

Universität Potsdam
Institut für Informatik
Lehrstuhl Maschinelles Lernen



Maschinelles Lernen Modelle, Version Spaces, Lernen

Christoph Sawade/Niels Landwehr

Silvia Makowski

Tobias Scheffer

Überblick

- Problemstellungen: Klassifikation und Regression
- Modelle und Hypothesenraum
- Verlustfunktionen und Regularisierer
- Unsicherheit, Wahrscheinlichkeiten, Bayes'sche Regel

Überblick

- Problemstellungen: Klassifikation und Regression
- Modelle und Hypothesenraum
- Verlustfunktionen und Regularisierer
- Unsicherheit, Wahrscheinlichkeiten, Bayes'sche Regel

Problemstellung Klassifikation

- Eingabe: Instanz (Objekt) $\mathbf{x} \in X$.
 - ◆ Objekte oft durch Vektor von *Attributen* repräsentiert ($X = \mathbb{R}^m$)
 - ◆ Instanz ist Belegung der Attribute.

- ◆ $\mathbf{x} = \begin{pmatrix} x_1 \\ \dots \\ x_m \end{pmatrix}$ Merkmalsvektor


- Ausgabe: Klasse $y \in Y$; endliche Menge Y .
 - ◆ Klasse wird auch als Zielattribut bezeichnet
 - ◆ y heißt auch (Klassen)Label





Klassifikation: Beispiel

- Eingabe: Instanz (Objekt) $\mathbf{x} \in X$.

X = Menge aller möglichen Kombinationen einer Menge von Medikamenten

Attribute	Instanz \mathbf{x}	Medikamenten- kombination
Medikament 1 enthalten?	$\begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ 0 \end{pmatrix}$	
⋮		
Medikament 6 enthalten?		

Belegung der Attribute,
Merkmalsvektor

- Ausgabe: $y \in Y = \{toxisch, ok\}$  / 

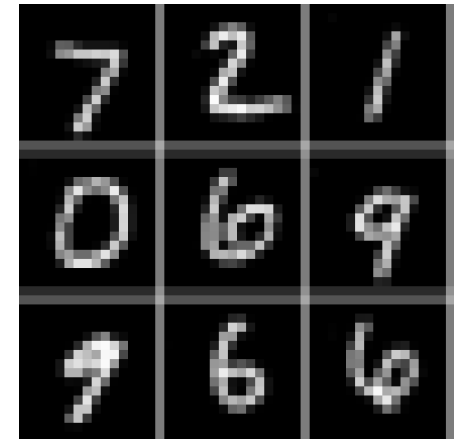


Klassifikation: Beispiel

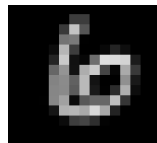
- Eingabe: Instanz (Objekt) $\mathbf{x} \in X$.

X = Menge aller 16x16 Pixel Bitmaps

Attribute	Instanz \mathbf{x}
Grauwert Pixel 1	$\begin{pmatrix} 0.1 \\ 0.3 \\ 0.45 \\ \dots \\ 0.65 \\ 0.87 \end{pmatrix}$ 256 Pixelwerte
...	
Grauwert Pixel 256	



- Ausgabe: $y \in Y = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$: erkannte Ziffer



→ Klassifikator → "6"

Klassifikation: Beispiel

- Eingabe: Instanz (Objekt) $\mathbf{x} \in X$.

X = Menge aller möglichen Email-Texte

Attribute

Instanz \mathbf{x}

Wort 1 kommt vor?

$\begin{pmatrix} 0 \\ 1 \\ 0 \\ \dots \\ 1 \\ 0 \end{pmatrix}$

Alternative

Beneficiary

Friend

...

Sterling

Zoo

...

Wort N kommt vor?

$N \approx 100000$

Email

Dear Beneficiary,

your Email address has been picked online in this years MICROSOFT CONSUMER AWARD as a Winner of One Hundred and Fifty Five Thousand Pounds Sterling...

- Ausgabe: $y \in Y = \{spam, ok\}$

Dear Beneficiary,
We are pleased to notify you that your Email address has been picked online in this second quarter's MICROSOFT CONSUMER AWARD (MCA) as a Winner of One Hundred and Fifty Five Thousand Pounds Sterling...



Klassifikator



„Spam“

Problemstellung Klassifikationslernen

- Idee: Klassifikator aus Daten lernen
- Eingabe Lernproblem: Trainingsdaten.

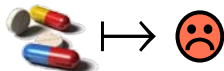
$$\diamond L = \langle (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \rangle$$

$$\diamond \mathbf{x}_i = \begin{pmatrix} x_{i1} \\ \dots \\ x_{im} \end{pmatrix}$$



- Ausgabe: Klassifikator (auch als Modell bezeichnet).

$$f : X \rightarrow Y$$



$$f(\mathbf{x}) = \begin{cases} \text{☹️}, & \text{wenn } x_1 = 1 \wedge x_3 = 0 \wedge x_6 = 1 \\ \text{😊}, & \text{sonst} \end{cases}$$

Problemstellung Klassifikationslernen

- Idee: Klassifikator aus Daten lernen
- Eingabe Lernproblem: Trainingsdaten.

$$\diamond L = \langle (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \rangle$$

$$\diamond \mathbf{x}_i = \begin{pmatrix} x_{i1} \\ \dots \\ x_{im} \end{pmatrix}$$

◆ \mathbf{x}_i Objektrepräsentation

◆ y_i Klassenlabel



$$\begin{pmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 1 \end{pmatrix}, ok$$

(\mathbf{x}_1, y_1)

$$\begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ 0 \end{pmatrix}, toxisch$$

(\mathbf{x}_2, y_2)

$$\begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 1 \end{pmatrix}, ok$$

(\mathbf{x}_3, y_3)

Problemstellung Klassifikationslernen

- Eingabe Lernproblem: Trainingsdaten.

