

# Support Vector Machine

---

Florian Abt

17. Mai 2018

Universität Hamburg - Informatik

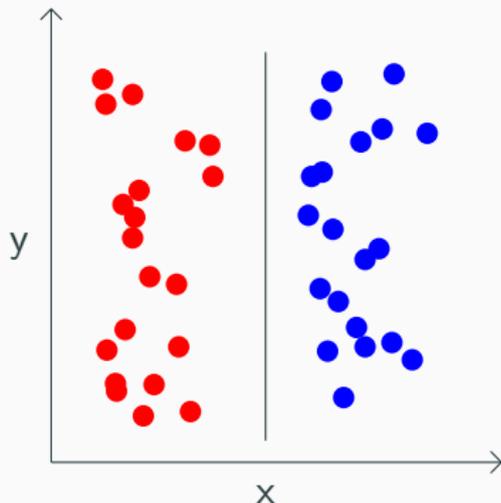
# Grundlegende Idee

---

# Grundlegende Idee - Das Problem

- Klassifikation von Elementen
- Verteilungen unbekannt
- Elemente befinden sich im  $n$ -dimensionalen Raum

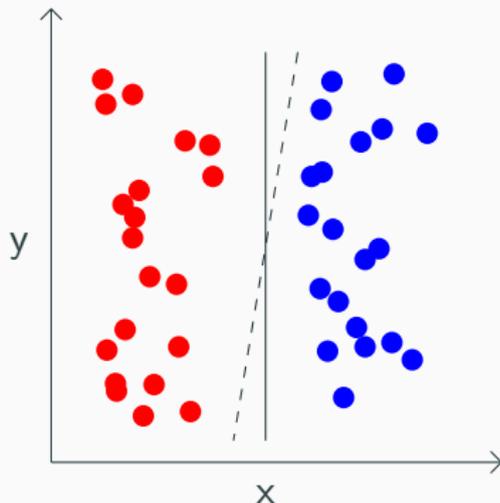
⇒ Platzieren einer Hyperebene zur Klassifikation  
(mit minimalem Trainings-Fehler)



# Grundlegende Idee - Das Problem

- Klassifikation von Elementen
- Verteilungen unbekannt
- Elemente befinden sich im  $n$ -dimensionalen Raum

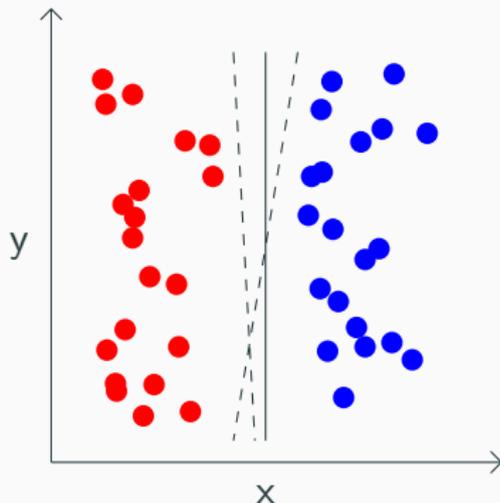
⇒ Platzieren einer Hyperebene zur Klassifikation  
(mit minimalem Trainings-Fehler)



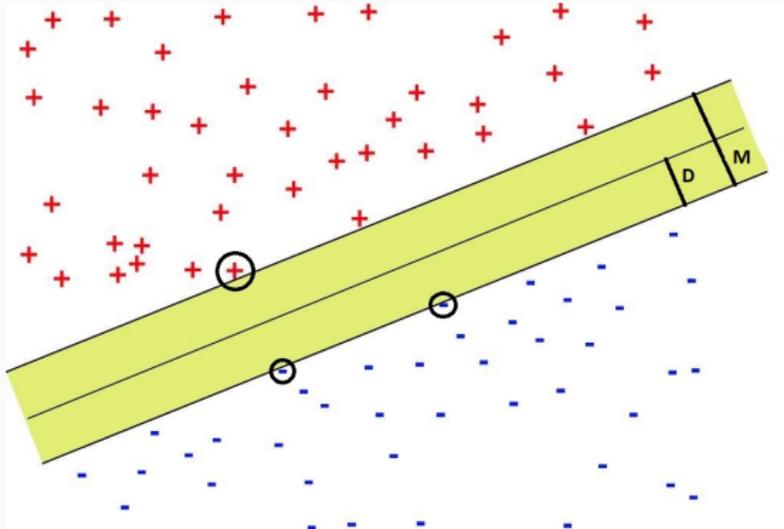
# Grundlegende Idee - Das Problem

- Klassifikation von Elementen
- Verteilungen unbekannt
- Elemente befinden sich im  $n$ -dimensionalen Raum

⇒ Platzieren einer Hyperebene zur Klassifikation  
(mit minimalem Trainings-Fehler)



# Grundlegende Idee - Vorgehen



- Maximaler Abstand  $\Rightarrow$  geringes Fehlerrisiko
- Stützvektoren zur Bestimmung der Hyperebene

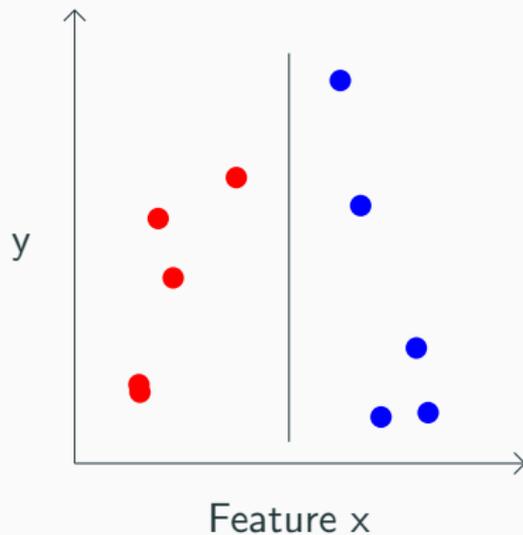
# Grundlegende Idee - Nicht-lineares Mapping

