

6.4. SVM - soft margins

$$(P) \quad \min_{a \in \mathbb{R}^p, b, \xi \in \mathbb{R}} \left\{ \frac{1}{2} \|a\|^2 + c \sum_{i=1}^n \xi_i \right\} \quad (c \text{ parameters})$$

$$\text{s.t.} \quad y_i (a^T x_i + b) \geq 1 - \xi_i$$

$$\xi_i \geq 0, \quad i=1, \dots, n$$

Consequences

$$\xi_i > 1 \Rightarrow \text{misclassification}$$

$$0 < \xi_i \leq 1 \Rightarrow \text{correctly classified, but lies in the margin}$$

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

$$\text{s.t.} \quad 0 \leq \lambda_i \leq c, \quad i=1, \dots, n$$

$$\sum_{i=1}^n \lambda_i y_i = 0$$

Investigating complementary slackness

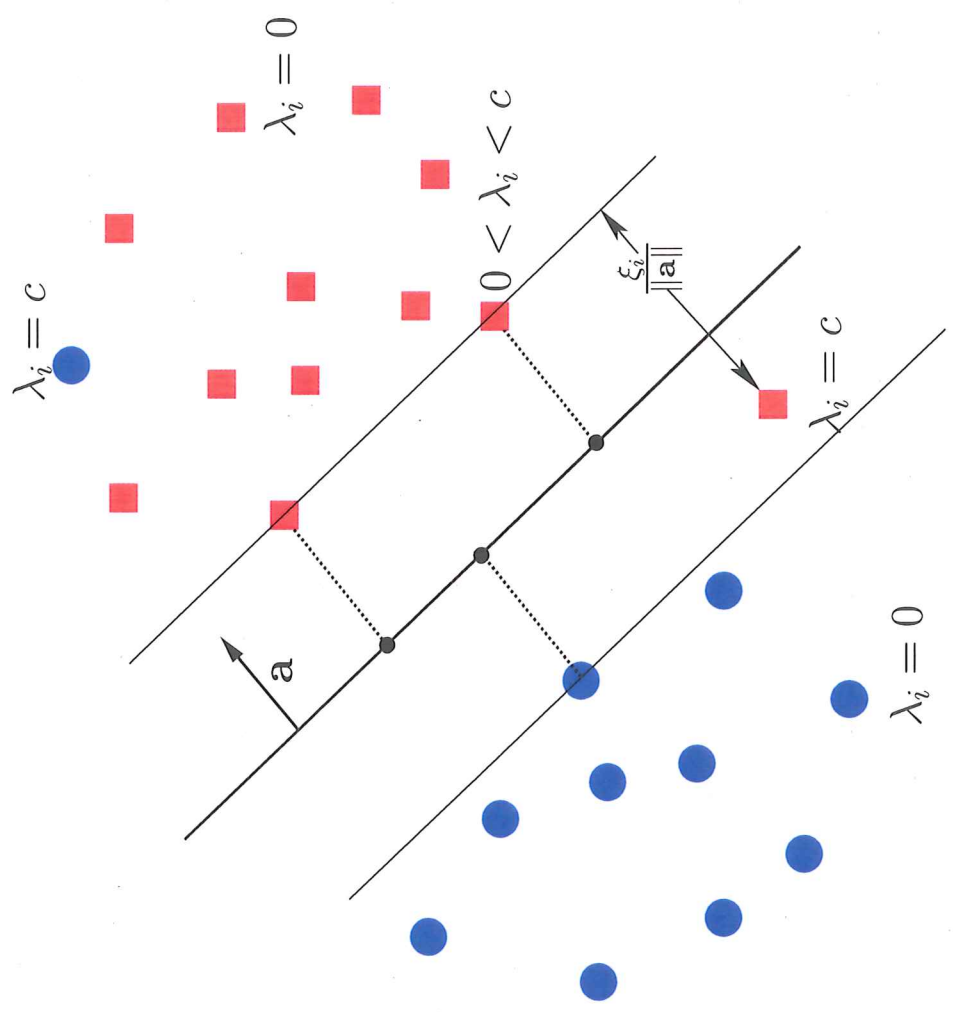
$$0 < \lambda_i < c \Rightarrow y_i (a^T x_i + b) = 1 \quad (\text{margin SVs})$$

