

Problem 1: Derive the weight update rule that maximize the conditional likelihood assuming that a dataset $D = \{(x_i, y_i)\}_{i=1}^n$ is given.

the objective function is to find the optimal w^*

$$\Rightarrow w^* = \arg \max_w \prod_{i=1}^n \sigma(w^T x_i)^{y_i} (1 - \sigma(w^T x_i))^{1-y_i}$$

but due to computational reason, we want to minimize the objective function rather than maximizing it.

$$\Rightarrow E(w) = \sum_{i=1}^N -y_i \ln(\sigma(w^T x_i)) - (1-y_i) \ln(1 - \sigma(w^T x_i))$$

$$= \sum_{i=1}^N -y_i \ln\left(\frac{1}{1+e^{-w^T x_i}}\right) - \ln(1 - \sigma(w^T x_i)) + y_i \ln(1 - \sigma(w^T x_i))$$

$$E(w) = \sum_{i=1}^N -y_i \ln\left(\frac{1}{1+e^{-w^T x_i}}\right) - \ln\left(1 - \frac{1}{1+e^{-w^T x_i}}\right) + y_i \ln\left(1 - \frac{1}{1+e^{-w^T x_i}}\right)$$

$$= \sum_{i=1}^N -y_i \ln(1+e^{-w^T x_i}) - \ln\left(\frac{1+e^{-w^T x_i}-1}{1+e^{-w^T x_i}}\right) + y_i \ln\left(\frac{1+e^{-w^T x_i}-1}{1+e^{-w^T x_i}}\right)$$

$$= \sum_{i=1}^N y_i \ln(1+e^{-w^T x_i}) - \ln\left(\frac{e^{-w^T x_i}}{1+e^{-w^T x_i}}\right) + y_i \ln\left(\frac{e^{-w^T x_i}}{1+e^{-w^T x_i}}\right)$$

$$= \sum_{i=1}^N y_i \ln(1+e^{-w^T x_i}) - \ln e^{-w^T x_i} + \ln(1+e^{-w^T x_i}) + y_i \ln(e^{-w^T x_i}) - y_i \ln(1+e^{-w^T x_i})$$

$$E(w) = \sum_{i=1}^N -\ln e^{-w^T x_i} + \ln(1+e^{-w^T x_i}) + y_i \ln(e^{-w^T x_i}) \quad (a^*)$$

$$\Rightarrow \frac{\delta E(w)}{\delta w} = \left[\sum_{i=1}^N -\ln e^{-w^T x_i} + \ln(1+e^{-w^T x_i}) + y_i \ln(e^{-w^T x_i}) \right]' \quad (a)$$

$$\frac{\delta \ln e^{-w^T x_i}}{\delta w} = \left(\frac{x_i}{e^{-w^T x_i}} \right) (e^{-w^T x_i})' \quad \text{chain rule} \quad \text{" ' ' is the derivative sign"}$$

$$= \left(\frac{1}{e^{-w^T x_i}} \right) (-w^T x_i)' e^{-w^T x_i} \quad \text{chain rule}$$

$$\frac{\delta \ln e^{-w^T x_i}}{\delta w} = -x_i \quad (1)$$

$$\frac{\delta (\ln(1+e^{-w^T x_i}))}{\delta w} = \left(\frac{1}{1+e^{-w^T x_i}} \right) (1+e^{-w^T x_i})' \quad \text{chain rule}$$

$$= \left(\frac{1}{1+e^{-w^T x_i}} \right) (-w^T x_i)' e^{-w^T x_i}$$

$$= \frac{-x_i e^{-w^T x_i}}{1+e^{-w^T x_i}} \quad (2)$$

put (1) & (2) into (a)

$$\begin{aligned}\Rightarrow \frac{\partial E(w)}{\partial w} &= \sum_{i=1}^N (-1)(-x) + \left(-\frac{x e^{-w^T x}}{1 + e^{-w^T x}} \right) + (y_i)(-x) \\ &= \sum_{i=1}^N x - x \frac{e^{-w^T x}}{1 + e^{-w^T x}} - x y_i \\ &= \sum_{i=1}^N x \left(1 - \frac{e^{-w^T x}}{1 + e^{-w^T x}} - y_i \right)\end{aligned}$$

$$\text{but } \frac{e^{-w^T x}}{1 + e^{-w^T x}} = 1 - \frac{1}{1 + e^{-w^T x}}$$

$$\begin{aligned}\Rightarrow \frac{\partial E(w)}{\partial w} &= \sum_{i=1}^N x \left[1 - \left(1 - \frac{1}{1 + e^{-w^T x}} \right) - y_i \right] \\ &= \sum_{i=1}^N x \left(\frac{1}{1 + e^{-w^T x}} - y_i \right)\end{aligned}$$

$$\frac{\partial E(w)}{\partial w} = \sum_{i=1}^N x (\sigma(w^T x) - y_i) \quad (c)$$

from a^* , w is a concave function of w ; hence, there is no local minima.

On top of that concave function is easy to optimize using gradient descent.