Eingaberaum



$\Rightarrow$   
Mapping

Feature-Raum



# Algorithmen

---

# Algorithmen - Übersicht

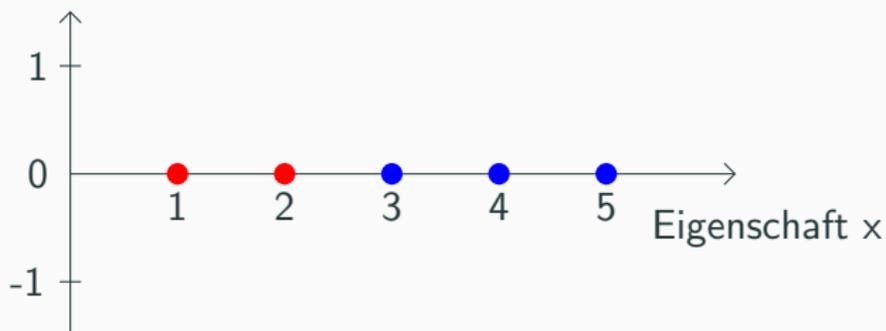
- Linear Maximum Margin Classifier
  - klare lineare Separierung
  - lineare Hyperebene
  - Maximaler Abstand Ebene - Elemente
  
- Linear Soft Margin Classifier
  - nicht klare lineare Separierung
  - lineare Hyperebene
  - bestimmte Anzahl Elemente falsch klassifiziert
  
- Nonlinear Classifier
  - keine lineare Separierung möglich
  - nicht-lineare Hyperebene  
bzw.
  - nicht-lineares Mapping erforderlich

# Algorithmen - Linear Maximum Margin Classifier

Beispiel:

$$(x_1, y_1), (x_2, y_2) \dots (x_n, y_n), \quad x \in \mathbb{R}^1, \quad y \in \{+1, -1\}$$

Klassifikation  $y$



# Algorithmen - Linear Maximum Margin Classifier

Annahme: Finde Hyperebene mit größtem Abstand  
⇒ geringstes Klassifizierungsrisiko

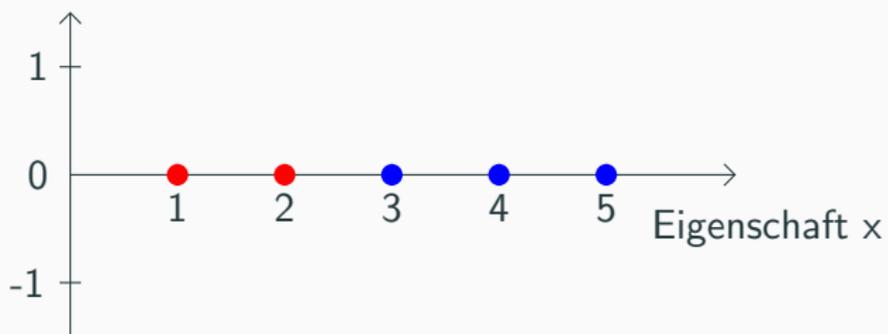
Entscheidungsfunktion:

$$d(x, w, b) = w^T \cdot x + b = \sum_{i=1}^n w_i \cdot x_i + b$$

Zu bestimmende Parameter der Entscheidungsfunktion  
(Ebenengleichung) sind  $w = [w_1 \quad w_2 \quad \dots \quad w_n]^T$  und  $b$

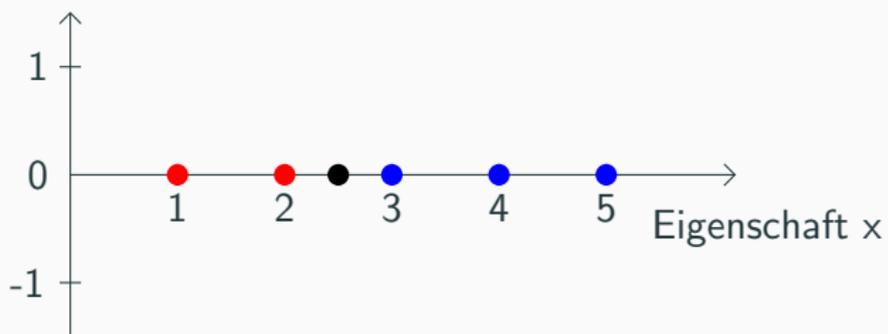
# Algorithmen - Linear Maximum Margin Classifier

Klassifikation  $y$



# Algorithmen - Linear Maximum Margin Classifier

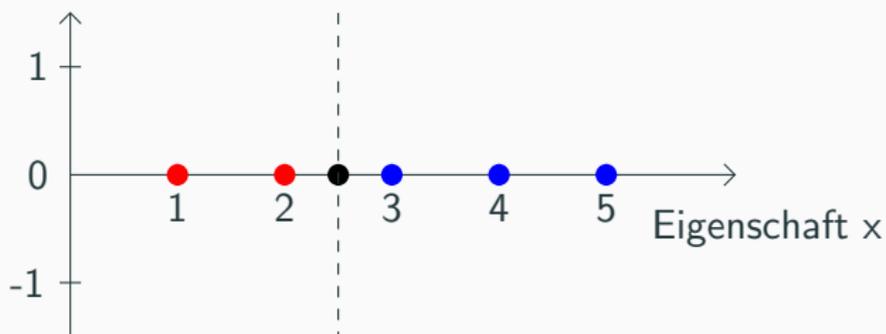
Klassifikation  $y$



# Algorithmen - Linear Maximum Margin Classifier

$d(x, w, b)$  beschreibt kanonische Hyperebene (Hypergerade)

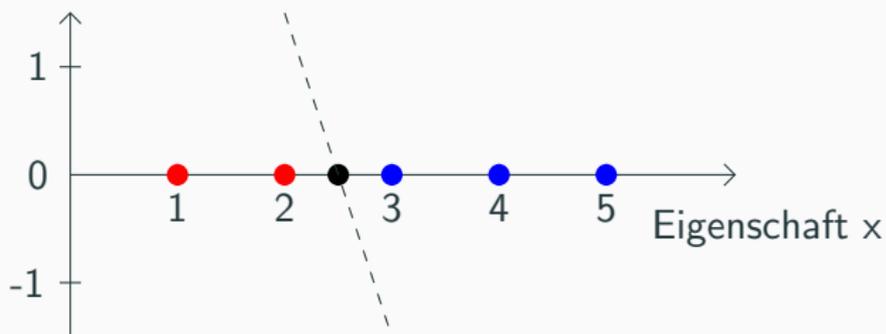
Klassifikation  $y$



# Algorithmen - Linear Maximum Margin Classifier

$d(x, w, b)$  beschreibt kanonische Hyperebene (Hypergerade)

