

Universität Potsdam
Institut für Informatik
Lehrstuhl Maschinelles Lernen



Maschinelles Lernen Modelle, Version Spaces, Lernen

Christoph Sawade/Niels Landwehr
Jules Rasetaharison
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

Klassifikation

- Eingabe: Instanz (Objekt) $\mathbf{x} \in X$.
 - ◆ Objekte oft durch Vektor von *Attributen* repräsentiert
 - ◆ 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}	
Medikament 1 enthalten?	$\begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ 0 \end{pmatrix}$	
⋮		
Medikament 6 enthalten?		

Belegung der Attribute, Merkmalsvektor

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

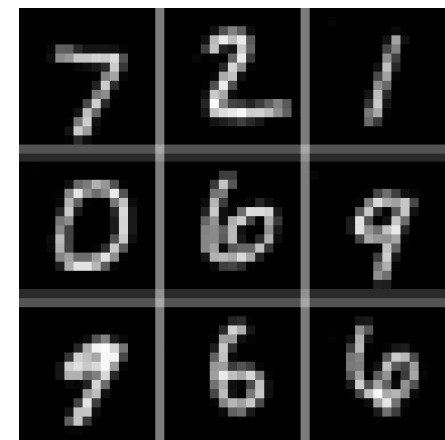


Klassifikation: Beispiel

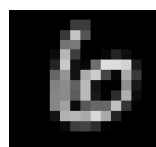
- Eingabe: Instanz (Objekt) $\mathbf{x} \in X$.

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



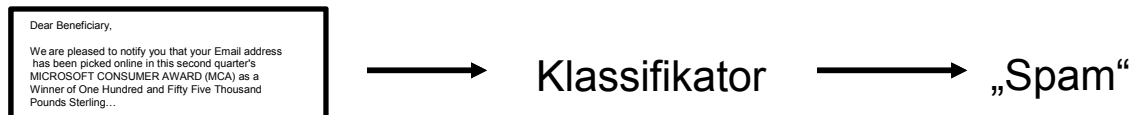
Klassifikation: Beispiel

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

X = Menge aller möglichen Email-Texte

Attribute	Instanz \mathbf{x}	Email	
Wort 1 kommt vor?	$\begin{pmatrix} 0 \\ 1 \\ 0 \\ \dots \\ 1 \\ 0 \end{pmatrix}$	<div style="border: 1px solid black; padding: 5px;">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...</div>	
...			Address
...			Beneficiary
...			Friend
...			...
Wort N kommt vor?	1	Sterling	
$N \approx 100000$	0	Science	

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



Klassifikationslernen

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

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

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

◆ \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)

Klassifikationslernen

- Eingabe Lernproblem: Trainingsdaten.

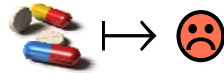
$$\diamond L = \langle (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \rangle$$

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



- Ausgabe: Klassifikator (auch als Modell bezeichnet).

$$f : X \rightarrow Y$$



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

Klassifikationslernen

- Eingabe Lernproblem: Trainingsdaten.

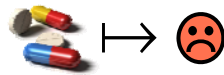
$$\diamond L = \langle (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \rangle$$

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



- Ausgabe: Klassifikator (auch als Modell bezeichnet).

$$f : X \rightarrow Y$$



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

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

Linearer Klassifikator mit
Parametervektor \mathbf{w} .

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

Klassifikationslernen

- Eingabe Lernproblem: Trainingsdaten.

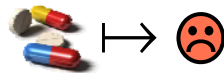
$$\diamond L = \langle (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \rangle$$

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



- 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

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

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

Regression

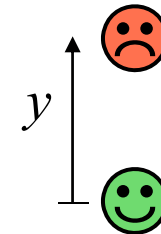
- Eingabe: Instanz (Objekt) $\mathbf{x} \in X$.
 - ◆ Objekte oft durch Attribut-Vektoren repräsentiert.
 - ◆ Instanz ist Belegung der Attribute.

- ◆ $\mathbf{x} = \begin{pmatrix} x_1 \\ \dots \\ x_m \end{pmatrix}$ 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 Ausgaberräumen.
- 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**

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

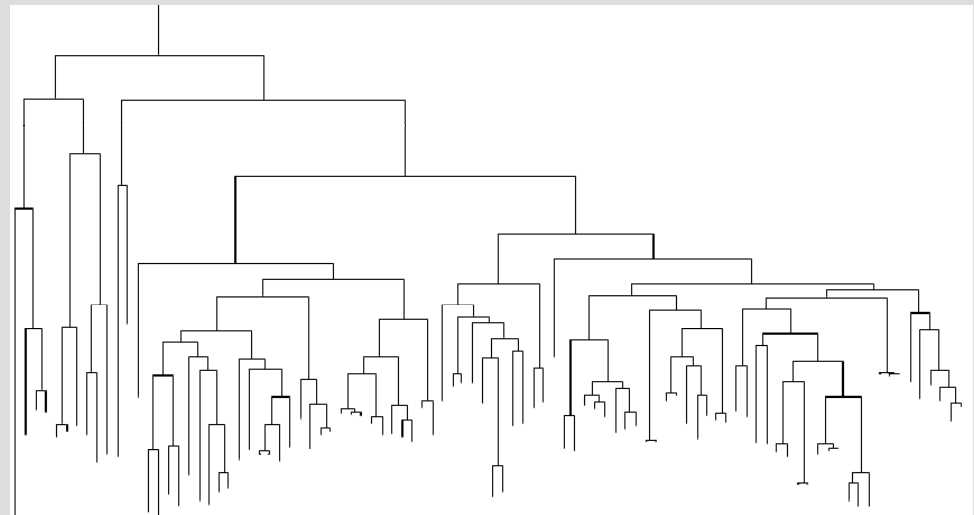
- Ordinale Regression.
- Präferenzlernen.
- Taxonomie-Klassifikation.
- Klassifikation und Regression mit strukturierten

