

Wissensentdeckung in Datenbanken

Support Vector Machine

Nico Piatkowski und Uwe Ligges

Informatik—Künstliche Intelligenz
Computergestützte Statistik
Technische Universität Dortmund

30.05.2017

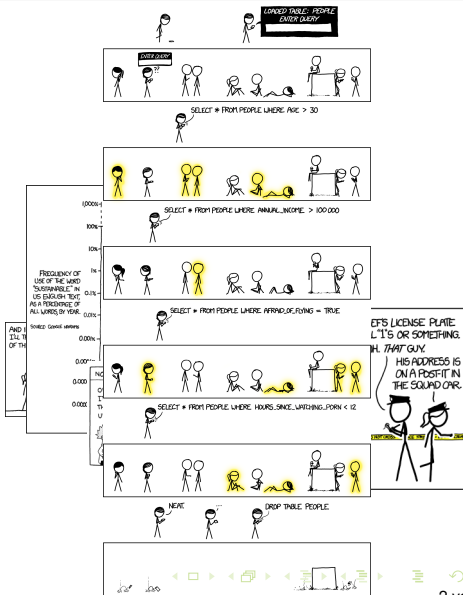
Überblick

Was bisher geschah...

- Modellklassen
- Verlustfunktionen
- Numerische Optimierung
- Regularisierung
- Überanpassung
- SQL, Häufige Mengen
- Klassifikation

Heute

- SVM



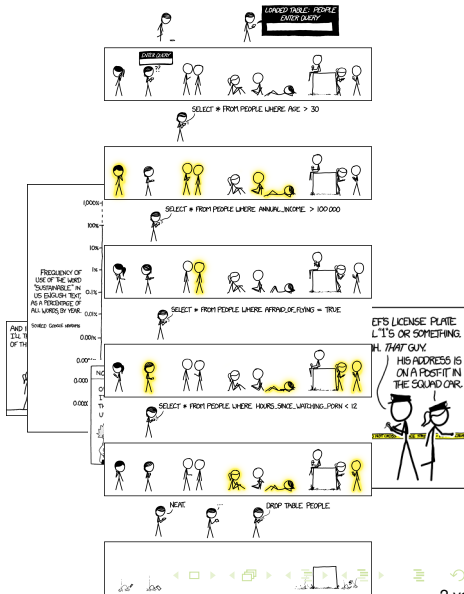
Überblick

Was bisher geschah...

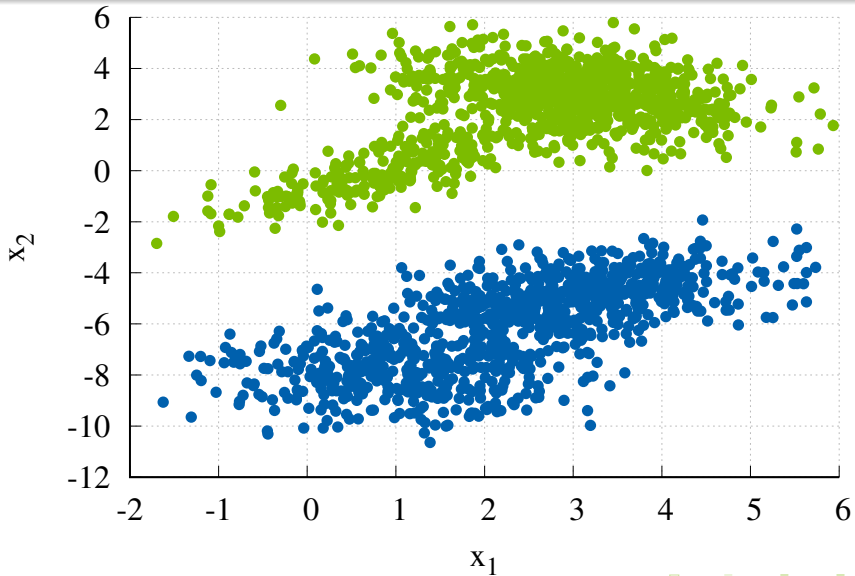
- Modellklassen
- Verlustfunktionen
- Numerische Optimierung
- Regularisierung
- Überanpassung
- SQL, Häufige Mengen
- Klassifikation