Klassifikation  $y$

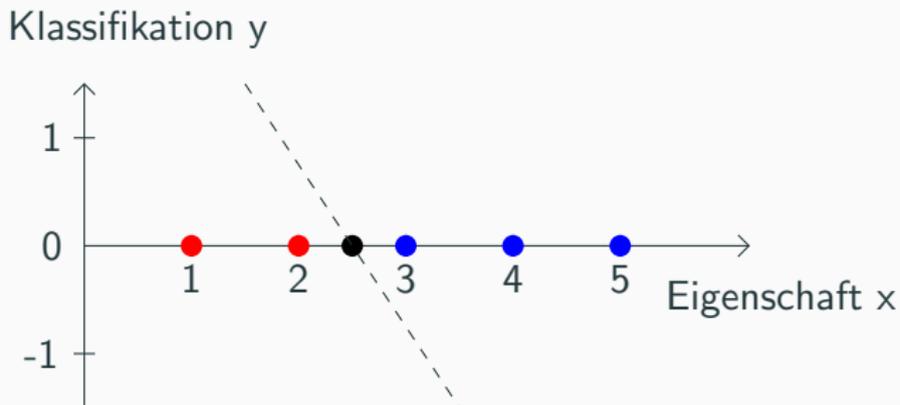


$d(x, k \cdot w, k \cdot b), \quad k \in \mathbb{N}$

beschreibt alle möglichen Hyperebenen

# Algorithmen - Linear Maximum Margin Classifier

$d(x, w, b)$  beschreibt kanonische Hyperebene (Hypergerade)



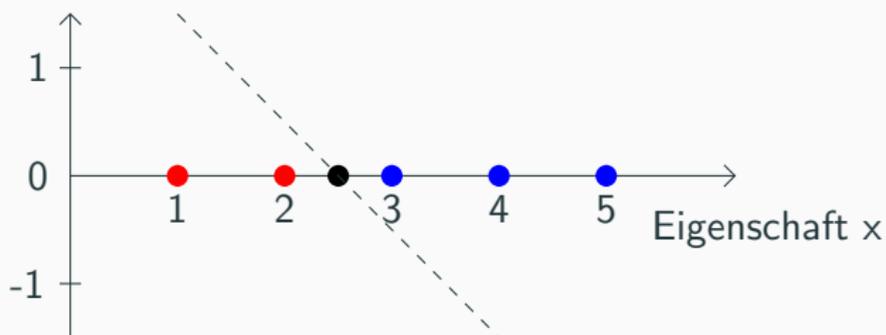
$d(x, k \cdot w, k \cdot b), \quad k \in \mathbb{N}$

beschreibt alle möglichen Hyperebenen

# Algorithmen - Linear Maximum Margin Classifier

$d(x, w, b)$  beschreibt kanonische Hyperebene (Hypergerade)

Klassifikation  $y$

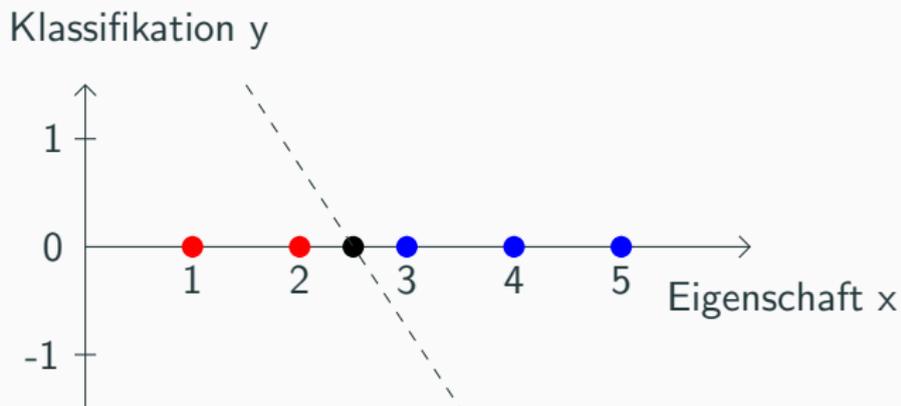


$d(x, k \cdot w, k \cdot b), \quad k \in \mathbb{N}$

beschreibt alle möglichen Hyperebenen

# Algorithmen - Linear Maximum Margin Classifier

$d(x, w, b)$  beschreibt kanonische Hyperebene (Hypergerade)



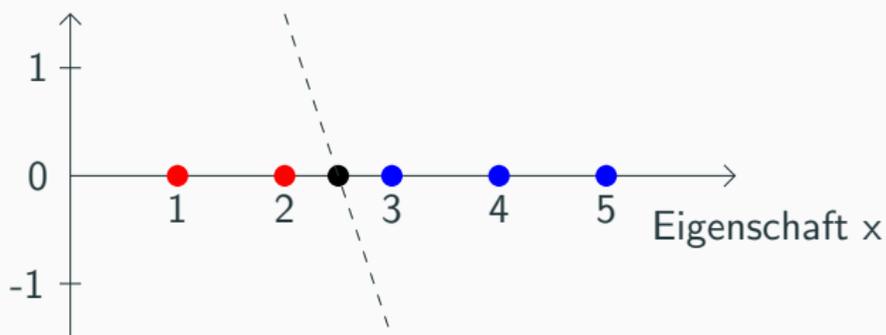
$d(x, k \cdot w, k \cdot b), \quad k \in \mathbb{N}$

beschreibt alle möglichen Hyperebenen

# Algorithmen - Linear Maximum Margin Classifier

$d(x, w, b)$  beschreibt kanonische Hyperebene (Hypergerade)

Klassifikation  $y$



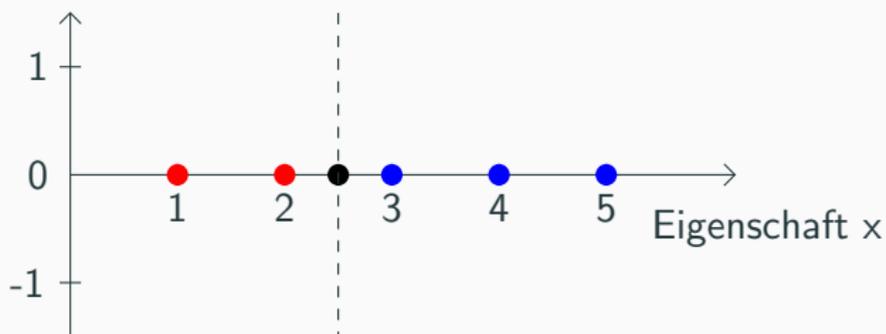
$d(x, k \cdot w, k \cdot b), \quad k \in \mathbb{N}$

beschreibt alle möglichen Hyperebenen

# Algorithmen - Linear Maximum Margin Classifier

$d(x, w, b)$  beschreibt kanonische Hyperebene (Hypergerade)