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.

- Kollaboratives Lernen

- Eingabe X und Ausgabe Y
strukturierte Räume

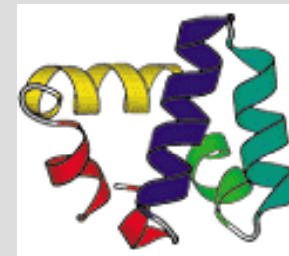
Beispiel:

Eingabe DNA

Ausgabe Proteinfaltung

Klassenlabel 3D Struktur

...AAGCTTGCACCTGCCGT...



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
[Add to Cart](#) [Add to Wish List](#)

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

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

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

Suche nach Hypothese

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

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

Trainingsdaten

Medikamente in der Kombination

Beispiel-Kombinationen

	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

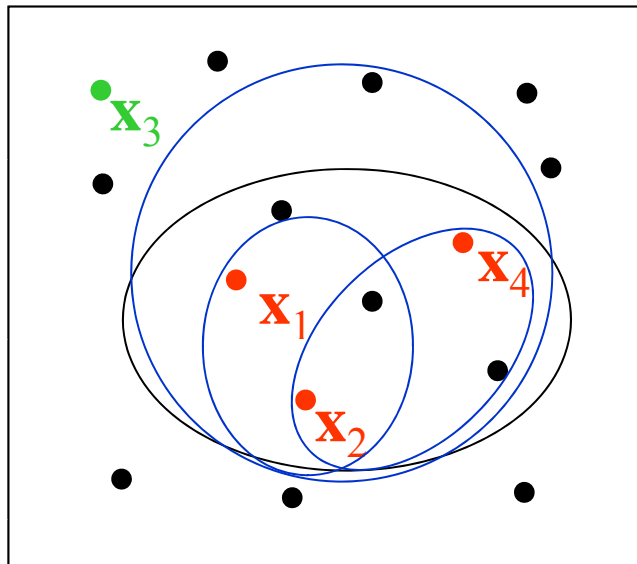
$$\forall i: f(\mathbf{x}_i) = y_i$$

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

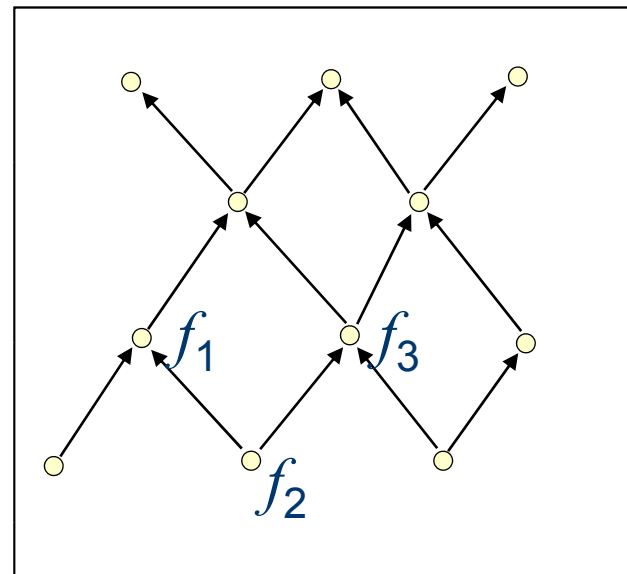
„Genereller-Als“-Ordnung

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

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 $x_2=1, x_3=1$

$f_2 \geq_g f_1$ $f_2 \geq_g f_3$ aber nicht $f_1 \leq_g / \geq_g f_3$

Version Space

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

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

- 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

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

- 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

