

Universität Potsdam
Institut für Informatik
Lehrstuhl Maschinelles Lernen



Textklassifikation, Informationsextraktion

Tobias Scheffer
Thomas Vanck

Textklassifikation, Informationsextraktion

PROSAR-AIDA

File Edit Options View Window ?

0 - [H:\Doku\Intern\Tabellen\Demos\Rechnungslesungl...

GLOBE LTD.

Globe Ltd. World Retail
Mars House
Leafield Way
Corsham, Wiltshire
SN13 9SW

Orders: 01483 8786545
Fax: 01483 8786425
order@world.co.uk

Taxpoint Date: 26/09/02
Invoice Number: 23598
Your Order: 68974
Please refer on all payments

INVOICE

Paradatec Ltd.
Oban House, Rope Yard
Wootton Bassett, Wiltshire
SN4 7BW

Pos.	Description	Qty.	Price	Value
Purchase order No. 4510425457				
01	4,000 Pcs Neon Light Bulb Material# 0124 Unit price 3,70			14,80
02	2,000 Pcs Heating Element NiChrome Material# 0453 Unit price 33,44			66,88
03	1,000 Pcs Hight Output LED Line (blue) Material# 0922 Unit price 12,45			12,45
04	8,000 Pcs Halogen Lamp Fixtures Chrome Material# 0765 Unit price 2,78			22,24
05	1,000 Pcs Transformer 12V Dual Purpose with Enhanced Screening Material# 0329 Unit price 22,95			22,95
06	2,000 Pcs Fuse Material# 0078 Unit price 0,75			1,50
Sub-total				140,82

Globe Ltd. World Retail, Mars House, Leafield Way, Corsham, Wiltshire SN13 9SW
VAT registration number 534 2342 28

Results (primary)

Search objects INVTABLE

	POS	ITEM	QUANTITY	PRICE	TOTAL
1	01	0124	4,000	3,70	14,80
2	02	0453	2,000	33,44	66,88
3	03	0922	1,000	12,45	12,45
4	04	0765	8,000	2,78	22,24
5	05	0329	1,000	22,95	22,95
6	06	0078	2,000	0,75	1,50

PROSAR-AIDA

Image #4
Process page?

OK Abbrechen

0: Page: Processing image file "H:\Doku\Intern\Tabellen\Demos\Rechnungslesung\Images\Rechnungsdemo_new.tif" #4

Pause 1543,618

Textklassifikation, Informationsextraktion

- Textklassifikator: Ordnet einen Text einer Menge von inhaltlichen Kategorien zu.
- Wird häufig aus Prozessdaten gelernt.
- Anwendungsbeispiel: Posteingangsverarbeitung.

- Informationsextraktion: Identifikation definierter Felder in Dokument.
- Wird häufig aus Prozessdaten gelernt.
- Anwendungsbeispiel: Automatisierung von Dokumentenverarbeitungsprozessen.

Textklassifikation: Repräsentation

- Nach Tokenisierung wird Text durch Vektor repräsentiert.
- Vektorraummodell:
 - ◆ Vektor der Worthäufigkeiten für probabilistische Modelle.
 - ◆ TFIDF-Repräsentation für lineare Verfahren.
- Wortreihenfolge bleibt unberücksichtigt.
- Vektorraummodell mit N-Grammen:
 - ◆ Wird für Spamererkennung verwendet.
 - ◆ Jedes N-Gramm durch Dimension repräsentiert, oder
 - ◆ Sparse Binary Polynomial Hashing oder
 - ◆ Orthogonal Sparse N-Grams.

Repräsentation

- Jedes N-Gramm durch Dimension repräsentiert,
 - ◆ Nur N-Gramme, die auch tatsächlich auftreten.
- Sparse Binary Polynomial Hashing
 - ◆ Fenster der Breite N wird über Text geschoben,
 - ◆ Jede Teilmenge von bis zu N Tokens, die in dieser Reihenfolge auftauchen und das am weitesten rechts stehende Token beinhalten, ist eine Dimension,
- Orthogonal Sparse Bigrams.
 - ◆ Fenster der Breite N wird über Text geschoben,
 - ◆ Jedes Paar aus einem beliebigen Token im Fenster und dem am rechten Fensterrand stehenden Token ist ein Merkmal.
- SBPH und OSB: N-Gramme mit Platzhaltern.

Klassifikator / Entscheidungsfunktion

- Für eine binäre Klassifikation ($y = +1$ oder -1) wird eine Entscheidungsfunktion $f(x)$ gelernt.
- Je größer $f(x)$, desto wahrscheinlicher ist, dass x zur Klasse gehört.
- Wenn $f(x) \geq \theta$, dann entscheide $h(x) = +1$, sonst $h(x) = -1$.
- Klassifikator $h(x)$, Entscheidungsfunktion $f(x)$.
- Der Wert für θ verschiebt „false positives“ zu „false negatives“.
- Optimaler Wert hängt von Kosten einer positiven oder negativen Fehlklassifikation ab.

