



Ostbayerische Technische Hochschule  
Amberg-Weiden

# Machine Learning

Prof. Dr. Fabian Brunner

<fa.brunner@oth-aw.de>

Amberg, 29. November 2021

## Thema heute: Support Vector Machines

- Grundidee
- Herleitung der Problemformulierung
- Vergleich mit Logistischer Regression
- Transformation in höherdimensionale Räume

## Problemstellung

- Binäres Klassifikationsproblem
- $p$  numerische Features („unabhängige Variablen“)
- Binäre Zielvariable („unabhängige Variable“) mit den Klassen  $-1$  und  $+1$ .
- $m$  Trainingsdatensätze

$$(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)}) ,$$

wobei  $\mathbf{x}^{(i)} \in \mathbb{R}^p$  und  $y^{(i)} \in \{-1, +1\}$ .

## Modell-Training

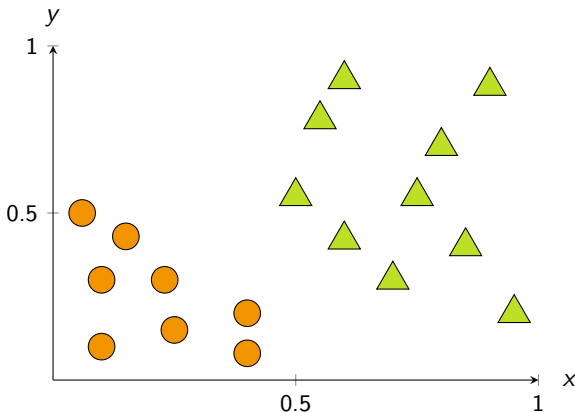
Bestimme anhand der Trainingsdaten eine Entscheidungsregel)

$$f : \mathbb{R}^p \rightarrow \{-1, +1\} .$$

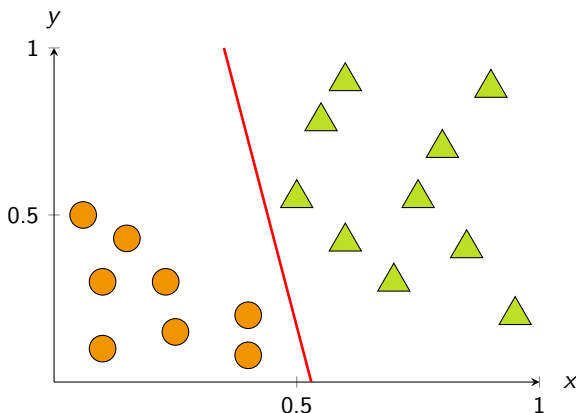
## Modellanwendung

Für einen Query Point  $\mathbf{x}_q$  prognostiziere das Label  $f(\mathbf{x}_q)$ .

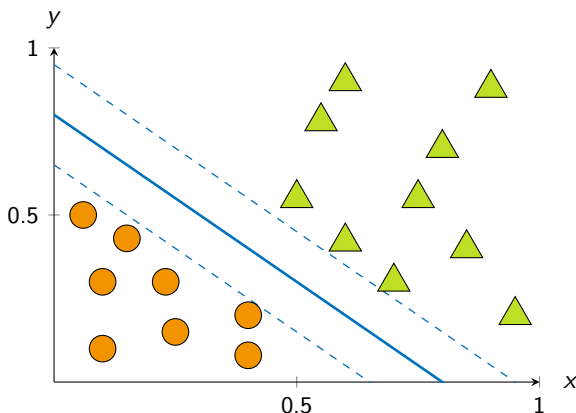
- Ziel bei SVMs ist die Bestimmung einer Decision Boundary mit möglichst großem Randbereich („margin“), in dem keine Samples liegen. Modelle mit solchen Entscheidungsgrenzen weisen tendenziell geringere Fehler beim Verallgemeinern auf unbekannte Daten auf.
- Man spricht bei SVM daher von einem „**large margin classifier**“.



- Ziel bei SVMs ist die Bestimmung eines Decision Boundary mit möglichst großem Randbereich („margin“), in dem keine Samples liegen. Modelle mit solchen Entscheidungsgrenzen weisen tendenziell geringere Fehler beim Verallgemeinern auf unbekannte Daten auf.
- Man spricht bei SVM daher von einem „**large margin classifier**“.



- Ziel bei SVMs ist die Bestimmung einer Decision Boundary mit möglichst großem Randbereich („margin“), in dem keine Samples liegen. Modelle mit solchen Entscheidungsgrenzen weisen tendenziell geringere Fehler beim Verallgemeinern auf unbekannte Daten auf.
- Man spricht bei SVM daher von einem „**large margin classifier**“.



# Das Euklidische Skalarprodukt

Gegeben seien zwei Vektoren

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}.$$

- Dann berechnet sich das Skalarprodukt von  $\mathbf{x}$  und  $\mathbf{y}$  durch

$$\mathbf{x} \cdot \mathbf{y} = \mathbf{x}^T \mathbf{y} = x_1 y_1 + x_2 y_2.$$

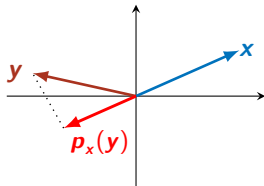
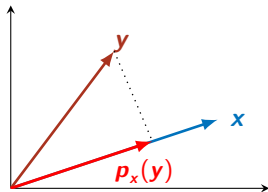
- Die Länge des Ortsvektors von  $\mathbf{0}$  zu  $\mathbf{x}$  lautet