Heute

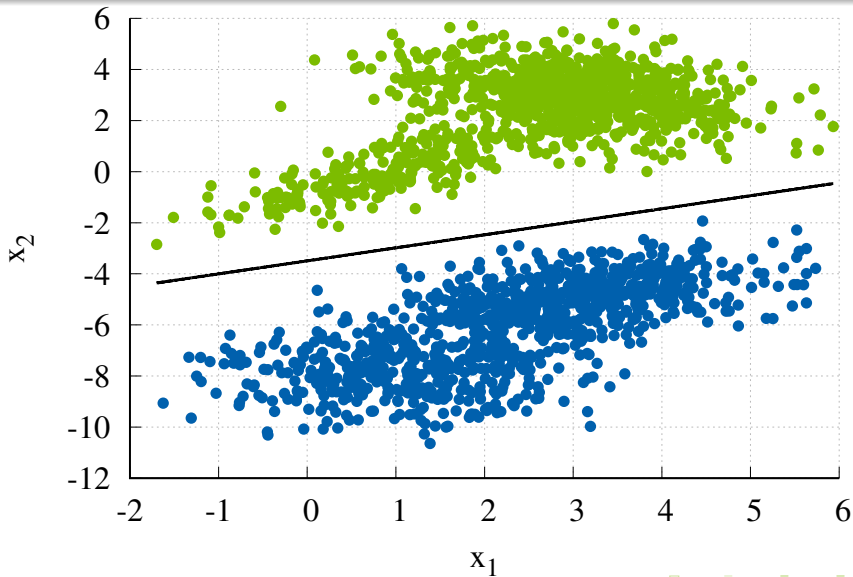
- SVM



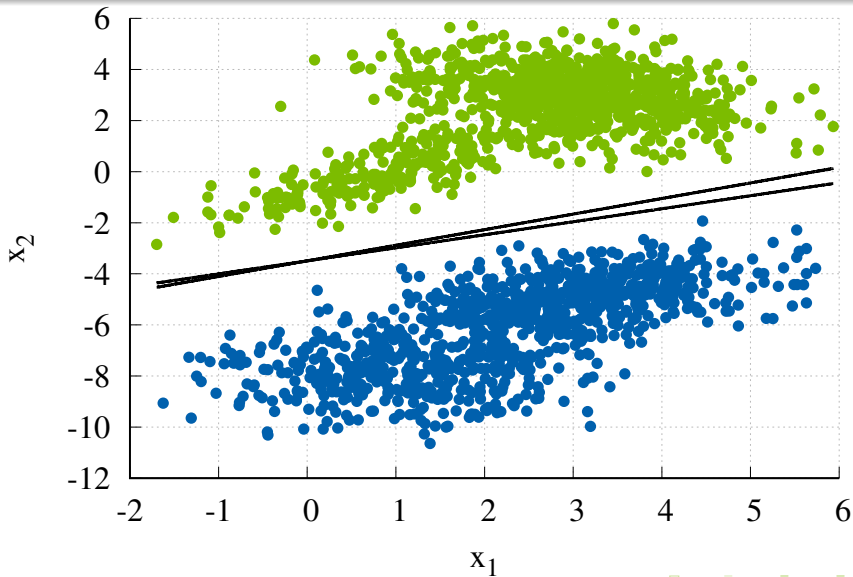
Klassifikation mit Hyperebenen



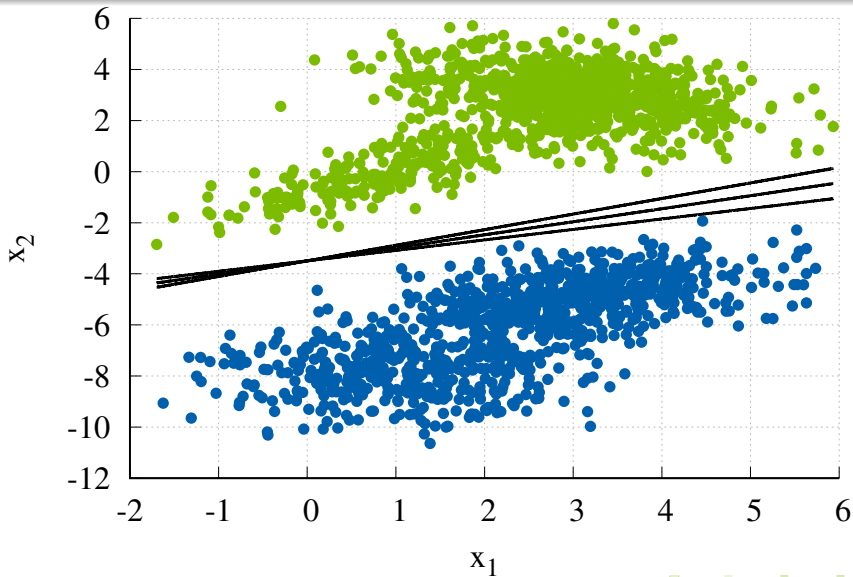
Klassifikation mit Hyperebenen



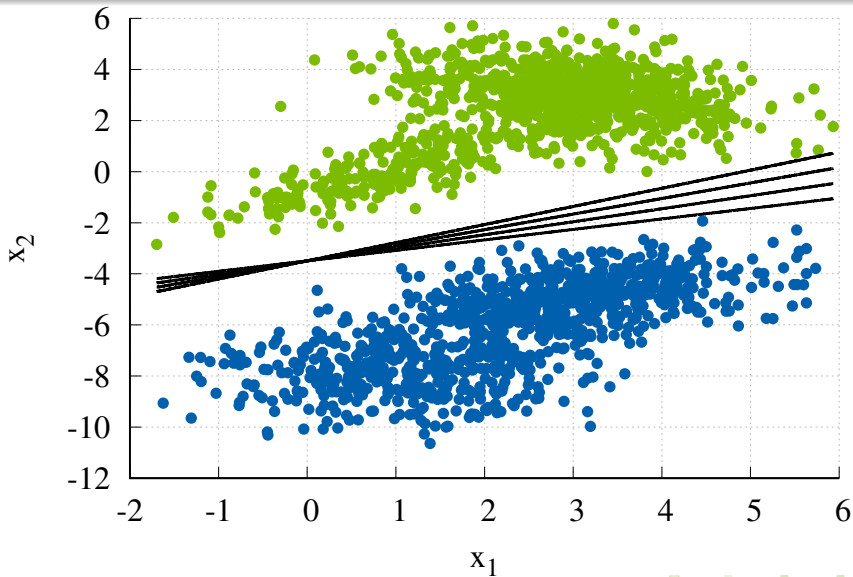
Klassifikation mit Hyperebenen



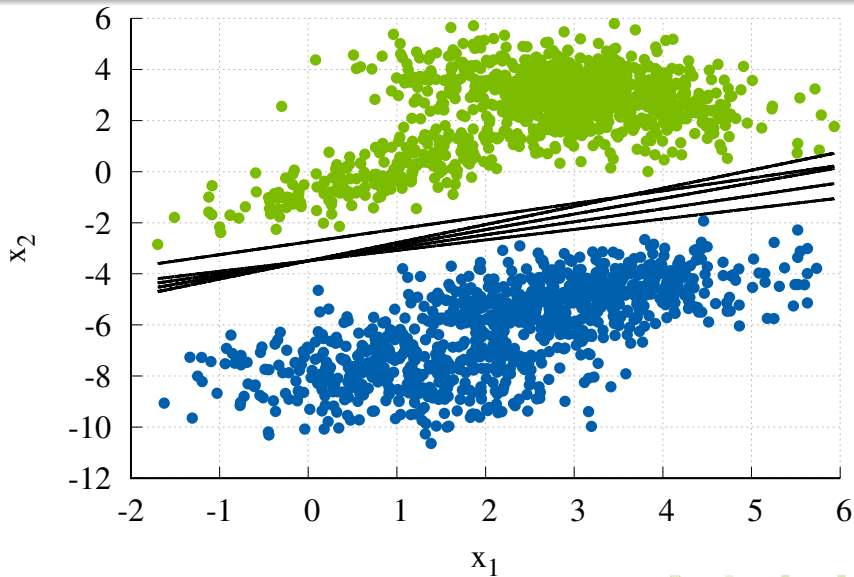
Klassifikation mit Hyperebenen



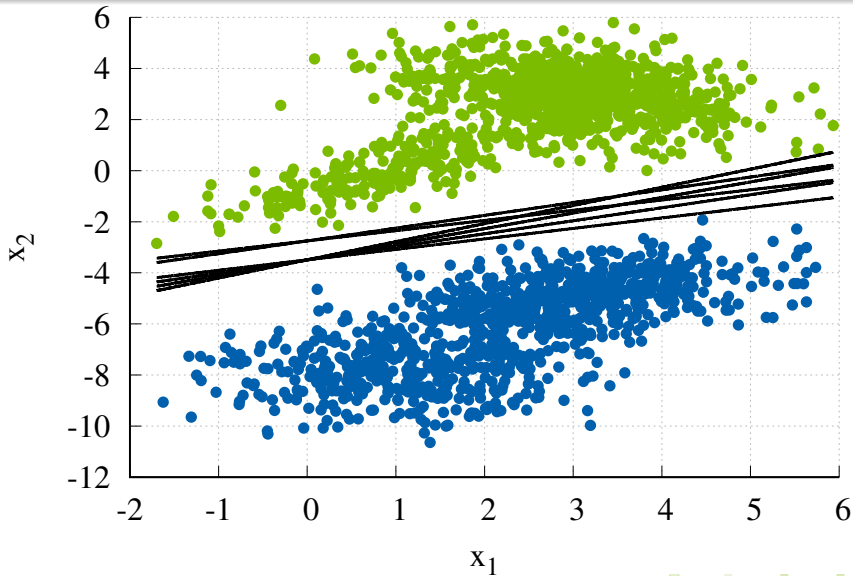
Klassifikation mit Hyperebenen



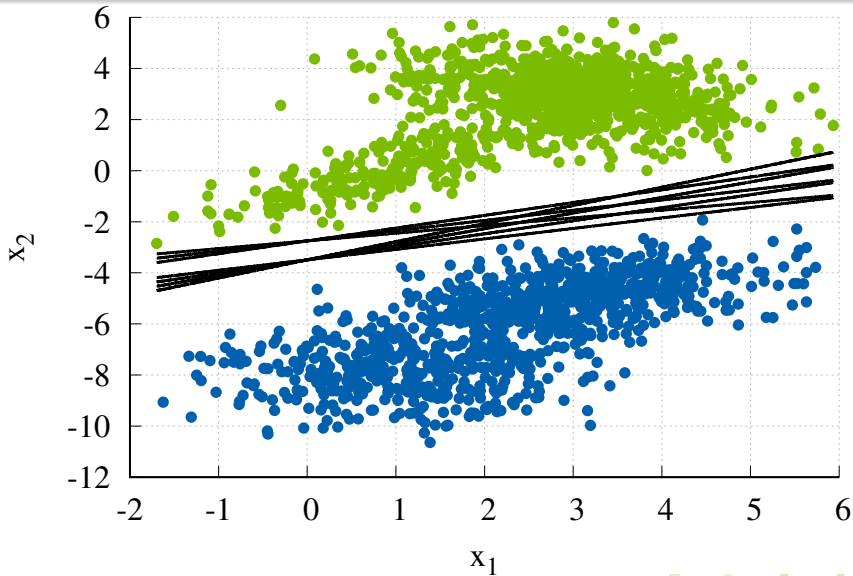
Klassifikation mit Hyperebenen



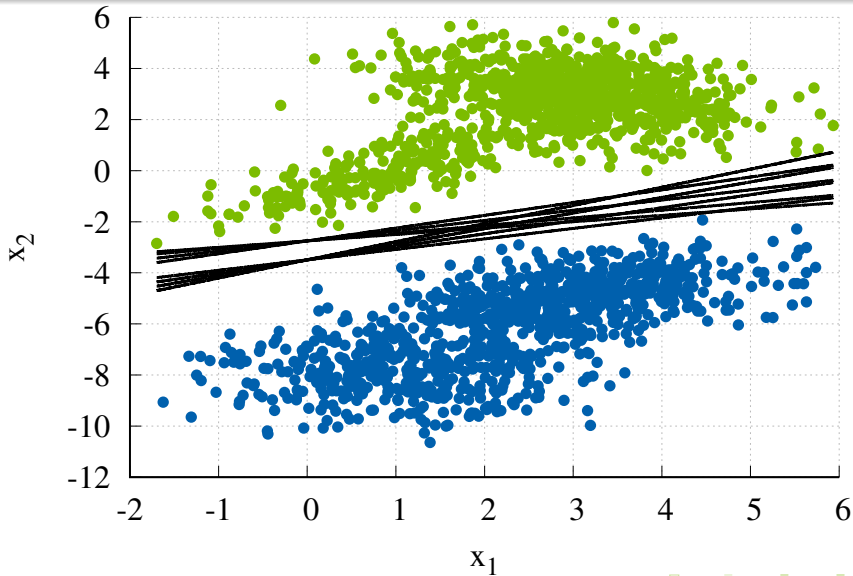
Klassifikation mit Hyperebenen



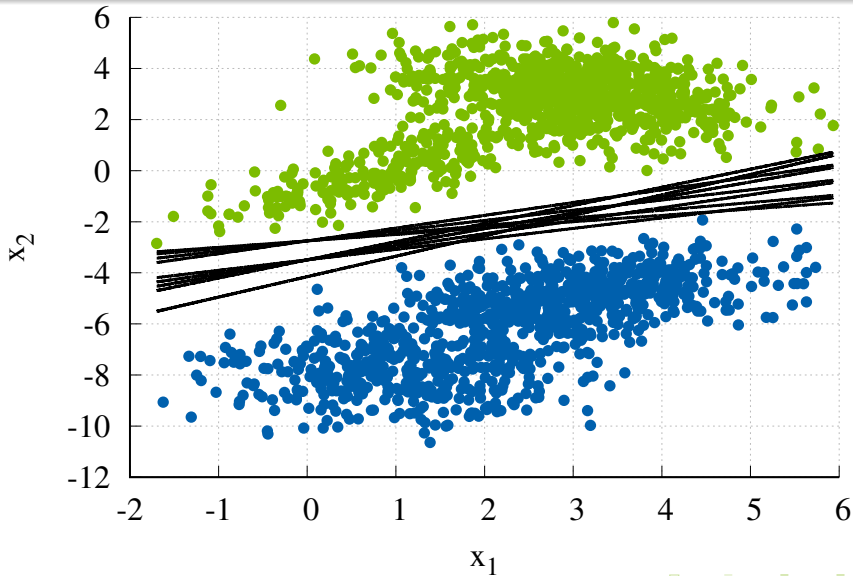
Klassifikation mit Hyperebenen



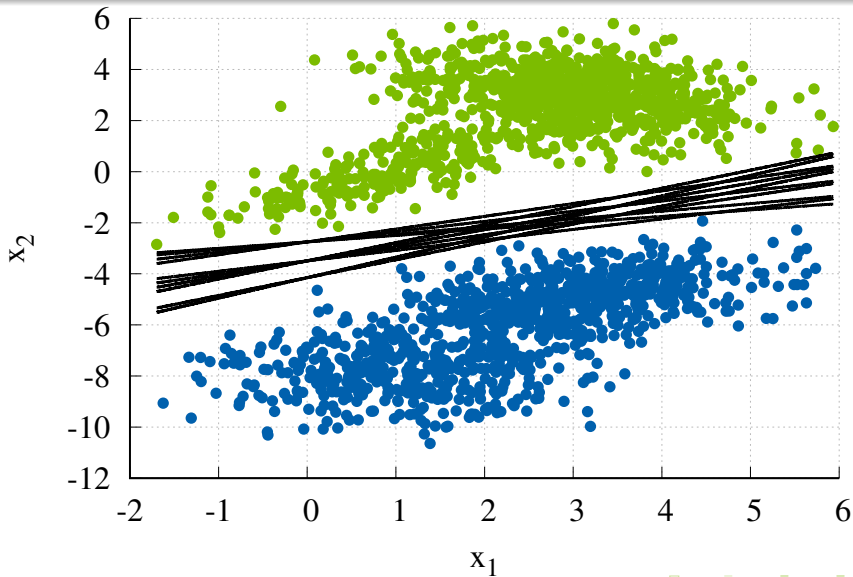
Klassifikation mit Hyperebenen



Klassifikation mit Hyperebenen



Klassifikation mit Hyperebenen



Hyperebenen

Hier: Daten \mathbf{x} aus \mathbb{R}^n , Klasse y aus $\{-1, 1\}$

- Geradengleichung: $y = f_{\beta_0, \beta}(\mathbf{x}) = \beta_0 + \langle \beta, \mathbf{x} \rangle$
- Wir wissen: $\langle \beta, \mathbf{x} \rangle = \|\beta\|_2 \|\mathbf{x}\|_2 \cos \angle(\beta, \mathbf{x})$
 $\Rightarrow \cos \angle(\beta, \mathbf{x})$ bestimmt das Vorzeichen von $\langle \beta, \mathbf{x} \rangle$
- $\text{sign}(\langle \beta, \mathbf{x} \rangle) = \text{sign}(\cos \angle(\beta, \mathbf{x}))$

$$\text{sign}(\beta_0 + \langle \beta, \mathbf{x} \rangle) = \begin{cases} +1 & \beta_0 > -\|\beta\|_2 \|\mathbf{x}\|_2 \cos \angle(\beta, \mathbf{x}) \\ -1 & \beta_0 < -\|\beta\|_2 \|\mathbf{x}\|_2 \cos \angle(\beta, \mathbf{x}) \\ 0 & \text{sonst} \end{cases}$$