Evaluation von Textklassifikatoren

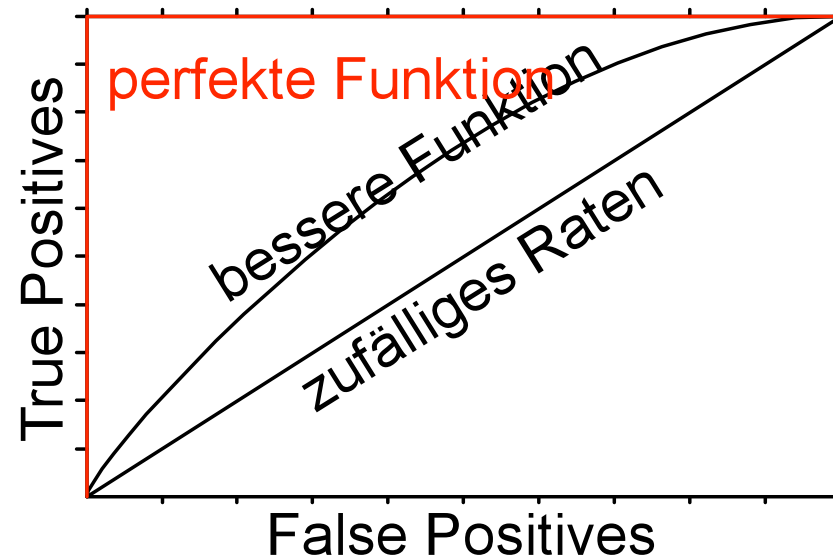
- Fehlklassifikationswahrscheinlichkeit
 - ◆ Häufig nicht aussagekräftig, weil $P(+1)$ sehr klein.
 - ◆ Wie gut sind 5% Fehler, wenn $P(+1)=3\%$?
 - ◆ Idee: Nicht Klassifikator bewerten, sondern Entscheidungsfunktion.
- Precision / Recall
 - ◆ Precision-Recall-Kurve bewertet Entscheidungsfunktion,
 - ◆ Jeder Wert für θ entspricht Punkt auf P-R-Kurve.
 - ◆ F-Measure: „Durchschnitt“ aus Precision und Recall.
- Receiver Operating Characteristic (ROC-Kurve)
 - ◆ Bewertet Entscheidungsfunktion,
 - ◆ Fläche unter ROC-Kurve = $P(\text{positives Beispiel hat höheren f-Wert als negatives Beispiel})$

ROC-Analyse

- Entscheidungsfunktion + Schwellwert = Klassifikator.
- Großer Schwellwert: Mehr positive Bsp falsch.
- Kleiner Schwellwert: Mehr negative Bsp falsch.
- Bewertung der Entscheidungsfunktion unabhängig vom konkreten Schwellwert.
- Receiver-Operating-Characteristic-Analyse
- Kommt aus der Radartechnik, im 2. Weltkrieg entwickelt.
- Bewertung der Qualität von Entscheidungsfunktionen.

ROC-Kurven

- Charakterisieren das Verhalten des Klassifikators für alle möglichen Schwellwerte.
- X-Achse: „False Positives“: Anzahl negativer Beispiele, die als positiv klassifiziert werden.
- Y-Achse: „True Positives“: Anzahl positiver Beispiele, die als positiv klassifiziert werden.



Bestimmen der ROC-Kurve von f

- Für alle positiven Beispiele x_p in Testmenge
 - ◆ Füge $f(x_p)$ in sortierte Liste L_p ein.
- Für alle negativen Beispiele x_n in Testmenge
 - ◆ Füge $f(x_n)$ in sortierte Liste L_n ein.
- $tp = fp = 0$
- Wiederhole solange L_p und L_n nicht leer sind.
 - ◆ Wenn $L_p \rightarrow \text{Element} \geq L_n \rightarrow \text{Element}$ dann $\text{increment}(tp)$ und $L_p = LP \rightarrow \text{Next}$.
 - ◆ Wenn $L_n \rightarrow \text{Element} \geq L_p \rightarrow \text{Element}$ dann $\text{increment}(fp)$ und $L_n = Ln \rightarrow \text{Next}$.
 - ◆ Zeichne neuen Punkt (fp, tp)

Flächeninhalt der ROC-Kurve

- Flächeninhalt AUC kann durch Integrieren (Summieren der Trapez-Flächeninhalte) bestimmt werden.
- p = zufällig gezogenes Positivbeispiel
- n = zufällig gezogenes Negativbeispiel
- Theorem: $AUC = P(f(p) > f(n))$.

- Beweisidee: ROC-Kurve mit nur einem Positiv- und einem Negativbeispiel (Flächeninhalt 0 oder 1);
Durchschnitt vieler solcher Kurven = AUC.

