

CS446: Machine Learning, Fall 2017, Homework 1

Name: Yiming Gao (yimingg2)

Worked individually

Logistic Regression: Deriving Gradient Descent Update Rules

Problem 1

Solution:

We have

$$\log \frac{P(y=1|\mathbf{x})}{P(y=0|\mathbf{x})} = \mathbf{w}^T \mathbf{x},$$

which is equivalent to

$$P(y=1|\mathbf{x}, \mathbf{w}) = \frac{1}{1 + e^{w_0 + \sum_{i=1}^d w_i x_i}} = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}},$$
$$P(y=0|\mathbf{x}, \mathbf{w}) = 1 - P(y=1|\mathbf{x}, \mathbf{w}) = \frac{e^{-\mathbf{w}^T \mathbf{x}}}{1 + e^{-\mathbf{w}^T \mathbf{x}}},$$

Problem 2

Find the derivative of the sigmoid function.

Solution:

We denote Sigmoid function as $\sigma(\mathbf{z})$ and derive its derivative as follows:

$$\begin{aligned} (1 + e^{-z})\sigma(\mathbf{z}) &= 1 \\ \Rightarrow -e^{-z}\sigma + (1 + e^{-z})\frac{d\sigma}{dz} &= 0 \\ \Rightarrow \frac{d}{dz}\sigma(\mathbf{z}) &= \sigma \cdot \frac{e^{-z}}{1 + e^{-z}} \\ &= \sigma \cdot \frac{(1 + e^{-z}) - 1}{1 + e^{-z}} \\ &= \sigma \cdot \left[1 - \frac{1}{1 + e^{-z}}\right] \\ &= \sigma \cdot (1 - \sigma) \end{aligned}$$

Problem 3, 4

Reference: <http://web.engr.oregonstate.edu/~xfern/classes/cs534/notes/logistic-regression-note.pdf>

Derive the Likelihood function of Logistic Regression.

Solution:

The log likelihood function is as follows:

$$\log p(D|M) = \sum_{i=1}^N \log p(\mathbf{x}_i, y_i) = \sum_{i=1}^N \log p(y_i|\mathbf{x}_i)p(\mathbf{x}_i).$$

Let $\sigma(\mathbf{x}, \mathbf{w}) = \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}}}$ denote the **sigmoid** function. Note that in logistic regression, we don't care about $p(\mathbf{x}_i)$ and only need to learn the $p(y|\mathbf{x})$. Thus we have

$$L(\mathbf{w}) = \sum_{i=1}^N \log p(y_i|\mathbf{x}_i) = \sum_{i=1}^N \log \sigma(\mathbf{x}_i, \mathbf{w})^{y_i} (1 - \sigma(\mathbf{x}_i, \mathbf{w}))^{1-y_i}$$

To maximize L with respect to \mathbf{w} , we look at each example:

$$L_i(\mathbf{w}) = \log \sigma(\mathbf{x}_i, \mathbf{w})^{y_i} (1 - \sigma(\mathbf{x}_i, \mathbf{w}))^{1-y_i} = y_i \log \sigma(\mathbf{x}_i, \mathbf{w}) + (1 - y_i) \log (1 - \sigma(\mathbf{x}_i, \mathbf{w}))$$

where $\sigma(\mathbf{x}, \mathbf{w}) = \frac{1}{1+e^{-\mathbf{w}^T \mathbf{x}}}$.

Taking gradient of L_i with respect to \mathbf{w} , we have

$$\begin{aligned} \nabla_{\mathbf{w}} L_i &= \frac{y_i}{\sigma(\mathbf{x}_i, \mathbf{w})} \nabla_{\mathbf{w}} \sigma - \frac{1 - y_i}{1 - \sigma(\mathbf{x}_i, \mathbf{w})} \nabla_{\mathbf{w}} \sigma \\ &= \frac{y_i}{\sigma} \sigma(1 - \sigma) \mathbf{x}_i - \frac{1 - y_i}{1 - \sigma} \sigma(1 - \sigma) \mathbf{x}_i \\ &= [y_i(1 - \sigma) - (1 - y_i)\sigma] \mathbf{x}_i \\ &= (y_i - \sigma(\mathbf{x}_i, \mathbf{w})) \mathbf{x}_i \end{aligned}$$

(1) So for a **single** training example, the update rule is:

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \eta(y_i - \sigma(\mathbf{x}_i, \mathbf{w})) \mathbf{x}_i$$

(2) Consider **all** training examples, we have

$$\nabla_{\mathbf{w}} L = \sum_{i=1}^N (y_i - \sigma(\mathbf{x}_i, \mathbf{w})) \mathbf{x}_i$$

The update rule for Gradient Descent is

$$\mathbf{w}_{k+1} = \mathbf{w}_k + \eta \nabla L L(\mathbf{w}_k)$$

where $\nabla L L(\mathbf{w}_k) = \sum_{i=1}^N (y_i - \sigma(\mathbf{x}_i, \mathbf{w}_k)) \mathbf{x}_i$ and η is the stepsize of Gradient Descent.