

## Vorlesung 16 Support vector machines

In den letzten Vorlesungen haben wir Regressionsprobleme behandelt. Dabei war die Response-Variable  $y \in \mathbb{R}^N$  kontinuierlich.

Heute:  $y \in \{-1, 1\}$  ist binär.

Wir wollen wieder ein Modell  $f_\Theta: \mathbb{R}^D \rightarrow \{-1, 1\}$  finden, das die Trainingsdaten  $(x_1, y_1), \dots, (x_n, y_n) \in \mathbb{R}^D \times \{-1, 1\}$  beschreibt.

Support vector machines sind ein Modell für dieses Setting:

$$f_\Theta(x) = \text{sgn}(\langle a, x \rangle + b), \quad \Theta = (a, b) \in \mathbb{R}^D \times \mathbb{R}.$$

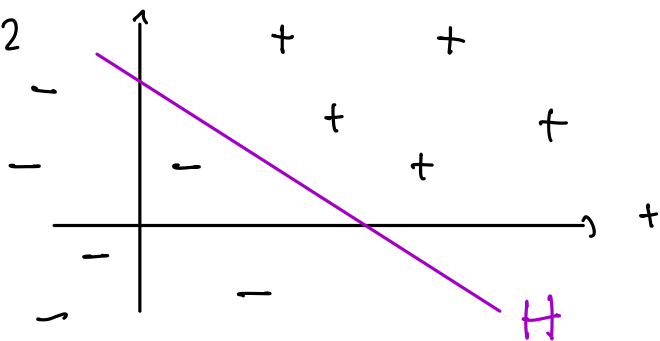
Beachte Wir behandeln hier ein deterministisches Modell.

Die Geometrie dieses Modells ist wie folgt.

Für alle  $\Theta = (a, b) \in \mathbb{R}^D \times \mathbb{R}$  haben wir die Hyperebene

$$H = \{x \in \mathbb{R}^D \mid \langle a, x \rangle + b = 0\}$$

Beispiel  $D=2$



Ziel Daten durch eine Hyperebene trennen.

Es bleibt zu klären, wie wir entsprechende Parameter für unsere Daten finden.

Dazu brauchen wir zwei Lemmata.

### Lemma 16.1

Sei  $x \in \mathbb{R}^D$  und  $\theta = (a, b) \in \mathbb{R}^D \times \mathbb{R}$ . Für alle  $y \in \{-1, 1\}$  gilt:

$$y = f_{\theta}(x) \iff y(\langle a, x \rangle + b) > 0.$$

### Beweis

Möglichkeit 1):  $y = f_{\theta}(x) = -1$ . Dann ist  $\langle a, x \rangle + b < 0$

$$\rightarrow y(\langle a, x \rangle + b) > 0$$

Möglichkeit 2)  $y = f_{\theta}(x) = 1$ . Dann ist  $\langle a, x \rangle + b > 0$

$$\rightarrow y(\langle a, x \rangle + b) > 0. \quad \square$$

### Lemma 16.2

Sei  $x_0 \in \mathbb{R}^D$  und  $H = \{x \in \mathbb{R}^D \mid \langle a, x \rangle + b = 0\}$ .

Sei  $\langle a, a \rangle = 1$ . Der Euklidische Abstand von  $x_0$  zu  $H$  ist

$$|\langle a, x_0 \rangle + b| = y_0(\langle a, x_0 \rangle + b), \quad y_0 = f_{\theta}(x_0).$$

### Beweis

Sei  $z \in H$  der Punkt auf  $H$ , der den Abstand zu  $x_0$  minimiert.

Dann können wir

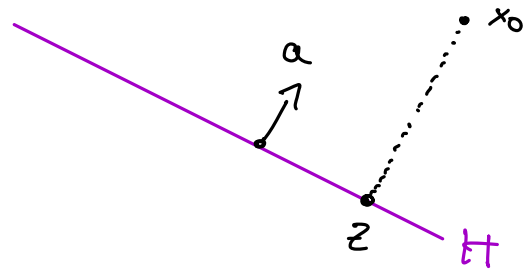
$$x_0 = z + \gamma a \quad \text{für } \gamma \in \mathbb{R}$$

schreiben, wobei  $|\gamma| = \|z - x_0\| = \text{Abstand von } x_0 \text{ zu } H$ .

