

11 Support Vektor Maschinen

11.1 Optimale Trennebene für linear separable Probleme

- betrachte (separables) Zweiklassenproblem
(auf das jedes Mehrklassenproblem zurückgeführt werden kann)
- Stichprobe: $\{(\vec{c}_n, y_n) \mid n = 1, \dots, N\}$, mit: $y_n \in \{+1, -1\}$
- gesucht ist Trennebene: $\vec{w}^T \vec{c} + b = 0$, sodaß gilt:

$$\vec{w}^T \vec{c}_n + b \geq 0 \text{ falls } y_n = +1 \quad (11.1)$$

$$\vec{w}^T \vec{c}_n + b < 0 \text{ falls } y_n = -1 \quad (11.2)$$

$$(11.3)$$

- mit Umskalierung erhalten wir:

$$y_n(\tilde{w}^T \vec{c}_n + \tilde{b}) \geq 1, \text{ für } n = 1, \dots, N \quad (11.4)$$

11.1 Optimale Trennebene für linear separable Probleme

1. wähle $\epsilon = \min_{n:y_n=-1} |\vec{w}^T \vec{c}_n + b|$

(beachte: für $y_n = +1$ könnte $\epsilon = 0$ werden)

2. sei $y_n = -1$ $\vec{w}^T \vec{c}_n + b < 0$

$$\vec{w}^T \vec{c}_n + b \leq -\epsilon$$

$$\vec{w}^T \vec{c}_n + b + \frac{\epsilon}{2} \leq -\frac{\epsilon}{2} \quad \left| +\frac{\epsilon}{2} \right.$$

$$\frac{2}{\epsilon} \vec{w}^T \vec{c}_n + \frac{2}{\epsilon} \left(b + \frac{\epsilon}{2} \right) \leq -1 \quad \left| \cdot \frac{2}{\epsilon} \right.$$

$$\vec{w}^T \vec{c}_n + \tilde{b} \leq -1,$$

$$\text{mit } \vec{w} := \frac{2}{\epsilon} \vec{w}^T, \tilde{b} := \frac{2}{\epsilon} \left(b + \frac{\epsilon}{2} \right)$$

$$y_n (\vec{w}^T \vec{c}_n + \tilde{b}) \geq 1$$

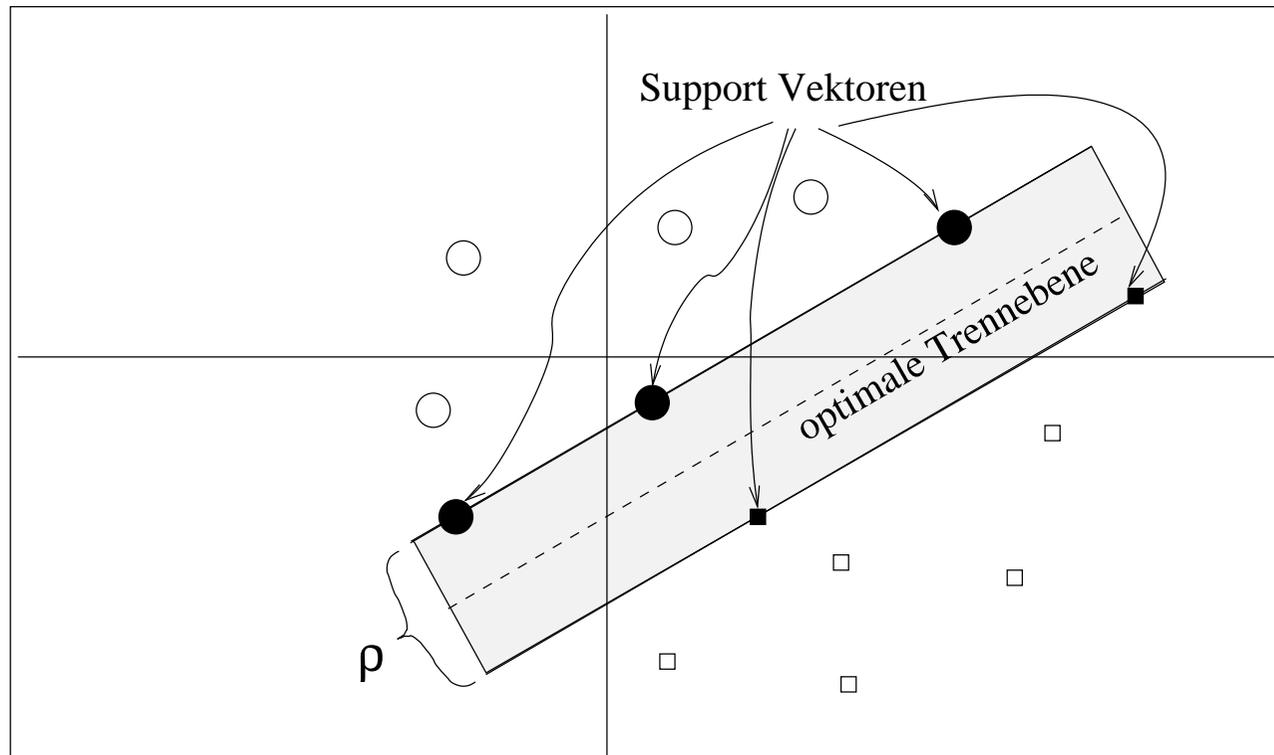
11.1 Optimale Trennebene für linear separable Probleme