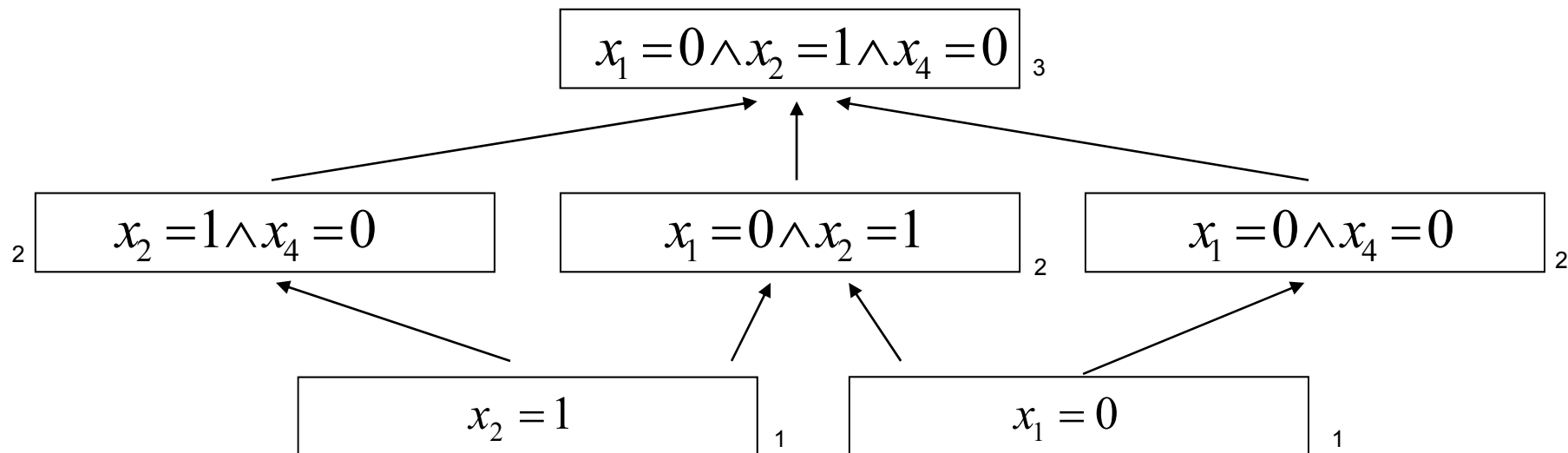
■ H: $f(\mathbf{x}) = \begin{cases} \text{☹️}, & \text{wenn } \bigwedge_{j \in J} x_j = v_j \\ \text{😊}, & \text{sonst} \end{cases}$

- 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.
- ◆ Jedes Element des Version Space könnte die Daten erzeugt haben / die korrekte Hypothese sein

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	☹️

Beispiel-Kombinationen

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

- Verlustfunktionen: Grad der Konsistenz mit Trainingsdaten
- A-Priori-Wahrscheinlichkeit (Prior) über Modelle.
- Wahrscheinlichstes Modell gegeben Daten.

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: Klassifikationsprobleme, 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

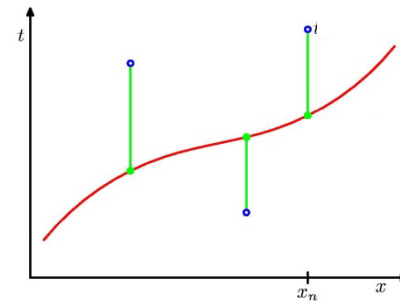
- Wie schlimm ist es, wenn Modell $f(\mathbf{x}_i)$ vorhersagt obwohl der echte Wert der Zielvariable y_i ist?
- Beispiel: diagnostische Klassifikationsprobleme, übersehene Erkrankungen (False Negatives) schlimmer als False Positives.
 - ◆ Kostenmatrix

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

Verlustfunktion

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

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

- ◆ mit $w_j \in \{-1, 0, +1\}$

w: Modellparameter

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

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

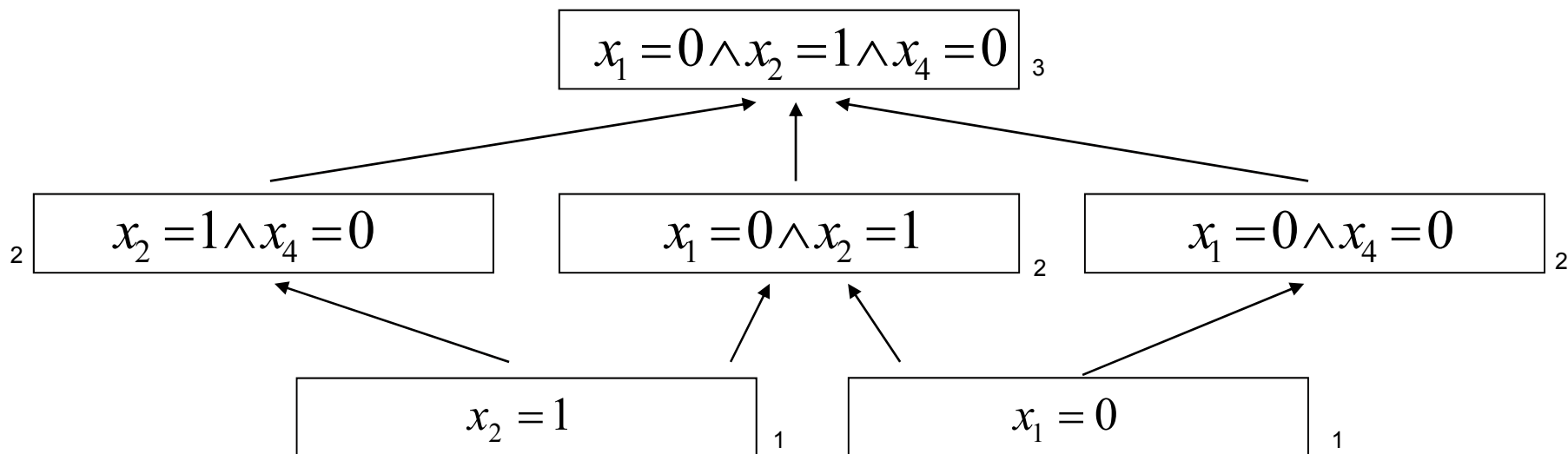
- ◆ Parameter λ steuert Stärke des Regularisierers