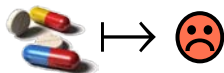
◆ $L = \langle (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \rangle$

◆ $\mathbf{x}_i = \begin{pmatrix} x_{i1} \\ \dots \\ x_{im} \end{pmatrix}$



- Ausgabe: Klassifikator (auch als Modell bezeichnet).

$f : X \rightarrow Y$



$$f(\mathbf{x}) = \begin{cases} \text{frowny face} & \text{wenn } x_1 = 1 \wedge x_3 = 0 \wedge x_6 = 1 \\ \text{smiley face} & \text{sonst} \end{cases}$$

$$f_{\mathbf{w}}(\mathbf{x}) = \begin{cases} \text{frowny face} & \text{wenn } \mathbf{w}^T \mathbf{x} + b \geq 0 \\ \text{smiley face} & \text{sonst} \end{cases}$$

Linearer Klassifikator mit
Parametervektor \mathbf{w} .

$$\mathbf{w}^T \mathbf{x} = \sum_{i=1}^m w_i x_i$$

Problemstellung Klassifikationslernen

- Eingabe Lernproblem: Trainingsdaten.

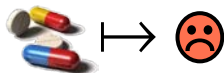
$$\diamond L = \langle (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \rangle$$

$$\diamond \mathbf{x}_i = \begin{pmatrix} x_{i1} \\ \dots \\ x_{im} \end{pmatrix}$$



- Ausgabe: Klassifikator (auch als Modell bezeichnet).

$$f : X \rightarrow Y$$



Verschiedene Klassen von Klassifikatoren

- Entscheidungsbäume.
- Generalisierte lineare Modelle (Kernel).
- ...
- Betrachtete Klassifikatoren wesentliches Unterscheidungsmerkmal zwischen Verfahren des ML

Problemstellung Klassifikationslernen

- Eingabe Lernproblem: Trainingsdaten.

- ◆ $L = \langle (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \rangle$

- ◆ $\mathbf{x}_i = \begin{pmatrix} x_{i1} \\ \dots \\ x_{im} \end{pmatrix}$



- Alternative Schreibweise:

- ◆ Trainingsinstanzen: Matrix $\mathbf{X} = (\mathbf{x}_1 \quad \dots \quad \mathbf{x}_N) = \begin{pmatrix} x_{11} & x_{N1} \\ \vdots & \ddots & \vdots \\ x_{1m} & x_{Nm} \end{pmatrix}$

- ◆ Trainingslabels: Vektor $\mathbf{y} = \begin{pmatrix} y_1 \\ \dots \\ y_N \end{pmatrix}$

Problemstellung: Regression

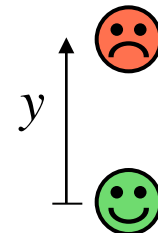
- Eingabe: Instanz (Objekt) $\mathbf{x} \in X$.
 - ◆ Objekte oft durch Attribut-Vektoren repräsentiert.
 - ◆ Instanz ist Belegung der Attribute.

- ◆ $\mathbf{x} = \begin{pmatrix} x_1 \\ \dots \\ x_m \end{pmatrix}$ Merkmalsvektor



Wie toxisch ist
Kombination?

- Ausgabe: kontinuierlicher Wert, $y \in \mathbb{R}$
 - ◆ z.B. *Toxizität*.



Regressionslernen

■ Eingabe: Trainingsdaten.

◆ $L = \langle (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \rangle$

◆ $\mathbf{x}_i = \begin{pmatrix} x_{i1} \\ \dots \\ x_{im} \end{pmatrix}$

◆ $y_i \in \mathbb{R}$

■ Ausgabe: Modell, Regressionsmodell.

◆ $f : X \rightarrow \mathbb{R}$

◆ Z.B. $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$



$$\begin{pmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 1 \end{pmatrix}, 0.05$$

$$(\mathbf{x}_1, y_1)$$



$$\begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ 0 \end{pmatrix}, 0.95$$

$$(\mathbf{x}_2, y_2)$$



$$\begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 1 \\ 1 \end{pmatrix}, 0.01$$

$$(\mathbf{x}_3, y_3)$$

Andere Lernprobleme

- Ordinale Regression.
- Präferenzlernen.
- Taxonomie-Klassifikation.
- Klassifikation und Regression mit strukturierten Ausgabebereichen.
- Kollaborative Vorhersage.
- ...

Andere Lernprobleme

- Ordinale Regression.