- $\|\mathbf{x}\|_2$ durch Daten bestimmt, β_0 und β werden gelernt

Hyperebenen

Hier: Daten \mathbf{x} aus \mathbb{R}^n , Klasse y aus $\{-1, 1\}$

- Geradengleichung: $y = f_{\beta_0, \beta}(\mathbf{x}) = \beta_0 + \langle \beta, \mathbf{x} \rangle$
- Wir wissen: $\langle \beta, \mathbf{x} \rangle = \|\beta\|_2 \|\mathbf{x}\|_2 \cos \angle(\beta, \mathbf{x})$
 $\Rightarrow \cos \angle(\beta, \mathbf{x})$ bestimmt das Vorzeichen von $\langle \beta, \mathbf{x} \rangle$
- $\text{sign}(\langle \beta, \mathbf{x} \rangle) = \text{sign}(\cos \angle(\beta, \mathbf{x}))$

$$\text{sign}(\beta_0 + \langle \beta, \mathbf{x} \rangle) = \begin{cases} +1 & \beta_0 > -\|\beta\|_2 \|\mathbf{x}\|_2 \cos \angle(\beta, \mathbf{x}) \\ -1 & \beta_0 < -\|\beta\|_2 \|\mathbf{x}\|_2 \cos \angle(\beta, \mathbf{x}) \\ 0 & \text{sonst} \end{cases}$$

- $\|\mathbf{x}\|_2$ durch Daten bestimmt, β_0 und β werden gelernt

Hyperebenen

Hier: Daten \mathbf{x} aus \mathbb{R}^n , Klasse y aus $\{-1, 1\}$

- Geradengleichung: $y = f_{\beta_0, \beta}(\mathbf{x}) = \beta_0 + \langle \beta, \mathbf{x} \rangle$
- Wir wissen: $\langle \beta, \mathbf{x} \rangle = \|\beta\|_2 \|\mathbf{x}\|_2 \cos \angle(\beta, \mathbf{x})$
 $\Rightarrow \cos \angle(\beta, \mathbf{x})$ bestimmt das Vorzeichen von $\langle \beta, \mathbf{x} \rangle$
- $\text{sign}(\langle \beta, \mathbf{x} \rangle) = \text{sign}(\cos \angle(\beta, \mathbf{x}))$

$$\text{sign}(\beta_0 + \langle \beta, \mathbf{x} \rangle) = \begin{cases} +1 & \beta_0 > -\|\beta\|_2 \|\mathbf{x}\|_2 \cos \angle(\beta, \mathbf{x}) \\ -1 & \beta_0 < -\|\beta\|_2 \|\mathbf{x}\|_2 \cos \angle(\beta, \mathbf{x}) \\ 0 & \text{sonst} \end{cases}$$