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

Optimierungsproblem: Beispiel

- $\hat{R}(\mathbf{w}, L) = \sum_i l(f(\mathbf{x}_i), y_i) + \lambda |\mathbf{w}|^2$
- Beste Hypothese für $\lambda = 0.1$?

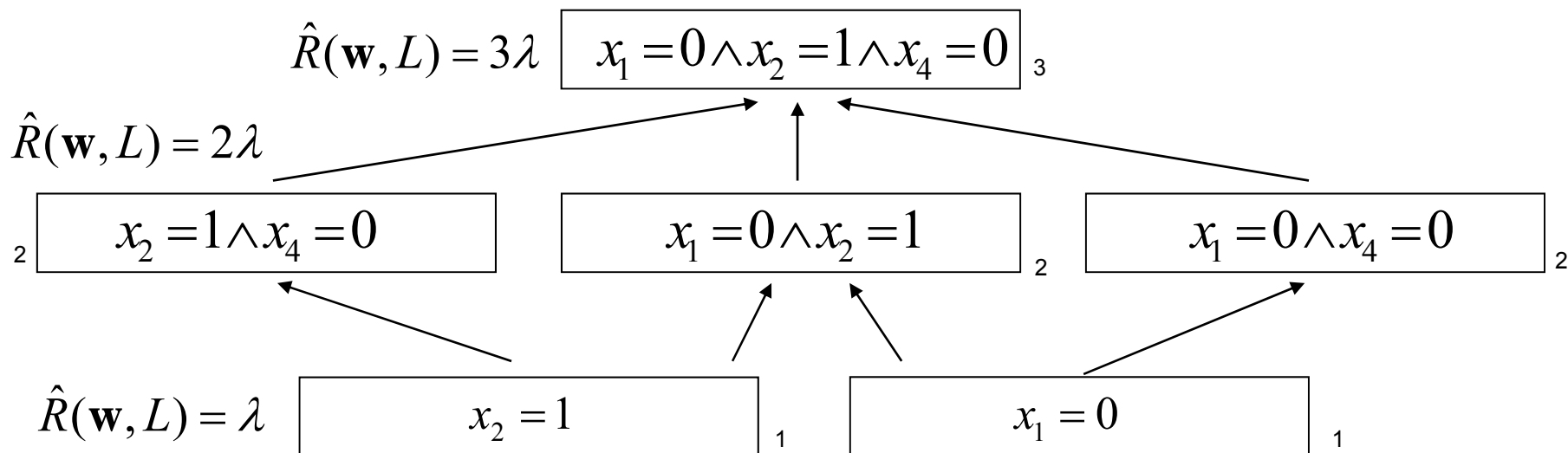
	x_1	x_2	x_3	x_4	x_5	x_6	y
\mathbf{x}_1	0	1	0	0	1	1	☹️
\mathbf{x}_2	0	1	1	0	1	1	☹️
\mathbf{x}_3	1	0	1	0	1	0	😊
\mathbf{x}_4	0	1	1	0	0	0	☹️



Optimierungsproblem: Beispiel

- $\hat{R}(\mathbf{w}, L) = \sum_i l(f(\mathbf{x}_i), y_i) + \lambda |\mathbf{w}|^2$
- 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; Lösung existiert manchmal nicht oder 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

Unsicherheit und Wahrscheinlichkeiten

- Viele Verfahren des maschinellen Lernens basieren auf probabilistischen Überlegungen
- Modellvorstellung beim Lernen:
 - ◆ Jemand hat echtes f^* nach A-Priori-Wahrscheinlichkeit („Prior“) $p(f)$ gezogen.
 - ◆ f^* ist nicht bekannt, aber $p(f)$ reflektiert Vorwissen (was sind wahrscheinliche Modelle?).
 - ◆ Trainingseingaben \mathbf{x}_i werden gezogen.
 - ◆ Klassenlabels y_i werden nach $p(y_i | \mathbf{x}_i, f^*)$ gezogen.
 - ◆ Fragestellung Lernen: Gegeben L und $p(f)$, was ist wahrscheinlichstes „echte“ Modell ?
→ Versuche, f^* (ungefähr) zu rekonstruieren

Wahrscheinlichkeitstheorie

- Zufallsexperiment: definierter Prozess, in dem ein Elementarereignis ω erzeugt wird.
- Ereignisraum Ω : Menge aller möglichen Elementarereignisse.
- Ereignis A : Teilmenge des Ereignisraums.
- Wahrscheinlichkeitsfunktion p : Funktion, die Wahrscheinlichkeitsmasse auf Ereignisse $A \subseteq \Omega$ verteilt.