3. nun sei $y_n = +1$

$$\begin{aligned} \vec{w}^T \vec{c}_n + b &\geq 0 \\ \vec{w}^T \vec{c}_n + b + \frac{\epsilon}{2} &\geq \frac{\epsilon}{2} && | +\frac{\epsilon}{2} \\ \frac{2}{\epsilon} \vec{w}^T \vec{c}_n + \frac{2}{\epsilon} (b + \frac{\epsilon}{2}) &\geq 1 && | \cdot \frac{2}{\epsilon} \\ \vec{\tilde{w}}^T \vec{c}_n + \tilde{b} &\geq 1 \\ y_n (\vec{\tilde{w}}^T \vec{c}_n + \tilde{b}) &\geq 1 \end{aligned}$$

11.1 Optimale Trennebene für linear separable Probleme

- **margin of separation** ρ ist die Breite des “Schlauches”, der die Stichprobenelemente \vec{c}_n beider Klassen trennt:

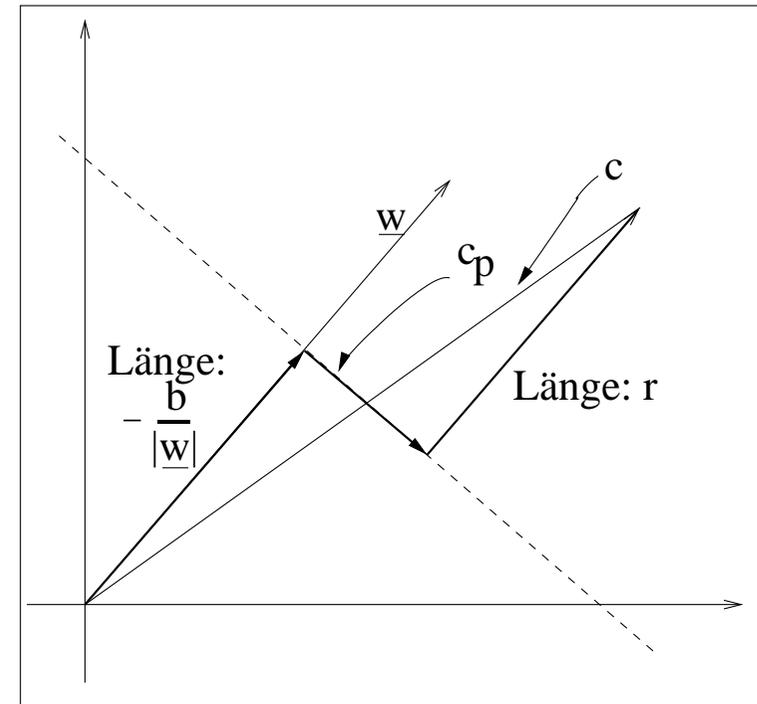


- die Trennebene mit maximalem ρ ist die **optimale Trennebene**
- die \vec{c}_n , die diesen Abstand einnehmen, sind die **Support Vektoren**
es gilt dann in (11.4) die Gleichheit

11.1 Optimale Trennebene für linear separable Probleme

- Abstand r von der Trennebene (\vec{w}, b) ,
mit der Funktion: $d(\vec{c}) := \vec{w}^T \vec{c} + b$ ausgedrückt

$$\begin{aligned}\vec{c} &= \vec{c}_p + \left(r - \frac{b}{|\vec{w}|} \right) \frac{\vec{w}}{|\vec{w}|} \\ d(\vec{c}) &= \vec{w}^T \left(\frac{r}{|\vec{w}|} - \frac{b}{|\vec{w}|^2} \right) \vec{w} + b \\ &= |\vec{w}|^2 \left(\frac{r}{|\vec{w}|} - \frac{b}{|\vec{w}|^2} \right) + b \\ &= r |\vec{w}| - b + b = r |\vec{w}| \\ r &= \frac{d(\vec{c})}{|\vec{w}|}\end{aligned}$$