Klassifikation  $y$



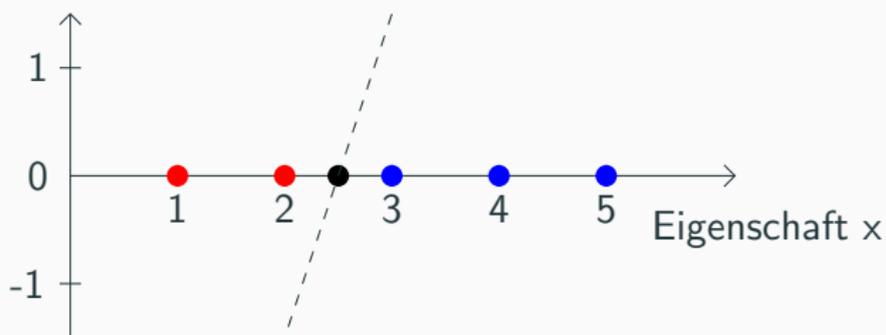
$d(x, k \cdot w, k \cdot b), \quad k \in \mathbb{N}$

beschreibt alle möglichen Hyperebenen

# Algorithmen - Linear Maximum Margin Classifier

$d(x, w, b)$  beschreibt kanonische Hyperebene (Hypergerade)

Klassifikation  $y$



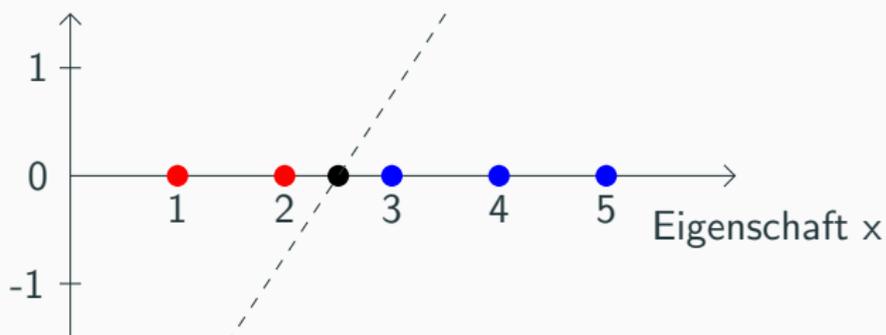
$d(x, k \cdot w, k \cdot b), \quad k \in \mathbb{N}$

beschreibt alle möglichen Hyperebenen

# Algorithmen - Linear Maximum Margin Classifier

$d(x, w, b)$  beschreibt kanonische Hyperebene (Hypergerade)

Klassifikation  $y$



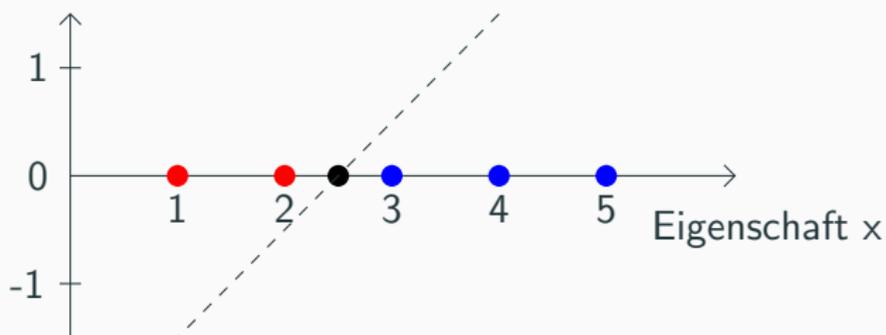
$d(x, k \cdot w, k \cdot b), \quad k \in \mathbb{N}$

beschreibt alle möglichen Hyperebenen

# Algorithmen - Linear Maximum Margin Classifier

$d(x, w, b)$  beschreibt kanonische Hyperebene (Hypergerade)

Klassifikation  $y$



$d(x, k \cdot w, k \cdot b), \quad k \in \mathbb{N}$

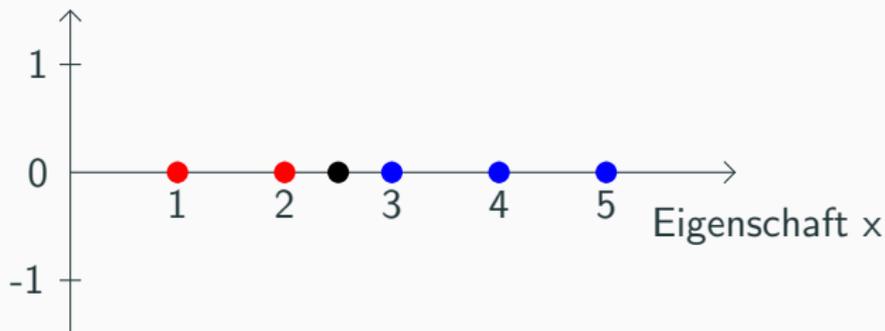
beschreibt alle möglichen Hyperebenen

# Algorithmen - Linear Maximum Margin Classifier

Indikatorfunktion  $i_F$  (Abbildung der Entscheidungsfunktion auf  $y \in \{1, -1\}$ )

$$i_F = \text{sign}(d(x_p, w, b)) := \begin{cases} 1, & \text{if } d(x_p, w, b) > 0 \quad (y_p = 1) \\ -1, & \text{if } d(x_p, w, b) < 0 \quad (y_p = -1) \end{cases}$$

Klassifikation  $y$

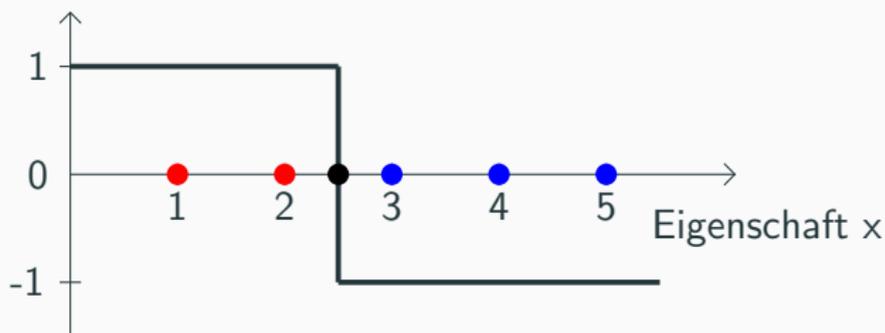


# Algorithmen - Linear Maximum Margin Classifier

Indikatorfunktion  $i_F$  (Abbildung der Entscheidungsfunktion auf  $y \in \{1, -1\}$ )

$$i_F = \text{sign}(d(x_p, w, b)) := \begin{cases} 1, & \text{if } d(x_p, w, b) > 0 \quad (y_p = 1) \\ -1, & \text{if } d(x_p, w, b) < 0 \quad (y_p = -1) \end{cases}$$