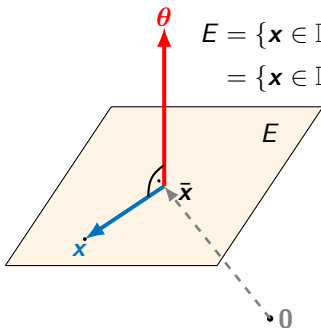
$$\|\mathbf{x}\| = \sqrt{\mathbf{x} \cdot \mathbf{x}} = \sqrt{x_1^2 + x_2^2}.$$

- Mit Hilfe des Skalarprodukts  $p := \mathbf{x} \cdot \mathbf{y}$  lautet die Projektion von  $\mathbf{y}$  auf  $\mathbf{x}$ :

$$\mathbf{p}_x(\mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\|^2} \cdot \mathbf{x} \Rightarrow |\mathbf{x} \cdot \mathbf{y}| = \|\mathbf{x}\| \cdot \|\mathbf{p}_x(\mathbf{y})\|.$$

- Ist  $\mathbf{x}$  normiert, gibt das Skalarprodukt  $\mathbf{x} \cdot \mathbf{y}$  also die (signierte) Länge der Projektion von  $\mathbf{y}$  auf  $\mathbf{x}$  an.





$$\begin{aligned} E &= \{x \in \mathbb{R}^3 : (x - \bar{x}) \cdot \theta = 0\} \\ &= \{x \in \mathbb{R}^3 : x \cdot \theta = \underbrace{\bar{x} \cdot \theta}_{=: b}\} \end{aligned}$$

Verständnisfragen:

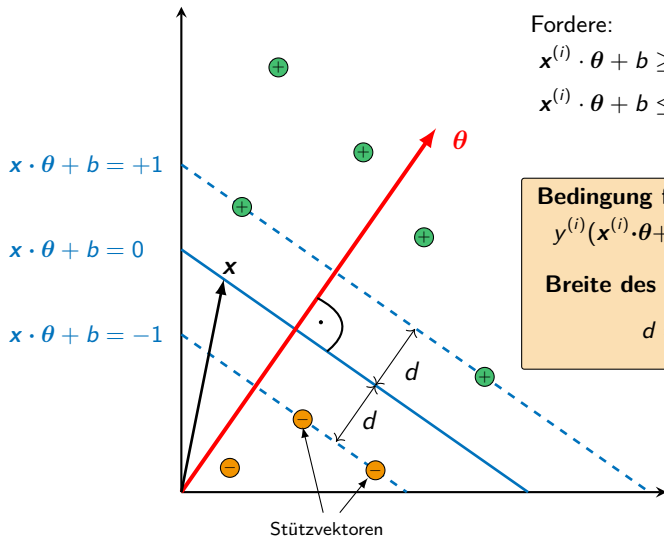
- Welchen Wert hat  $b$ , wenn die Ebene durch den Nullpunkt verläuft?
- Welche relative Lage haben zwei Ebenen mit den Darstellungen

$$E_1 = \{x \in \mathbb{R}^3 : x \cdot \theta = b_1\},$$

$$E_2 = \{x \in \mathbb{R}^3 : x \cdot \theta = b_2\} ?$$

- Wie berechnet man den Abstand von  $E$  zum Nullpunkt?





Fordere:

$$\mathbf{x}^{(i)} \cdot \boldsymbol{\theta} + b \geq 1 \quad \text{für } y^{(i)} = +1$$

$$\mathbf{x}^{(i)} \cdot \boldsymbol{\theta} + b \leq -1 \quad \text{für } y^{(i)} = -1$$

**Bedingung für Margin:**

$$y^{(i)}(\mathbf{x}^{(i)} \cdot \boldsymbol{\theta} + b) \geq 1 \quad \text{für alle } \mathbf{x}^{(i)}, y^{(i)}$$

**Breite des Margins:**

$$d = \frac{1}{\|\boldsymbol{\theta}\|}$$

# SVM - linear separierbarer Fall

Um einen möglichst breiten Margin zu bekommen, ist also  $\frac{1}{\|\theta\|}$  unter den Nebenbedingungen  $y^{(i)}(\mathbf{x}^{(i)} \cdot \theta - b) \geq 1$  zu maximieren. Alternativ kann auch  $\|\theta\|^2$  minimiert werden. Es ergibt sich ein restringiertes Optimierungsproblem:

## Optimierungsproblem im Fall linear separierbarer Samples

$$\min_{\theta \in \mathbb{R}^p, b \in \mathbb{R}} \frac{1}{2} \|\theta\|^2 \quad (\text{P})$$

$$\text{u.d.N. } y^{(i)}(\theta \cdot \mathbf{x}^{(i)} + b) \geq 1, \quad i = 1, \dots, m$$

- Es handelt sich um ein (konvexes) quadratisches Problem mit linearen Ungleichungs-Nebenbedingungen.
- Es ist nur dann lösbar, wenn die Daten linear trennbar sind. Andernfalls können die Nebenbedingungen nicht erfüllt werden.
- Ausweg (s. später): modifiziere das Problem, damit auch Fehler (z.B. Punkte innerhalb des Margins oder auf der falschen Seite) toleriert werden können. Dies wird durch Einführung von Schlupfvariablen erreicht.