- $\|\mathbf{x}\|_2$ durch Daten bestimmt, β_0 und β werden gelernt

Hyperebenen

Hier: Daten \mathbf{x} aus \mathbb{R}^n , Klasse y aus $\{-1, 1\}$

- Geradengleichung: $y = f_{\beta_0, \beta}(\mathbf{x}) = \beta_0 + \langle \beta, \mathbf{x} \rangle$
- Wir wissen: $\langle \beta, \mathbf{x} \rangle = \|\beta\|_2 \|\mathbf{x}\|_2 \cos \angle(\beta, \mathbf{x})$
 $\Rightarrow \cos \angle(\beta, \mathbf{x})$ bestimmt das Vorzeichen von $\langle \beta, \mathbf{x} \rangle$
- $\text{sign}(\langle \beta, \mathbf{x} \rangle) = \text{sign}(\cos \angle(\beta, \mathbf{x}))$

$$\text{sign}(\beta_0 + \langle \beta, \mathbf{x} \rangle) = \begin{cases} +1 & \beta_0 > -\|\beta\|_2 \|\mathbf{x}\|_2 \cos \angle(\beta, \mathbf{x}) \\ -1 & \beta_0 < -\|\beta\|_2 \|\mathbf{x}\|_2 \cos \angle(\beta, \mathbf{x}) \\ 0 & \text{sonst} \end{cases}$$

- $\|\mathbf{x}\|_2$ durch Daten bestimmt, β_0 und β werden gelernt

Hyperebenenbestimmung als Optimierungsproblem (I)

Datensatz $\mathcal{D} = \{(y^1, \mathbf{x}^1), (y^2, \mathbf{x}^2), \dots, (y^N, \mathbf{x}^N)\}$

Idealisiert:

$$\min_{\beta_0, \beta} \frac{1}{2} \|\beta\|_2^2$$

$$\text{s.t. } y^i = \text{sign}(\beta_0 + \langle \beta, \mathbf{x}^i \rangle) \quad \forall 1 \leq i \leq N$$

Pragmatisch:

$$\min_{\beta_0, \beta} \frac{1}{2} \|\beta\|_2^2$$

$$\text{s.t. } y^i (\beta_0 + \langle \beta, \mathbf{x}^i \rangle) \geq 1 \quad \forall 1 \leq i \leq N$$

Hyperebenenbestimmung als Optimierungsproblem (I)

Datensatz $\mathcal{D} = \{(y^1, \mathbf{x}^1), (y^2, \mathbf{x}^2), \dots, (y^N, \mathbf{x}^N)\}$