Klassifikation  $y$

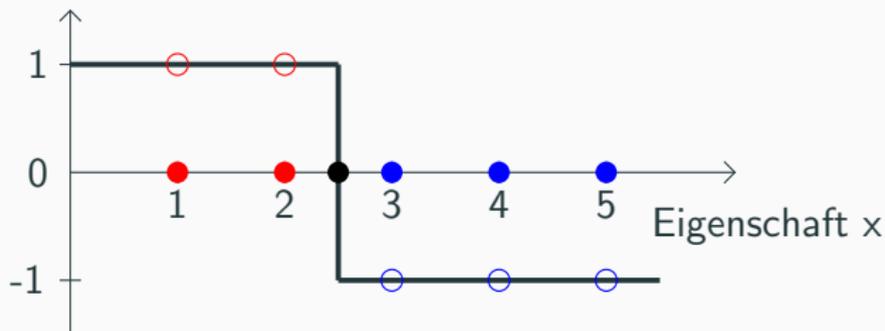


# Algorithmen - Linear Maximum Margin Classifier

Indikatorfunktion  $i_F$  (Abbildung der Entscheidungsfunktion auf  $y \in \{1, -1\}$ )

$$i_F = \text{sign}(d(x_p, w, b)) := \begin{cases} 1, & \text{if } d(x_p, w, b) > 0 \quad (y_p = 1) \\ -1, & \text{if } d(x_p, w, b) < 0 \quad (y_p = -1) \end{cases}$$

Klassifikation  $y$

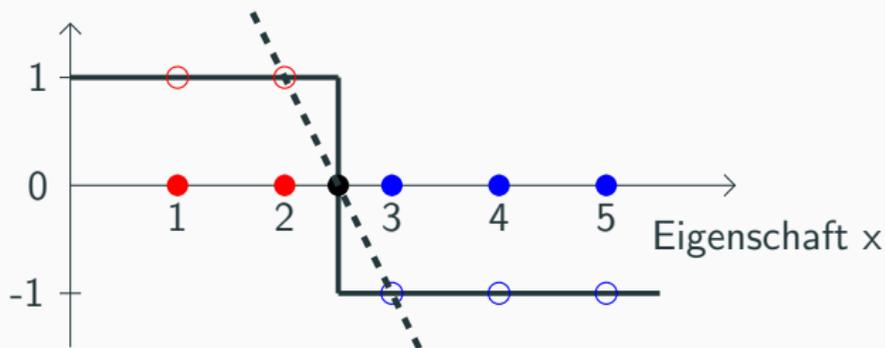


# Algorithmen - Linear Maximum Margin Classifier

Indikatorfunktion  $i_F$  (Abbildung der Entscheidungsfunktion auf  $y \in \{1, -1\}$ )

$$i_F = \text{sign}(d(x_p, w, b)) := \begin{cases} 1, & \text{if } d(x_p, w, b) > 0 \quad (y_p = 1) \\ -1, & \text{if } d(x_p, w, b) < 0 \quad (y_p = -1) \end{cases}$$

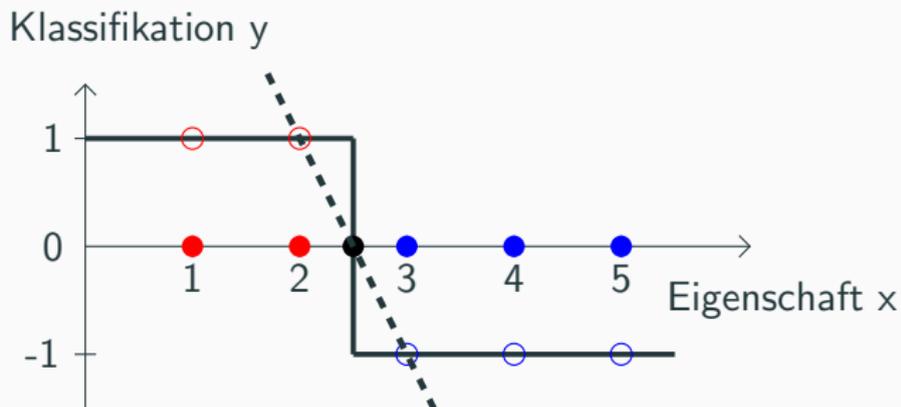
Klassifikation  $y$



# Algorithmen - Linear Maximum Margin Classifier

Indikatorfunktion  $i_F$  (Abbildung der Entscheidungsfunktion auf  $y \in \{1, -1\}$ )

$$i_F = \text{sign}(d(x_p, w, b)) := \begin{cases} 1, & \text{if } d(x_p, w, b) > 0 \quad (y_p = 1) \\ -1, & \text{if } d(x_p, w, b) < 0 \quad (y_p = -1) \end{cases}$$



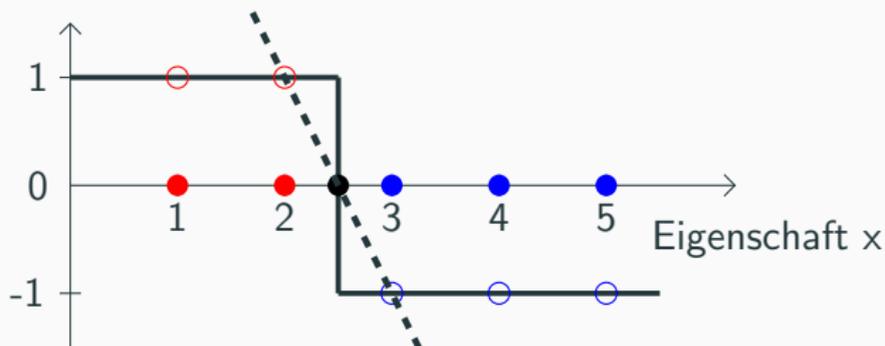
Für Stützvektoren gilt:  $d(x, w, b) = i_F = 1$