Beachte:  $x_0 - z = \gamma a$  ist orthogonal zu  $H$ .

Daraus folgt:

$$\begin{aligned} \langle a, x_0 \rangle + b &= \langle a, z + \gamma a \rangle + b \\ &= \langle a, z \rangle + \gamma \langle a, a \rangle + b = \gamma, \quad \text{da } z \in H. \end{aligned}$$



$$\Rightarrow |\langle a, x_0 \rangle + b| = |r| = \text{Abstand von } x_0 \text{ zu } H.$$

Eine geeignete Parameterwahl ist es nun den Abstand der Input-Daten  $x_1, \dots, x_n \in \mathbb{R}^D$  zur Hyperebene  $H$  zu maximieren unter den Bedingungen, dass  $f_\Theta(x_i) = y_i$ . Dies führt zu folgendem Optimierungsproblem:

$$\begin{aligned} \|a\|^2 = \langle a, a \rangle = 1 \quad & \xrightarrow{\quad} \max_{\Theta = (a, b) \in \mathbb{R}^D \times \mathbb{R}} \\ (\alpha) \quad & \text{s.t.: } y_k (\langle a, x_k \rangle + b) \geq r. \quad \text{für } k=1, \dots, n, \\ & r \geq 0. \end{aligned}$$

Abstand zu  $H$ .

↑ Nebenbedingungen

Das Optimierungsproblem  $(\alpha)$  wird oft in einer äquivalenten Form beschrieben: Sei

$$a' = \frac{a}{r}, \quad b' = \frac{b}{r}$$

Dann erhalten wir die Constraints

$$y_k (\langle a', x \rangle + b') \geq 1.$$

und es gilt:

$$\|a'\| = \frac{\|a\|}{r} = \frac{1}{r}$$

Definition 16-3.

Die Hard-Margin-SVM ist durch folgendes Optimierungsproblem gegeben:

$$\begin{aligned} \min_{\Theta = (a, b) \in \mathbb{R}^D \times \mathbb{R}} \quad & \|a\|^2 \\ \text{s.t.} \quad & y_k (\langle a, x_k \rangle + b) \geq 1 \quad \text{für } k=1, \dots, n. \end{aligned}$$

Das Problem mit Hard-Margin SVM ist, dass Ausreißer in den Daten dafür sorgen, dass das Optimierungsproblem keine Lösung hat. Soft-Margin SVM löst dieses Problem wie folgt:

### Definition 16.4

Soft-Margin SVM ist durch folgendes Optimierungsproblem gegeben.

$$\begin{aligned} \min_{a, b, \xi} \quad & \|a\|^2 + C \sum_{k=1}^n \xi_k \\ \text{s.t.} \quad & y_k (\langle a, x_k \rangle + b) \geq 1 - \xi_k, \quad k=1, \dots, n \\ & \xi_k \geq 0, \end{aligned}$$

wobei  $C > 0$  ein Regularisierungsparameter ist.

### Proposition 16.5

Der Hinge-Loss ist  $\ell(y, \hat{y}) = \max\{0, 1 - y \cdot \hat{y}\}$

Soft-Margin SVM ist durch ERM bzgl. eines regularisierten Hinge-Loss gegeben.

Beweis

Übung.

Hard-Margin-SVM und Soft-Margin-SVM werden

Primal SVMs genannt. Eine andere Formulierung ist die Dual SVM.

Wir betrachten dazu die Lagrange-Funktion

$$\mathcal{L}(a, b, \xi, \alpha, \beta) = \|a\|^2 + C \sum_{k=1}^n \xi_k - \sum_{k=1}^n \alpha_k (y_k (\langle a, x_k \rangle + b) - (1 - \xi_k)) - \sum_{k=1}^n \beta_k \xi_k.$$

xx) ist dann äquivalent zu

$$\min_{a, b, \xi} \max_{\alpha, \beta} \mathcal{L}(a, b, \xi, \alpha, \beta).$$

$$\text{s.t. } \alpha_k, \beta_k \geq 0, \quad k=1, \dots, n.$$

Die KKT-Bedingungen sagen, dass ein optimaler Wert folgendes erfüllen muss:

$$\frac{\partial \mathcal{L}}{\partial a} = \frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial \xi} = 0.$$