Wahrscheinlichkeitstheorie

- Gültige Wahrscheinlichkeitsfunktion p (Kolmogorow-Axiome)
 - ◆ Wahrscheinlichkeit von Ereignis $A \subseteq \Omega$: $0 \leq p(A) \leq 1$
 - ◆ Sicheres Ereignis: $p(\Omega) = 1$, und $p(\emptyset) = 0$
 - ◆ Für die Wahrscheinlichkeit zweier inkompatibler Ereignisse $A \subseteq \Omega, B \subseteq \Omega$ (d.h. $A \cap B = \emptyset$) gilt:

$$p(A \cup B) = p(A) + p(B)$$

Wahrscheinlichkeitstheorie: Beispiel

■ Würfeln

- ◆ Ereignisraum $\Omega = \{1, 2, 3, 4, 5, 6\}$
- ◆ Elementarereignisse haben Wsk $p(\{\omega\}) = 1/6$
- ◆ Ereignis gerade Zahl: $A = \{2, 4, 6\}$
- ◆ Wahrscheinlichkeit des Ereignisses: $p(A) = 1/2$

Zufallsvariable

- Zufallsvariable X : Abbildung eines elementaren Ereignisses auf einen numerischen Wert

$$X : \Omega \rightarrow \mathbb{R}$$

$$X : \omega \mapsto x$$

- Wahrscheinlichkeit, dass ZV einen Wert annimmt

$$p(X = x) = p(\{\omega \in \Omega \mid X(\omega) = x\})$$

- Bei kontinuierlichen ZV oft $p(X = x) = 0$, man betrachtet dann auch Verteilungsfunktion

$$p(X \leq x) = p(\{\omega \in \Omega \mid X(\omega) \leq x\})$$

Zufallsvariable: Beispiel

- Würfeln mit 2 Würfeln

- ◆ Ereignisraum $\Omega = \{(\omega_1, \omega_2) \mid \omega_i \in \{1, 2, 3, 4, 5, 6\}\}$
- ◆ Elementarereignisse haben Wsk $p(\{(\omega_1, \omega_2)\}) = 1/36$
- ◆ Zufallsvariable: Summe der beide Würfel

$$X((\omega_1, \omega_2)) = \omega_1 + \omega_2$$

- ◆ Wahrscheinlichkeit für Wert der ZV:

$$p(X = 5) = ?$$

Zufallsvariable: Beispiel

■ Würfeln mit 2 Würfeln

- ◆ Ereignisraum $\Omega = \{(\omega_1, \omega_2) \mid \omega_i \in \{1, 2, 3, 4, 5, 6\}\}$
- ◆ Elementarereignisse haben Wsk $p(\{(\omega_1, \omega_2)\}) = 1/36$
- ◆ Zufallsvariable: Summe der beide Würfel

$$X((\omega_1, \omega_2)) = \omega_1 + \omega_2$$

- ◆ Wahrscheinlichkeit für Wert der ZV:

$$\begin{aligned} p(X = 5) &= p(\{(1, 4), (2, 3), (3, 2), (4, 1)\}) \\ &= 4 / 36 \end{aligned}$$

Gemeinsame und bedingte Wahrscheinlichkeiten

- Oft gibt es mehrere Zufallsvariablen X, Y, Z, \dots
- Gemeinsame Wahrscheinlichkeit:

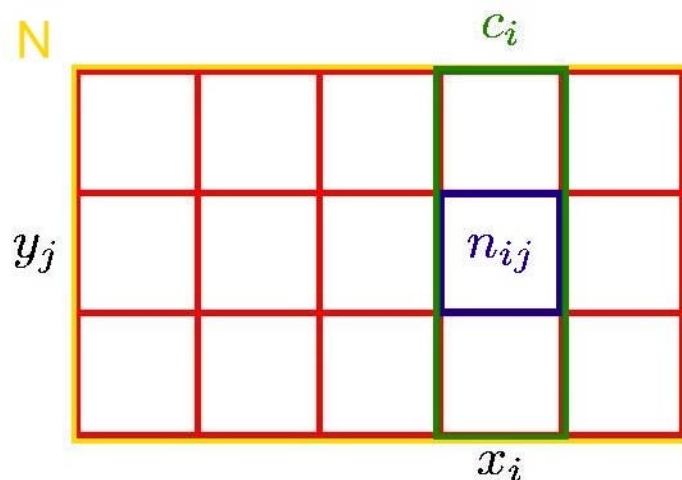
$$p(X = x, Y = y) = p(\{\omega \in \Omega \mid X(\omega) = x \wedge Y(\omega) = y\})$$