- Anstelle des oben hergeleiteten Optimierungsproblems wird beim Modelltraining üblicherweise das dazu **duale Problem** gelöst, da es eine vorteilhafte Struktur im Hinblick auf die Verwendung nichtlinearer Kernels (s. später) aufweist.
- Seien  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  und  $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $i = 1, \dots, l$  stetig differenzierbare konvexe Funktionen. Betrachte das folgende Optimierungsproblem mit Nebenbedingungen:

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \\ \text{u.d.N. } g_i(\mathbf{x}) = 0 \text{ für } i = 1, \dots, l. \end{aligned}$$

- Das Lagrange-Funktional lautet

$$\mathcal{L}(\mathbf{x}, \lambda_1, \dots, \lambda_l) := f(\mathbf{x}) + \sum_{i=1}^l \lambda_i g_i(\mathbf{x})$$

- Das Optimierungsproblem ist äquivalent zum Problem

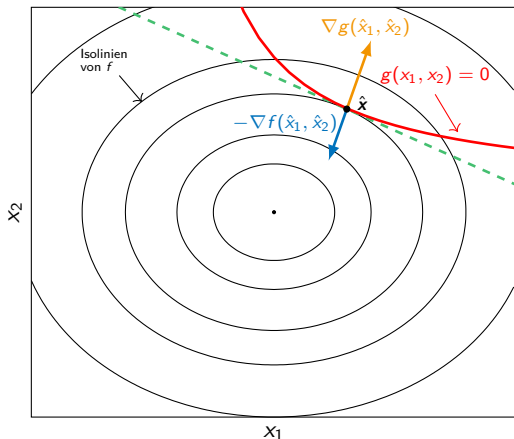
$$\min_{\mathbf{x} \in \mathbb{R}^n} \max_{\boldsymbol{\lambda} \in \mathbb{R}^l} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) .$$

## Methode der Lagrange-Multiplikatoren

Ist  $\hat{\mathbf{x}} \in \mathbb{R}^n$  eine Lösung des obigen (konvexen) Optimierungsproblems und sind bestimmte Regularitätsvoraussetzungen (z.B. Slater-Bedingung) erfüllt, dann gibt es Zahlen  $\lambda_1, \dots, \lambda_l \in \mathbb{R}$  (**Lagrange-Multiplikatoren**), sodass

$$\begin{aligned}\nabla f(\mathbf{x}) &= \sum_{i=1}^l \lambda_i \nabla g_i(\mathbf{x}) , \\ g_i(\mathbf{x}) &= 0 \text{ für } i = 1, \dots, l .\end{aligned}$$

- Die obigen Bedingungen sind gleichbedeutend damit, dass  $\nabla \mathcal{L}(\mathbf{x}, \lambda_1, \dots, \lambda_l) = \mathbf{0}$  gilt.
- Es handelt sich allgemein nur um notwendige Bedingungen. Sind sie in einem Punkt erfüllt, dann heißt er **kritischer Punkt** und ist damit ein „Kandidat“ für ein Minimum.



Im Minimum  $\hat{x}$  gilt:

$$\nabla f(\hat{x}_1, \hat{x}_2) = \lambda \nabla g(\hat{x}_1, \hat{x}_2)$$

für ein  $\lambda \in \mathbb{R}$ , d.h. die beiden Vektoren sind kollinear.

Wir erweitern unser Problem nun um Ungleichungsnebenbedingungen:

$$\begin{array}{ll} \min_{\mathbf{x} \in \mathbb{R}^n} & f(\mathbf{x}) \\ \text{u.d.N.} & g_i(\mathbf{x}) = 0 \text{ für } i = 1, \dots, l, \\ & h_i(\mathbf{x}) \leq 0 \text{ für } i = 1, \dots, k. \end{array}$$

Das Lagrange-Funktional dazu lautet:

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) := f(\mathbf{x}) + \sum_{i=1}^l \lambda_i g_i(\mathbf{x}) + \sum_{i=1}^k \alpha_i h_i(\mathbf{x}) .$$

## Primales Problem

Mit Hilfe des Lagrange-Funktional kann man das obige Optimierungsproblem in folgendes äquivalentes Optimierungsproblem umformulieren:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \max_{\boldsymbol{\lambda}, \boldsymbol{\alpha}; \boldsymbol{\alpha} \geq 0} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\alpha}) .$$

Durch Vertauschung von max und min erhält man das sog. **duale Problem**

## Duales Problem

$$\max_{\lambda, \alpha; \alpha \geq 0} \min_{\mathbf{x} \in \mathbb{R}^n} \mathcal{L}(\mathbf{x}, \lambda, \alpha)$$

- Es gilt stets  $d^* \leq p^*$ , wobei  $d^*$  die Lösung des dualen und  $p^*$  die Lösung des primalen Problems bezeichnen.
- Sei  $d^*$  die Lösung des dualen Problems und es gelte  $d^* = \mathcal{L}(\mathbf{x}^*, \lambda^*, \alpha^*)$ . Dann kann  $h_i(\mathbf{x}^*) < 0$  nur gelten, falls  $\alpha_i = 0$ .
- Unter bestimmten Voraussetzungen (die bei SVM vorliegen!), gilt auch  $d^* = p^*$ . Statt des primalen Problems kann man also auch das duale Problem lösen.

Ist  $f$  konvex und sind die Funktionen  $h_i$  und  $f_i$  linear, so gilt  $d^* = p^*$  und die folgenden KKT-Bedingungen sind notwendige und hinreichende Bedingungen, damit ein Punkt  $\hat{\mathbf{x}}$  ein Optimum ist:

## KKT-Bedingungen

$$\begin{aligned}\frac{\partial \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\alpha})}{\partial x_i} \Big|_{\mathbf{x}=\hat{\mathbf{x}}} &= 0 , \\ \alpha_i h_i(\hat{\mathbf{x}}) &= 0 , \\ \alpha_i &\geq 0 , \\ h_i(\hat{\mathbf{x}}) &\leq 0 , \\ g_i(\hat{\mathbf{x}}) &= 0 ,\end{aligned}$$



Durch Anwenden der KKT-Bedingungen erhält man das zu (P) duale Problem:

## Duales Problem

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^m} \quad & \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)} \\ \text{u.d.N.} \quad & \alpha_i \geq 0 \quad \text{für } i = 1, \dots, m, \\ & \sum_{i=1}^m \alpha_i y^{(i)} = 0. \end{aligned}$$

- Es handelt sich um ein (konvexes) quadratisches Optimierungsproblem mit Nebenbedingungen.
- Man kann zeigen, dass es eindeutig lösbar ist, sofern die Daten linear separierbar sind.
- Für die Lösung solcher Probleme existieren numerische Löser.