- Präferenzlernen

- T

- Mischung aus Klassifikation und Regression

- • Endliche, diskrete Labels

- • Ordnung

Evaluation

Very Good

Good

Average

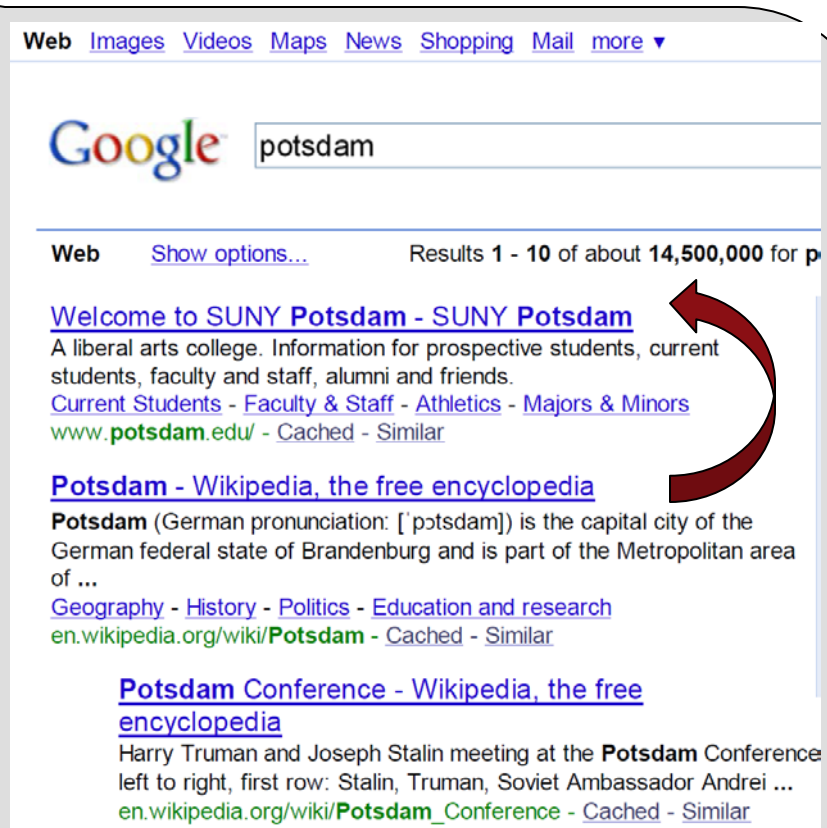
Bad

Very Bad

Andere Lernprobleme

- Ordinale Regression.
- Präferenzlernen.
- Taxonomische Klassifikation.
- ...

- • Keine direkten Klassen beobachtet, sondern nur Präferenzen
- • z.B Reihenfolge von Suchresultaten aus Clickstreams lernen



The screenshot shows a Google search interface with the search term 'potsdam'. The search results are displayed under the 'Web' tab. The first result is 'Welcome to SUNY Potsdam - SUNY Potsdam', which is a liberal arts college. The second result is 'Potsdam - Wikipedia, the free encyclopedia', which is the capital city of the German federal state of Brandenburg. The third result is 'Potsdam Conference - Wikipedia, the free encyclopedia', which is a meeting between Harry Truman and Joseph Stalin. A red arrow points from the first result to the second result.

Web Images Videos Maps News Shopping Mail more ▾

Google potsdam

Web Show options... Results 1 - 10 of about 14,500,000 for p

[Welcome to SUNY Potsdam - SUNY Potsdam](#)
A liberal arts college. Information for prospective students, current students, faculty and staff, alumni and friends.
[Current Students](#) - [Faculty & Staff](#) - [Athletics](#) - [Majors & Minors](#)
[www.potsdam.edu/](#) - [Cached](#) - [Similar](#)

[Potsdam - Wikipedia, the free encyclopedia](#)
Potsdam (German pronunciation: [ˈpɔtsdam]) is the capital city of the German federal state of Brandenburg and is part of the Metropolitan area of ...
[Geography](#) - [History](#) - [Politics](#) - [Education and research](#)
[en.wikipedia.org/wiki/Potsdam](#) - [Cached](#) - [Similar](#)

[Potsdam Conference - Wikipedia, the free encyclopedia](#)
Harry Truman and Joseph Stalin meeting at the **Potsdam** Conference left to right, first row: Stalin, Truman, Soviet Ambassador Andrei ...
[en.wikipedia.org/wiki/Potsdam_Conference](#) - [Cached](#) - [Similar](#)

Andere Lernprobleme

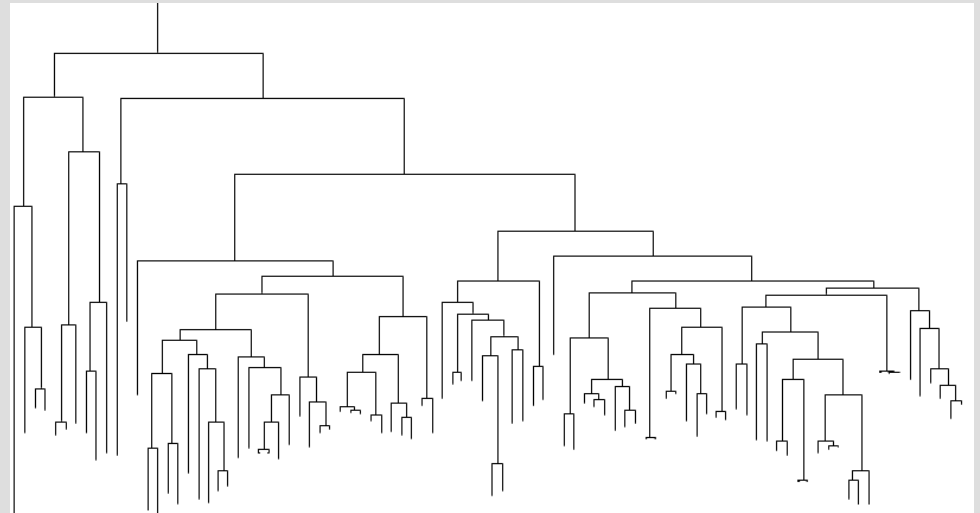
- Ordinale Regression.
- Präferenzlernen.
- Taxonomie-Klassifikation.
- Klassifikation und Regression mit strukturierten

Hierarchie von Klassen

- Ein Objekt hat mehrere
- Klassenlabels

Panther ist...

- >Tier
- >Säugetier
- >Katze
- >Panther



Andere Lernprobleme

- Ordinale Regression.
- Präferenzlernen.
- Taxonomie-Klassifikation.
- Klassifikation und Regression mit strukturierten Ausgaberräumen.