- für jeden Support Vektoren \vec{c} gilt

$$d(\vec{c}) = \pm 1 \quad (11.5)$$

$$\tilde{r} = \frac{\pm 1}{|\vec{w}|} \quad (11.6)$$

$$\Rightarrow \rho = 2 \cdot \tilde{r} = \frac{2}{|\vec{w}|} \quad (11.7)$$

- Folgerung:

- der Vektor \vec{w} der optimalen Trennebene ist derjenige Vektor minimaler Norm, der die Bedingung (11.4) erfüllt
- die optimale Trennebene ist möglichst “flach” ($|\vec{w}|$ minimale)

- Bemerkung/Definition:

1. für jede Entscheidungsfunktion $d(\vec{c}; \vec{w}, b)$, die korrekt die gesamte Stichprobe trennt, ist auch $d(\vec{c}; \lambda \vec{w}, \lambda b)$, $\lambda > 0$, eine korrekte Trennebene
wir betrachten gerade diejenigen Trennebenen, für die die Support Vektoren c gilt: $d(\vec{c}) = \pm 1$
diese nennen wir **kanonische Trennebenen**

11.1 Optimale Trennebene für linear separable Probleme

2. von allen kanonischen Trennebenen wählen wir diejenige, mit minimalem $|\vec{w}|$

11.2 Strukturelle Risikominimierung

Satz (Vapnik) für eine Menge $d(\vec{c}; \vec{w}, b)$ von kanonischen Trennebenen in \mathbb{R}^d

- für eine Stichprobe der Größe N , sodaß R der Radius der kleinsten Hyperkugel ist, die alle Datenpunkte enthält,
- mit $|\vec{w}| \leq A$, gilt

die VC-Dimension ist $h \leq \min(R^2 A^2, d) + 1$

- alle kanonischen Trennebenen liefern $R_{\text{emp}}[d(\vec{c}; \vec{w}, b)] = 0$
- von all diesen hat diejenige mit minimalem $|\vec{w}|$ die kleinste VC-Dimension
- formal definieren wir also die geschachtelte Menge von Funktionsklassen:

$$H_l := \{d(\vec{c}; \vec{w}, b) \mid d(\vec{c}; \vec{w}, b) \text{ ist kanonische Trennebene mit } |\vec{w}| \leq A_l\}$$

mit $A_1 \leq A_2 \leq \dots \leq A_l$

- und die optimale kanonische Trennebene ist daher Lösung für SRM-Prinzip

11.2 Strukturelle Risikominimierung

Bemerkung es könnte (natürlich) passieren, daß eine Trennebene, die die Stichprobe nicht korrekt trenne, d.h. $R_{\text{emp}}[d(\vec{c}; \vec{w}, b)] > 0$ liefert, noch kleineres $|\vec{w}|$ und daher kleinere VC-Dimension hat,
und u.U. noch kleineres strukturelles Risiko resultiert

11.3 Optimierung

Primäres Optimierungsproblem

- gesucht sind die optimalen Werte für \vec{w} und b , die
 - $\Phi(\vec{w}) = \frac{1}{2}\vec{w}^T \vec{w}$ minimieren und
 - $y_n(\vec{w}^T \vec{c}_n + b) \geq 1$, für $n = 1, \dots, N$ erfüllen
- Minimierung der Lagrange-Funktion:

$$J(\vec{w}, b, \vec{\alpha}) = \frac{1}{2}\vec{w}^T \vec{w} - \sum_{n=1}^N \alpha_n [y_n(\vec{w}^T \vec{c}_n + b) - 1]$$

mit den Lagrange-Multiplikatoren $\alpha_n \geq 0$

- leichtere Lösung mit identischem Optimum mit:

Duales Optimierungsproblem

- primäres und duales Problem haben identische Lösung
- gesucht sind die optimalen Werte für α_n , die

$$- Q(\vec{\alpha}) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{j=1}^I \alpha_n \alpha_j y_n y_j \vec{c}_n^T \vec{c}_j \text{ minimieren und}$$

$$- \sum_{n=1}^N \alpha_n y_n = 0 \text{ sowie}$$

$$- \alpha_n \geq 0 \text{ erfüllen}$$

- das duale Problem enthält \vec{w} und b **nicht** ! Optimum für

$$- \vec{w} = \sum_{n=1}^N \alpha_n y_n \vec{c}_n$$

wobei nur für Support Vektoren \vec{c}_n gelten kann: $\alpha_n \neq 0$

(d.h. es gehen ausschließlich Support Vektoren in die Lösung ein)

$$- b = y_n - \vec{w}^T \vec{c}_n, \text{ für ein beliebiges } i \text{ mit } \alpha_n \neq 0$$