$$\lambda_i = c \Rightarrow y_i (a^T x_i + b) = 1 - \xi_i, \quad \xi_i \geq 0$$

$$\lambda_i = 0 \Rightarrow y_i (a^T x_i + b) \geq 1 \quad (\text{non SVs})$$

(margin error)

Hence, the solution is sparse, most of the points have $\lambda_i = 0$.



6.5. The SMO-Algorithm

Sequential Minimal Optimization to solve the dual problem.

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

Assume λ is a feasible point, i.e., λ satisfies the constraints.

Idea: "cyclic coordinatewise optimization" does not work, since, e.g.,

$$\lambda_1 y_1 = - \sum_{i=2}^n \lambda_i y_i \quad \text{or} \quad \lambda_1 = -y_1 \sum_{i=2}^n \lambda_i y_i$$

Hence each λ_j is determined by fixing $\lambda_i, i \neq j$.

Idea: update at least two λ_j simultaneously \rightarrow SMO alg.

(SMO) -

repeat

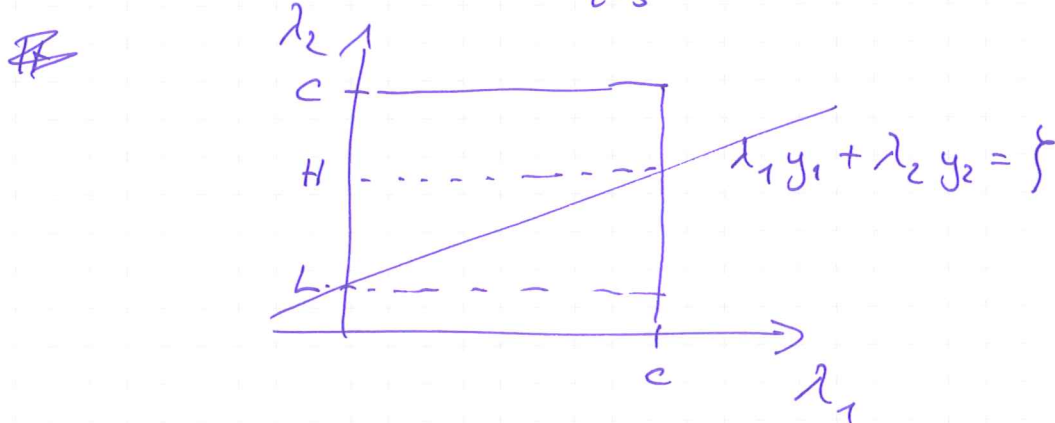
1. Select a pair (i, j) to update next (the one which promises most progress)
 2. Optimize $W(\lambda)$ w.r.t. λ_i and λ_j while holding $\lambda_k, k \neq i, j$, fixed
- until convergence

Check KKT within a tolerance limit

$\epsilon = 0.01$ or 0.001 , to verify convergence.

Optimize $W(\lambda)$ w.r.t. λ_1, λ_2 with $\lambda_3, \dots, \lambda_n$ fixed,
 λ feasible.

It holds $\lambda_1 y_1 + \lambda_2 y_2 = -\sum_{i=3}^n \lambda_i y_i = \xi$, ξ fixed.



Derive $0 \leq \lambda_1, \lambda_2 \leq c$ (*)
 $L \leq \lambda_2 \leq H$

Moreover $\lambda_1 = (\xi - \lambda_2 y_2) y_1$ (**)

Hence: $W(\lambda_1, \dots, \lambda_n) = W((\xi - \lambda_2 y_2) y_1, \underbrace{\lambda_2, \lambda_3, \dots, \lambda_n}_{\text{fixed}})$

is a quadratic function of λ_2 .