# Algorithmen - Linear Maximum Margin Classifier

Indikatorfunktion  $i_F$  (Abbildung der Entscheidungsfunktion auf  $y \in \{1, -1\}$ )

$$i_F = \text{sign}(d(x_p, w, b)) := \begin{cases} 1, & \text{if } d(x_p, w, b) > 0 \quad (y_p = 1) \\ -1, & \text{if } d(x_p, w, b) < 0 \quad (y_p = -1) \end{cases}$$

Klassifikation  $y$



Für Stützvektoren gilt:  $d(x, w, b) = i_F = 1$

Entscheidungsgrenze:  $d(x, w, b) = 0$  (im Eingaberaum)

Kanonische Hyperebene:

- für alle Trainingselemente  $x$  gilt:  $|d| \geq |i_F|$
- $R^1$  Eingaberaum  $\Rightarrow$  nur 1 kanonische Hyperebene
- $R^n$  Eingaberaum  $\Rightarrow$   $n$  kanonische Hyperebenen

Optimale kanonische Hyperebenen im  $R^n$

$\Rightarrow$  Maximaler Abstand  $D$  zu den Stützvektoren

Abstand Punkt  $p$  zur Hyperebene:

$$D = \frac{|(w x_p) \pm b|}{\|w\|} = \frac{|w_1 x_{1p} + w_2 x_{2p} + \dots + w_n x_{np} \pm b|}{\sqrt{w_1^2 + w_2^2 + \dots + w_n^2}} \quad (1)$$

## Algorithmen - Linear Maximum Margin Classifier

Abstand Punkt  $p$  zur Hyperebene:

$$D = \frac{|(wx_p) \pm b|}{\|w\|} = \frac{|w_1x_{1p} + w_2x_{2p} + \dots + w_nx_{np} \pm b|}{\sqrt{w_1^2 + w_2^2 + \dots + w_n^2}} \quad (1)$$

Da Punkt  $p$  ein Stützvektor gilt:

$$d(x_p, w, b) = w^T x_p + b = i_F = 1 \quad (2)$$

# Algorithmen - Linear Maximum Margin Classifier

Abstand Punkt  $p$  zur Hyperebene:

$$D = \frac{|(w x_p) \pm b|}{\|w\|} = \frac{|w_1 x_{1p} + w_2 x_{2p} + \dots + w_n x_{np} \pm b|}{\sqrt{w_1^2 + w_2^2 + \dots + w_n^2}} \quad (1)$$

Da Punkt  $p$  ein Stützvektor gilt:

$$d(x_p, w, b) = w^T x_p + b = i_F = 1 \quad (2)$$

Es folgt:

$$D = \frac{1}{\sqrt{w_1^2 + w_2^2 + \dots + w_n^2}} \quad (3)$$

# Algorithmen - Linear Maximum Margin Classifier

Abstand Punkt  $p$  zur Hyperebene:

$$D = \frac{|(w x_p) \pm b|}{\|w\|} = \frac{|w_1 x_{1p} + w_2 x_{2p} + \dots + w_n x_{np} \pm b|}{\sqrt{w_1^2 + w_2^2 + \dots + w_n^2}} \quad (1)$$

Da Punkt  $p$  ein Stützvektor gilt:

$$d(x_p, w, b) = w^T x_p + b = i_F = 1 \quad (2)$$

Es folgt:

$$D = \frac{1}{\sqrt{w_1^2 + w_2^2 + \dots + w_n^2}} \quad (3)$$

Die **Minimierung von  $w^T w$**  (Äquivalent zu Minimierung von  $\sqrt{w^T w}$ ) unter der Bedingung  $y_i [w^T x_i + b] \geq 1$  für jedes Trainingselement, **maximiert** den Abstand **D**.

Nicht-lineares Optimierungsproblem mit Ungleichungsbedingung:

*Minimiere  $f(x)$*

*unter der Nebenbedingung  $g(x) \leq 0$*

Nicht-lineares Optimierungsproblem mit Ungleichungsbedingung:

$$\begin{aligned} & \text{Minimiere } f(x) \\ & \text{unter der Nebenbedingung } g(x) \leq 0 \end{aligned}$$

Lagrange-Dualität mit Lagrange-Multiplikatoren  $\alpha_i$ :

$$L(w, b, \alpha) = \underbrace{\frac{1}{2} w^T w}_{f(x)=M^{-1}=(2D)^{-1}} - \sum_{i=1}^l \alpha_i \underbrace{\{y_i[w^T x_i + b] - 1\}}_{g(x)=y_i[w^T x_i + b] - 1 \geq 0} \quad (4)$$

Nicht-lineares Optimierungsproblem mit Ungleichungsbedingung:

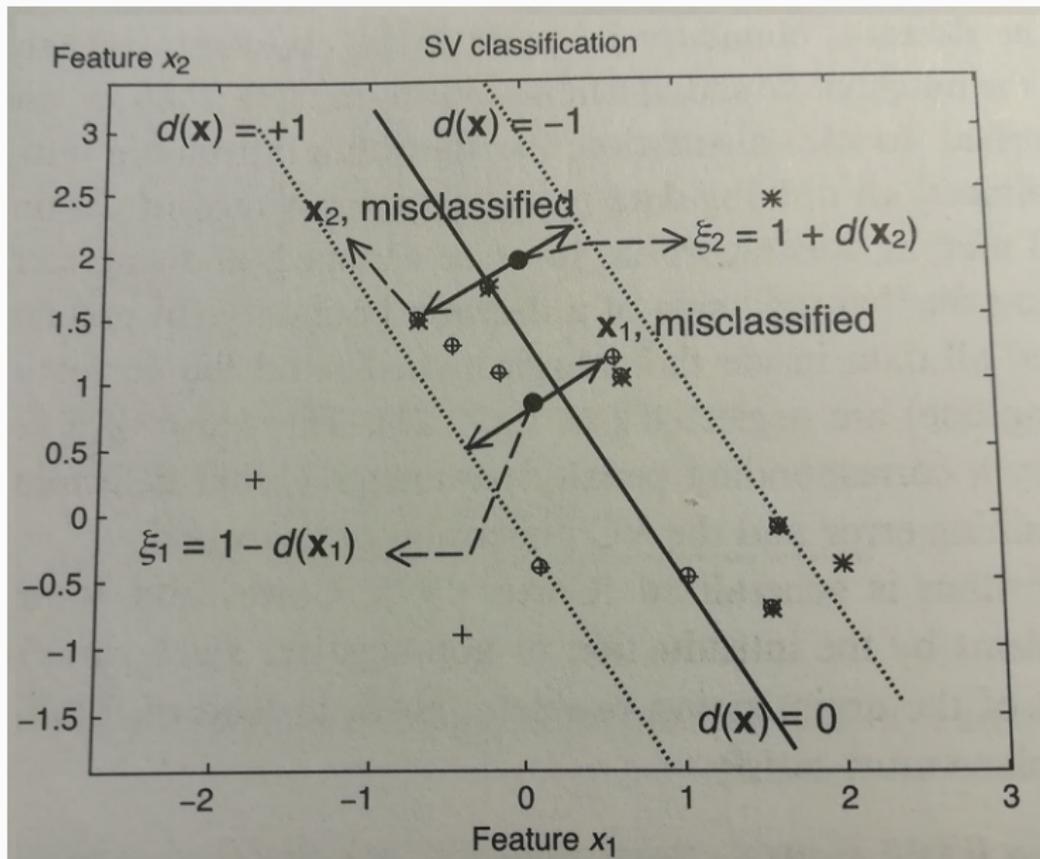
$$\begin{aligned} & \text{Minimiere } f(x) \\ & \text{unter der Nebenbedingung } g(x) \leq 0 \end{aligned}$$