- Bedingte Wahrscheinlichkeit

$$p(X = x \mid Y = y) = \frac{p(X = x, Y = y)}{p(Y = y)}$$

Gemeinsame und bedingte Wahrscheinlichkeiten

„Zufälligen Punkt auf Fläche auswählen“



Gemeinsame Wahrscheinlichkeit

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$$

Randwahrscheinlichkeit

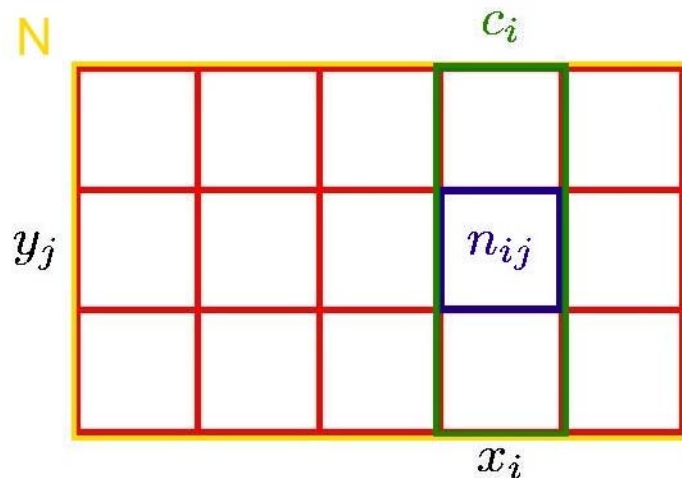
$$p(X = x_i) = \frac{c_i}{N}$$

Bedingte Wahrscheinlichkeit

$$\begin{aligned} p(Y = y_j | X = x_i) &= \frac{n_{ij}}{c_i} \\ &= \frac{p(Y = y_j, X = x_i)}{p(X = x_i)} \end{aligned}$$

Wichtige Rechenregeln für Wahrscheinlichkeiten

„Zufälligen Punkt auf Fläche auswählen“



Summenregel

$$p(X = x_i) = \sum_j p(X = x_i, Y = y_j)$$

$$p(X = x_i) = \frac{c_i}{N} = \frac{1}{N} \sum_j n_{ij}$$
$$= \sum_j p(X = x_i, Y = y_j)$$

Produktregel

$$p(Y = y_j, X = x_i) = p(Y = y_j | X = x_i) p(X = x_i)$$

$$p(Y = y_j, X = x_i) = \frac{n_{ij}}{N} = \frac{n_{ij}}{c_i} \frac{c_i}{N} = p(Y = y_j | X = x_i) p(X = x_i)$$

Wichtige Rechenregeln für Wahrscheinlichkeiten

- Summenregel $p(X) = \sum_Y p(X, Y)$
- Produktregel $p(X, Y) = p(Y|X)p(X)$

Unabhängigkeit von Zufallsvariablen

- Variablen X, Y unabhängig gdw

$$p(X = x_i, Y = y_j) = p(X = x_i)p(Y = y_j)$$

- Äquivalente Definition Unabhängigkeit:

$$p(X = x_i | Y = y_j) = p(X = x_i) \quad \text{und} \quad p(Y = y_j | X = x_i) = p(Y = y_j)$$

- Bedingte Unabhängigkeit

$$p(X = x_i, Y = y_i | Z = z_k) = p(X = x_i | Z = z_k)p(Y = y_i | Z = z_k)$$

Bayessche Regel

- Bayessche Regel:

$$p(X | Y) = \frac{p(Y | X)p(X)}{p(Y)}$$

- Beweis einfach:

$$p(X | Y) = \frac{p(X, Y)}{p(Y)} = \frac{p(Y | X)p(X)}{p(Y)}$$

Definition bedingte Wahrscheinlichkeit

Produktregel

- Wichtige Grundeinsicht für das maschinelle Lernen: Erlaubt den Rückschluss auf *Modellwahrscheinlichkeiten* gegeben *Wahrscheinlichkeiten von Beobachtungen*

Bayessche Regel

- Modellwahrscheinlichkeit gegeben Daten und Vorwissen

$$p(\text{Modell} | \text{Daten}) = \frac{p(\text{Daten} | \text{Modell}) p(\text{Modell})}{p(\text{Daten})}$$

$p(\text{Daten})$ konstant,
unabhängig von *Modell*

$$\propto p(\text{Daten} | \text{Modell}) p(\text{Modell})$$

Likelihood: wie gut erklärt
Modell die Daten?

Prior: wie wahrscheinlich
ist Modell a priori?

Maximum-A-Posteriori-Hypothese

- Wahrscheinlichstes Modell gegeben die Daten.