Precision / Recall

- Alternative zur ROC-Analyse.
- Stammt aus dem Information Retrieval.
- $$\text{Precision} = \frac{\text{True positives}}{\text{True} + \text{false positives}}$$
- $$\text{Recall} = \frac{\text{True positives}}{\text{True positive} + \text{false negatives}}$$
- Precision: $P(\text{richtig} \mid \text{als positiv erkannt})$
- Recall: $P(\text{als positiv erkannt} \mid \text{ist positiv})$

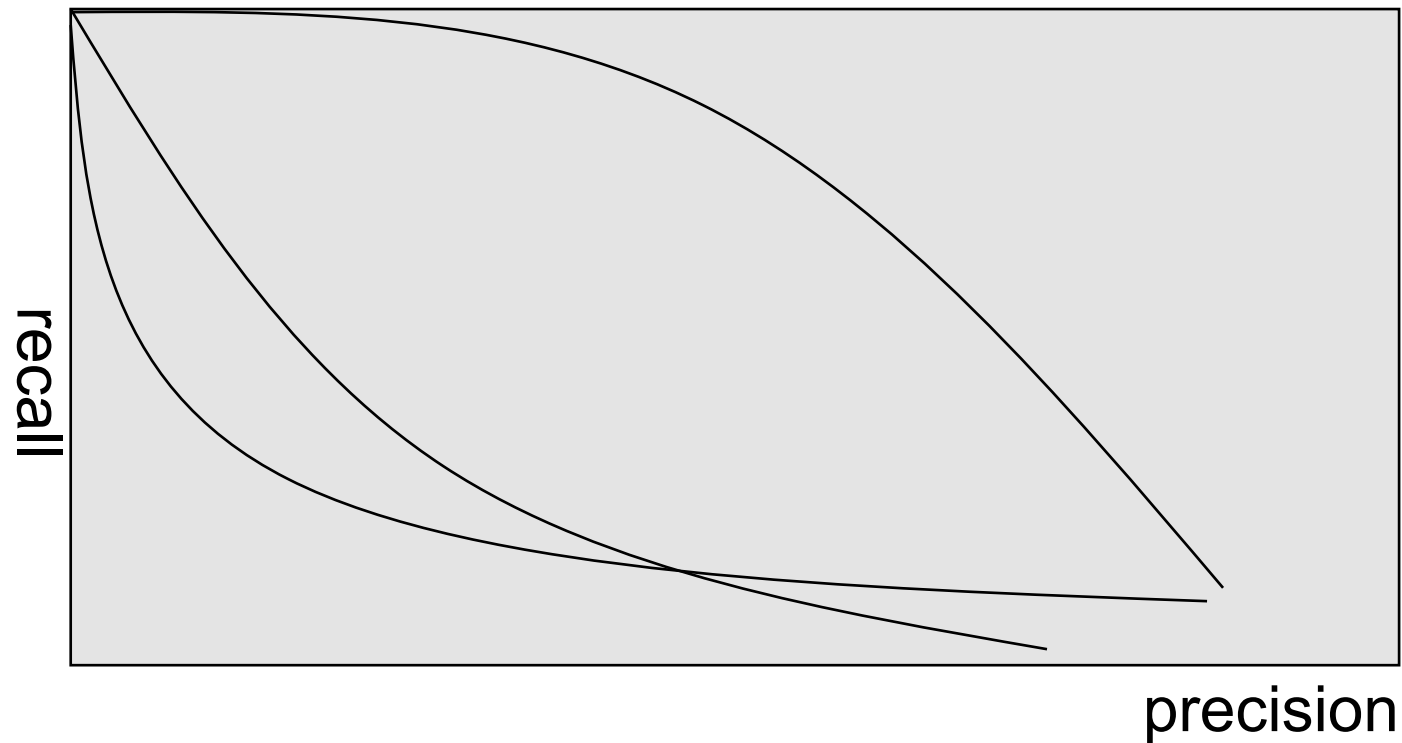
Precision / Recall

- Zusammenfassungen der Kurve in einer Zahl:
 - ◆ F-Measure: Maximum über alle (p,r)-Paare auf der Kurve:

$$\text{F-measure} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

- ◆ Precision-Recall-Breakeven-Point: Derjenige Wert PRBEP für den gilt:
Precision(θ) = Recall(θ) = PRBEP.

Precision / Recall Trade-Off



- Precision-/Recall-Kurven
- Welcher Klassifikator ist der Beste / Schlechteste

Fehlerschätzung

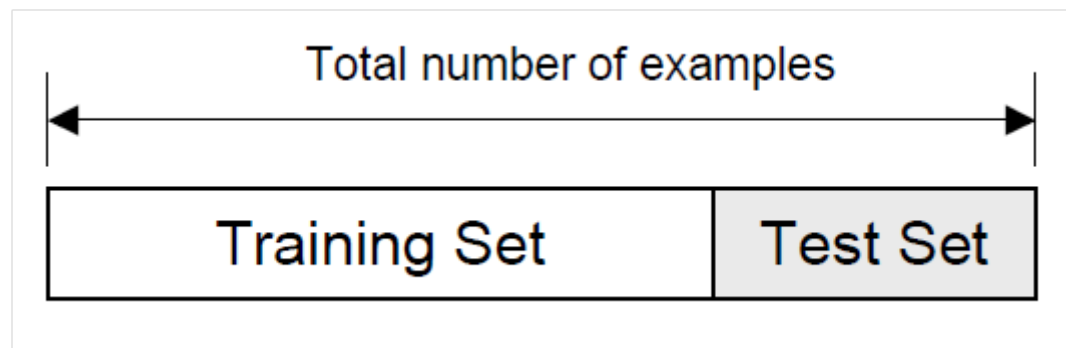
- Wir haben Klassifikator h gelernt, und wollen nun wissen wie gut er ist.
- Echter Fehler E ist unbekannt.
- Verwende Fehlerschätzer \hat{E}
- Bias = $\mathbb{E}[\hat{E}] - E = \begin{cases} < 0 & \text{optimistisch} \\ > 0 & \text{pessimistisch} \\ = 0 & \text{erwartungstreu} \end{cases}$
- $Var[\hat{E}] = \begin{cases} \text{groß bei wenigen Daten} \Rightarrow \text{Schätzung unsicher} \\ \text{klein bei vielen Daten} \Rightarrow \text{Schätzung sicher} \end{cases}$