Sei jetzt  $u := \left( \frac{1}{2} \alpha_k y_k \right)_{k=1}^n$ ,  $v = (\langle a, x_k \rangle + b)_{k=1}^n$ ,

$e = (1, \dots, 1)^T$ . Dann:

$$\mathcal{L} = \|a\|^2 + C \sum_{k=1}^n \xi_k - \sum_{k=1}^n \alpha_k (y_k (\langle a, x_k \rangle + b) - (1 - \xi_k)) - \sum_{k=1}^n \beta_k \xi_k.$$

$$= \langle a, a \rangle + C \cdot \langle e, \xi \rangle - 2 \langle u, v \rangle - \langle \alpha + \beta, \xi \rangle + \langle \alpha, e \rangle$$

Damit folgt:

$$\text{xxx)} \quad \begin{cases} \frac{\partial \mathcal{L}}{\partial a} = 2a - 2 \sum_{k=1}^n u_k x_k = 0 & \Rightarrow a = \sum_{k=1}^n u_k x_k. \\ \frac{\partial \mathcal{L}}{\partial b} = 2 \sum_{k=1}^n u_k = 0 & \Rightarrow \sum_{k=1}^n u_k = 0 \end{cases}$$

$$\left| \frac{\partial \mathcal{L}}{\partial \xi} = C e - (\alpha + \beta) = 0 \Rightarrow \alpha + \beta = C e \Rightarrow \alpha_k \leq C, \quad k=1, \dots, n.$$

Bemerkung Der Name "Support Vector Machine" kommt von der Gleichung  $\alpha = \sum_{k=1}^n u_k x_k = \frac{1}{2} \sum_{k=1}^n \alpha_k y_k x_k$ .  
Die  $x_k$  mit  $\alpha_k \neq 0$  heißen "Support Vectors".

Wir setzen die Optimalitätsbedingungen (KKT) in  $\mathcal{L}$  ein:

$$\begin{aligned} \mathcal{L} &= \langle \alpha, \alpha \rangle + C \cdot \langle e, \xi \rangle - \lambda \langle u, v \rangle - \langle \alpha + \beta, \xi \rangle + \langle \alpha, e \rangle \\ &= \langle \alpha, \alpha \rangle - \lambda \langle u, v \rangle + \langle \alpha, e \rangle \quad (\text{weil } C e - (\alpha + \beta) = 0) \\ &= \left\langle \sum_{k=1}^n u_k x_k, \sum_{l=1}^n u_l x_l \right\rangle - \lambda \langle u, v \rangle + \langle \alpha, e \rangle \\ &= \sum_{k=1}^n \sum_{l=1}^n u_k u_l \langle x_k, x_l \rangle - \lambda \langle u, v \rangle + \langle \alpha, e \rangle \\ &= u^T G u - \lambda \langle u, v \rangle + \langle \alpha, e \rangle, \end{aligned}$$

wobei  $G = (\langle x_k, x_l \rangle)_{k,l=1}^n \in \mathbb{R}^{n \times n}$ .

Weiterhin:

$$\begin{aligned} \langle u, v \rangle &= \langle u, (\langle \alpha, x_k \rangle + b)_{k=1}^n \rangle \\ &= \sum_{k=1}^n u_k \langle \alpha, x_k \rangle + b \underbrace{\sum_{k=1}^n u_k}_{=0} \\ &= \sum_{k=1}^n u_k \left\langle \sum_{l=1}^n u_l x_l, x_k \right\rangle \\ &= \sum_{k=1}^n \sum_{l=1}^n u_k u_l \langle x_k, x_l \rangle = u^T G u \end{aligned}$$

$$\Rightarrow \mathcal{L} = -u^T G u + \sum_{k=1}^n \alpha_k$$

Definition 16.6

Die Dual SVM ist gegeben durch das Optimierungsproblem

$$\begin{array}{ll} \max_{\alpha} & -u^T G u + \sum_{k=1}^n \alpha_k \\ \text{s.t.} & \sum_{k=1}^n \alpha_k y_k = 0, \quad 0 < \alpha_k \leq C \end{array}$$

wobei  $u = (\sum_{k=1}^n \alpha_k y_k)$  und  $G = (\langle x_u, x_\ell \rangle)_{u, \ell=1}^n$ .