Lagrange-Dualität mit Lagrange-Multiplikatoren  $\alpha_i$ :

$$L(w, b, \alpha) = \underbrace{\frac{1}{2} w^T w}_{f(x)=M^{-1}=(2D)^{-1}} - \sum_{i=1}^l \alpha_i \underbrace{\{y_i[w^T x_i + b] - 1\}}_{g(x)=y_i[w^T x_i + b] - 1 \geq 0} \quad (4)$$

Der Sattelpunkt von  $L(w, b, \alpha)$  ist das gesuchte Minimum von  $f(x)$  und damit der größte Abstand  $D$ .

# Algorithmen - Linear Soft Margin Classifier



## Algorithmen - Linear Soft Margin Classifier

Erweiterung der Bedingung  $y_i[w^T x_i + b] \geq 1$  mit lockeren Variablen  $\xi_i$  ( $i = 0, l$ ):

$$y_i[w^T x_i + b] \geq 1 - \xi_i, \quad i = 1, l \quad (5)$$

$$\xi_i \geq 0 \quad (6)$$

## Algorithmen - Linear Soft Margin Classifier

Erweiterung der Bedingung  $y_i[w^T x_i + b] \geq 1$  mit lockeren Variablen  $\xi_i$  ( $i = 0, l$ ):

$$y_i[w^T x_i + b] \geq 1 - \xi_i, \quad i = 1, l \quad (5)$$

$$\xi_i \geq 0 \quad (6)$$

Einführung eines frei wählbaren Strafparameters  $C$  (Überlappungsfehler) in der zu minimierenden Funktion:

$$J(w, \xi) = \frac{1}{2} w^T w + C \left( \sum_{i=1}^l \xi_i \right) \quad (7)$$

Lagrange-Dualität mit Lagrange-Multiplikatoren  $\alpha_i$  und  $\beta_i$

$$\begin{aligned} L(w, b, \xi, \alpha, \beta) = & \frac{1}{2} w^T w + C \left( \sum_{i=1}^l \xi \right) \\ & - \sum_{i=1}^l \alpha_i \{ [w^T x_i + b] - 1 + \xi \} \\ & - \sum_{i=1}^l \beta_i \xi_i \end{aligned} \quad (8)$$

Der Sattelpunkt von  $(w, b, \xi, \alpha, \beta)$  ist das gesuchte Minimum von  $f(x)$  und damit der größte Abstand D.

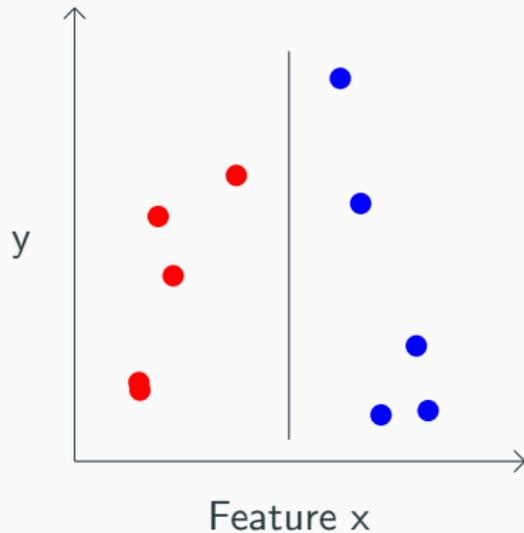
# Algorithmen - Nonlinear Classifier

Eingaberaum



$\Rightarrow$   
Mapping

Feature-Raum



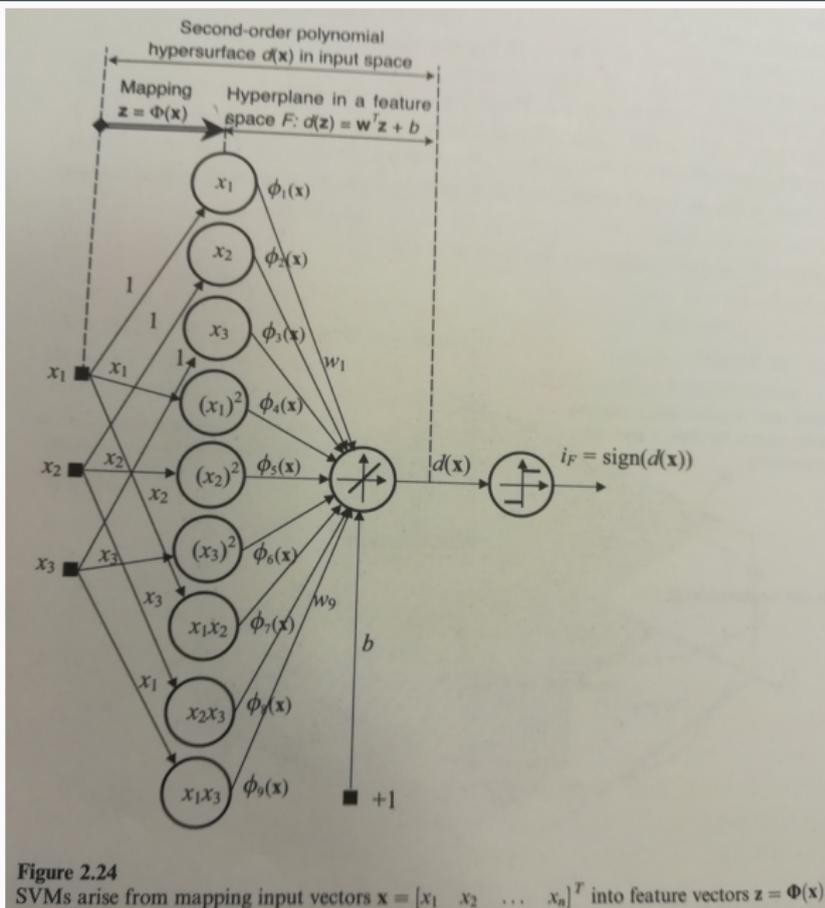
Erhöhung der Dimension!