Fehlerschätzung

- Problem: Fehlerschätzung nur auf den Trainingsdaten
- Empirischen Fehler: Messe wie viele Datenpunkte falsch klassifiziert werden.
- Es kann sein, dass $\mathbb{E}[\hat{E}] = 0$ (Overfitting)
 - ◆ Dann gilt für den Bias = $\mathbb{E}[\hat{E}] - E < 0$
 - ◆ D.h. wir sind sehr optimistisch, was unseren Fehler angeht, da wir ihn geringer schätzen als den tatsächlichen echten Fehler.
- Deshalb: Fehlerschätzung auf Trainingsdaten schlecht!

Holdout-Testing

- Besser: Verwende Testdaten, die nicht zum Trainieren des Klassifikators benutzt wurden.
- Training-und-Test:
 - ◆ Verwende 80% der Daten zum Trainieren
 - ◆ Verwende 20% der Daten zum Testen
- Auf 20% der Daten wird dann \hat{E} , ROC-Kurve, PR-Kurve ermittelt.



Holdout-Testing

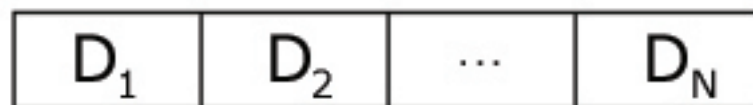
- Algorithmus Training-and-Test
 - ◆ Aufteilen der Datenbank in Trainingsmenge und Testmenge.
 - ◆ h_1 = Lernalgorithmus (Trainingsmenge)
 - ◆ Bestimme \hat{E} anhand der Testmenge.
 - ◆ h = Lernalgorithmus (alle Beispiele)
 - ◆ Liefere Hypothese h zusammen mit Fehlerschätzer

$$\hat{E} \pm z_{\delta/2} \sqrt{\frac{\hat{E}(1-\hat{E})}{20\% \cdot m}}$$

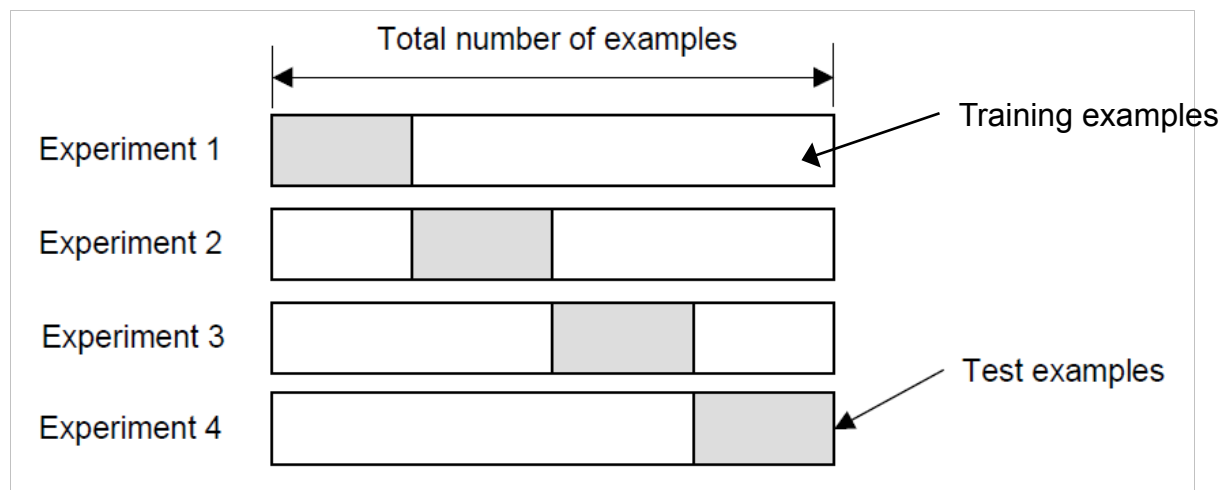
- Training-and-Test ist für große Datenbanken gut anwendbar.
- Problematisch für kleine Datenbanken

N-Fold Cross-Validation

- Idee: Teile Daten D in N gleichgroße Teilstücke



- Trainiere auf $D \setminus D_i$ und teste dann auf D_i
- Iteriere wie folgt:



N-Fold Cross-Validation

- Algorithmus N-CV (Beispiele D)
 - ◆ Bilde N etwa gleich große Blöcke D_1, \dots, D_N von Beispielen, die zusammen D ergeben. $\hat{E}=0$.
 - ◆ For $i = 1 \dots N$
 - ★ $h = \text{Lernalgorithmus}(D \setminus D_i)$
 - ★ $\hat{E} = \hat{E} + \text{empirischer Fehler von } h \text{ auf } D_i$
 - ◆ $\hat{E} = \hat{E}/N$
 - ◆ $h = \text{Lernalgorithmus}(D)$
 - ◆ Liefere Hypothese h mit Fehlerschätzer

$$\hat{E} \pm z_{\delta/2} \cdot \sqrt{\frac{S_N}{N}}$$

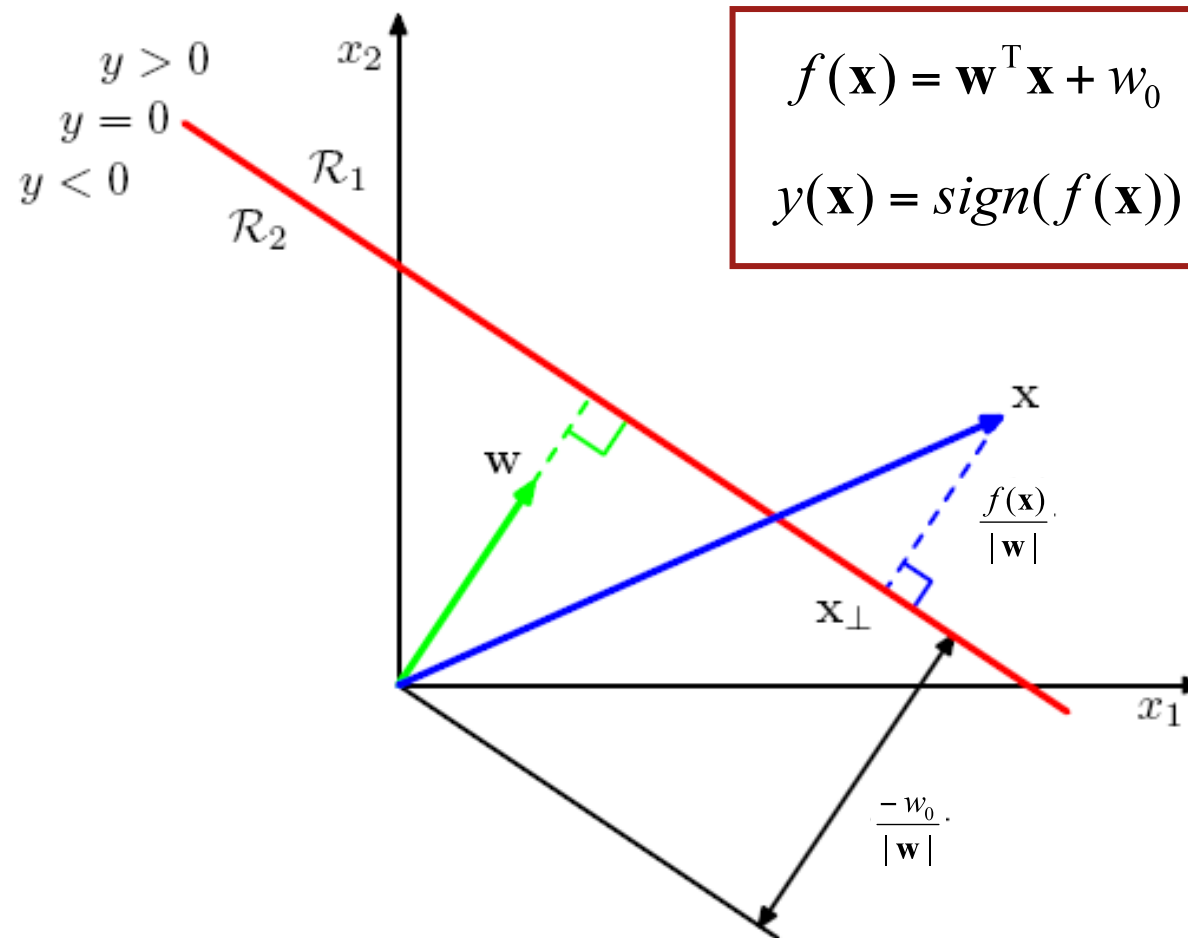
- Wenn $|D| = N$, heisst das Verfahren Leave-one-Out.
- Nur leicht pessimistischer Schätzer.

Was man über Hyperebenen wissen sollte...

- Wir beschreiben eine Hyperebene durch $\langle \mathbf{w}, \mathbf{x} \rangle = 0$, wobei $\mathbf{x}, \mathbf{w} \in \mathbb{R}^d$
 1. Was ist \mathbf{w} ?
 2. Wann gilt $\langle \mathbf{w}, \mathbf{x} \rangle > 0$ bzw. $\langle \mathbf{w}, \mathbf{x} \rangle < 0$

- Die Ebene ist ein Model, mit dem klassifiziert werden kann
 1. Wie bestimmt man eine derartige Ebene für eine Menge $D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$
 2. Wie viele gibt es?
 3. Welche ist die beste (für ungesehene Daten)?

