

Logistic Regression - MLE

Given n pairs (\vec{x}_i, y_i) where \vec{x}_i is $d \times 1$
and y_i is 0 or 1

We aim to maximize the conditional likelihood of all y_i 's, given x_i 's

$$\vec{w}_{MLE} = \arg \max_{\vec{w}} \prod_{i=1}^n P(y_i | \vec{x}_i, \vec{w})$$

$$\log l(\vec{w}) = \log \left[\prod_{i=1}^n P(y_i | \vec{x}_i, \vec{w}) \right]$$

$$= \sum_{i=1}^n \log \left[P(y_i | \vec{x}_i, \vec{w}) \right]$$

$$= \sum_{i=1}^n \left[y_i \log \left[P(y_i=1 | \vec{x}_i, \vec{w}) \right] + (1-y_i) \log \left[P(y_i=0 | \vec{x}_i, \vec{w}) \right] \right]$$

Exactly one of the two terms is always zero because $y_i = 0$ or 1

$$= \sum_{i=1}^n \left[y_i \left[\log \left[\frac{P(y_i=1 | \vec{x}_i, \vec{w})}{P(y_i=0 | \vec{x}_i, \vec{w})} \right] \right] + \log \left[P(y_i=0 | \vec{x}_i, \vec{w}) \right] \right]$$

$$= \sum_{i=1}^n \left[y_i \left[\log \left[\frac{e^{\vec{w}^T \vec{x}_i} / (1 + e^{\vec{w}^T \vec{x}_i})}{(1) / (1 + e^{\vec{w}^T \vec{x}_i})} \right] \right] + \log \left[\frac{1}{1 + e^{\vec{w}^T \vec{x}_i}} \right] \right]$$

$$= \sum_{i=1}^n \left[y_i \left[\vec{w}^T \vec{x}_i \right] - \log \left[1 + e^{\vec{w}^T \vec{x}_i} \right] \right]$$

$$\nabla_{\vec{w}} [\log \ell] = \begin{bmatrix} \frac{\partial \log \ell(\vec{w})}{\partial w_0} \\ \vdots \\ \frac{\partial \log \ell(\vec{w})}{\partial w_n} \end{bmatrix}$$

$$\frac{\partial \log \ell(\vec{w})}{\partial w_j} = \sum_{i=1}^n \left[y_i x_i^j - \frac{1 \times e^{\vec{w}^T \vec{x}_i}}{1 + e^{\vec{w}^T \vec{x}_i}} \times x_i^j \right] \text{ where } x_i^j \text{ is } j^{\text{th}} \text{ component of } x_i$$

$$= \sum_{i=1}^n \left[x_i^j \left[y_i - \frac{e^{\vec{w}^T \vec{x}_i}}{1 + e^{\vec{w}^T \vec{x}_i}} \right] \right]$$

This gradient can be used to perform gradient descent on \vec{w} .

L₂-regularized logistic regression

In addition to maximized log-conditional-likelihood, we penalize the L₂ norm of \vec{w}

$$\therefore \text{log}l(\vec{w}) = \log \left[\prod_{i=1}^n P(y_i | \vec{x}_i, \vec{w}) \right] - \lambda \|\vec{w}\|_2^2$$

The MLE steps follow exactly the same pattern as ~~those~~ for plain vanilla logistic regression, except the $-\lambda \|\vec{w}\|_2^2$ term. Eventually, component-wise gradient is the following:

$$\frac{\partial \text{log}l}{\partial w_j} = \sum_{i=1}^n \left[x_i^j \left[y_i - \frac{e^{\vec{w}^T \vec{x}_i}}{1 + e^{\vec{w}^T \vec{x}_i}} \right] \right] - 2\lambda w_j$$

This gradient can be used for gradient descent to estimate \vec{w} .