

## 7. Machine Learning

### 7.1 Supervised Learning

Given  $(x_i, y_i), i=1, \dots, n$  training examples/samples

$x_i \in X$  : input variables, feature variables

$y_i \in Y$  : output  $u$ , target variables

$\{(x_i, y_i) \mid i=1, \dots, n\}$  training set.

Supervised learning : determine a fct.

$$h: X \rightarrow Y \quad (\text{a "hypothesis"})$$

so that  $h(x)$  is a "good" predictor of  $y$ .

$y$  continuous  $\Rightarrow$  regression problem

$y$  discrete  $\Rightarrow$  classification problem

#### 7.1.1 Linear regression

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

Assumption / hypothesis

$$y_i = \alpha_0 + x_{i1} \alpha_1 + x_{i2} \alpha_2 + \dots + x_{ip} \alpha_p + \epsilon_i$$

$\epsilon_i$  : random error

$$= (1, x_i^T) \alpha + \epsilon_i$$

Hence, learn  $h_w(x) = (1, x^T) \alpha, \alpha = (\alpha_0, \alpha_1, \dots, \alpha_p)^T$   
parameter

Set  $n \times (p+1)$   

$$X = \begin{pmatrix} 1 & x_1^T \\ \vdots & \vdots \\ 1 & x_n^T \end{pmatrix} \quad y = (y_1, \dots, y_n)^T$$

$$\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^T$$

Then  $y = X\beta + \varepsilon$

Problem: Find the best  $\beta$  by solving

$$\min_{\beta \in \mathbb{R}^{p+1}} \|y - X\beta\|$$

Solution: (i) project  $y$  onto  $\mathcal{J}_n(X) : \hat{y}$

(ii) find  $\beta$  such that  $\hat{y} = X\beta$

(i)  $X(X^T X)^{-1} X^T$  is an orthogonal proj. onto  $\mathcal{J}_n(X)$ , provided  $(X^T X)^{-1}$  exists.

$$\hat{y} = X(X^T X)^{-1} X^T y = \arg \min_{z \in \mathcal{J}_n(X)} \|y - z\|$$

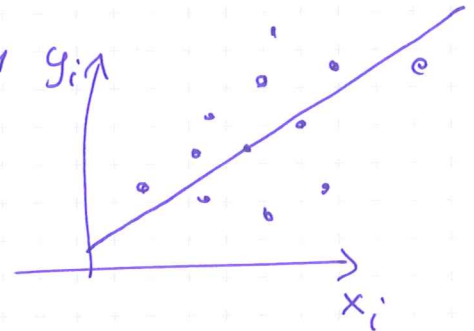
$$\begin{aligned} \text{(ii)} \quad \hat{y} = X\beta &\Rightarrow X^T X(X^T X)^{-1} X^T y = X^T X\beta \\ &\Rightarrow X^T y = X^T X\beta \quad (\text{normal equations}) \\ &\Rightarrow \beta^* = (X^T X)^{-1} X^T y \\ \text{and } X\beta^* &= X(X^T X)^{-1} X^T y = \hat{y} \end{aligned}$$

In summary:  $\beta^* = (X^T X)^{-1} X^T y$  is a solution.

Note: the inverse  $(X^T X)^{-1}$  must exist. If not, replace  $(X^T X)^{-1}$  by the so called Moore-Penrose inverse  $(X^T X)^+$ .

Example 7.1. (1-dim. regression)

$$y_i = \vartheta_0 + \vartheta_1 x_i + \varepsilon_i, \quad i=1, \dots, 14$$



Solution:

$$\vartheta_1^* = \frac{\sigma_{xy}}{\sigma_x^2}$$

$$\vartheta_0^* = \bar{y} - \vartheta_1^* \bar{x}$$

$$\sigma_{xy} = \frac{1}{n} \sum x_i y_i - \bar{x} \bar{y}$$

$$\sigma_x^2 = \frac{1}{n} \sum x_i^2 - \bar{x}^2$$

$$\bar{x} = \frac{1}{n} \sum x_i, \quad \bar{y} = \frac{1}{n} \sum y_i$$

7.1.2. Logistic Regression

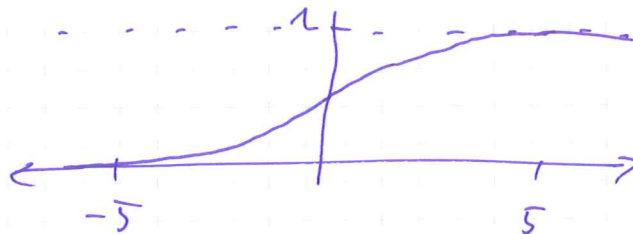
Classification approach. In the following regression with output 0, 1.