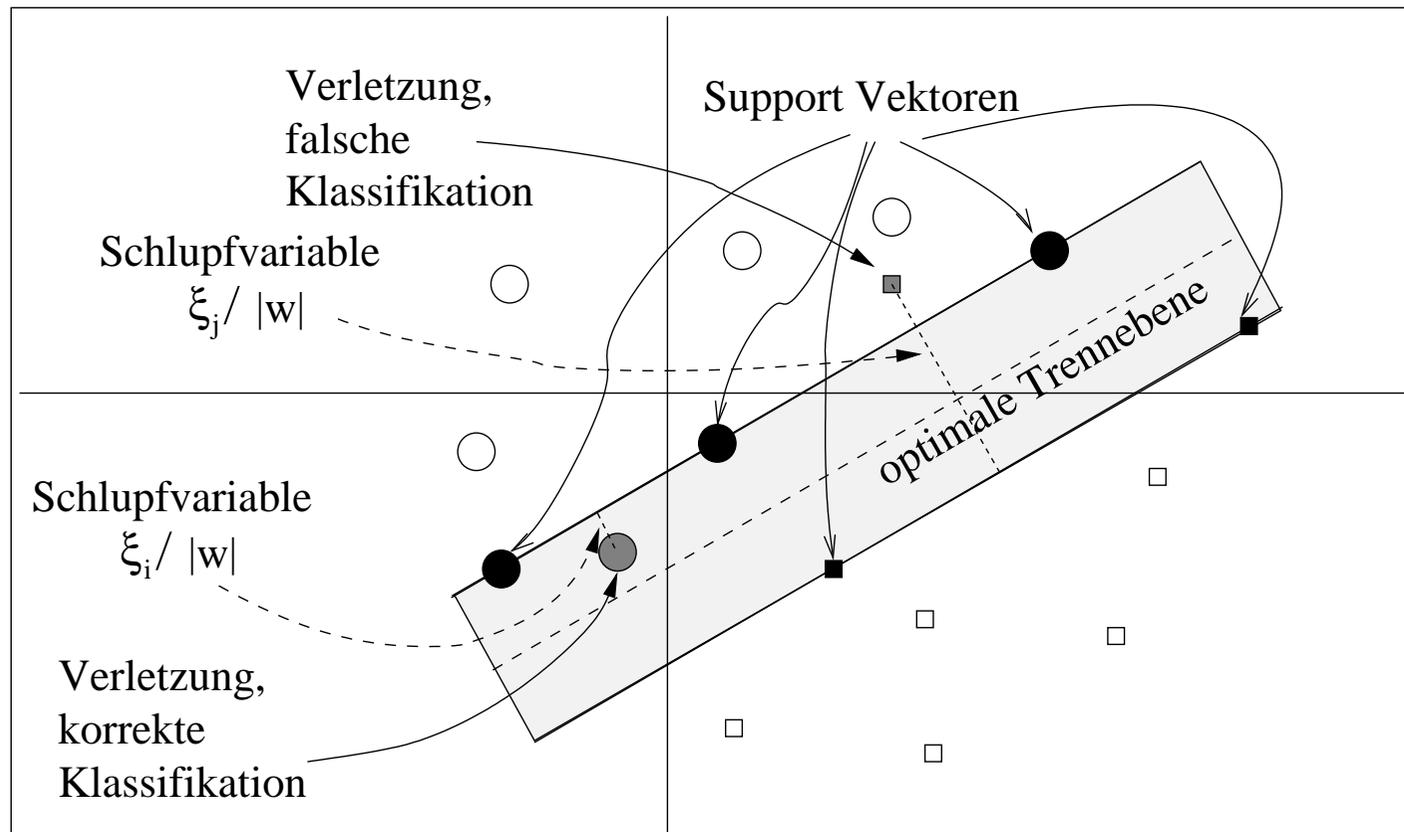
- liefert konvexes quadratisches Optimierungsproblem (nicht trivial)