Lineare Klassifikatoren



Lineare Klassifikatoren

- Umformulierung mit zusätzlichem, konstanten Eingabeattribut $x_0=1$:

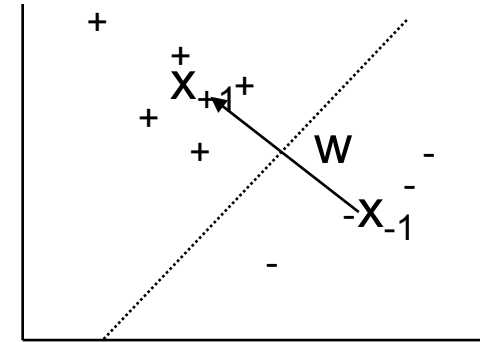
- ◆ $f(\mathbf{x}) = \mathbf{w}_{(1..n)}^T \mathbf{x}_{(1..n)} + w_0$

$$= (w_1 \quad \dots \quad w_n) \begin{pmatrix} x_1 \\ \dots \\ x_n \end{pmatrix} + w_0 = (w_0 \quad w_1 \quad \dots \quad w_n) \begin{pmatrix} x_0 \\ x_1 \\ \dots \\ x_n \end{pmatrix}$$
$$= \mathbf{w}_{(0..n)}^T \mathbf{x}_{(0..n)}$$

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

$$y(\mathbf{x}) = \text{sign}(f(\mathbf{x}))$$

Rocchio



- $\bar{\mathbf{x}}_{-1}$: Mittelpunkt der neg. Beispiele
- $\bar{\mathbf{x}}_{+1}$: Mittelpunkt der pos. Beispiele
- Trennebene: Normalenvektor = $(\mathbf{x}_{+1} - \mathbf{x}_{-1})$

$$f(\mathbf{x}) = (\bar{\mathbf{x}}_{+1} - \bar{\mathbf{x}}_{-1})\mathbf{x} + w_0$$

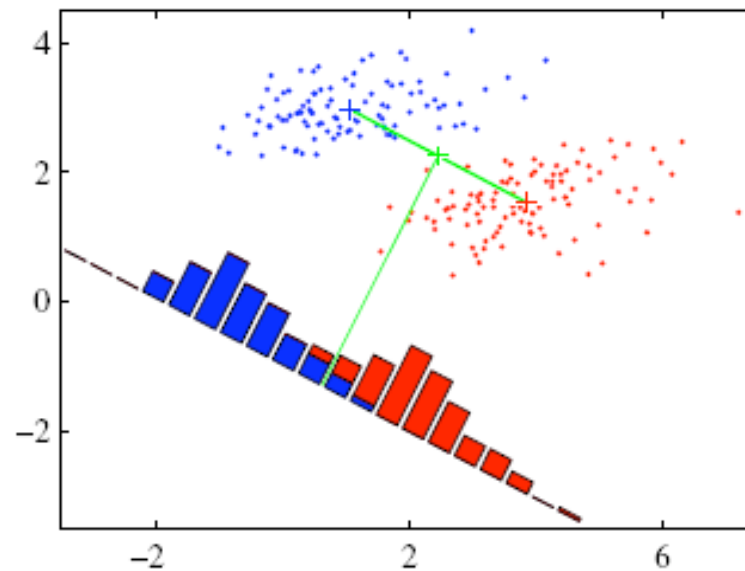
- Bestimmung von w_0 : Mittelpunkt $(\mathbf{x}_{+1} + \mathbf{x}_{-1})/2$ muss auf der Ebene liegen.

$$f((\bar{\mathbf{x}}_{-1} + \bar{\mathbf{x}}_{+1})/2) = (\bar{\mathbf{x}}_{+1} - \bar{\mathbf{x}}_{-1})(\bar{\mathbf{x}}_{-1} + \bar{\mathbf{x}}_{+1})/2 + w_0 = 0$$

$$\Leftrightarrow w_0 =$$

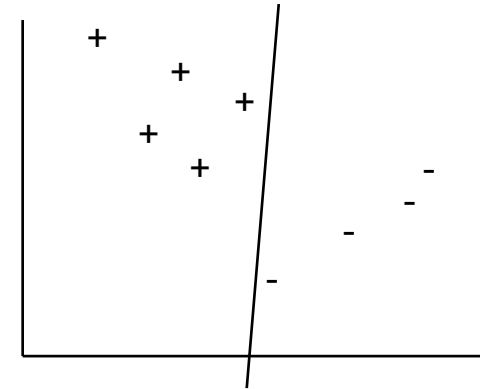
Rocchio

- Trennebenen hat maximalen Abstand von den Mittelpunkten der Klassen.
- Trainingsbeispiele können falsch klassifiziert werden.
- Differenz der Mittelwerte kann schlechter Normalenvektor für Diskrimination sein.