- ◆ $f_{MAP} = \arg \max_{f_w} p(f_w | L)$

$$= \arg \max_{f_w} \frac{p(L|f_w)p(f_w)}{p(L)}$$

Anwendung
Bayes'sche Regel

$$= \arg \max_{f_w} p(L|f_w)p(f_w)$$

$$= \arg \min_{f_w} -\log P(L|f_w) - \log p(f_w)$$

Log-Likelihood

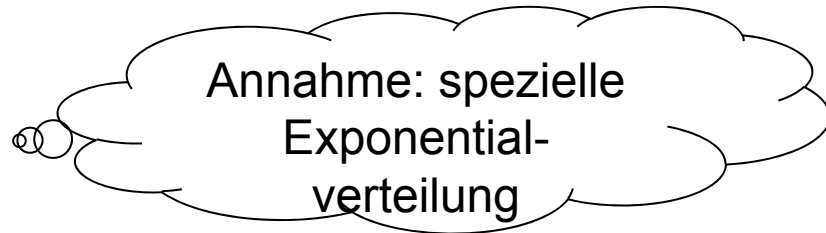
Log-Prior

⇒ Optimierungskriterium bestehend aus log-likelihood und log-prior

w Parameter des Modells $f_w(\mathbf{x})$

Log-Likelihood

- Wie wahrscheinlich sind die Daten gegeben das Modell?
 - ◆ $-\log p(L|f_w) = -\log p(y_1, \dots, y_N | f_w, \mathbf{x}_1, \dots, \mathbf{x}_N)$
- Annahme: Datenpunkte sind unabhängig gezogen.
 - ◆ $-\log p(y_1, \dots, y_N | f_w, \mathbf{x}_1, \dots, \mathbf{x}_N)$
 $= -\log \prod_i p(y_i | f_w, \mathbf{x}_i)$
 $= \sum_i -\log p(y_i | f_w, \mathbf{x}_i)$



Log-Likelihood

- Wie wahrscheinlich sind die Daten gegeben das Modell?
 - ◆ $-\log p(L|f_w) = -\log p(y_1, \dots, y_N | f_w, \mathbf{x}_1, \dots, \mathbf{x}_N)$
- Annahme: Datenpunkte sind unabhängig gezogen.

- ◆ $-\log p(y_1, \dots, y_N | f_w, \mathbf{x}_1, \dots, \mathbf{x}_N)$

$$= -\log \prod_i p(y_i | f_w, \mathbf{x}_i)$$

$$= \sum_i -\log p(y_i | f_w, \mathbf{x}_i)$$

Annahme: spezielle Exponentialverteilung

$$p(y_i | f_w, \mathbf{x}_i) = \frac{1}{Z} e^{-l(f(\mathbf{x}_i), y_i)}$$

Verlustfunktion

Normalisierer

$$l(f(\mathbf{x}_i), y_i) = \begin{cases} f(\mathbf{x}_i) = +1 & \begin{matrix} y_i = +1 & y_i = -1 \\ 0 & c \end{matrix} \\ f(\mathbf{x}_i) = -1 & \begin{matrix} c & 0 \end{matrix} \end{cases}$$

$$c, \bar{c} \neq 0$$

Log-Likelihood

- Wie wahrscheinlich sind die Daten gegeben das Modell?
 - ◆ $-\log p(L|f_w) = -\log p(y_1, \dots, y_N | f_w, \mathbf{x}_1, \dots, \mathbf{x}_N)$
- Annahme: Datenpunkte sind unabhängig gezogen.