Idealisiert:

$$\min_{\beta_0, \boldsymbol{\beta}} \frac{1}{2} \|\boldsymbol{\beta}\|_2^2$$

$$\text{s.t. } y^i = \text{sign}(\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}^i \rangle) \quad \forall 1 \leq i \leq N$$

Pragmatisch:

$$\min_{\beta_0, \boldsymbol{\beta}} \frac{1}{2} \|\boldsymbol{\beta}\|_2^2$$

$$\text{s.t. } y^i (\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}^i \rangle) \geq 1 \quad \forall 1 \leq i \leq N$$

Hyperebenenbestimmung als Optimierungsproblem (I)

Datensatz $\mathcal{D} = \{(y^1, \mathbf{x}^1), (y^2, \mathbf{x}^2), \dots, (y^N, \mathbf{x}^N)\}$

Idealisiert:

$$\min_{\beta_0, \boldsymbol{\beta}} \frac{1}{2} \|\boldsymbol{\beta}\|_2^2$$

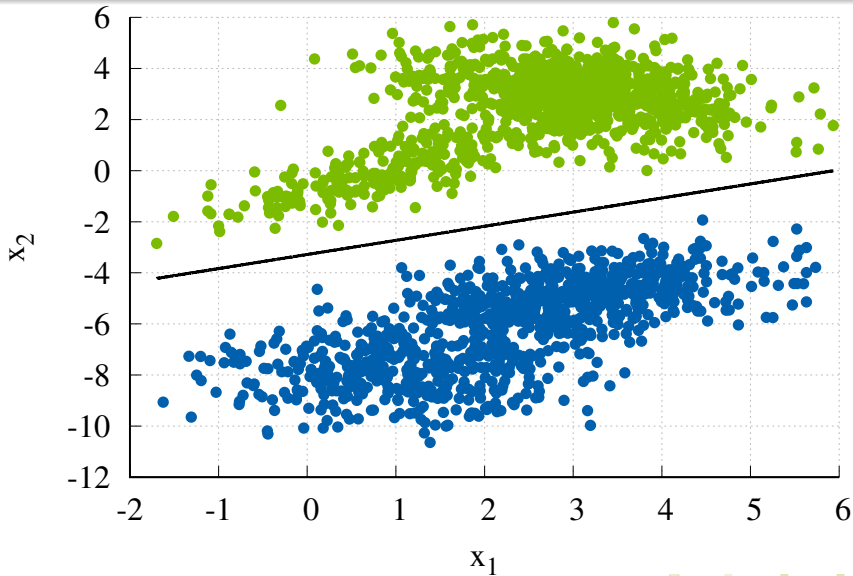
$$\text{s.t. } y^i = \text{sign}(\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}^i \rangle) \quad \forall 1 \leq i \leq N$$

Pragmatisch:

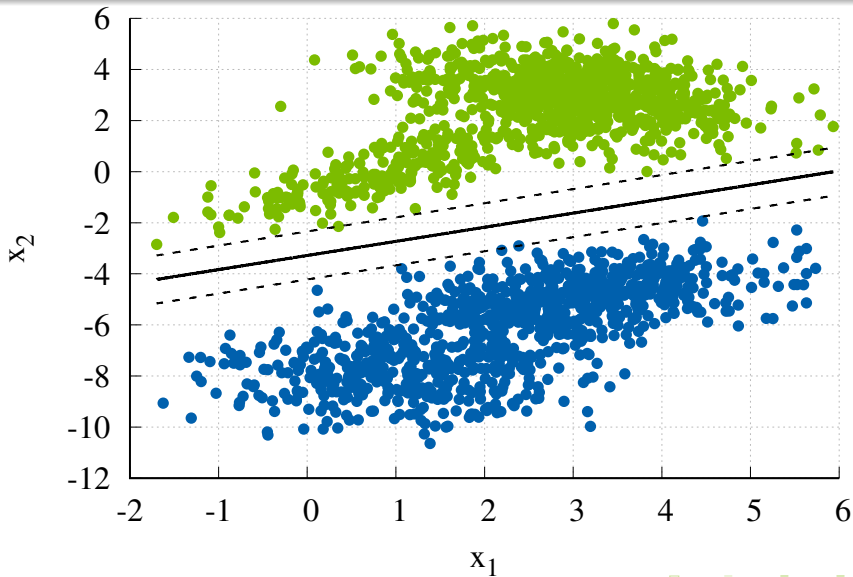
$$\min_{\beta_0, \boldsymbol{\beta}} \frac{1}{2} \|\boldsymbol{\beta}\|_2^2$$

$$\text{s.t. } y^i (\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}^i \rangle) \geq 1 \quad \forall 1 \leq i \leq N$$

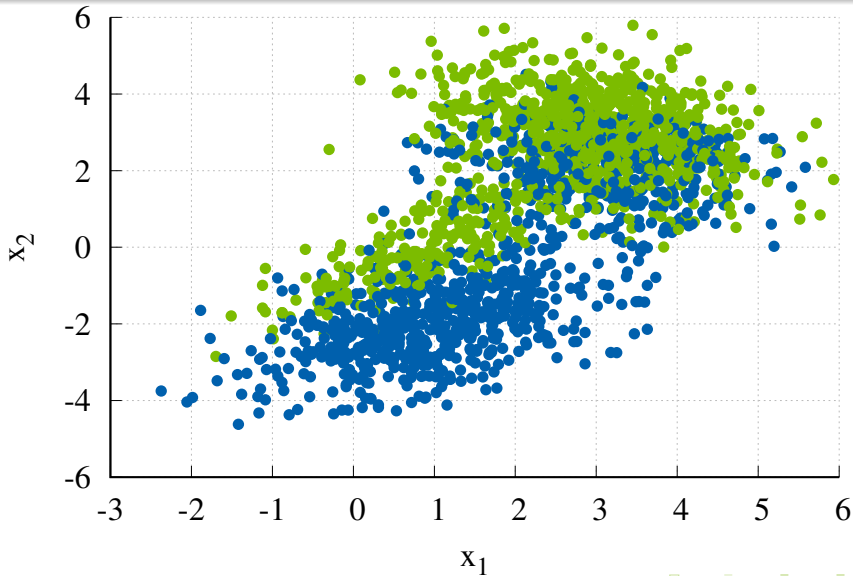
Klassifikation mit Hyperebenen



Klassifikation mit Hyperebenen



Klassifikation mit Hyperebenen (nicht trennbar)





