

7.1.1. Linear Regression

Training set (x_i, y_i) , $x_i \in \mathbb{R}^p$, $y_i \in \mathbb{R}$, $i=1, \dots, n$

Hypothesis:

$$y_i = \vartheta_0 + x_{i1} \vartheta_1 + \dots + x_{ip} \vartheta_p + \varepsilon_i, \quad i=1, \dots, n$$

$$\underline{y} = X \underline{\vartheta} + \underline{\varepsilon}, \quad x_{ij}, y_i \in \mathbb{R}, \quad X = \begin{pmatrix} 1 & x_1^T \\ \vdots & \vdots \\ 1 & x_n^T \end{pmatrix}$$

$$\underline{\hat{\vartheta}} = (X^T X)^{-1} X^T \underline{y} \quad (\text{or } \underline{\hat{\vartheta}} = (X^T X)^+ X^T \underline{y})$$

New observation

$$\hat{y} = (1, x^T) \underline{\hat{\vartheta}} \quad \text{best predictor.}$$

7.1.2. Logistic Regression

Training set (x_i, y_i) , $x_i \in \mathbb{R}^p$, $y_i \in \{0, 1\}$

Example: Spam mail

$x_i \in \mathbb{R}^p$ feature vector of some email

$$y_i = \begin{cases} 0, & \text{spam} \\ 1, & \text{no spam} \end{cases}$$

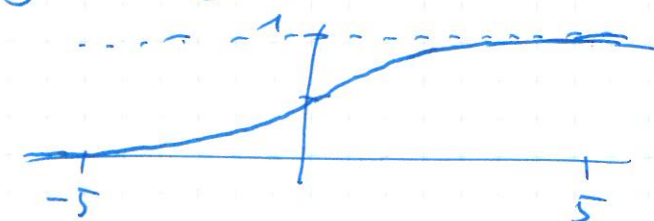
Hypothesis:

$$h_{\underline{\vartheta}}(\underline{x}) = g(\underline{\vartheta}^T \underline{x}) = \frac{1}{1 + e^{-\underline{\vartheta}^T \underline{x}}}$$

$g(z) = \frac{1}{1 + e^{-z}}$ is called logistic or sigmoid function.

It holds $g'(z) = g(z)(1 - g(z))$, $z \in \mathbb{R}$

Plot:



Given training set

$$(x_i, y_i), \quad i=1, \dots, n, \quad \underline{x}_i \in \mathbb{R}^p, \quad y_i \in \{0, 1\}$$

Set $x_{i0} = 1$, s.t.

$$\underline{w}^T \underline{x}_i = w_0 + \sum_{j=1}^p w_j x_{ij}, \quad \underline{x}_i = (1, x_{i1}, \dots, x_{ip})$$

Objective: find the best \underline{w} !

Probabilistic interpretation:

Assume

$$P(y=1 | \underline{x}, \underline{w}) = \underline{h}_{\underline{w}}(\underline{x})$$

$$P(y=0 | \underline{x}, \underline{w}) = 1 - \underline{h}_{\underline{w}}(\underline{x})$$

Hence

$$p(y | \underline{x}, \underline{w}) = (\underline{h}_{\underline{w}}(\underline{x}))^y (1 - \underline{h}_{\underline{w}}(\underline{x}))^{1-y}, \quad y \in \{0, 1\}$$

Assume n independent training samples, set

$$X = \begin{pmatrix} 1 & \underline{x}_1^T \\ \vdots & \vdots \\ 1 & \underline{x}_n^T \end{pmatrix} = (x_{ij})_{\substack{1 \leq i \leq n \\ 0 \leq j \leq p}}$$

Likelihood fct.

$$L(\underline{w}) = p(y | X, \underline{w}) = \prod_{i=1}^n p(y_i | \underline{x}_i, \underline{w})$$

$$= \prod_{i=1}^n (\underline{h}_{\underline{w}}(\underline{x}_i))^{y_i} (1 - \underline{h}_{\underline{w}}(\underline{x}_i))^{1-y_i}$$

log-Likelihood fct.

$$\ell(\underline{w}) = \log L(\underline{w}) = \sum_{i=1}^n y_i \log \underline{h}_{\underline{w}}(\underline{x}_i) + (1-y_i) \log(1 - \underline{h}_{\underline{w}}(\underline{x}_i)) (*)$$

Objective: $\max_{\underline{w} \in \mathbb{R}^{p+1}} \ell(\underline{w})$

Algorithm to compute the optimum: gradient ascent

$$\underline{w}^{(k+1)} = \underline{w}^{(k)} + \alpha \nabla_{\underline{w}} \ell(\underline{w}^{(k)}), \alpha: \text{"learning parameter"}$$

Compute $\nabla_{\underline{w}}$ for each addend i in $(*)$, set $(x_i, y_i) = (x, y)$

$$\begin{aligned} & \frac{\partial}{\partial w_j} [y \log h_{\underline{w}}(x) + (1-y) \log (1-h_{\underline{w}}(x))] \\ &= \left(y \frac{1}{g(\underline{w}^T \underline{x})} - (1-y) \frac{1}{1-g(\underline{w}^T \underline{x})} \right) \frac{\partial}{\partial w_j} g(\underline{w}^T \underline{x}) \\ &= \left(\dots \right) g(\underline{w}^T \underline{x})(1-g(\underline{w}^T \underline{x})) \frac{\partial}{\partial w_j} \underline{w}^T \underline{x} \\ &= (y(1-g(\underline{w}^T \underline{x})) - (1-y)g(\underline{w}^T \underline{x})) x_j, \quad \underline{x} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_p \end{pmatrix} \\ &= (y - h_{\underline{w}}(x)) x_j, \quad j=0,1,\dots,p \end{aligned}$$

Hence

$$\frac{\partial}{\partial w_j} \ell(\underline{w}) = \sum_{i=1}^n (y_i - h_{\underline{w}}(x_i)) x_{ij}, \quad j=0,1,\dots,p$$

Alternatively: Newton's method

$$\underline{w}^{(k+1)} = \underline{w}^{(k)} - H^{-1} \nabla_{\underline{w}} \ell(\underline{w}^{(k)})$$

Denote \underline{w}^* the optimal \underline{w} , i.e., $\underline{w}^* = \arg \max_{\underline{w} \in \mathbb{R}^p} \ell(\underline{w})$.

Let \underline{x} be a future observation

$$\text{Soft decision rule: } y = h_{\underline{w}^*}(\underline{x}) = \frac{1}{1 + e^{-\underline{w}^{*T} \underline{x}}}$$

$$y \in (0,1)$$

7.1.3. The perceptron learning algorithm

"Soft" decision in logistic regression is a value $y \in (0, 1)$.

Force hard decision, values to be 0 or 1 by

$$g(z) = \begin{cases} 1, & z \geq 0 \\ 0, & z < 0 \end{cases}$$

$$h_{\underline{w}}(\underline{x}) = g(\underline{w}^T \underline{x})$$

Use the same update rule as above

$$\underline{w}_j^{(k+1)} = \underline{w}_j^{(k)} + \alpha \sum_{i=1}^n (y_i - h_{\underline{w}_j^{(k)}}(\underline{x}_i)) \underline{x}_{ij}, \quad j=0, 1, \dots, p.$$

This is the perceptron learning algorithm.

A rough model for how neurons in the brain work.

Lacking: meaningful probabilistic interpretation as MLE.