- **Kette** → **Struktur**

Eingabe X und Ausgabe Y
strukturierte Räume

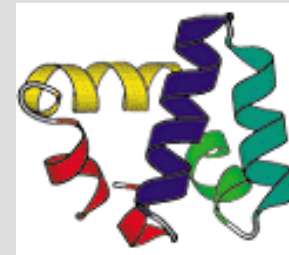
Beispiel:

Eingabe DNA

Ausgabe Proteinfaltung

Klassenlabel 3D Struktur

...AAGCTTGCCTGCGT...



Ausnutzen von Relationen
zwischen Objekten

Beispiel: Produktempfehlungen

Vorhersage interessanter Produkte

Was hat der Nutzer/haben
ähnliche Nutzer vorher gekauft?

amazon.com [Help](#) | [Close window](#)

Recommended for You

 **High Performance Web Sites:
Essential Knowledge for
Front-End Engineers**
by Steve Souders (Author)
Our Price: \$19.79
Used & new from \$16.24

I own it
 Not interested

Because you purchased...

**Programming Collective Intelligence: Building
Smart Web 2.0 Applications** (Paperback)
by Toby Segaran (Author)

This was a gift
 Don't use for
recommendations

■ **Kollaborative Vorhersage.**

■ ...

Überblick

- Lernprobleme: Klassifikation und Regression
- Modelle und Hypothesenraum
- Verlustfunktionen und Regularisierer
- Unsicherheit, Wahrscheinlichkeiten, Bayes'sche Regel

Klassifikationslernen

- Eingabe: Trainingsdaten.

- ◆ $L = \langle (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \rangle$

- ◆ $\mathbf{x}_i = \begin{pmatrix} x_{i1} \\ \dots \\ x_{im} \end{pmatrix}$



- Ausgabe: Klassifikator.

- ◆ $f : X \rightarrow Y$

$$f(\mathbf{x}) = \begin{cases} \text{frowny face}, & \text{wenn } x_1 = 1 \wedge x_3 = 0 \wedge x_6 = 1 \\ \text{smiley face}, & \text{sonst} \end{cases}$$

- Wie Klassifikator lernen aus Trainingsdaten?

- ◆ Ansatz: Klassifikator, der Trainingsdaten (Beobachtungen) erklärt

- ◆ Suchproblem im Raum aller (betrachteten) Klassifikatoren

Hypothesenraum

- Hypothesenraum, Modellraum H :
 - ◆ Menge der Klassifikationsmodelle, die Lernverfahren in Betracht zieht.
 - ◆ Hypothesenraum ist einer der Freiheitsgrade beim maschinellen Lernen, viele Räume gebräuchlich.
 - ◆ Hypothesenraum heisst auch *Language Bias*
- Beispiel:
 - ◆ Alle möglichen Konjunktionen von Bedingungen

$$f(\mathbf{x}) = \begin{cases} \text{☹️}, & \text{wenn } \bigwedge_{j \in J} x_j = v_j \\ \text{😊}, & \text{sonst} \end{cases} \quad J \subseteq \{1, \dots, m\}, v_j \in \{0, 1\}$$

- ◆ Wie groß ist Hypothesenraum (m binäre Attribute)?

Suche nach Hypothese

- Suche nach Klassifikator für „Kombination toxisch“.
- Hypothesenraum:

$$f(\mathbf{x}) = \begin{cases} \text{☹️}, & \text{wenn } \bigwedge_{j \in J} x_j = v_j \\ \text{😊}, & \text{sonst} \end{cases}$$

Beispiel-
Kombinationen

Trainingsdaten

Medikamente in der Kombination

	x_1	x_2	x_3	x_4	x_5	x_6	y
\mathbf{x}_1	0	1	0	0	1	1	☹️
\mathbf{x}_2	0	1	1	0	1	1	☹️
\mathbf{x}_3	1	0	1	0	1	0	😊
\mathbf{x}_4	0	1	1	0	0	0	☹️

- Ansatz: Hypothese sollte konsistent sein mit Trainingsdaten

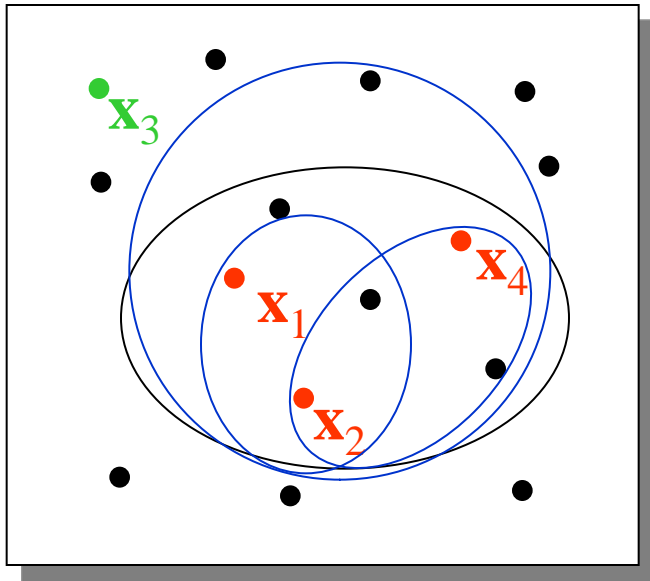
$$\forall i: f(\mathbf{x}_i) = y_i$$

- Identifizieren aller solchen Hypothesen?
- Nutze Struktur auf dem Hypothesenraum (generell/speziell)