- Für die Entscheidungsfunktion kann man zeigen:

$$\theta \cdot \mathbf{x}_q + b = \sum_{i=1}^m \alpha_i y^{(i)} \mathbf{x}^{(i)} \cdot \mathbf{x}_q + b ,$$

d.h. die zum dualen Problem gehörende Entscheidungsregel lautet:

## Entscheidungsregel bei Lösung des dualen Problems

$$\sum_{i=1}^m \alpha_i y^{(i)} \mathbf{x}^{(i)} \cdot \mathbf{x}_q + b \geq 0 \quad \text{Zuordnung zur Klasse } +1 ,$$

$$\sum_{i=1}^m \alpha_i y^{(i)} \mathbf{x}^{(i)} \cdot \mathbf{x}_q + b < 0 \quad \text{Zuordnung zur Klasse } -1 .$$

- Die Entscheidungsregel hängt, im Gegensatz zum Primalen Problem, von den Trainingsdatenpunkten ab.
- Die resultierenden Gewichte  $\alpha_i$  sind allerdings nur für die Stützvektoren ungleich Null, d.h. obige Summe enthält nur sehr wenige Summanden.

## Minimierungsproblem

$$\max_{\alpha \in \mathbb{R}^m} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}$$

+ Nebenbedingungen

## Entscheidungsregel

$$\sum_{i=1}^m \alpha_i y^{(i)} \mathbf{x}^{(i)} \cdot \mathbf{x}_q + b \geq 0 \quad \rightarrow \text{Zuordnung zur Klasse } +1 ,$$

$$\sum_{i=1}^m \alpha_i y^{(i)} \mathbf{x}^{(i)} \cdot \mathbf{x}_q + b < 0 \quad \rightarrow \text{Zuordnung zur Klasse } -1 .$$

# Struktur des dualen Problems

## Minimierungsproblem

$$\max_{\alpha \in \mathbb{R}^m} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \mathbf{x}^{(i)} \cdot \mathbf{x}^{(j)}$$

+ Nebenbedingungen

## Entscheidungsregel

$$\sum_{i=1}^m \alpha_i y^{(i)} \mathbf{x}^{(i)} \cdot \mathbf{x}_q + b \geq 0 \quad \rightarrow \text{Zuordnung zur Klasse } +1 ,$$

$$\sum_{i=1}^m \alpha_i y^{(i)} \mathbf{x}^{(i)} \cdot \mathbf{x}_q + b < 0 \quad \rightarrow \text{Zuordnung zur Klasse } -1 .$$

## Struktur des dualen Problems

- Das Minimierungsproblem hängt nur von Skalarprodukten zwischen den Feature-Vektoren ab.
- Die Entscheidungsregel hängt nur von Skalarprodukten der Feature-Vektoren mit dem Query Point ab.

Häufig sind die Daten nicht exakt linear separierbar. Dann besitzt das obige Problem keine Lösung. Durch Einführung von **Schlupfvariablen** wird das Verfahren robust gegen Trainings-Fehler.

## Optimierungsproblem im Fall nicht linear separierbarer Samples

$$\min_{\theta \in \mathbb{R}^p, \xi \in \mathbb{R}^m, b \in \mathbb{R}} \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^m \xi_i$$
$$\text{u.d.N. } y^{(i)}(\theta \cdot x^{(i)} + b) \geq 1 - \xi_i, \quad \xi_i \geq 0 \quad i = 1, \dots, m.$$

- Das obige Problem ist ein quadratisches Optimierungsproblem mit Ungleichungsnebenbedingungen.
- Jede Nebenbedingung kann erfüllt werden, wenn  $\xi_i$  hinreichend groß ist.
- Ist  $\xi_i > 1$ , so liegt  $x^{(i)}$  auf der falschen Seite des Decision Boundary.
- Ist  $0 \leq \xi_i \leq 1$ , so liegt  $x^{(i)}$  im Margin (bzw. auf den Ebenen, die diesen abgrenzen)

Bisher hatten wir SVM als restringiertes Optimierungsproblem formuliert:

$$\min_{\theta \in \mathbb{R}^p, \xi \in \mathbb{R}^m, b \in \mathbb{R}} \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^m \xi_i$$

unter den Nebenbed.  $y^{(i)}(\theta \cdot \mathbf{x}^{(i)} + b) \geq 1 - \xi_i$ ,  $\xi_i \geq 0$ ,  $i = 1, \dots, m$ .