\Rightarrow update rule $w = w - (\text{learning rate})(\text{partial derivative of } E(w) \text{ w.r.t } w)$

$$\Rightarrow w_j = w_j - \alpha \nabla_w E(w)$$

from above derivative, we found that $\nabla_w E(w) = \sum_{i=1}^N x (\sigma(w^T x) - y_i)$

$$\therefore w_j = w_j - \alpha \sum_{i=1}^N x (\sigma(w^T x) - y_i)$$

Problem 2:

1. Compute $\frac{\delta \sigma(a)}{\delta w}$ when $a = w^T x$ where $x, w \in \mathbb{R}^m$

$$\sigma(a) = \frac{1}{1+e^{-a}} = (1+e^{-a})^{-1} = (1+e^{-w^T x})^{-1}$$

$$\Rightarrow \frac{\delta (1+e^{-w^T x})^{-1}}{\delta w} = \frac{(-1)(1+e^{-w^T x})^{-2} \delta (1+e^{-w^T x})}{\delta w} \quad \text{chain rule}$$

$$= (-1)(-w^T x)(e^{-w^T x})(1+e^{-w^T x})^{-2}$$

$$= x e^{-w^T x} (1+e^{-w^T x})^{-2}$$

$$\Rightarrow \frac{\delta (1+e^{-w^T x})^{-1}}{\delta w} = \frac{x e^{-w^T x}}{(1+e^{-w^T x})^2}$$

2. For logistic regression, the probability of y , given input value x and its estimated weight vector w , can be presented as a posterior probability of y , $P(y|x, w)$ using bayes' theorem. If $y=1 \Rightarrow P(y=1|x, w) = \frac{1}{1+e^{-w^T x}}$ (a)

$$\text{conversely, if } y=0 \Rightarrow P(y=0|x, w) = 1 - \frac{1}{1+e^{-w^T x}}$$

$$P(y=0|x) = \frac{1+e^{-w^T x} - 1}{1+e^{-w^T x}} = \frac{e^{-w^T x}}{1+e^{-w^T x}} \quad (b)$$

$P(y=1|x, w)$ means the probability of y or output label = 1, given the input feature x and its estimated model parameter w .

$P(y=0|x, w)$ means the probability of y or output label = 0, given the input feature x and its estimated model parameter w .

Using the Bernoulli distribution, $P^k(1-p)^{1-k}$ $k \in \{0, 1\}$, we could combine (a) & (b) together, then we get

$$P(y|x, w) = \sigma(w^T x)^y (1 - \sigma(w^T x))^{1-y}$$

$$\text{if we generate } N \text{ samples, } P(y|x, w) = \prod_{i=1}^N \sigma(w^T x_i)^{y_i} (1 - \sigma(w^T x_i))^{1-y_i}$$

3. Show the loss function for logistic regression and explain how we learn w

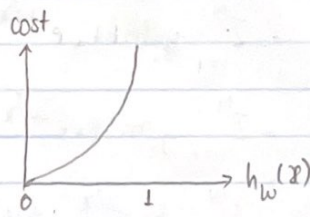
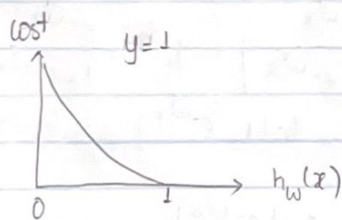
In logistic regression, we focus on doing the binary classification; hence, $y \in \{0, 1\}$. We constraint the prediction to some value of 0, 1. Therefore, we use the sigmoid function $g(z) = \frac{1}{1+e^{(-z)}}$, hypothesis: $h_w(x) = \frac{1}{1+e^{-w^T x}}$

The hypothesis function returns the probability that $y=1$ given x , parameterized by w . It is written as $h(x) = P(y=1|x, w)$. The decision boundary is:

$$\begin{cases} \text{if } w^T x \geq 0 \rightarrow h(x) \geq 0.5 \\ \text{elif } w^T x < 0 \rightarrow h(x) < 0.5 \end{cases}$$

The loss function is in place because we want to create a punishment when predicting 1 while it is actually 0 or 0 when it is actually 1.

$$\text{cost}(h_w(x), y) = \begin{cases} -\log(h_w(x)) & \text{if } y=1 & (1) \\ -\log(1-h_w(x)) & \text{if } y=0 & (2) \end{cases}$$



we could also add the cost function together. Now it will be written as

$$\text{cost}(h_w(x), y) = -y \log(h_w(x)) - (1-y) \log(1-h_w(x))$$

Hence, the cost function of the model is the summation of all the training data:

$$E(w) = \sum_{i=1}^N -y_i \log(h_w(x_i)) - (1-y_i) \log(1-h_w(x_i))$$

$$E(w) = \sum_{i=1}^N -y_i \log\left(\frac{1}{e^{-w^T x_i}}\right) - (1-y_i) \log\left(1 - \frac{1}{e^{-w^T x_i}}\right)$$

Problem 2.3 continue~

There is no closed-form equation to compute the value of w , that minimize the cost function. However, the cost function $E(w)$ is convex. Hence, gradient Descent is guaranteed to find global minimum.

The gradient descent update rule $w_{\text{new}} = w_{\text{old}} - d \nabla_w E(w_{\text{old}})$

meaning, to compute w , we use the previous w value, subtracts with d (the learning rate) times $\nabla_w E(w_{\text{old}})$ partial derivative of $E(w)$ w.r.t w .

derive from problem (1) $w_{\text{new}} = w_j = w_j - d \sum_{i=1}^N x_i (\sigma(w^T x) - y_i)$

for each instance, it computes the prediction error and multiplies it by the j^{th} feature value $j \in \{0, 1, 2, \dots, m\}$. Once, we get the gradient vector, containing all the partial derivative, we can use it in Batch Gradient Descent algorithm.