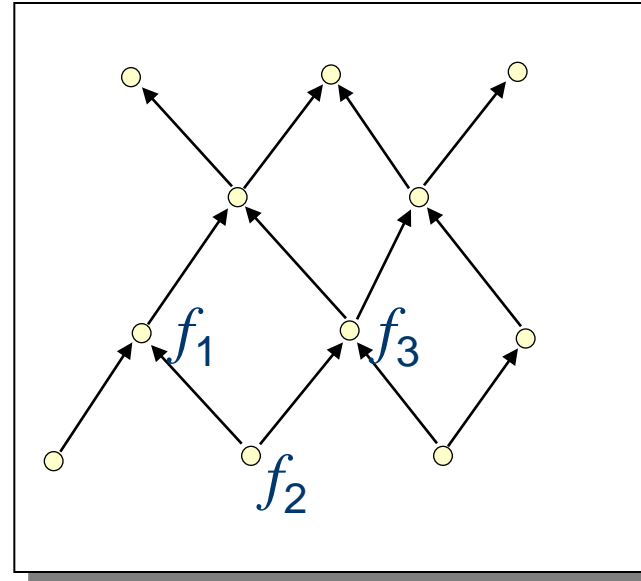
„Genereller-Als“-Ordnung

$$f_g \geq_g f_s \quad \text{gdw.} \quad f_s(\mathbf{x}) = \text{☹} \Rightarrow f_g(\mathbf{x}) = \text{☹} \quad \forall \mathbf{x} \in X$$

Grundmenge X



Hypothesen H



spezifisch

Immer ☺

generell

Immer ☹

$f_1 = \text{☹}$, wenn $x_2=1, x_6=1$ $f_2 = \text{☹}$, wenn $x_2=1$ $f_3 = \text{☹}$, wenn $x_2=1, x_3=1$

$$f_2 \geq_g f_1 \quad f_2 \geq_g f_3 \quad \text{aber nicht } f_1 \stackrel{\leq_g}{\geq_g} f_3$$

Version Space

- Menge aller Hypothesen, die mit den Trainingsdaten konsistent sind, nennen wir den Version Space:

$$VS_{H,L} = \{f \in H \mid \forall (\mathbf{x}_i, y_i) \in L : f(\mathbf{x}_i) = y_i\}$$

- Version Space begrenzt durch generellste/speziellste Hypothesen, die Daten erklären

$$G = \{f \in VS_{H,L} \mid \neg \exists f' \in VS_{H,L} : f' >_g f\}$$

Generellste konsistente Hypothesen

$$S = \{f \in VS_{H,L} \mid \neg \exists f' \in VS_{H,L} : f >_g f'\}$$

Speziellste konsistente Hypothesen

Version Space: alles „zwischen“ G und S (keine unendlichen Ketten)

$$VS_{H,L} = \{f \in H \mid \exists f_g \in G, f_s \in S : f_g \geq_g f \geq_g f_s\}$$

Version Space: Beobachtungen

$$VS_{H,L} = \{f \in H \mid \forall (\mathbf{x}_i, y_i) \in L : f(\mathbf{x}_i) = y_i\}$$

- Version Space wird kleiner, je mehr Daten vorhanden
- Version Space leer: Trainingsmenge widersprüchlich (es existiert keine Hypothese in H , die Daten erklärt)
- Version Space einelementig:
 - ◆ Richtiges Modell gefunden,
 - ◆ Oder richtiges Modell ist nicht im Hypothesenraum.
- Mehrere Elemente im Version Space:
 - ◆ Noch nicht fertig.

Version Space: Brute Force Konstruktion

- Initialisiere V auf Menge aller Hypothesen.
- Für alle Trainingsbeispiele (\mathbf{x}_i, y_i) :
 - ◆ Lösche alle Hypothesen f aus V , die mit \mathbf{x}_i inkonsistent sind, also $f(\mathbf{x}_i) \neq y_i$.
- V ist jetzt der Version Space
- Bessere Verfahren unter Benutzung von G und S (keine Details)

Beispiel: Version Space

■ H: $f(\mathbf{x}) = \begin{cases} \text{☹️}, & \text{wenn } \bigwedge_{j \in J} x_j = v_j \\ \text{😊}, & \text{sonst} \end{cases}$

- Welche Hypothesen sind im Version Space?

Medikamente in der Kombination

	x_1	x_2	x_3	x_4	x_5	x_6	y
\mathbf{x}_1	0	1	0	0	1	1	☹️
\mathbf{x}_2	0	1	1	0	1	1	☹️
\mathbf{x}_3	1	0	1	0	1	0	😊
\mathbf{x}_4	0	1	1	0	0	0	☹️

$L = \langle (\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), (\mathbf{x}_3, y_3), (\mathbf{x}_4, y_4) \rangle$