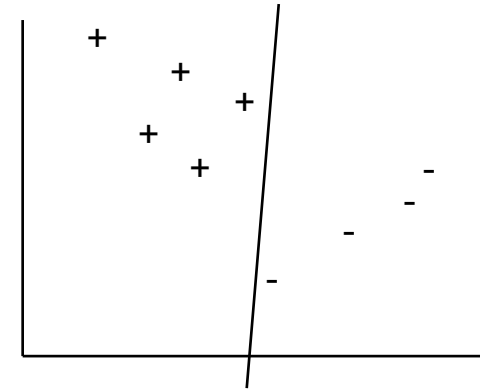
Perzeptron

- Lineares Modell:
 - ◆ $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
- Ziel:
 - ◆ Für alle Beispiele $(\mathbf{x}_i, +1)$:
 $f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i > 0$
 - ◆ Für alle Beispiele $(\mathbf{x}_i, -1)$:
 $f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i < 0$



Perzeptron

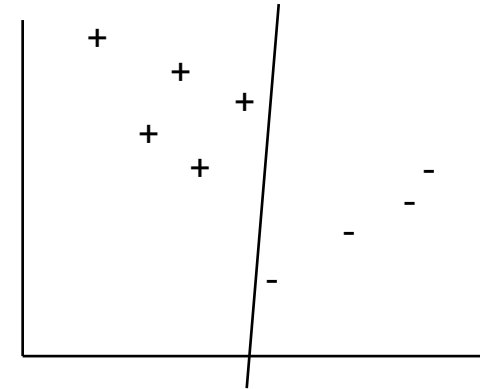
- Lineares Modell:
 - ◆ $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
- Ziel:
 - ◆ Für alle Beispiele $(\mathbf{x}_i, +1)$:
$$f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i > 0$$
 - ◆ Für alle Beispiele $(\mathbf{x}_i, -1)$:
$$f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i < 0$$
- Ziel: Für alle Beispiele:
 - ◆ $y_i f(\mathbf{x}_i) = y_i \mathbf{w}^T \mathbf{x}_i > 0$
 - ◆ = Beispiel liegt auf der richtigen Seite der Ebene.



Perzeptron

- Lineares Modell:
 - ◆ $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
- Ziel: Für alle Beispiele:
 - ◆ $y_i f(\mathbf{x}_i) = y_i \mathbf{w}^T \mathbf{x}_i > 0$
- Perzeptron-Optimierungskriterium:

- ◆
$$J_P(\mathbf{w}) = \sum_{(\mathbf{x}_i, y_i) \in L} \min \{y_i \mathbf{w}^T \mathbf{x}_i, 0\}$$



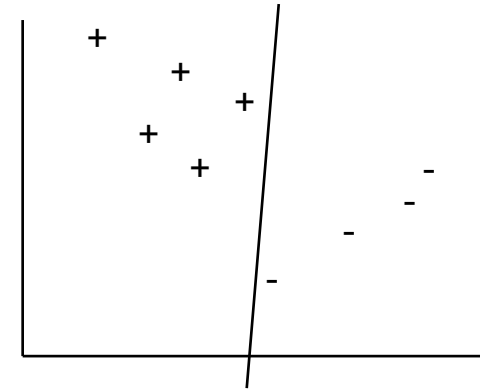
Perzeptron

- Lineares Modell:
 - ◆ $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
- Ziel: Für alle Beispiele:
 - ◆ $y_i f(\mathbf{x}_i) = y_i \mathbf{w}^T \mathbf{x}_i > 0$
- Perzeptron-Optimierungskriterium:

- ◆
$$J_P(\mathbf{w}) = \sum_{(\mathbf{x}_i, y_i) \in L} \min \{y_i \mathbf{w}^T \mathbf{x}_i, 0\}$$

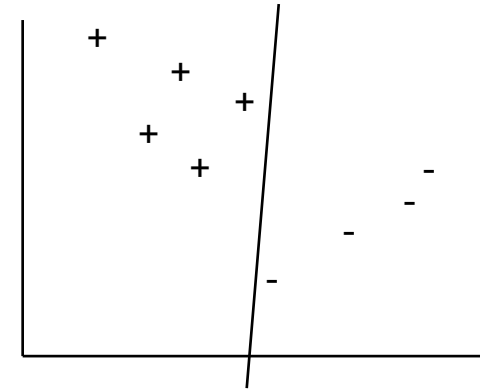
- Subgradient für Beispiel (\mathbf{x}_i, y_i) :

