

1. MLE minimizes KL divergence to the empirical distribution (Exercise 2.15 of Murphy's book) (1 point)

$$KL(P||Q) = \sum_{k=1}^K P_k \log \frac{q_k}{p_k} = \sum_{k=1}^K P_k \log p_k - \sum_{k=1}^K P_k \log q_k$$

Since P is empirical distribution, then P_k are the same. So fix $P_k = P$
Then it's minimize $KP \log P - P \sum_{k=1}^K \log q_k$. So we need to maximize $\sum_{k=1}^K \log q_k$, which is actually MLE.

2. Centering and ridge regression (Exercise 7.3 of Murphy's book) (1 point)

$$J = (y^T - w^T x^T - w_0 \mathbf{1}^T)(y - xw - w_0 \mathbf{1}) + \lambda w^T w = y^T y + w^T x^T x w + w_0 \mathbf{1}^T w_0 \mathbf{1} + w^T x^T w_0 \mathbf{1} + w_0 \mathbf{1}^T x w - y^T x w - w^T x^T y - y^T w_0 \mathbf{1} - w_0 \mathbf{1}^T y + \lambda w^T w$$

$$\frac{dJ}{dw} = 2x^T x w + 2x^T w_0 \mathbf{1} - x^T y - x^T y + 2\lambda w = 0 \quad x^T w_0 \mathbf{1} = w_0 \sum_{i=1}^n x_i = 0$$

$$x^T x w + x^T w_0 \mathbf{1} - x^T y + \lambda w = 0 \quad (x^T x + \lambda I) w = x^T y, \quad w = (x^T x + \lambda I)^{-1} x^T y$$

$$\frac{dJ}{dw_0} = 2n w_0 + 2w^T (\sum_{i=1}^n x_i + x_n) - 2(y + \dots + y_n) = 0 \quad \text{Since } \bar{x} = 0,$$

$$\text{it's } n w_0 = \sum_{i=1}^n y_i, \quad w_0 = \bar{y},$$

$$\frac{d^2 J}{dw^2} = x^T w + \lambda, \quad \frac{d^2 J}{dw_0^2} = 2n. \quad \frac{d^2 J}{dw dw_0} = 0 \Rightarrow \text{Here } J \text{ is the local minimum}$$

3. Symmetric version of ℓ_2 regularized multinomial logistic regression (Exercise 8.5 of Murphy's book) (1 point)

$$J = \sum_{i=1}^N \log \frac{\exp(w_{k0} + w_k^T x_i)}{\sum_{c=1}^C \exp(w_{c0} + w_c^T x_i)} - \lambda \sum_{c=1}^C \|w_c\|_2^2$$

$$\textcircled{1} \frac{dJ}{dw_c} = -2\lambda w_c + \sum_{i=1}^N x_i - \sum_{i=1}^N \frac{x_i \exp(w_{k0} + w_k^T x_i)}{\sum_{c=1}^C \exp(w_{c0} + w_c^T x_i)}, \quad \frac{dJ}{dw_c} \big|_{1 \leq k \leq C-1} = -2\lambda w_k - \sum_{i=1}^N \frac{x_i \exp(w_{k0} + w_k^T x_i)}{\sum_{c=1}^C \exp(w_{c0} + w_c^T x_i)}$$

$$\textcircled{2} \sum_{c=1}^C \frac{dJ}{dw_c} = -2\lambda \left(\sum_{c=1}^C w_c \right) + \sum_{i=1}^N x_i - \sum_{i=1}^N x_i \frac{\sum_{c=1}^C \exp(w_{c0} + w_c^T x_i)}{\sum_{c=1}^C \exp(w_{c0} + w_c^T x_i)} = -2\lambda \left(\sum_{c=1}^C w_c \right). \text{ At the optimization, } \frac{dJ}{dw_c} = 0, \quad \frac{dJ}{dw_c} = 0 = -2\lambda \sum_{c=1}^C w_c, \quad \sum_{c=1}^C w_c = 0$$

4. Elementary properties of ℓ_2 regularized logistic regression (Exercise 8.6 of Murphy's book) (1 point)

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

a. It's a convex function. So just one local minimum, False

b. False, It doesn't encourage reducing a subset of the weights to zero

c. True. If $\lambda = 0$, the model is easy to be overfitted.

d. False. If λ is small, then it can cause overfitting, which produce large $J(w, D_{train})$

e. False. Although the increase of λ eclipse the overfitting, which increase $J(w, D_{test})$

But if λ is too large, then the model is inflexible and the likelihood function decreases.

$$5. \frac{d(x^T w)}{dw} = x, \quad X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nn} \end{bmatrix} \quad x^T w = \begin{bmatrix} x_{11} w_1 + x_{12} w_2 + \dots + x_{1n} w_n \\ x_{12} w_1 + x_{22} w_2 + \dots + x_{n2} w_n \\ \vdots \\ x_{n1} w_1 + x_{n2} w_2 + \dots + x_{nn} w_n \end{bmatrix}$$

$$\frac{d(x^T w)}{dw} = \begin{bmatrix} x_{11} & x_{12} & x_{12} & \dots & x_{1n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & \dots & x_{nn} \end{bmatrix} = X$$

$$\frac{d(y^T X w)}{dw} = X^T y \quad y^T X = [y_1 x_{11} + \dots + y_n x_{n1}, y_1 x_{12} + \dots + y_n x_{n2}, \dots, y_1 x_{1n} + \dots + y_n x_{nn}]$$

$$y^T X w = w_1 (y_1 x_{11} + \dots + y_n x_{n1}) + \dots + w_n (y_1 x_{1n} + \dots + y_n x_{nn})$$

$$\frac{d(y^T X w)}{dw} = \begin{bmatrix} y_1 x_{11} + \dots + y_n x_{n1} \\ \vdots \\ y_1 x_{1n} + \dots + y_n x_{nn} \end{bmatrix} = X^T y$$

$$\frac{d(w^T X w)}{dw} = (X + X^T) w$$

$$w^T X w = w_1 (w_1 x_{11} + \dots + w_n x_{n1}) + \dots + w_n (w_1 x_{1n} + \dots + w_n x_{nn})$$

$$\frac{d(w^T X w)}{dw} = \begin{pmatrix} x_{11} + \dots + x_{n1} + x_{11} + \dots + x_{1n} \\ x_{12} + \dots + x_{n2} + x_{21} + \dots + x_{2n} \\ \vdots \\ x_{1n} + \dots + x_{nn} + x_{n1} + \dots + x_{nn} \end{pmatrix} = X w + X^T w = (X + X^T) w$$