Beispiel: Version Space

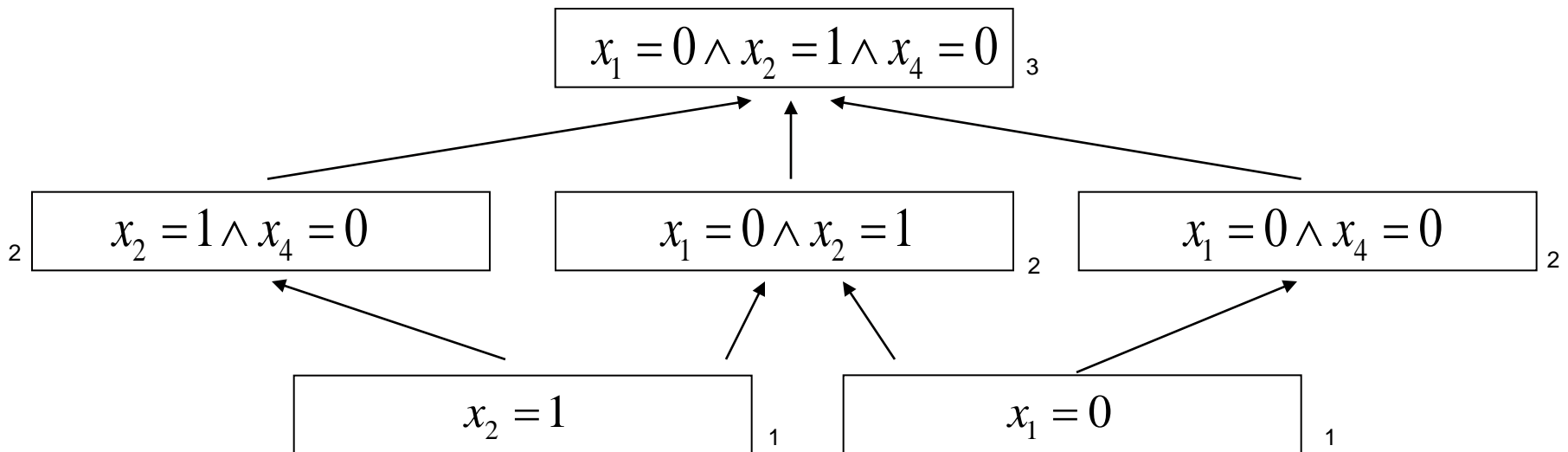
■ H: $f(\mathbf{x}) = \begin{cases} \text{☹️}, & \text{wenn } \bigwedge_{j \in J} x_j = v_j \\ \text{😊️}, & \text{sonst} \end{cases}$

- Welche Hypothesen sind im Version Space?

Medikamente in der Kombination

		x_1	x_2	x_3	x_4	x_5	x_6	y
Beispiel-Kombinationen	\mathbf{x}_1	0	1	0	0	1	1	☹️
	\mathbf{x}_2	0	1	1	0	1	1	☹️
	\mathbf{x}_3	1	0	1	0	1	0	😊️
	\mathbf{x}_4	0	1	1	0	0	0	☹️

$$L = \langle (\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), (\mathbf{x}_3, y_3), (\mathbf{x}_4, y_4) \rangle$$



Version Space

- H: $f(\mathbf{x}) = \begin{cases} \text{☹️}, & \text{wenn } \bigwedge_j x_j = v_j \\ \text{😊}, & \text{sonst} \end{cases}$

- Welche Hypothesen sind im Version Space?

- Problem:

- ◆ Alle Elemente des Version Space erklären die Daten gleichermaßen gut.
- ◆ Version Space nicht robust bei fehlerhaften Daten

Medikamente in der Kombination

	x_1	x_2	x_3	x_4	x_5	x_6	y
\mathbf{x}_1	0	1	0	0	1	1	☹️
\mathbf{x}_2	0	1	1	0	1	1	☹️
\mathbf{x}_3	1	0	1	0	1	0	😊
\mathbf{x}_4	0	1	1	0	0	0	☹️

Unsicherheit

- In der Praxis erreicht man niemals Gewissheit darüber, ein korrektes Modell gefunden zu haben.
- Version Space-Ansatz problematisch
 - ◆ Der Hypothesenraum ist meist unendlich groß.
 - ◆ Der Version Space ist dann meist auch unendlich groß, oder leer.

Alternative/zusätzliche Konzepte

- Lernen als *Optimierungsproblem*
 - ◆ Verlustfunktionen: Grad der Konsistenz mit Trainingsdaten
 - ◆ A-priori Verteilung über Modelle, Regularisierer

Überblick

- Lernprobleme: Klassifikation und Regression
- Modelle und Hypothesenraum
- **Verlustfunktionen und Regularisierer**
- Unsicherheit, Wahrscheinlichkeiten, Bayes'sche Regel

Verlustfunktion, Optimierungskriterium

- Alternative zu Version Spaces: Lernprobleme werden als Optimierungsprobleme formuliert.
 - ◆ *Verlustfunktion* misst, wie gut Modell zu Trainingsdaten passt
 - ◆ *Regularisierungsfunktion* misst, ob das Modell nach unserem Vorwissen *wahrscheinlich* ist.
 - ◆ *Optimierungskriterium* ist Summe aus Verlust und Regularisierer.
 - ◆ Suche Minimum des Optimierungskriteriums
- Insgesamt wahrscheinlichstes Modell, gegeben Trainingsdaten und Vorwissen.

Verlustfunktion

- Wie schlimm ist es, wenn Modell $f(\mathbf{x}_i)$ vorhersagt obwohl der echte Wert der Zielvariable y_i ist?

$$l(f(\mathbf{x}_i), y_i)$$

- Verlust auf den ganzen Trainingsdaten L :

$$\sum_{i=1}^N l(f(\mathbf{x}_i), y_i)$$

- Beispiel: Binäres Klassifikationsproblem mit positiver Klasse (+1) und negativer Klasse (-1). False Positives und False Negatives gleich schlimm.

