

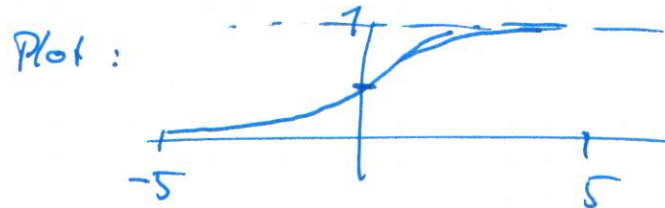
7.1.2. Logistic Regression

Classification like regression with finitely many values of y . In the following only $y \in \{0, 1\}$.

Hypothesis:

$$h_w(x) = g(w^T x) = \frac{1}{1 + e^{-w^T x}} \in [0, 1]$$

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



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

Given training (x_i, y_i) , $i = 1, \dots, m$

$$x_i \in \mathbb{R}^p, y_i \in \{0, 1\}$$

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

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

(i) $X(X^T X)^{-1} X^T$ is an orth. proj. onto $\text{Im}(X)$
provided $(X^T X)^{-1}$ exists.

$$\hat{y} = X(X^T X)^{-1} X^T y = \arg \min_{z \in \text{Im}(X)} \|y - z\|$$

(ii) $\hat{y} = X\alpha$

$$\Rightarrow \cancel{X^T X} (X^T X)^{-1} X^T y = X^T X \alpha$$

$$\Rightarrow X^T y = X^T X \alpha \quad (\text{normal equations})$$

and ~~α~~ *

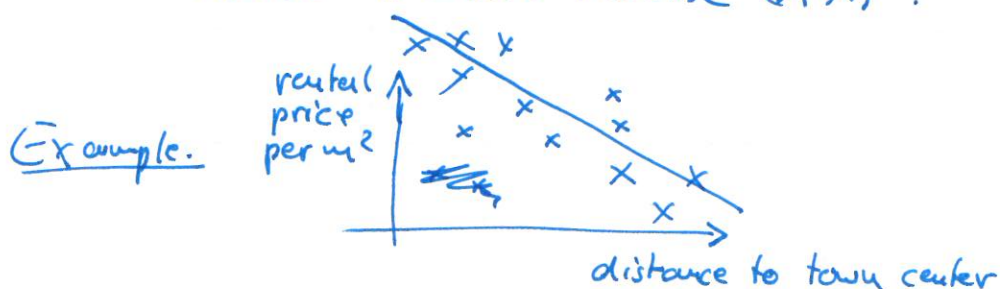
$$\Rightarrow \alpha^* = (X^T X)^{-1} X^T y$$

$$\text{and } X\alpha^* = X(X^T X)^{-1} X^T y = \hat{y}$$

In summary:

$$\alpha^* = (X^T X)^{-1} X^T y \text{ is a solution of (1).}$$

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



7.1.1. Linear Regression

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

Assumption / Hypothesis

$$y_i = \alpha_0 + x_{i1} \alpha_1 + \dots + x_{ip} \alpha_p + \varepsilon_i$$

ε_i : random vector

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

Hence, learn $h_\alpha(x) = (1, x^T) \alpha$,

$\alpha = (\alpha_0, \alpha_1, \dots, \alpha_p)^T$ parameter

Set $X = \begin{pmatrix} 1 & x_1^T \\ \vdots & \vdots \\ 1 & x_n^T \end{pmatrix}$, $y = (y_1, \dots, y_n)^T$
 $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^T$

Then $y = X \alpha + \varepsilon$

$\begin{matrix} n & n \times (p+1) & \leftarrow p+1 & \leftarrow n \\ \downarrow & \downarrow & & \\ y & X & \alpha & \varepsilon \end{matrix}$

Problem: Find the "best" α

$$\min_{\alpha \in \mathbb{R}^{p+1}} \|y - X\alpha\| \quad (1)$$

Solution: (i) project y onto $\text{Im}(X)$: \hat{y}

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

Example 6.6. (Polynomial kernel)

$$K(x, y) = (x^T y + c)^d$$

$$x, y \in \mathbb{R}^p, c \in \mathbb{R}, d \in \mathbb{N}, d \geq 2.$$

Feature space of $\dim \binom{p+d}{d}$

containing all monomials of degree $\leq d$. \square

7. Machine Learning

7.1 Supervised Learning

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

$x_i \in \mathcal{X}$: input variables, feature variables

$y_i \in \mathcal{Y}$: output variables, target variables

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

Supervised learning problem: determine a function

$$h: \mathcal{X} \rightarrow \mathcal{Y}$$

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

If \mathcal{Y} continuous: regression problem

\mathcal{Y} discrete: classification problem

Def. 6.4. Kernel $K(x, y)$ is called valid if there is a feature function ϕ s.t.

