

Machine Learning, Lecture 6: Logistic regression

S. Nõmm

¹Department of Computer Science, Tallinn University of Technology

12.03.2015

Generative approach *versus* Discriminative approach

- ▶ *Generative* approach - create a model of the form $p(y, \mathbf{x})$ and then derive $p(y | \mathbf{x})$.
- ▶ *Discriminative* approach - fit the model of the form $p(y | \mathbf{x})$ directly.

Logistic regression

- ▶ Linear regression model $p(y | \mathbf{x}; \boldsymbol{\theta}) = \mathcal{N}(y | \mu(\mathbf{x}))$
 - ▶ Replace Gaussian distribution for y with a Bernoulli distribution (more appropriate for the binary response)

$$p(y | \mathbf{x}, \boldsymbol{\theta}) = \text{Ber}(y | \mu(\mathbf{x}))$$

where $\mu(\mathbf{x}) = \mathbb{E}[y | \mathbf{x}] = p(y = 1 | \mathbf{x})$.

- ▶ Ensure that $0 \leq \mu(\mathbf{x}) \leq 1$ by

$$\mu(\mathbf{x}) = \text{sigm}(\boldsymbol{\theta}^T \mathbf{x})$$

where $\text{sigm}(\eta)$ is the *sigmoid* or *logistic* or *logit* function:

$$\mu(\mathbf{x}) = \frac{1}{1 + e^{-\eta}} = \frac{e^{\eta}}{e^{\eta} + 1}$$



$$p(y | \mathbf{x}, \boldsymbol{\theta}) = \text{Ber}(y | \text{sigm}(\boldsymbol{\theta}^T \mathbf{x}))$$

Some important properties

- ▶ For the logistic function

$$g(\eta) = \frac{1}{1 + e^{-\eta}}$$

$$g(\eta) = 0.5 \quad \text{if } \eta = 0$$

$$g(\eta) > 0.5 \quad \text{if } \eta > 0$$

$$g(\eta) < 0.5 \quad \text{if } \eta < 0$$

- ▶ Derivative of the logistic function

$$g'(\eta) = g(\eta)(1 - g(\eta))$$

Probabilistic interpretation

- ▶ Let us compute the probabilities of $y = 1$ and $y = 0$

$$P(y = 1 \mid \mathbf{x}, \boldsymbol{\theta}) = \text{sigm}(\boldsymbol{\theta}^T \mathbf{x})$$

$$P(y = 0 \mid \mathbf{x}, \boldsymbol{\theta}) = 1 - \text{sigm}(\boldsymbol{\theta}^T \mathbf{x})$$

Could you write this statement in a more compact form?

$$P(y \mid \mathbf{x}, \boldsymbol{\theta}) = ?$$

- ▶ The meaning of $\boldsymbol{\theta}^T \mathbf{x}$

$$g(\boldsymbol{\theta}^T \mathbf{x}) = \frac{e^{\boldsymbol{\theta}^T \mathbf{x}}}{1 + e^{\boldsymbol{\theta}^T \mathbf{x}}}$$

after the straight but tedious calculations one gets

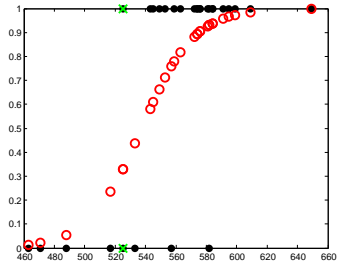
$$\boldsymbol{\theta}^T \mathbf{x} = \log \frac{g(\boldsymbol{\theta}^T \mathbf{x})}{1 - g(\boldsymbol{\theta}^T \mathbf{x})}$$

here and after referred as *log-odds*, probability of event occurring is divided by the probability of not occurring.

Example

Denote x_i to be the SAT score of the student i and y_i is whether they passed or failed a class.

$$p(y_i = 1 \mid x_i; \mathbf{w}) = \text{sigm}(\omega_0 + \omega_1 x_i)$$



Likelihood

- ▶ Likelihood of the parameters (probability of the entire data set)

$$\mathcal{L}(\boldsymbol{\theta}) = P(Y | \mathbf{X}; \boldsymbol{\theta}) = \prod_{i=1}^m (\text{sigm}(\boldsymbol{\theta}^T \mathbf{x}_i))^{y_i} (1 - \text{sigm}(\boldsymbol{\theta}^T \mathbf{x}_i))^{1-y_i}$$

- ▶ We use log- likelihood which leads:

$$\begin{aligned} \ell(\boldsymbol{\theta}) &= \log \mathcal{L}(\boldsymbol{\theta}) \\ &= \log \prod_{i=1}^m (\text{sigm}(\boldsymbol{\theta}^T \mathbf{x}_i))^{y_i} (1 - \text{sigm}(\boldsymbol{\theta}^T \mathbf{x}_i))^{1-y_i} \\ &= \sum_{i=1}^m (y_i \log \text{sigm}(\boldsymbol{\theta}^T \mathbf{x}_i) + (1 - y_i) \log(1 - \text{sigm}(\boldsymbol{\theta}^T \mathbf{x}_i))) \end{aligned}$$

Likelihood maximization

- ▶ Gradient descent to minimize the negative log-likelihood.
Update step:

$$\theta_j^{k+1} = \theta_j^k - \alpha \frac{\partial}{\partial \theta_j^k} \ell(\boldsymbol{\theta})$$

- ▶ Gradient ascent to maximize log likelihood. Update step:

$$\theta_j^{k+1} = \theta_j^k + \alpha \frac{\partial}{\partial \theta_j^k} \ell(\boldsymbol{\theta})$$

- ▶ By derivation the log -likelihood one gets the gradient ascend update for the logistic regression:

$$\theta_j^{k+1} = \theta_j^k + \alpha \sum_{i=1}^m (y_i - \text{sigm}(\boldsymbol{\theta}^T \mathbf{x}_i)) x_{i,j}$$

simultaneously for each θ_j , $j = 0, \dots, n$.