Klassifikation als Optimierungsproblem (II)

$$\min_{\beta_0, \beta, \zeta} \frac{1}{2} \|\beta\|_2^2 + C \sum_{i=1}^N \zeta_i$$

$$\text{s.t. } y^i(\beta_0 + \langle \beta, x^i \rangle) \geq 1 - \zeta_i, \quad \forall 1 \leq i \leq N$$

$$\zeta_i \geq 0, \quad \forall 1 \leq i \leq N$$

Also suchen wir die kleinsten ζ_i mit

$$\zeta_i \geq 1 - y^i(\beta_0 + \langle \beta, x^i \rangle) \quad \text{und} \quad \zeta_i \geq 0$$

für jedes $1 \leq i \leq N$ und damit

$$\zeta_i^* = \max\{0, 1 - y^i(\beta_0 + \langle \beta, x^i \rangle)\}$$

Klassifikation als Optimierungsproblem (II)

$$\min_{\beta_0, \boldsymbol{\beta}, \zeta} \frac{1}{2} \|\boldsymbol{\beta}\|_2^2 + C \sum_{i=1}^N \zeta_i$$

$$\text{s.t. } y^i (\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}^i \rangle) \geq 1 - \zeta_i, \quad \forall 1 \leq i \leq N$$

$$\zeta_i \geq 0, \quad \forall 1 \leq i \leq N$$

Also suchen wir die kleinsten ζ_i mit

$$\zeta_i \geq 1 - y^i (\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}^i \rangle) \quad \text{und} \quad \zeta_i \geq 0$$

für jedes $1 \leq i \leq N$ und damit

$$\zeta_i^* = \max\{0, 1 - y^i (\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}^i \rangle)\}$$

Klassifikation als Optimierungsproblem (III)

$$\zeta_i^* = \max\{0, 1 - y^i(\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}^i \rangle)\}$$

Einsetzen in Zielfunktion:

$$\min_{\beta_0, \boldsymbol{\beta}} \frac{1}{2} \|\boldsymbol{\beta}\|_2^2 + C \sum_{i=1}^N \max\{0, 1 - y^i(\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}^i \rangle)\}$$

Minimierung der l_2 -regularisierten Hinge-Verlustfunktion

$$\text{Hinge}(\beta_0, \boldsymbol{\beta}; \mathcal{D}) + \lambda R_2(\boldsymbol{\beta}) = \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, y) \in \mathcal{D}} \max\{0, 1 - y f_{\beta_0, \boldsymbol{\beta}}(\mathbf{x})\} + \lambda \|\boldsymbol{\beta}\|_2^2$$

Klassifikation als Optimierungsproblem (III)

$$\zeta_i^* = \max\{0, 1 - y^i(\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}^i \rangle)\}$$

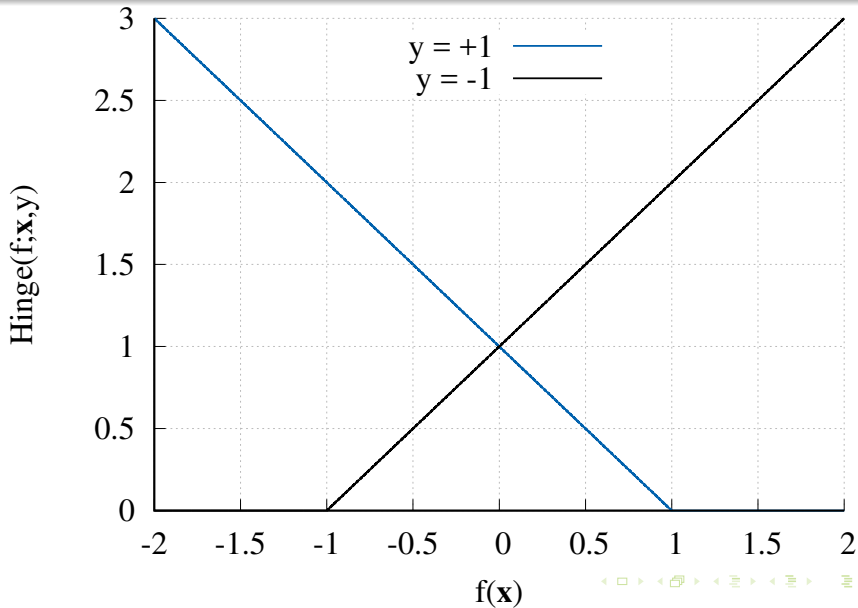
Einsetzen in Zielfunktion:

$$\min_{\beta_0, \boldsymbol{\beta}} \frac{1}{2} \|\boldsymbol{\beta}\|_2^2 + C \sum_{i=1}^N \max\{0, 1 - y^i(\beta_0 + \langle \boldsymbol{\beta}, \mathbf{x}^i \rangle)\}$$

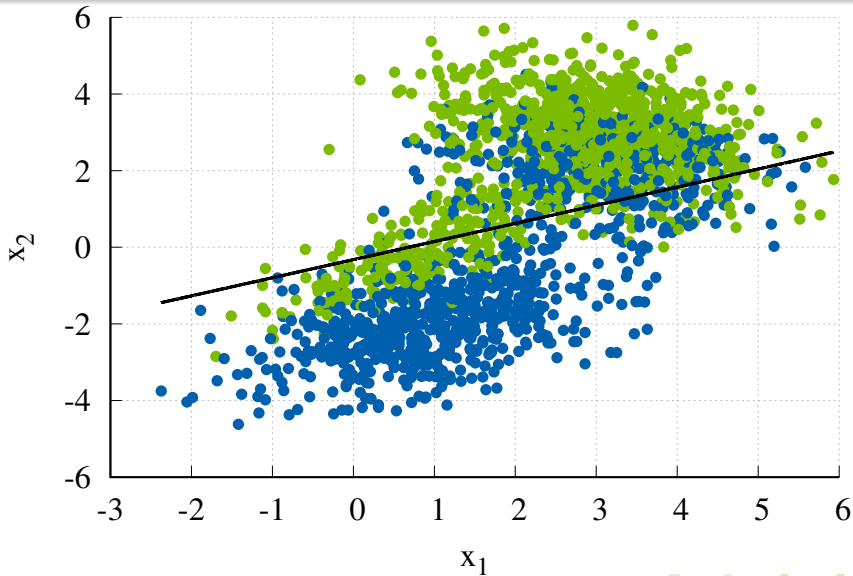
Minimierung der l_2 -regularisierten Hinge-Verlustfunktion

$$\text{Hinge}(\beta_0, \boldsymbol{\beta}; \mathcal{D}) + \lambda R_2(\boldsymbol{\beta}) = \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, y) \in \mathcal{D}} \max\{0, 1 - y f_{\beta_0, \boldsymbol{\beta}}(\mathbf{x})\} + \lambda \|\boldsymbol{\beta}\|_2^2$$