$$K(x, y) = \langle \phi(x), \phi(y) \rangle \text{ for all } x, y \in \mathbb{R}^p. \quad \square$$

Th. 6.5 (Mercer)

Given $K : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$. K is valid if and only if for any $x_1, \dots, x_n \in \mathbb{R}^p$ the kernel matrix

$$(K(x_i, x_j))_{i,j=1,\dots,n} \text{ is u.p.d.} \quad \square$$

Proof. only " \Rightarrow "

$$K \text{ valid} \Rightarrow \exists \phi : K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle = \langle \phi(x_j), \phi(x_i) \rangle$$

hence symmetry \checkmark . Moreover

$$\begin{aligned} & z^T (K(x_i, x_j))_{i,j=1,\dots,n} z \\ &= z^T (\langle \phi(x_i), \phi(x_j) \rangle)_{i,j=1,\dots,n} z = \sum_{k,l} z_k z_l \langle \phi(x_k), \phi(x_l) \rangle \\ &= \left\langle \sum_k z_k \phi(x_k), \sum_l z_l \phi(x_l) \right\rangle \geq 0 \quad \square \end{aligned}$$

Needed: inner product in the feature space
 $\{\phi(x) \mid x \in \mathbb{R}^p\}$.

Example 6.2.

$$x, y \in \mathbb{R}^p, \quad K(x, y) = \langle x, y \rangle^2 = \left(\sum_{i=1}^p x_i y_i \right)^2$$

Question: Is there some ϕ s.t. $\langle x, y \rangle^2$ is
 an inner product in the feature space?

$$p=2: \quad x = (x_1, x_2)^T, \quad y = (y_1, y_2)^T$$

$$\text{Use } \phi(x) = (x_1^2, x_2^2, x_1 x_2, x_2 x_1)^T : \mathbb{R}^2 \rightarrow \mathbb{R}^4$$

$$\begin{aligned} \langle \phi(x), \phi(y) \rangle &= x_1^2 y_1^2 + x_2^2 y_2^2 + x_1 x_2 y_1 y_2 + x_2 x_1 y_2 y_1 \\ &= x_1^2 y_1^2 + x_2^2 y_2^2 + 2x_1 x_2 y_1 y_2 \\ &= (x_1 y_1 + x_2 y_2)^2 = \langle x, y \rangle^2 \quad \square \end{aligned}$$

Example 6.3. (Gaussian kernel)

$$x, y \in \mathbb{R}^p, \quad K(x, y) = \exp\left(-\frac{\|x-y\|^2}{2\sigma^2}\right)$$

Question: \exists feature fct. ϕ and a feature space with
 $\langle \cdot, \cdot \rangle$?
 \perp

6.6. Kernels

raw data ~~→~~ (attributes)

→ transformed data (features)

$x_i \rightarrow \phi(x_i)$, ϕ feature mapping

$$\begin{aligned} \textcircled{1) \quad} \max_{\lambda} \quad g(\lambda) &= \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \lambda_i \lambda_j x_i^T x_j \\ \text{s.t.} \quad &0 \leq \lambda \leq c \\ &\sum_{i=1}^n \lambda_i y_i = 0 \end{aligned}$$

$g(\lambda)$ only depends on the inner products
 $x_i^T x_j = \langle x_i, x_j \rangle$

Substitute x_i by $\phi(x_i)$ and use some inner product
 $\langle \cdot, \cdot \rangle$. Replace $x_i^T x_j$ by $K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$

Remark: $K(x, y)$ is often easier to compute
 than $\phi(x)$ itself.

Intuition: If $\phi(x), \phi(y)$ are close, $\langle \phi(x), \phi(y) \rangle$ is large.

$\phi(x) \perp \phi(y)$ then $\langle \phi(x), \phi(y) \rangle = 0$. Hence,
 $K(x, y)$ measures similarity of x and y .