- ◆ Zero-One Loss: $l(f(\mathbf{x}_i), y_i) = \begin{cases} 0, & \text{wenn } f(\mathbf{x}_i) = y_i \\ 1, & \text{sonst} \end{cases}$

Verlustfunktion

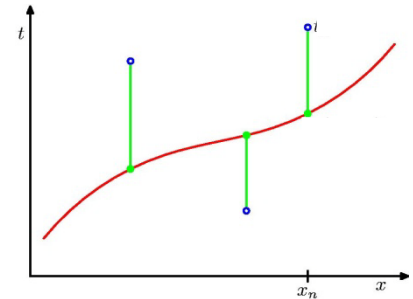
- Beispiel: diagnostische Klassifikationsprobleme, übersehene Erkrankungen (False Negatives) schlimmer als False Positives.
 - ◆ Kostenmatrix

$$l(f(\mathbf{x}_i), y_i) = \begin{cases} f(\mathbf{x}_i) = +1 & \begin{array}{|c|c|} \hline y_i = +1 & y_i = -1 \\ \hline 0 & c_{FP} \\ \hline \end{array} \\ f(\mathbf{x}_i) = -1 & \begin{array}{|c|c|} \hline y_i = +1 & y_i = -1 \\ \hline c_{FN} & 0 \\ \hline \end{array} \end{cases}$$

Verlustfunktion

- Beispiel Verlustfunktion Regression: Vorhersage möglichst dicht an echtem Wert des Zielattributes
 - ◆ Quadratischer Fehler

$$l(f(\mathbf{x}_i), y_i) = (f(\mathbf{x}_i) - y_i)^2$$



Verlustfunktion

- Wie schlimm ist es, wenn Modell $f(\mathbf{x}_i)$ vorhersagt obwohl der echte Wert der Zielvariable y_i ist?
 - ◆ Verlust $l(f(\mathbf{x}_i), y_i)$.
- Verlustfunktion ist aus der jeweiligen Anwendung heraus motiviert.



Regularisierer

- Verlustfunktion drückt aus, wie gut Modell zu Daten passt
- Regularisierer:
 - ◆ drückt Annahme darüber aus, ob Modell *a priori* wahrscheinlich ist.
 - ◆ Unabhängig von den Trainingsdaten.
 - ◆ Je höher der Regularisierungsterm für ein Modell, desto unwahrscheinlicher
- Häufig wird die Annahme ausgedrückt, dass wenige der Attribute für ein gutes Modell ausreichen.
 - ◆ Anzahl der Attribute, L_0 -Regularisierung
 - ◆ Betrag der Attribut-Gewichtungen, L_1 -Regularisierung
 - ◆ Quadrat der Attribut-Gewichtungen, L_2 -Regularisierung.

Regularisierer: Beispiel

- Hypothesenraum: Konjunktion von Bedingungen

- ◆ $f(\mathbf{x}) = \begin{cases} \text{☹️}, & \text{wenn } x_1 = 1 \wedge x_3 = 1 \wedge x_7 = 1 \\ \text{😊}, & \text{sonst} \end{cases}$

- Lineares Modell: Lässt sich schreiben als

- ◆ $f(\mathbf{x}) = \begin{cases} \text{☹️}, & \text{wenn } x_1 + x_3 + x_7 \geq 3 \\ \text{😊}, & \text{sonst} \end{cases}$

- Allgemein: äquivalente Darstellung ist

- ◆ $f_{\mathbf{w}}(\mathbf{x}) = \begin{cases} \text{☹️}, & \text{wenn } \sum_{j=1}^m w_j x_j \geq b \\ \text{😊}, & \text{sonst} \end{cases}$
 $= \begin{cases} \text{☹️}, & \text{wenn } \mathbf{w}^T \mathbf{x} \geq b \\ \text{😊}, & \text{sonst} \end{cases}$

w: Modellparameter

$$w_j \in \{-1, 0, +1\}$$

$w_j \in \{-1, +1\}$ falls Attribut in logischer Bedingung vorkommt

Regularisierer: Beispiel

- Linearer Klassifikator

- ◆ $f_{\mathbf{w}}(\mathbf{x}) = \begin{cases} \text{☹️}, & \text{wenn } \mathbf{w}^T \mathbf{x} \geq b \\ \text{😊}, & \text{sonst} \end{cases}$

- L_2 -Regularisierung:

- ◆ $\lambda |\mathbf{w}|^2 \quad |\mathbf{w}|^2 = \sum_i w_i^2$

- ◆ Addiert λ für jedes von null verschiedene Gewicht.

- Optimierungskriterium: Verlust+Regularisierer

- ◆ $\hat{R}(\mathbf{w}, L) = \sum_i l(f_{\mathbf{w}}(\mathbf{x}_i), y_i) + \lambda |\mathbf{w}|^2$

- ◆ Parameter λ steuert Stärke des Regularisierers

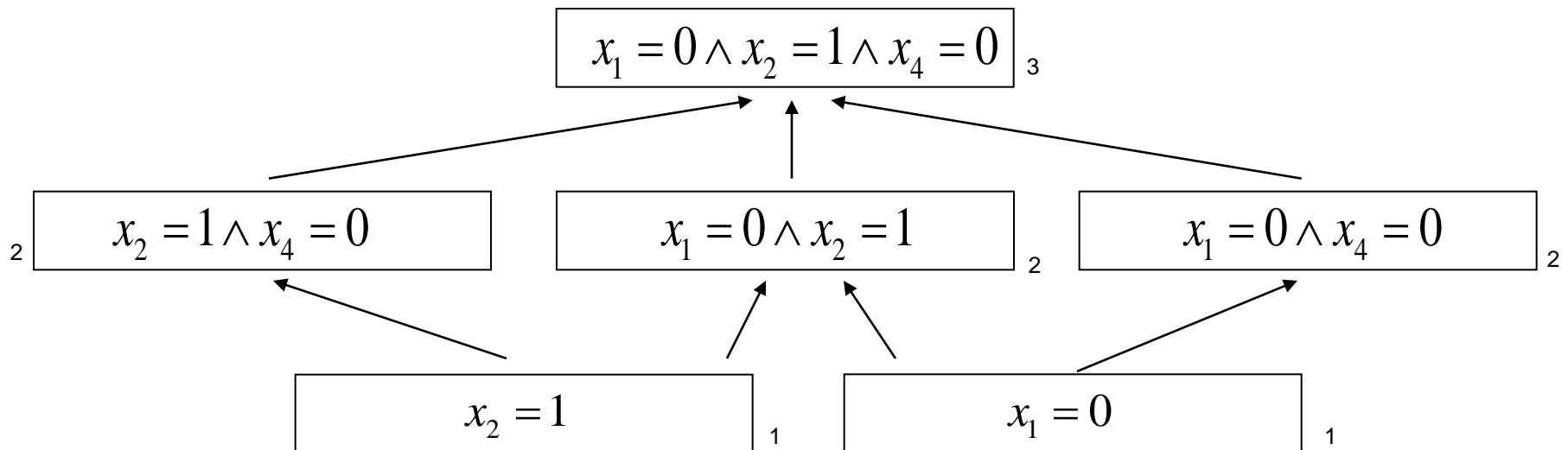
- Durch den Regularisierer implementierte Präferenz des Lernalgorithmus wird auch *Inductive Bias* genannt.