## Umformulierung der Ungleichungs-Nebenbedingungen

Wir erweitern nun  $\mathbf{x}^{(i)}$  und  $\theta$  um eine Komponente, indem wir  $x_0^{(i)} := 1$  und  $\theta_0 := b$  setzen und definieren  $f_\theta(\mathbf{x}) := \theta \cdot \mathbf{x}$ . Damit lauten die **Ungleichungs-Nebenbedingungen**

$$y^{(i)} f_\theta(\mathbf{x}^{(i)}) \geq 1 - \xi_i \quad \wedge \quad \xi_i \geq 0.$$

Diese lassen sich wie folgt als äquivalente **Gleichungen** schreiben:

$$\xi_i = \max(0, 1 - y^{(i)} f_\theta(\mathbf{x}^{(i)})).$$

Durch die Einführung der Maximum-Funktion ist ein Minimierungsproblem ohne

## SVM als unrestringiertes Optimierungsproblem

$$\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^{p+1}} J(\theta) := \frac{1}{2} \sum_{i=1}^p \theta_i^2 + C \sum_{i=1}^m \max(0, 1 - y^{(i)} f_{\theta}(\mathbf{x}^{(i)})) .$$

Nachdem die Parameter  $\theta \in \mathbb{R}^{p+1}$  durch Lösen dieses Optimierungsproblems bestimmt worden sind, wird die Klassenzuordnung eines neuen Datenpunkts  $\mathbf{x}_q$  wie folgt definiert:

## Klassenzuordnung bei SVM

Falls  $\theta \cdot \mathbf{x}_q \geq 0 \rightarrow$  Zuordnung zur Klasse  $y = +1$

Falls  $\theta \cdot \mathbf{x}_q < 0 \rightarrow$  Zuordnung zur Klasse  $y = -1$

Der Parameter  $C$  in der zu minimierenden Zielfunktion

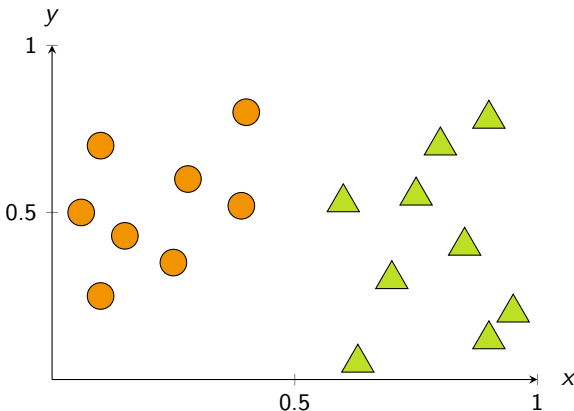
$$\frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^m \xi_i$$

kann als Regularisierungsparameter aufgefasst werden:

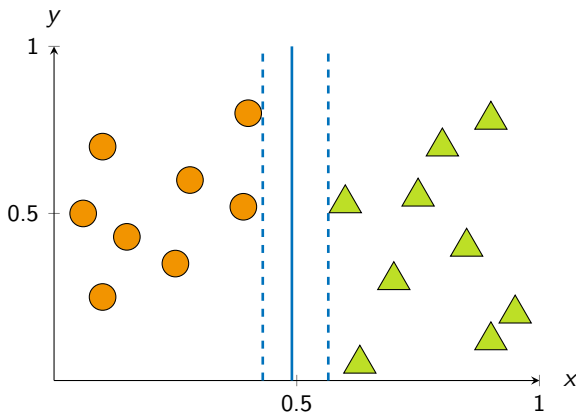
- Ist  $C$  klein, so erhalten die Nebenbedingungen weniger Gewicht und es werden Trainingsfehler zugelassen („soft margin“). Im Gegenzug wird versucht, die Parameter so zu bestimmen, dass der Margin möglichst breit wird.
- Ist  $C$  groß, werden die Parameter so bestimmt, dass die Nebenbedingungen möglichst strikt eingehalten werden. Dafür wird der Margin tendenziell schmaler.
- $C \rightarrow \infty$  erzwingt alle Nebenbedingungen („hard margin“)
- $C$  klein: geringere Varianz, höhere Verzerrung
- $C$  groß: höhere Varianz, geringere Verzerrung
- $C$  ist ein Hyperparameter und muss, z.B. mit Hilfe von Kreuzvalidierung, festgelegt werden.



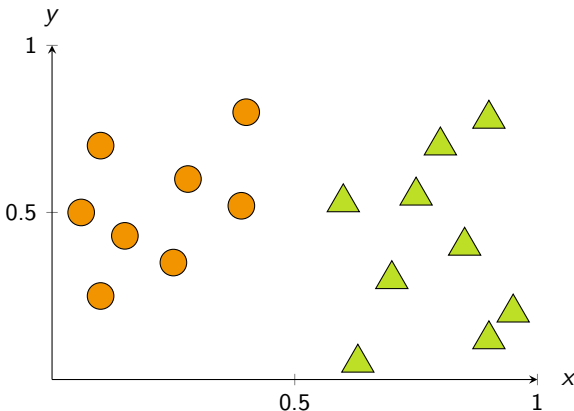
- Je größer  $C$  ist, desto stärker wird das Modell an die Daten angepasst, d.h. die Gerade wird so bestimmt, dass möglichst viele Datenpunkte korrekt klassifiziert werden. Ausreißer haben dann ggf. großen Einfluss.
- Für kleineres  $C$  ist das Verfahren „robuster“ gegenüber Ausreißern und einzelnen Fehlern (entspricht stärkerer Regularisierung).



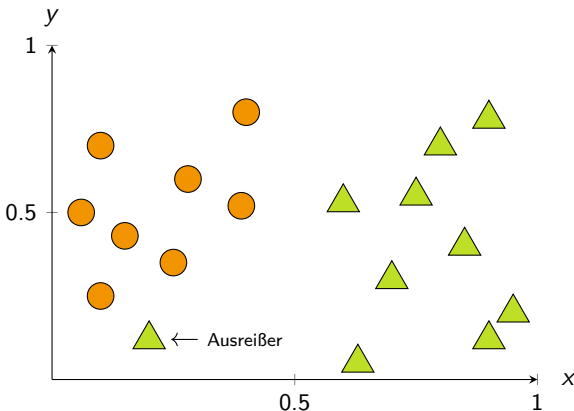
- Je größer  $C$  ist, desto stärker wird das Modell an die Daten angepasst, d.h. die Gerade wird so bestimmt, dass möglichst viele Datenpunkte korrekt klassifiziert werden. Ausreißer haben dann ggf. großen Einfluss.
- Für kleineres  $C$  ist das Verfahren „robuster“ gegenüber Ausreißern und einzelnen Fehlern (entspricht stärkerer Regularisierung).