Hypothesis:

$$h_{\vartheta}(x) = g(\vartheta^T x) = \frac{1}{1 + e^{-\vartheta^T x}} \in (0, 1)$$

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

Plot



It holds  $g'(z) = g(z)(1 - g(z))$

Given training set  $(x_i, y_i)$ ,  $i=1, \dots, n$ ,  $x_i \in \mathbb{R}^p$ ,  $y_i \in \{0, 1\}$

As before set  $x_{i0} = 1$ , s.t.

$$\mathcal{L}^T x_i = \mathcal{L}_0 + \sum_{j=1}^p \mathcal{L}_j x_{ij}, \quad x_i = (1, x_{i1}, \dots, x_{ip})$$

Probabilistic interpretation:

Assume:  $P(y=1 | \mathbb{R}^p x, \mathcal{L}) = \text{hse}(x)$

$$P(y=0 | x, \mathcal{L}) = 1 - \text{hse}(x)$$

Write:  $p(y | x, \mathcal{L}) = (\text{hse}(x))^y (1 - \text{hse}(x))^{1-y}$ ,  $y \in \{0, 1\}$

Assume  $n$  independent training samples

$$X = \begin{pmatrix} 1 & x_1^T \\ \vdots & \vdots \\ 1 & x_n^T \end{pmatrix} \in \mathbb{R}^{n \times (p+1)}, \quad X = (x_{ij})_{\substack{1 \leq i \leq n \\ 1 \leq j \leq p+1}}$$

Likelihood

$$L(\mathcal{L}) = p(y | X, \mathcal{L}) = \prod_{i=1}^n p(y_i | x_i, \mathcal{L})$$

$$= \prod_{i=1}^n (\text{hse}(x_i))^{y_i} (1 - \text{hse}(x_i))^{1-y_i}$$

Log likelihood

$$l(\mathcal{L}) = \log L(\mathcal{L}) = \sum_{i=1}^n \left[ y_i \log \text{hse}(x_i) + (1-y_i) \log (1 - \text{hse}(x_i)) \right] \quad (*)$$

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

An algorithm herefor: gradient ascent with updates

$$\mathcal{D}^{(k+1)} = \mathcal{D}^{(k)} + \alpha \nabla_{\mathcal{D}} \ell(\mathcal{D}^{(k)}), \quad \alpha \text{ learning parameter}$$

Compute  $\nabla$  for each addend in  $(*)$ . Set  $(x_i, y_i) = (x, y)$

$$\begin{aligned} & \frac{\partial}{\partial \mathcal{D}_j} \left[ y \log h_{\mathcal{D}}(x) + (1-y) \log (1-h_{\mathcal{D}}(x)) \right] \\ &= \left( y \frac{1}{g(\mathcal{D}^T x)} - (1-y) \frac{1}{1-g(\mathcal{D}^T x)} \right) \frac{\partial}{\partial \mathcal{D}_j} g(\mathcal{D}^T x) \\ &= (\dots) g(\mathcal{D}^T x) (1-g(\mathcal{D}^T x)) \frac{\partial}{\partial \mathcal{D}_j} \mathcal{D}^T x \\ &= \left( y(1-g(\mathcal{D}^T x)) - (1-y)g(\mathcal{D}^T x) \right) x_j \\ &= (y - h_{\mathcal{D}}(x)) x_j \end{aligned}$$

Hence

$$\frac{\partial}{\partial \mathcal{D}_j} \ell(\mathcal{D}) = \sum_{i=1}^n (y_i - h_{\mathcal{D}}(x_i)) x_{ij}$$

Update rule:

$$\mathcal{D}_j := \mathcal{D}_j + \alpha \sum_{i=1}^n (y_i - h_{\mathcal{D}}(x_i)) x_{ij}$$

Alternativ: Newton's method

$$\mathcal{D} := \mathcal{D} - H^{-1} \nabla_{\mathcal{D}} \ell(\mathcal{D})$$

$H$ : Hessian matrix

### 7.1.3 The perceptron learning algorithm

Logistic regression has values in  $(0,1)$  as output.

(soft decision). Now, force values to be 0 or 1 by

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

$$h_w(x) = g(w^T x)$$

Use the update rule, analogously to the above

$$w_j := w_j + \alpha (y_i - h_w(x_i)) x_{ij}$$

This is the perceptron learning alg.

Was taken as a rough model for how neurons in the brain work.

Lacking: meaningful interpretation, no MLE.

## 7.2. Reinforcement Learning

By now: training sets

Now: no training set, instead: reward function  
which indicates if the learning is doing well or not.

Applications: robot leg coordination, autonomous flying,  
cell phone routing, factory control, ...

Basis for an analytical description are

### 7.2.1. Markov Decision Processes (MDP)

MDP is a tuple  $(S, A, \{P_{sa}\}, \gamma, R)$

- $S$  set of states
- $A$  set of actions
- $P_{sa}$  transition probabilities

$P_{sa}$  is a prob. distr. over ~~the~~  $S$ ,  
transition to take if the present state is  $s$   
and action  $a$  is taken.

- $\gamma \in [0, 1]$  discount factor
- $R : S \times A \rightarrow \mathbb{R}$  reward function, often  $R : S \rightarrow \mathbb{R}$  only.

Dynamics of the process:

$$s_0 \xrightarrow{a_0} s_1 \xrightarrow{a_1} s_2 \xrightarrow{a_2} s_3 \rightarrow \dots$$

Total payoff  $R(s_0, a_0) + \gamma R(s_1, a_1) + \gamma^2 R(s_2, a_2) + \dots$

or  $R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \gamma^3 R(s_3) + \dots$

Goal:  $\mathbb{E} \max_{\text{over actions in } A} E[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots]$