11.4 Optimale Trennebene für nicht-linear separable Probleme

- erlaube Verletzung der optimalen Trennung um ξ_n :

$$y_n(\vec{w}^T \vec{c}_n + b) \geq 1 - \xi_n, \text{ für } n = 1, \dots, N \quad (11.8)$$

mit den **Schlupfvariablen (slack variables)** $\xi_n \geq 0$



11.4 Optimale Trennebene für nicht-linear separable Probleme

- neues Optimierungsziel:
optimale Trennebene, die den mittleren Klassifikationsfehler minimiert
(d.h. mittlere Anzahl ξ_n , mit $\frac{\xi_n}{|\vec{w}|} > 1$)
führt auf NP-vollständiges Problem, daher:
- als Approximation suche optimalen Werte für \vec{w} , die

- $\Phi(\vec{w}, \vec{\xi}) = \frac{1}{2} \vec{w}^T \vec{w} + C \sum_{n=1}^N \xi_n$ minimieren und

- $y_n(\vec{w}^T \vec{c}_n + b) \geq 1 - \xi_n$, für $n = 1, \dots, N$ sowie

- $\xi_n \geq 0$ erfüllen

wobei C ein Kontrollparameter ist

- Optimierung über ähnliches duales Problem wie im separablen Fall

11.5 Nicht-lineare Einbettung

- linear separable Probleme sind nicht sehr spannend (und real)
- durch nicht-lineare Einbettung (= Abbildung) in einen hoch-dimensionalen Merkmalsraum wird jedoch ein nicht-linear separables Problem mit hoher Wahrscheinlichkeit linear separabel
(vgl. Polynom-Klassifikator!)

- wir betrachten $M > N$ nicht-lineare Funktionen $\phi_n : \mathbb{R}^N \rightarrow \mathbb{R}$
Diese liefern zusammen:

$$\vec{\phi} : \mathbb{R}^N \rightarrow \mathbb{R}^M, \quad \vec{\phi}(\vec{c}) = (\phi_1(\vec{c}), \dots, \phi_M(\vec{c}))$$

- Optimierung wie bisher, aber in \mathbb{R}^M

11.5 Nicht-lineare Einbettung

- zur Klassifikation müssen wir auswerten:

$$\vec{w}^T \vec{\phi}(\vec{c}) + b = \left(\sum_{n=1}^N \alpha_n y_n \vec{\phi}^T(\vec{c}_n) \right) \vec{\phi}(\vec{c}) + b = \sum_{n=1}^N \alpha_n y_n \underbrace{\left(\vec{\phi}^T(\vec{c}_n) \vec{\phi}(\vec{c}) \right)}_{K(\vec{c}_n, \vec{c})} + b \quad (11.9)$$

(wegen $\vec{w} = \sum_{n=1}^N \alpha_n y_n \vec{c}_n$) mit dem (symmetrischen) **Kernoperator** $K(\vec{c}_n, \vec{c})$

11.5 Nicht-lineare Einbettung

- was haben wir gewonnen?

für gewisse Klassen von Kernoperator können wir die Operation im ursprünglichen Eingaberaum durchführen, und müssen nicht explizit in den hoch-dimensionalen Merkmalsraum gehen

(muß also doch nicht in \mathbb{R}^M erfolgen)

- Beispiel:

$$\begin{aligned}\vec{\phi}(c_1, c_2) &= (c_1^2, \sqrt{2} c_1 c_2, c_2^2) && \Rightarrow K(\vec{c}, \vec{c}_n) = (\vec{c}^T \vec{c}_n)^2 \\ \vec{\phi}(c_1, c_2) &= (c_1^2, c_2^2, \sqrt{2} c_1 c_2, \sqrt{2} c_1, \sqrt{2} c_2, 1) && \Rightarrow K(\vec{c}, \vec{c}_n) = (\vec{c}^T \vec{c}_n + 1)^2\end{aligned}$$

- auch das (duale) Optimierungsproblem – im linear separablen wie nicht-linear separablen Fall – kann damit im ursprünglichen Eingaberaum bearbeitet werden !

11.5 Nicht-lineare Einbettung

Wichtige Kernoperatoren

Kernoperator	Klassifikator	Bemerkung
$(\vec{c}^T \vec{c}_n + 1)^G$	Polynomklassifikator	Polynomgrad G apriori vorgegeben
$e^{-\frac{1}{2}\left(\frac{\vec{c}-\vec{c}_n}{\sigma}\right)^2}$	Radiale Basisfunktionen	gemeinsame, vorgegebene Breite σ^2
$\tanh(\beta_0 \vec{c}^T \vec{c}_n + \beta_1)$	zweilagiges Perzeptron	nur für einige Werte von β_0 und β_1

Bemerkungen

- die Dimensionalität des Merkmalsraumes wird durch die Optimierung als Anzahl der Support Vektoren vorgegeben
(nicht durch den Designer)