It can be written as (remember: $y_i \in \{-1, 1\}$)

$$\delta_2 \lambda_2^2 + \delta_1 \lambda_2 + \delta_0 \quad , \quad \delta_0, \delta_1, \delta_2 \text{ appropriate}$$

Determine the max. by differentiation.

$$2\delta_2 \lambda_2 + \delta_1 \stackrel{!}{=} 0 \quad \lambda_2 = -\frac{\delta_1}{2\delta_2}$$

with optimum solution $\lambda_2^{(r)}$ (r: row)

The final solution of (*) is

$$\lambda_2^{(c)} = \begin{cases} H, & \text{if } \lambda_2^{(r)} \geq H \\ \lambda_2^{(r)}, & \text{if } L \leq \lambda_2^{(r)} \leq H \\ L, & \text{if } \lambda_2^{(r)} < L \end{cases} \quad (c: \text{clipped})$$

λ_1 is computed from (**).

Still to clarify:

- o What is the best choice of the next pair (i, j) to update.
- o How to update the coefficients $\gamma_0, \gamma_1, \gamma_2$ in the run of SMO.
- o The algorithm converges, however, the right choice of (i, j) in each step accelerates the rate of convergence.
- o Osuna, Freund, Girosi (1997):
generalization of SMO algorithm.

6.6. Kernels

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

ϕ is called feature mapping.

Aim: achieve better separability.

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

$$\text{s.t.} \quad 0 \leq \lambda_i \leq c$$

$$\sum_{i=1}^n \lambda_i y_i = 0$$

$g(\alpha)$ only depends on the inner products $x_i^T x_j$.

Substitute x_i by $\phi(x_i)$ and use some inner product

$\langle \cdot, \cdot \rangle$ - Replace $x_i^T x_j$ by

$$\langle \phi(x_i), \phi(x_j) \rangle = K(x_i, x_j)$$

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.

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

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

$$p=2: \quad x = (x_1, x_2)^T, \quad z = (z_1, z_2)^T$$

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

$$\begin{aligned} \langle \phi(x), \phi(z) \rangle &= x_1^2 z_1^2 + x_2^2 z_2^2 + x_1 x_2 z_1 z_2 + x_2 x_1 z_2 z_1 \\ &= x_1^2 z_1^2 + x_2^2 z_2^2 + 2 x_1 x_2 z_1 z_2 \\ &= (x_1 z_1 + x_2 z_2)^2 = \langle x, z \rangle^2 \end{aligned}$$

Example. 6.3. (Gaussian Kernel)

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

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

Def. 6.4. Kernel $K(x, z)$ is called valid, if there exists a feature function ϕ such that $K(x, z) = \langle \phi(x), \phi(z) \rangle$ for all $x, z \in \mathbb{R}^p$.

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 $\{x_1, \dots, x_n\}$ the kernel matrix $(K(x_i, x_j))_{i,j=1, \dots, n}$ is u.u.d. \perp

Proof. \Rightarrow

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

Moreover,

$$z^T (K(x_i, x_j))_{i,j} z = z^T (\langle \phi(x_i), \phi(x_j) \rangle)_{i,j} z$$

$$= \sum_{k,e} z_k z_e \langle \phi(x_k), \phi(x_e) \rangle$$

$$= \langle \sum_k z_k \phi(x_k), \sum_e z_e \phi(x_e) \rangle \geq 0 \quad \square$$

Example 6 (Polynomial kernel)

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

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

Determine $\phi(x) \rightarrow \mathbb{R}^{\dots}$. \perp

Kernels can also be constructed over
infinite dimensional spaces, e.g.,
function space or probability distributions.
Provides a lot of modeling power.
The solution of the opt. problem is still a
convex problem with linear constraints.