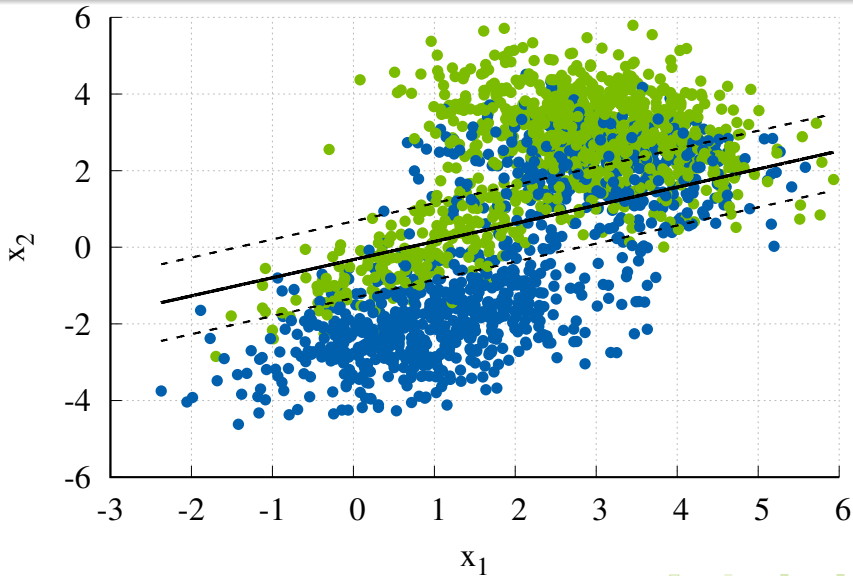
Hinge Verlustfunktion



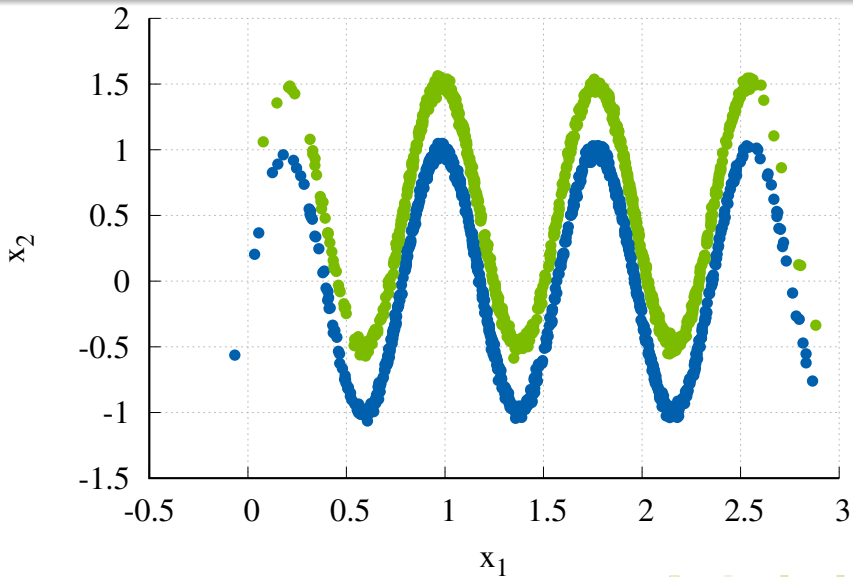
Klassifikation mit Hyperebenen (nicht trennbar)



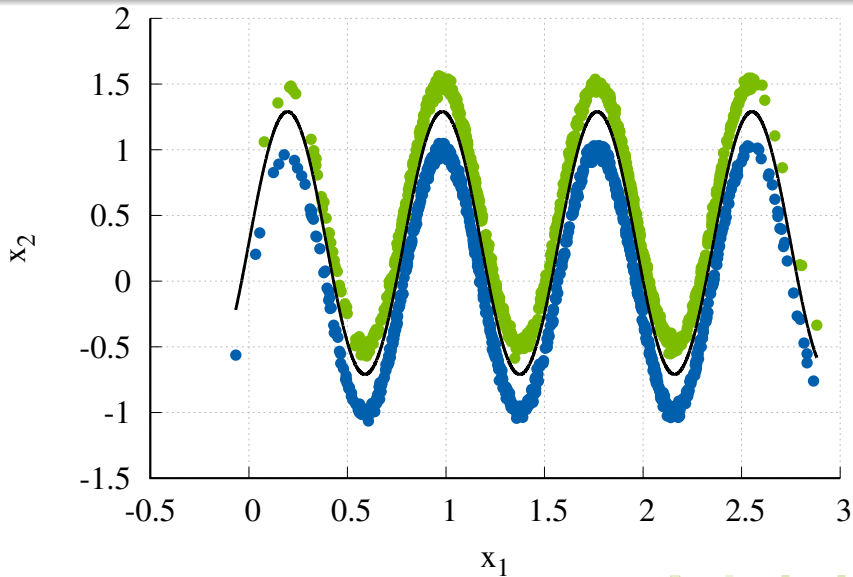
Klassifikation mit Hyperebenen (nicht trennbar)



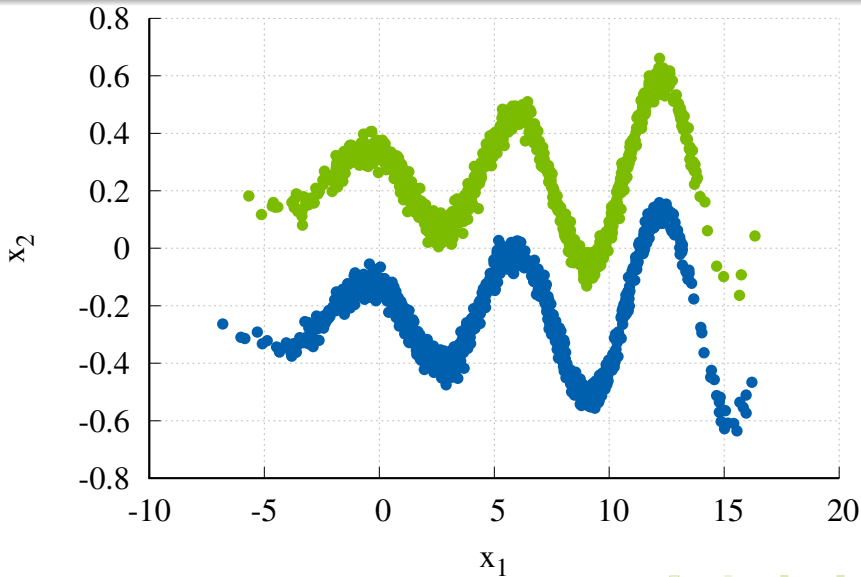
Nicht-linear trennbare Daten



Nicht-linear trennbare Daten



Nicht-linear trennbare Daten



Nicht-lineare Trennbarkeit durch Merkmalstransformation

- Trennung nicht-linearer Daten durch Transformation
 $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^d$
- ϕ kann jede beliebige Funktion sein
- Bild von ϕ kann höher- oder niedrig-dimensionaler als \mathcal{X} sein
- ϕ soll möglichst einfach zu berechnen sein (!!)
- Daten sollen im Bildraum von ϕ möglichst linear trennbar sein

Wo kommt ϕ her??