Erweiterung Eingaberaum auf Feature-Raum

$$R^n \rightarrow R^f \quad \text{mit Mapping} \quad \Phi(x) \quad (9)$$

für Vektor  $x = [x_1, x_2, x_3]$

# Algorithmen - Nonlinear Classifier



Erweiterung Eingaberaum auf Feature-Raum

$$R^n \rightarrow R^f \quad \text{mit Mapping } \Phi(x) \quad (9)$$

für Vektor  $x = [x_1, x_2, x_3]$

- Lineare Hyperebene im Feature-Raum
- Nicht-lineare Hyperoberfläche im Eingaberaum

Erweiterung Eingaberaum auf Feature-Raum

$$R^n \rightarrow R^f \quad \text{mit Mapping} \quad \Phi(x) \quad (9)$$

für Vektor  $x = [x_1, x_2, x_3]$

- Lineare Hyperebene im Feature-Raum
- Nicht-lineare Hyperoberfläche im Eingaberaum

Probleme:

- Finden von  $\Phi(x)$
- Explosion der Dimension

Beispiel:  $\Phi(x)$

# Algorithmen - Nonlinear Classifier

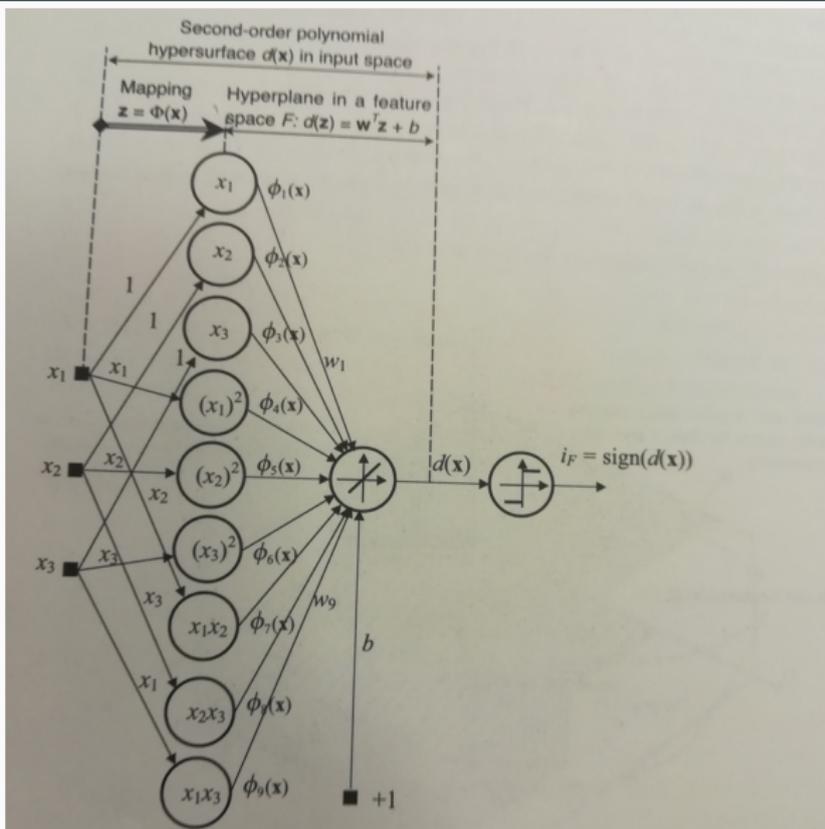


Figure 2.24

SVMs arise from mapping input vectors  $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_n]^T$  into feature vectors  $\mathbf{z} = \Phi(\mathbf{x})$ .

Beispiel:  $\Phi(x)$

Beispiel:  $\Phi(x)$

für Eingaberaum  $R^{256}$

ergibt sich Feature-Raum  $R^{33.152}$

Beispiel:  $\Phi(x)$

für Eingaberaum  $R^{256}$

ergibt sich Feature-Raum  $R^{33.152}$

$$L(w, b, \alpha) = \underbrace{\frac{1}{2} w^T w}_{\text{Funktion in } R^n} - \sum_{i=1}^l \alpha_i \cdot \text{Nebenbedingung}$$

Beispiel:  $\Phi(x)$

für Eingaberaum  $R^{256}$

ergibt sich Feature-Raum  $R^{33.152}$

$$L(w, b, \alpha) = \underbrace{\frac{1}{2} w^T w}_{\text{Funktion in } R^n} - \sum_{i=1}^l \alpha_i \cdot \text{Nebenbedingung}$$

Kernel Funktionen  $K(x_i, x_j)$  als Funktion im Eingaberaum:

$$K(x_i, x_j) = \Phi^T(x_i)\Phi(x_j)$$

# Anwendung

---

### SemEval-2018 task 7 Semantic Relation Extraction and Classification in Scientific Papers

- Klassifikation von semantischen Beziehungen
- Analyse Abstracts wissenschaftlicher Paper
- 6 verschiedene Klassen
- 2 Teil Task
  - Task a mit bereits extrahierten gelabelten semantischen Beziehungen
  - Task b mit Extraktion der semantischen Beziehungen

Problem:

1. Wörter in Vektoren konvertieren
  - word2Vec
  - elmo
2. Welche Feature zum Platzieren im Feature-Raum?

Zwei frameworks für SVM's:

- LIBSVM <https://www.csie.ntu.edu.tw/~cjlin/libsvm/>
  - Datenformat: <label> <index1>:<value1>  
<index2>:<value2> ...
- SVM<sup>light</sup> <http://svmlight.joachims.org/>

**Ende**

---

Danke für eure Aufmerksamkeit

Gibt es noch Fragen?

# Literatur

---

-  Kecman, V. (2001). Learning and soft computing: support vector machines, neural networks, and fuzzy logic models. MIT press.
-  <http://www.statsoft.com/Textbook/Support-Vector-Machines> (Stand: 09.05.2018)
-  <https://www.csie.ntu.edu.tw/~cjlin/libsvm/> (Stand:16.05.2018)
-  <http://svmlight.joachims.org/> (Stand:16.05.2018)