- Je größer  $C$  ist, desto stärker wird das Modell an die Daten angepasst, d.h. die Gerade wird so bestimmt, dass möglichst viele Datenpunkte korrekt klassifiziert werden. Ausreißer haben dann ggf. großen Einfluss.
- Für kleineres  $C$  ist das Verfahren „robuster“ gegenüber Ausreißern und einzelnen Fehlern (entspricht stärkerer Regularisierung).

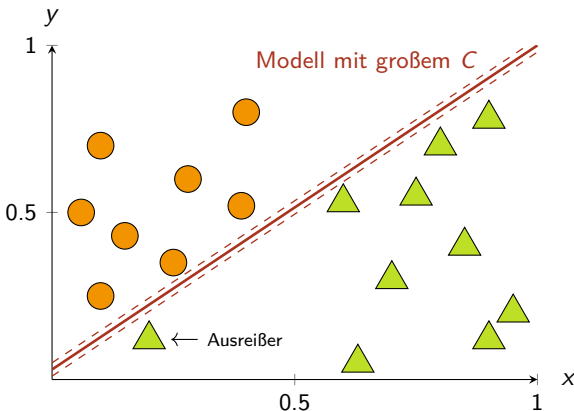


- Je größer  $C$  ist, desto stärker wird das Modell an die Daten angepasst, d.h. die Gerade wird so bestimmt, dass möglichst viele Datenpunkte korrekt klassifiziert werden. Ausreißer haben dann ggf. großen Einfluss.
- Für kleineres  $C$  ist das Verfahren „robuster“ gegenüber Ausreißern und einzelnen Fehlern (entspricht stärkerer Regularisierung).



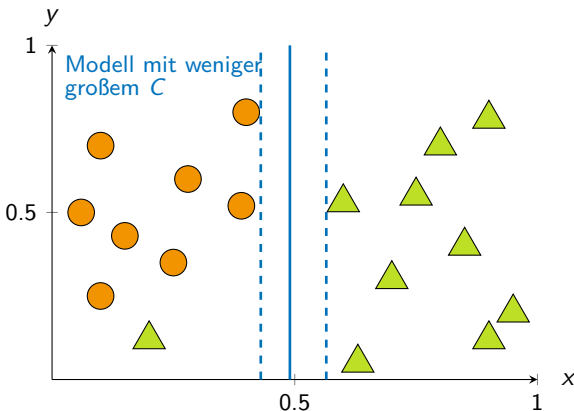
# Large Margin classification bei Ausreißern

- Je größer  $C$  ist, desto stärker wird das Modell an die Daten angepasst, d.h. die Gerade wird so bestimmt, dass möglichst viele Datenpunkte korrekt klassifiziert werden. Ausreißer haben dann ggf. großen Einfluss.
- Für kleineres  $C$  ist das Verfahren „robuster“ gegenüber Ausreißern und einzelnen Fehlern (entspricht stärkerer Regularisierung).



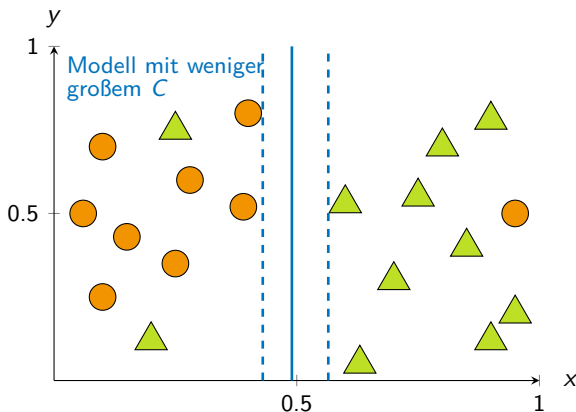
# Large Margin classification bei Ausreißern

- Je größer  $C$  ist, desto stärker wird das Modell an die Daten angepasst, d.h. die Gerade wird so bestimmt, dass möglichst viele Datenpunkte korrekt klassifiziert werden. Ausreißer haben dann ggf. großen Einfluss.
- Für kleineres  $C$  ist das Verfahren „robuster“ gegenüber Ausreißern und einzelnen Fehlern (entspricht stärkerer Regularisierung).



# Large Margin classification bei Ausreißern

- Je größer  $C$  ist, desto stärker wird das Modell an die Daten angepasst, d.h. die Gerade wird so bestimmt, dass möglichst viele Datenpunkte korrekt klassifiziert werden. Ausreißer haben dann ggf. großen Einfluss.
- Für kleineres  $C$  ist das Verfahren „robuster“ gegenüber Ausreißern und einzelnen Fehlern (entspricht stärkerer Regularisierung).