Oft: "Abstand" im neuen Raum einfacher zu berechnen als ϕ

Nicht-lineare Trennbarkeit durch Merkmalstransformation

- Trennung nicht-linearer Daten durch Transformation
 $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^d$
- ϕ kann jede beliebige Funktion sein
- Bild von ϕ kann höher- oder niedrig-dimensionaler als \mathcal{X} sein
- ϕ soll möglichst einfach zu berechnen sein (!!)
- Daten sollen im Bildraum von ϕ möglichst linear trennbar sein

Wo kommt ϕ her??

Oft: "Abstand" im neuen Raum einfacher zu berechnen als ϕ

Optimierung mit Nebenbedingungen

Lagrange-Dualität:

$$\min_{\beta} f(\beta)$$

$$\text{s.t. } g(\beta) \geq c$$

≡

$$\min_{\beta} \max_{\alpha} f(\beta) - \alpha(g(\beta) - c)$$

$$\text{s.t. } \alpha \geq 0$$

Optimierung mit Nebenbedingungen

Lagrange-Dualität:

$$\min_{\beta} f(\beta)$$

$$\text{s.t. } g(\beta) \geq c$$

≡

$$\min_{\beta} \max_{\alpha} f(\beta) - \alpha(g(\beta) - c)$$

$$\text{s.t. } \alpha \geq 0$$



Duales Problem und Stützvektoren

Lagrange-Dual der SVM:

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^N \alpha_i - \frac{1}{2} \alpha_i \alpha_j \sum_{i=1}^N \sum_{j=1}^N y^i y^j \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle \\ \text{s.t.} \quad & \sum_{i=1}^N \alpha_i y^i = 0 \\ & 0 \leq \alpha_i \leq C, \forall 1 \leq i \leq N \end{aligned}$$

Neues Problem hat “nur” N Parameter (ist das gut?)

Original Gewichte: $\beta = \sum_{i=1}^N \alpha_i y^i \mathbf{x}^i$

Datenpunkte mit $\alpha_i > 0$ heißen Stützvektor (Support-Vector)



Duales Problem und Stützvektoren

Lagrange-Dual der SVM:

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^N \alpha_i - \frac{1}{2} \alpha_i \alpha_j \sum_{i=1}^N \sum_{j=1}^N y^i y^j \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle \\ \text{s.t.} \quad & \sum_{i=1}^N \alpha_i y^i = 0 \\ & 0 \leq \alpha_i \leq C, \forall 1 \leq i \leq N \end{aligned}$$

Neues Problem hat “nur” N Parameter (ist das gut?)

Original Gewichte: $\beta = \sum_{i=1}^N \alpha_i y^i \mathbf{x}^i$

Datenpunkte mit $\alpha_i > 0$ heißen Stützvektor (Support-Vector)



Duales Problem und Stützvektoren

Lagrange-Dual der SVM:

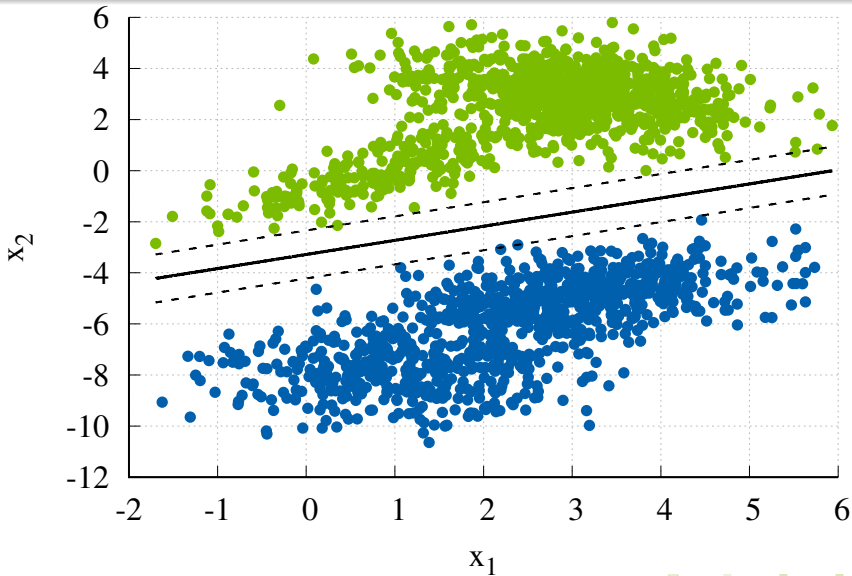
$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^N \alpha_i - \frac{1}{2} \alpha_i \alpha_j \sum_{i=1}^N \sum_{j=1}^N y^i y^j \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle \\ \text{s.t.} \quad & \sum_{i=1}^N \alpha_i y^i = 0 \\ & 0 \leq \alpha_i \leq C, \forall 1 \leq i \leq N \end{aligned}$$

Neues Problem hat “nur” N Parameter (ist das gut?)

Original Gewichte: $\beta = \sum_{i=1}^N \alpha_i y^i \mathbf{x}^i$

Datenpunkte mit $\alpha_i > 0$ heißen Stützvektor (Support-Vector)

Klassifikation mit Hyperebenen



Kernel “Trick”

Anstatt $\phi(\mathbf{x})$ für jeden \mathbf{x} auszurechnen,

- Finde Funktion K , mit $K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$
- K heißt *Kernel* oder *Kernfunktion*
- $\phi(\mathbf{x})$ (und $\phi(\mathbf{x}')$) muss nie explizit berechnet werden!
- ϕ muss nicht bekannt sein!!
- $K_{\gamma, \delta, \dots}(\mathbf{x}, \mathbf{x}')$ kann Parameter haben

Kernel “Trick”

Anstatt $\phi(\mathbf{x})$ für jeden \mathbf{x} auszurechnen,

- Finde Funktion K , mit $K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$
- K heißt *Kernel* oder *Kernfunktion*
- $\phi(\mathbf{x})$ (und $\phi(\mathbf{x}')$) muss nie explizit berechnet werden!
- ϕ muss nicht bekannt sein!!
- $K_{\gamma, \delta, \dots}(\mathbf{x}, \mathbf{x}')$ kann Parameter haben

Kernfunktionen

$$K_{\text{Linear}}(\mathbf{x}, \mathbf{x}') = \langle \mathbf{x}, \mathbf{x}' \rangle$$

$$K_{\text{Poly},k,c}(\mathbf{x}, \mathbf{x}') = (\langle \mathbf{x}, \mathbf{x}' \rangle + c)^k$$

$$K_{\text{Gauss},\gamma}(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{\gamma} \|\mathbf{x} - \mathbf{x}'\|_2^2\right)$$

- Graph Kernel (common subgraphs, Bäume, Pfade)
- String Kernel (gemeinsame Zeichenketten, gemeinsame Parsebäume)