Optimierungsproblem: Beispiel

- $\hat{R}(\mathbf{w}, L) = \sum_i l(f_{\mathbf{w}}(\mathbf{x}_i), y_i) + \lambda |\mathbf{w}|^2$

- Beste Hypothese für $\lambda = 0.1$?

	x_1	x_2	x_3	x_4	x_5	x_6	y
\mathbf{x}_1	0	1	0	0	1	1	☹️
\mathbf{x}_2	0	1	1	0	1	1	☹️
\mathbf{x}_3	1	0	1	0	1	0	😊
\mathbf{x}_4	0	1	1	0	0	0	☹️

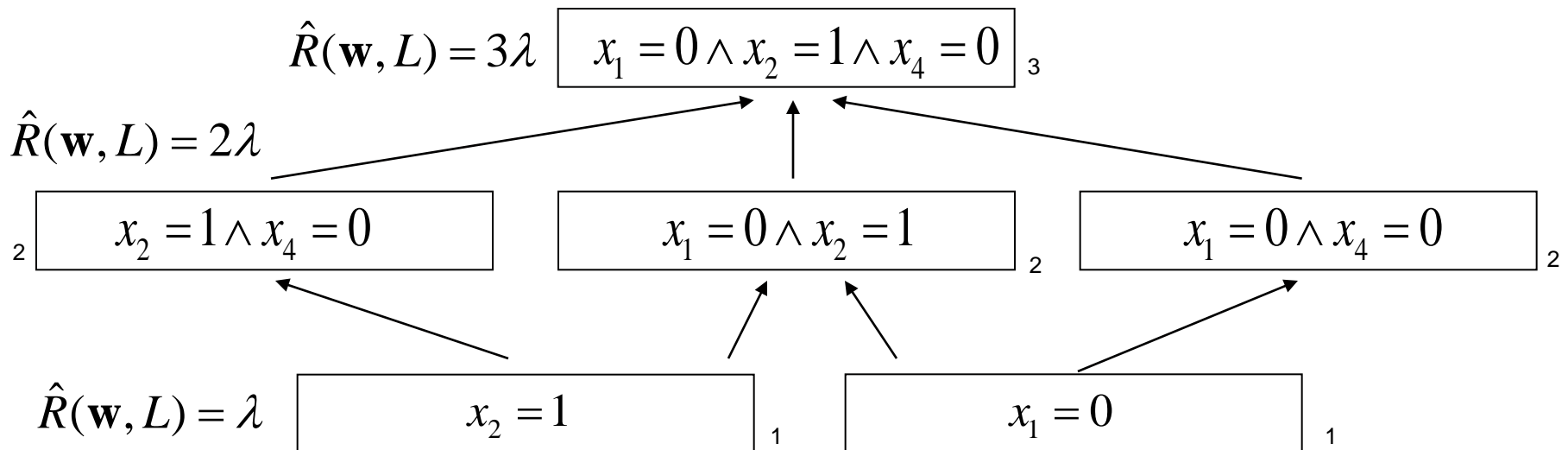


Optimierungsproblem: Beispiel

- $\hat{R}(\mathbf{w}, L) = \sum_i l(f_{\mathbf{w}}(\mathbf{x}_i), y_i) + \lambda \|\mathbf{w}\|^2$

- Beste Hypothese für $\lambda = 0.1$?

	x_1	x_2	x_3	x_4	x_5	x_6	y
\mathbf{x}_1	0	1	0	0	1	1	☹️
\mathbf{x}_2	0	1	1	0	1	1	☹️
\mathbf{x}_3	1	0	1	0	1	0	😊
\mathbf{x}_4	0	1	1	0	0	0	☹️



Optimierungsproblem

- Einstellung von λ ?
- Rechtfertigung für Optimierungskriterium?
- Mehrere Rechtfertigungen und Herleitungen.
 - ◆ **Wahrscheinlichste Hypothese (MAP-Hypothese).**
 - ◆ Hypothese, die Daten am stärksten komprimiert (*Minimum Description Length*).
 - ◆ Niedrige obere Schranke für Fehler auf zukünftigen Daten abhängig von $|\mathbf{w}|$. (*SRM*).
- Lernen ohne Regularisierung ist *ill-posed* Problem; keine eindeutige Lösung, oder Lösung hängt extrem stark von minimalen Änderungen in den Daten ab.

Überblick

- Lernprobleme: Klassifikation und Regression
- Modelle und Hypothesenraum
- Verlustfunktionen und Regularisierer
- Unsicherheit, Wahrscheinlichkeiten, Bayes'sche Regel