$$-\log p(y_1, \dots, y_N | f_w, \mathbf{x}_1, \dots, \mathbf{x}_N)$$

$$= -\log \prod_i p(y_i | f_w, \mathbf{x}_i)$$

$$= \sum_i -\log p(y_i | f_w, \mathbf{x}_i)$$

$$= \sum_i -\log \frac{e^{-l(f(\mathbf{x}_i), y_i)}}{Z}$$

$$= \sum_i l(f_w(\mathbf{x}_i), y_i) + \text{const}$$

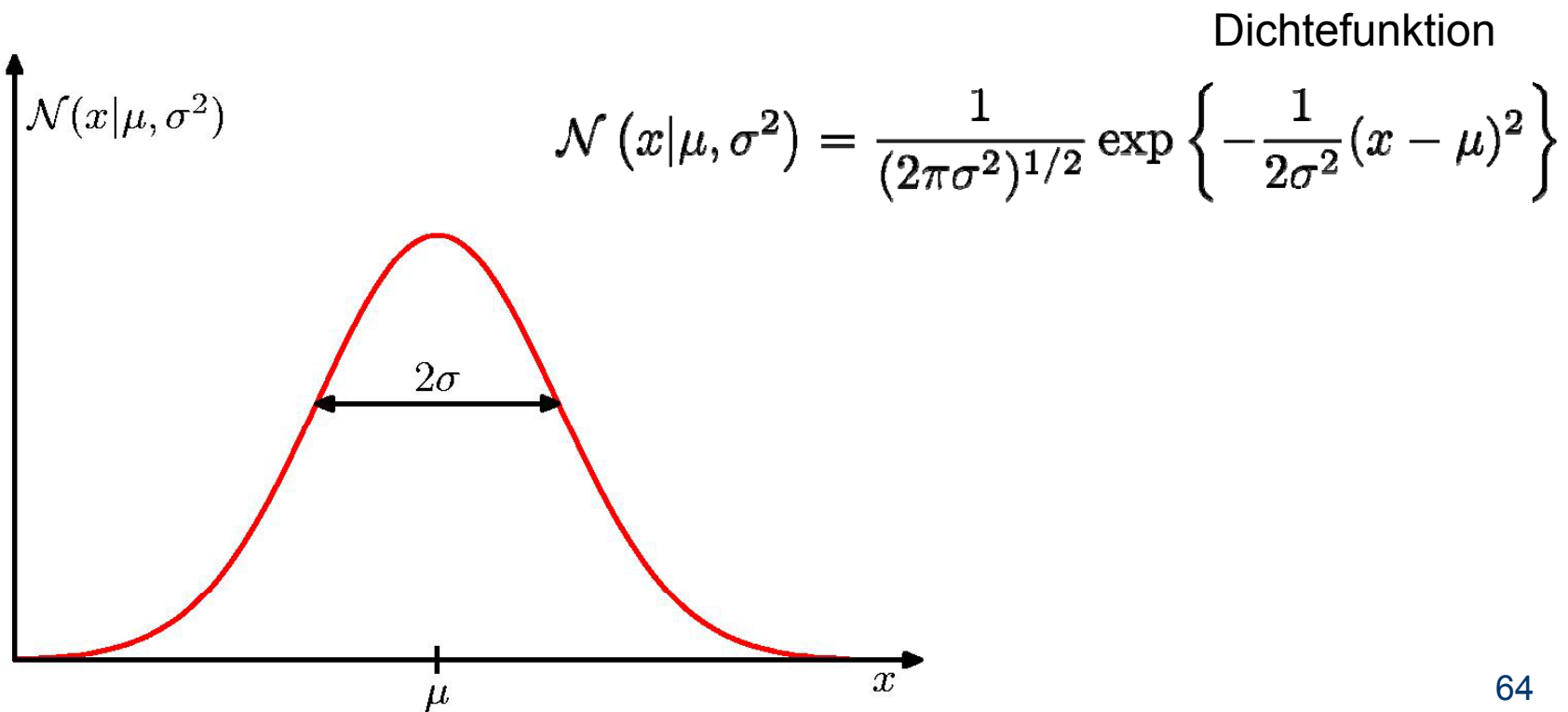
Konstanter Faktor
(unabhängig von f)

- Negative Log-Likelihood entspricht Verlustterm!

Prior: Gaußverteilung, Normalverteilung

Verteilung für reelle (oder vektor-wertige) Zufallsvariable x

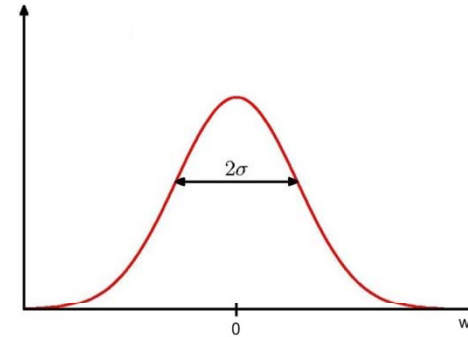
Streut um Mittelwert μ mit Varianz σ^2



A-Priori-Wahrscheinlichkeit (Prior)

- Annahme: Modellparameter normalverteilt mit Mittelwert $\mu = 0$

- ◆
$$p(f_{\mathbf{w}}) = \mathcal{N}(\mathbf{w} | 0, \sigma^2)$$
$$= \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}|\mathbf{w}|^2}$$



- Negativer Log-Prior:

- ◆
$$-\log p(f_{\mathbf{w}}) = \frac{1}{2\sigma^2} |\mathbf{w}|^2 + const$$

Konstanter Faktor
(unabhängig von f)

- Negativer Log-Prior = Regularisierer!

A-Posteriori-Wahrscheinlichkeit (Posterior)

- Wahrscheinlichstes Modell gegeben Vorwissen und Daten.

$$\begin{aligned}\diamond f_{MAP} &= \arg \max_{f_{\mathbf{w}}} p(f_{\mathbf{w}}|L) \\ &= \arg \min_{f_{\mathbf{w}}} -\log p(L|f_{\mathbf{w}}) - \log p(f_{\mathbf{w}}) \\ &= \arg \min_{f_{\mathbf{w}}} \sum_i l(f_{\mathbf{w}}(\mathbf{x}_i), y_i) + \lambda \|\mathbf{w}\|^2\end{aligned}$$

(Folie 60 zusammen mit 63 & 65, $\lambda = \frac{1}{2\sigma^2}$)

- ArgMin über regularisierte Verlustfunktion!
- Rechtfertigung für Optimierungskriterium?
 - ◆ **Wahrscheinlichste Hypothese (MAP-Hypothese).**

Zusammenfassung

- Klassifikation, Regression, Modelle für Vorhersage.
- Parametrisierung von Modellen (Parametervektor \mathbf{w})
- Hypothesenraum: Raum aller Modelle.
- Version Space: Modelle, die mit Trainingsdaten konsistent sind.
- Lernen (Parameterschätzung) aus Trainingsdaten.
 - ◆ Optimierungskriterium: Verlust über Trainingsbeispiele plus Regularisierungsterm.
 - ◆ Minimum des Kriteriums: wahrscheinlichstes Modell gegeben Daten und Prior.

Übungen

- Erstes Übungsblatt: Ausgabe morgen, Besprechung kommende Woche
- Sie können für einzelne Aufgaben votieren
- Sie müssen für 2/3 aller Aufgaben votieren