{1 + e^{-y_i x_i^T w}} \cdot e^{-y_i x_i^T w} \\
 \frac{\partial^2 \log(1 + e^{-y_i x_i^T w})}{\partial w^2} &= -y_i x_i \cdot \left(\frac{y_i x_i e^{-y_i x_i^T w} (1 + e^{-y_i x_i^T w}) - (-y_i x_i) e^{-y_i x_i^T w}}{(1 + e^{-y_i x_i^T w})^2} \right) \\
 &= \frac{y_i^2 x_i x_i^T e^{-y_i x_i^T w}}{(1 + e^{-y_i x_i^T w})^2} \\
 \frac{\partial^2 J(w)}{\partial w^2} &= \frac{1}{|D|} \sum_{i \in D} \frac{y_i^2 x_i x_i^T e^{-y_i x_i^T w}}{(1 + e^{-y_i x_i^T w})^2} + 2\lambda I, \quad y_i \in \{-1, 1\}
 \end{aligned}$$

We have $\alpha_i = \frac{y_i^2 e^{-y_i x_i^T w}}{(1 + e^{-y_i x_i^T w})^2} > 0 \therefore \frac{\partial^2 J(w)}{\partial w^2} = \frac{1}{|D|} \sum_{i \in D} \alpha_i x_i x_i^T + 2\lambda I$

$x_i x_i^T$ is positive semidefinite $\therefore J(w)$ is convex. Q.E.D.

b, F, L_2 regularized won't let w be sparse since its "smooth" property

c, T, if the training data is linearly separable, the greater w is, the greater likelihood is.

d, F, when λ increases, since we minimize $J(w)$, weights become smaller, the likelihood decreases.

e, F, at first it can reduce overfitting and increase $L(w, D_{test})$ but later $L(w, D_{test})$ decrease since weights get smaller.

$$T_5.$$

$$X^T w = [x_1, x_2, \dots, x_d] \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} = w_1 x_1 + w_2 x_2 + \dots + w_d x_d$$

$$= \begin{bmatrix} w_1 x_{11} + w_2 x_{21} + \dots + w_d x_{d1} \\ \vdots \\ w_1 x_{1h} + w_2 x_{2h} + \dots + w_d x_{dh} \end{bmatrix}$$

$$\therefore \frac{d(X^T w)}{dw} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1h} \\ & \vdots & & \\ & & \ddots & \\ x_{d1} & x_{d2} & \cdots & x_{dh} \end{bmatrix} = X$$

$$y^T X = [y_1, \dots, y_d] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} = y_1 x_1 + \dots + y_d x_d$$

$$= [y_1 x_{11} + y_2 x_{21} + \dots + y_d x_{d1}, \dots, y_1 x_{1d} + y_2 x_{2d} + \dots + y_d x_{dd}]$$

$$y^T X w = w_1(y_1 x_{11} + \dots + y_d x_{d1}) + \dots + w_d(y_1 x_{1d} + \dots + y_d x_{dd})$$

$$\frac{d(y^T X w)}{dw} = \begin{bmatrix} y_1 x_{11} + \dots + y_d x_{d1} \\ \vdots \\ y_1 x_{1d} + \dots + y_d x_{dd} \end{bmatrix} = (y^T X)^T = X^T y$$

$$W^T X = [w_1, w_2, \dots, w_d] \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} = w_1 x_1 + w_2 x_2 + \dots + w_d x_d$$

$$= [w_1 x_{11} + w_2 x_{21} + \dots + w_d x_{d1}, \dots, w_1 x_{1d} + w_2 x_{2d} + \dots + w_d x_{dd}]$$

$$W^T X W = w_1^2 x_{11} + w_1 (w_2 x_{21} + \dots + w_d x_{d1}) + \dots + w_d (w_1 x_{1d} + w_2 x_{2d} + \dots + w_{d-1} x_{(d-1)d}) + w_d^2 x_{dd}$$

$$\therefore \frac{d(W^T X W)}{d w} = \begin{bmatrix} 2w_1 x_{11} + w_2 x_{21} + \dots + w_d x_{d1} + w_2 x_{12} + \dots + w_d x_{1d} \\ w_1 x_{21} + 2w_2 x_{22} + w_1 x_{12} + \dots + w_d x_{d2} + \dots + w_d x_{2d} \\ \vdots \\ w_1 x_{d1} + \dots + w_d x_{1d} + \dots + w_{d-1} x_{(d-1)d} + 2w_d x_{dd} \end{bmatrix}$$

$$= \left(\begin{bmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ \vdots & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ x_{d1} & x_{d2} & \dots & x_{dd} \end{bmatrix} + \begin{bmatrix} x_{11} & x_{21} & \dots & x_{d1} \\ \vdots & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ x_{1d} & x_{2d} & \dots & x_{dd} \end{bmatrix} \right) \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$

$$= (X + X^T) W$$