

## 6.5. The SMO Algorithm

Sequential Minimal Optimization

- Select a pair  $(i, j)$  of  $\lambda$ 's to update
- Optimize the obj. fun. w.r.t.  $(\lambda_i, \lambda_j)$
- Explicit solution ~~for~~  $(\lambda_i^*, \lambda_j^*)$
- Iterate until convergence

## 6.6. Kernels

Instead of applying SVM to the raw data ("attributes")  $x_i$  apply it to transformed data ("features")  $\phi(x_i)$ .  
 $\phi$  is called feature mapping.

Aim: achieve better separability.

$$(D) \quad \max_{\lambda} g(\lambda) = \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \lambda_i \lambda_j x_i^T x_j$$

$g(\lambda)$  only depends on the inner product  $x_i^T x_j$ .

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.

If  $\phi(x) \perp \phi(y)$  then  $\langle \phi(x), \phi(y) \rangle = 0$ . Hence,

$K(x, y)$  measures how similar  $x$  and  $y$  are.

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

Example 6.2.

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

Question: Is there some  $\phi$  such that  $\langle x, y \rangle^2$   
is an inner product in the feature space.

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

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

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

Example 6.3. (Gaussian kernel)

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

Question:  $\exists$  feature mapping  $\phi$  and a feature space  
with  $\langle \cdot, \cdot \rangle$ ?

Def. 6.4. Kernel  $K(x, y)$  is called valid if there exists a feature function  $\phi$  such that  

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

Theorem 6.5. (Mercer)

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

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

Proof. only " $\Rightarrow$ ":

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

hence symmetry. Moreover

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



Example 6.6. (Polynomial kernel)

$$K(x, y) = (x^T y + c)^d, \quad x, y \in \mathbb{R}^p, c \in \mathbb{R}, d \in \mathbb{N}, d \geq 2$$

Feature space of dim  $\binom{p+d}{d}$  containing all monomials of degree  $\leq d$ .  $\perp$

Ex determine  $\phi$

## 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 variables, target variables

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

Supervised learning problem: determine a

function  $h: X \rightarrow Y$  (a "hypothesis")

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

If  $Y$  is continuous: regression problem

$Y$  is discrete: classification problem

7.1.1 Linear Regression

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

$$y_i = \alpha_0 + x_{i1} \alpha_1 + \dots + x_{ip} \alpha_p + \epsilon_i, \epsilon_i: \text{random error}$$

$$= (1, \underline{x}_i^T) \underline{\alpha} + \epsilon_i$$

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

Set  $X = \begin{pmatrix} 1 & \underline{x}_1^T \\ \vdots & \vdots \\ 1 & \underline{x}_n^T \end{pmatrix}$ ,  $\underline{y} = (y_1, \dots, y_n)^T$ ,  $\underline{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^T$

Then  $\underline{y} = X \underline{\alpha} + \underline{\epsilon}$

Problem: find the best  $\underline{\alpha}$  by solving

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

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

(ii) find  $\underline{\alpha}$  s.t.  $\hat{\underline{y}} = X \underline{\alpha}$

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

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

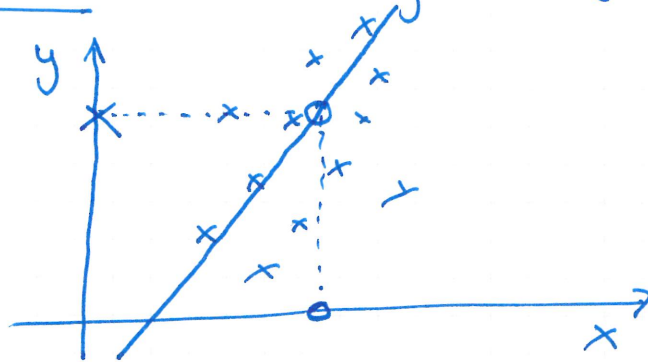
(ii)  $\hat{\underline{y}} = X \underline{\alpha} \Rightarrow X^T X (X^T X)^{-1} X^T \underline{y} = X^T X \underline{\alpha}$   
 ~~$\Rightarrow \underline{\alpha} = (X^T X)^{-1} X^T \underline{y}$~~   
 $\Rightarrow \underline{\alpha}^* = (X^T X)^{-1} X^T \underline{y}$

is a solution.

In summary  $\mathcal{Q}^* = (X^T X)^{-1} X^T y$  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. Linear Regression (1-dim)



Solution

$$\mathcal{Q}_1^* = \frac{\sigma_{xy}}{\sigma_x^2}, \quad \mathcal{Q}_0^* = \bar{y} - \mathcal{Q}_1^* \bar{x}$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

$$\sigma_{xy} = \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y}, \quad \sigma_x^2 = \sigma_{xx} \quad \perp$$