## Ansatz bei der Logistischen Regression

Bei der Logistischen Regression hat die Modellfunktion die folgende Gestalt:

$$f_{\theta}(\mathbf{x}) = g(\theta^T \mathbf{x}) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_p x_p) .$$

wobei

$$g(x) = \frac{1}{1 + e^{-x}} .$$

Die Parameter werden durch Minimierung der Cross Entropy bestimmt:

## Cross Entropy als Kostenfunktional bei der Logistischen Regression

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \left( f_{\theta}(\mathbf{x}^{(i)}) \right) + (1 - y^{(i)}) \log \left( 1 - f_{\theta}(\mathbf{x}^{(i)}) \right)$$



## Regularisierte Logistische Regression

$$J_{LR}(\theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(f_{\theta}(\mathbf{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\theta}(\mathbf{x}^{(i)})) + \frac{\lambda}{2m} \sum_{i=1}^p \theta_i^2.$$

## SVM

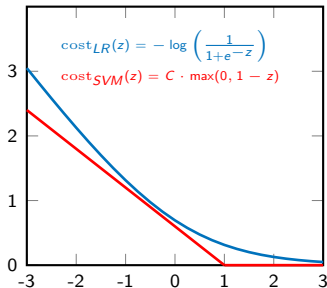
$$J_{SVM}(\theta) = \frac{1}{2} \sum_{i=1}^p \theta_i^2 + C \sum_{i=1}^m \max(0, 1 - y^{(i)} f_{\theta}(\mathbf{x}^{(i)}))$$

- Der Term zur Beschreibung der Breite des Margins entspricht einem Regularisierungsterm (Varianzreduzierung).
- Klassifikationsfehler auf dem Trainingsdatensatz werden mit Kosten versehen (Reduzierung der Verzerrung)
- Worin unterscheiden sich die Terme, die die Klassifikationsfehler bestrafen?

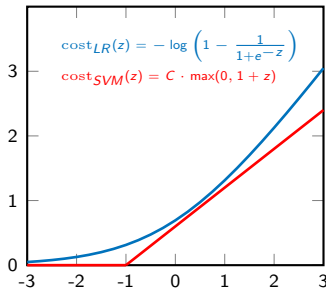
## Kosten eines Samples bei der logistischen Regression

$$\begin{aligned} & -y \log(f_{\theta}(\mathbf{x})) - (1 - y) \log(1 - f_{\theta}(\mathbf{x})) \\ &= -y \log\left(\frac{1}{1 + e^{-\theta^T \mathbf{x}}}\right) - (1 - y) \log\left(1 - \frac{1}{1 + e^{-\theta^T \mathbf{x}}}\right) \end{aligned}$$

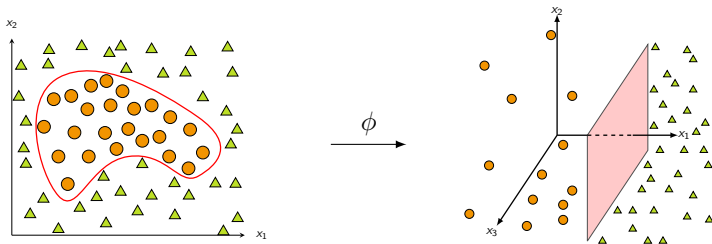
$y = 1$



$y = 0$  bzw.  $y = -1$



Idee: Nichtlineare Transformation des Problems



- Nichtlineare Transformation:  $\phi : \mathbb{R}^p \rightarrow \mathbb{R}^P$  mit  $P \geq p$ .
- Fitte ein Modell der Form  $f_{\theta}(\mathbf{x}) = \theta \cdot \phi(\mathbf{x}) + b$ .
- Die transformierten Daten  $\phi(\mathbf{x}^{(1)}), \dots, \phi(\mathbf{x}^{(m)})$  sind ggf. linear trennbar.
- Das transformierte Problem ist weiterhin linear in  $\theta$ .
- Beispiel: Übergang von der univariaten linearen Regression zur polynomialen Regression:

$$\phi : \mathbb{R} \rightarrow \mathbb{R}^{d+1}, \quad x \mapsto \phi(x) = (1, x, x^2, x^3, \dots, x^d)^T.$$

## Minimierungsproblem

$$\max_{\alpha \in \mathbb{R}^m} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \phi(\mathbf{x}^{(i)}) \cdot \phi(\mathbf{x}^{(j)})$$

+ Nebenbedingungen

## Entscheidungsregel

$$\sum_{i=1}^m \alpha_i y^{(i)} \phi(\mathbf{x}^{(i)}) \cdot \phi(\mathbf{x}_q) + b \geq 0 \quad \rightarrow \text{Zuordnung zur Klasse } +1$$

## Minimierungsproblem

$$\max_{\alpha \in \mathbb{R}^m} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \phi(\mathbf{x}^{(i)}) \cdot \phi(\mathbf{x}^{(j)})$$

+ Nebenbedingungen

## Entscheidungsregel

$$\sum_{i=1}^m \alpha_i y^{(i)} \phi(\mathbf{x}^{(i)}) \cdot \phi(\mathbf{x}_q) + b \geq 0 \quad \rightarrow \text{Zuordnung zur Klasse } +1$$

## Kernel Trick

- Definiere den **Kernel**  $K(\mathbf{x}, \mathbf{z}) := \phi(\mathbf{x}) \cdot \phi(\mathbf{z})$ .
- Wenn man  $K$  kennt, kann der SVM-Klassifikator trainiert werden, ohne dass  $\phi$  explizit berechnet werden muss. Dies ist häufig effizienter.
- Es muss kein nichtlineares Problem gelöst werden. Die mathematischen Eigenschaften (Konvexität) des SVM-Klassifikators bleiben erhalten.

