

Problem 3(a).

Given that, the maximum likelihood estimation:

$$\min_W F(W) \text{ where } F(W) = \frac{1}{n} \sum_{i=1}^n -\log [Pr(y=y_i | x=X_i; W)] + \frac{\eta}{2} \|W\|_F^2$$

The gradient of $F(W) \rightarrow \frac{\partial F(W)}{\partial W}$:

Each data point is independent. So the probability of all the data is:

$$\begin{aligned} F(W) &= \prod_{i=1}^n Pr(y=y^{(i)} | x=X^{(i)}) \\ &= \prod_{i=1}^n \kappa(\theta^T X^{(i)})^{y^{(i)}} [1 - \kappa(\theta^T X^{(i)})]^{(1-y^{(i)})} \end{aligned}$$

$$\text{Here, } \kappa = \frac{\exp(z_k - \max_j z_j)}{\sum_{j=1}^C \exp(z_j - \max_j z_j)} \text{ where } z = \theta^T x$$

where, C = number of possible classes.

and $k \in \{1, 2, \dots, C\}$

Derivative of gradient for one datapoint (x, y) without regularization

$$\begin{aligned} \frac{\partial F(w)}{\partial w} &= \frac{\partial}{\partial w} y \log \mu(\theta^T x) + \frac{\partial}{\partial w} (1-y) \log [1-\mu(\theta^T x)] \\ &= \left[\frac{y}{\mu(\theta^T x)} - \frac{1-y}{1-\mu(\theta^T x)} \right] \frac{\partial}{\partial w} \mu(\theta^T x) \quad (\text{derivative of } \log f(x)) \\ &= \left[\frac{y}{\mu(\theta^T x)} - \frac{1-y}{1-\mu(\theta^T x)} \right] \mu(\theta^T x) [1-\mu(\theta^T x)] \quad [\text{chain rule + derivative of sigma}] \\ &= \left[\frac{y - \mu(\theta^T x)}{\mu(\theta^T x) [1-\mu(\theta^T x)]} \right] \mu(\theta^T x) [1-\mu(\theta^T x)] x_j \quad [\text{algebraic manipulation}] \end{aligned}$$

$$= [y - \mu(\theta^T x)] x_j \quad [\text{cancelling term}]$$

The gradient of theta is simply the sum of this term for each training datapoint.

Finally, the gradient of $F(w)$ with regularizer

$$\frac{\partial F(w)}{\partial w} = -\frac{\eta}{2} w_i^{(t)} + \sum_j x_i^j [y^j - \hat{p}(y^j=1|x^j, w^{(t)})]$$

The gradient descent rule for w :

- First we have to initialize $w^{(1)} \in \mathbb{R}^D$ randomly

- Then we have to iterate through the following convergence:

$$w_i^o(t+1) \leftarrow w_i^o(t) + \alpha \left\{ -\lambda w_i^o(t) + \sum_j x_i^j [y_j^o - \hat{p}(y_j^o = 1 | x^j w^{(t)})] \right\}$$

new value previous value learning rate gradient at previous value

The updates give larger weights to those examples on which the current model makes larger mistakes

Comparison with least mean square algorithm:

$$w(n+1) = w(n) + \frac{1}{2} \mu [-\nabla J(n)]$$

cost function

LMS algorithm is based on the idea of gradient descent to search for the minimum error with a cost function equal to the mean squared error,

where gradient descent doesn't force the use of

any particular cost function, it hunts for the minimum cost solution