Betrachte die Transformation

$$\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3, \quad \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \end{pmatrix}.$$

Dann gilt

$$\begin{aligned} \phi(\mathbf{x}) \cdot \phi(\mathbf{z}) &= (x_1^2, x_2^2, \sqrt{2}x_1x_2) \cdot \begin{pmatrix} z_1^2 \\ z_2^2 \\ \sqrt{2}z_1z_2 \end{pmatrix} \\ &= x_1^2z_1^2 + x_2^2z_2^2 + 2x_1x_2z_1z_2 \\ &= (x_1z_1 + x_2z_2)^2 \\ &= (\mathbf{x} \cdot \mathbf{z})^2. \end{aligned}$$

Die direkte Berechnung des Kernels  $K(\mathbf{x}, \mathbf{z}) = (\mathbf{x} \cdot \mathbf{z})^2$  ist effizienter als die Transformation und anschließende Berechnung des Skalarprodukts.

## Linearer Kernel

$$K(\mathbf{x}, \mathbf{z}) = \mathbf{x} \cdot \mathbf{z}$$

Entspricht der linearen SVM ohne Verwendung eines Kernels.

## Polynomialer Kernel

$$K(\mathbf{x}, \mathbf{z}) = (1 + \mathbf{x} \cdot \mathbf{z})^d \text{ für } d > 0 .$$

Der polynimiale Kernel enthält alle polynomialen Terme bis zum Grad  $d$ .

## Gauß'scher Kernel

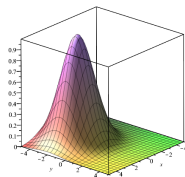
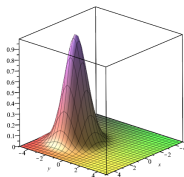
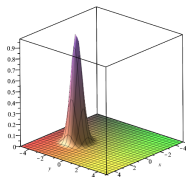
$$K(\mathbf{x}, \mathbf{z}) = \exp \left( -\frac{\|\mathbf{x} - \mathbf{z}\|^2}{2\sigma^2} \right) \text{ für } \sigma > 0 .$$

Der Gauß'sche Kernel führt auf einen unendlichdimensionalen Feature Space.

Entscheidungsfunktion:

$$f(\mathbf{x}) = \sum_{i=1}^m \alpha_i y^{(i)} \underbrace{K(\mathbf{x}^{(i)}, \mathbf{x})}_{=: f_i(\mathbf{x})} + b .$$

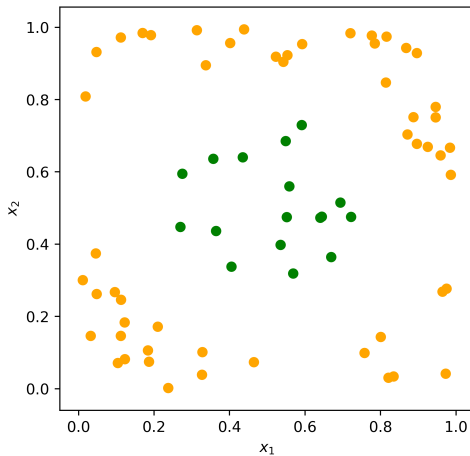
- Die Funktion  $f_i(\mathbf{x}) = K(\mathbf{x}^{(i)}, \mathbf{x})$  ist eine radiale Basisfunktion (RBF) um den Punkt  $\mathbf{x}^{(i)}$ . Man erhält also für jedes Item des Trainingsdatensatzes eine Basisfunktion ( $m$  sollte nicht zu groß sein!)
- Für  $\mathbf{x} \approx \mathbf{x}^{(i)}$  ist  $f_i(\mathbf{x}) \approx 1$ . Falls  $\mathbf{x}$  weit von  $\mathbf{x}^{(i)}$  entfernt ist, ist  $f_i(\mathbf{x}) \approx 0$ . Der Parameter  $\sigma > 0$  gibt an, wie eng  $f_i$  um  $\mathbf{x}^{(i)}$  zentriert ist.



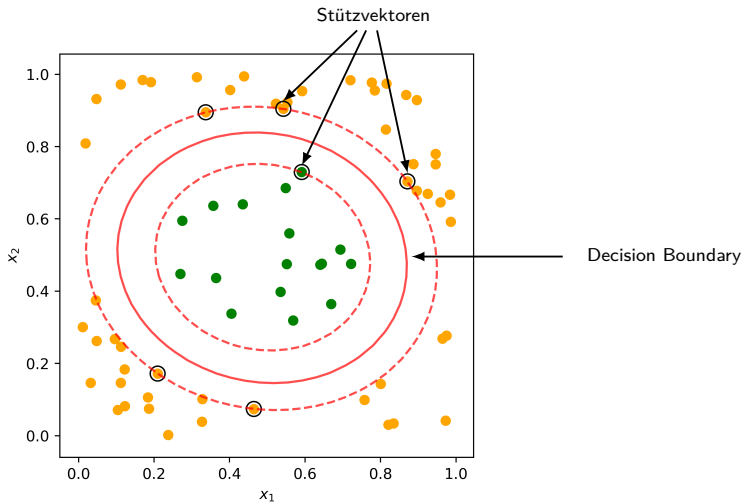
**Abbildung:** Gauß-Kernel für  $\sigma = 0.5$  (links)  $\sigma = 1$  (Mitte) und  $\sigma = 1.5$  (rechts)



# SVM-Klassifikator mit Gauß-Kernel - Beispiel



# SVM-Klassifikator mit Gauß-Kernel - Beispiel



- SVM als Large-Margin-Klassifikator
- Hard margin vs. Soft margin
- Formulierung des primalen und dualen Optimierungsproblems (für den linear separierbaren Fall)
- Struktur des dualen Optimierungsproblems
- Einführung von Schlupfvariablen zur Behandlung des nicht linear separierbaren Falls, Regularisierung
- Vergleich mit logistischer Regression
- Transformation in höherdimensionale Feature-Räume unter Verwendung von Kernels