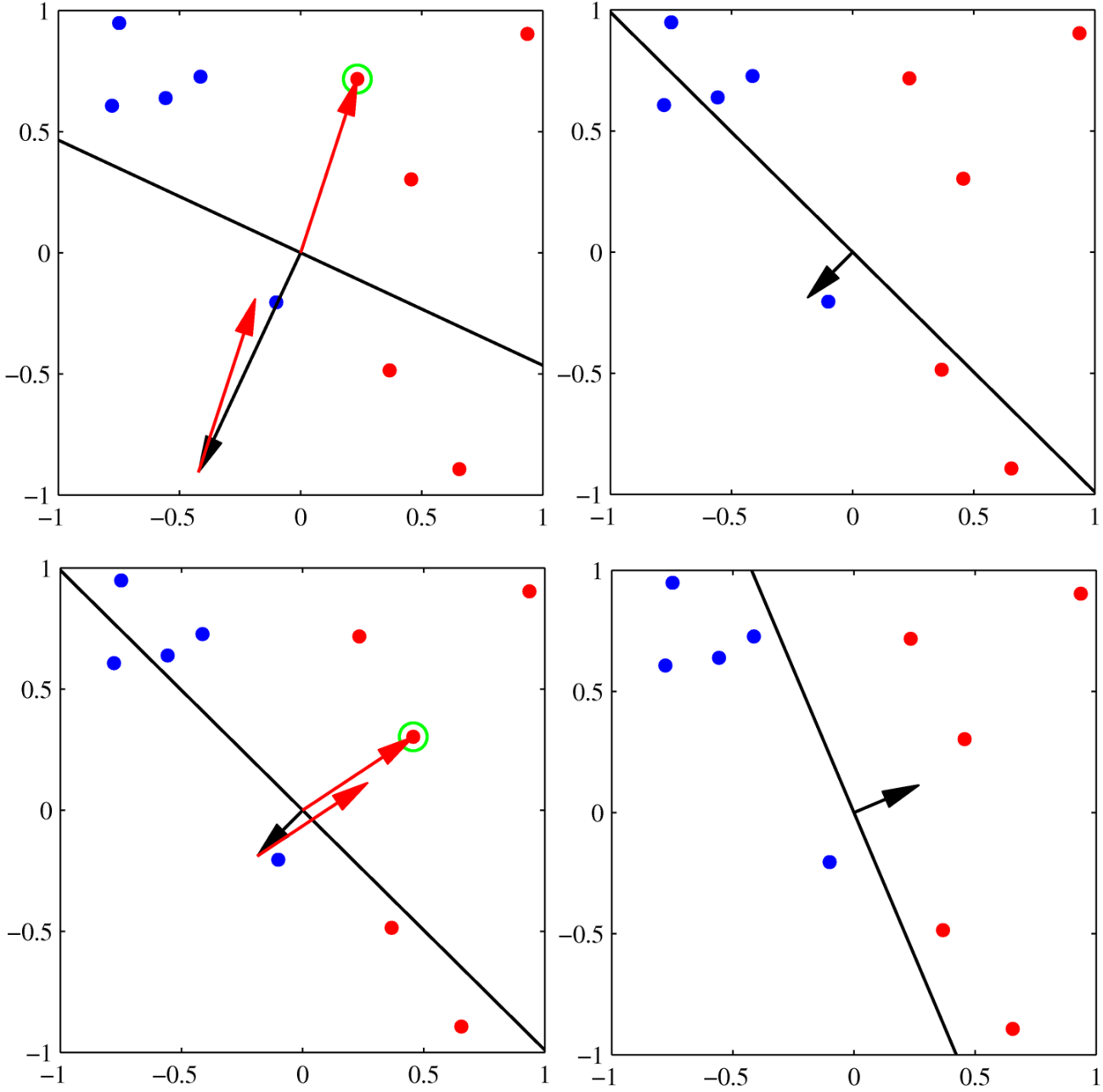
- ◆
$$\nabla_i J_P(\mathbf{w}) = \begin{cases} 0, & \text{wenn } y_i \mathbf{w}^T \mathbf{x}_i > 0 \\ y_i \mathbf{x}_i & \text{sonst} \end{cases}$$



Perzeptron

- Lineares Modell:
 - ◆ $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
- Perzeptron-Optimierungskriterium:
 - ◆ $J_P(\mathbf{w}) = \sum_{(\mathbf{x}_i, y_i) \in L} \min \{y_i \mathbf{w}^T \mathbf{x}_i, 0\}$
- Subgradient für Beispiel (\mathbf{x}_i, y_i) :
 - ◆ $\nabla_i J_P(\mathbf{w}) = \begin{cases} 0, & \text{wenn } y_i \mathbf{w}^T \mathbf{x}_i > 0 \\ y_i \mathbf{x}_i & \text{sonst} \end{cases}$
- Gradientenaufstieg: Wiederhole,
für alle Beispiele mit $y_i \mathbf{w}^T \mathbf{x}_i \leq 0$
 - ◆ $\mathbf{w} \leftarrow \mathbf{w} + y_i \mathbf{x}_i$





Perzeptron-Algorithmus

- Lineares Modell:
 - ◆ $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
- Perzeptron-Trainingsalgorithmus:
- Solange noch Beispiele (\mathbf{x}_i, y_i) mit der Hypothese inkonsistent sind, iteriere über alle Beispiele
 - ◆ Wenn $y_i \mathbf{w}^T \mathbf{x}_i \leq 0$ dann $\mathbf{w} \leftarrow \mathbf{w} + y_i \mathbf{x}_i$

Perzeptron-Algorithmus

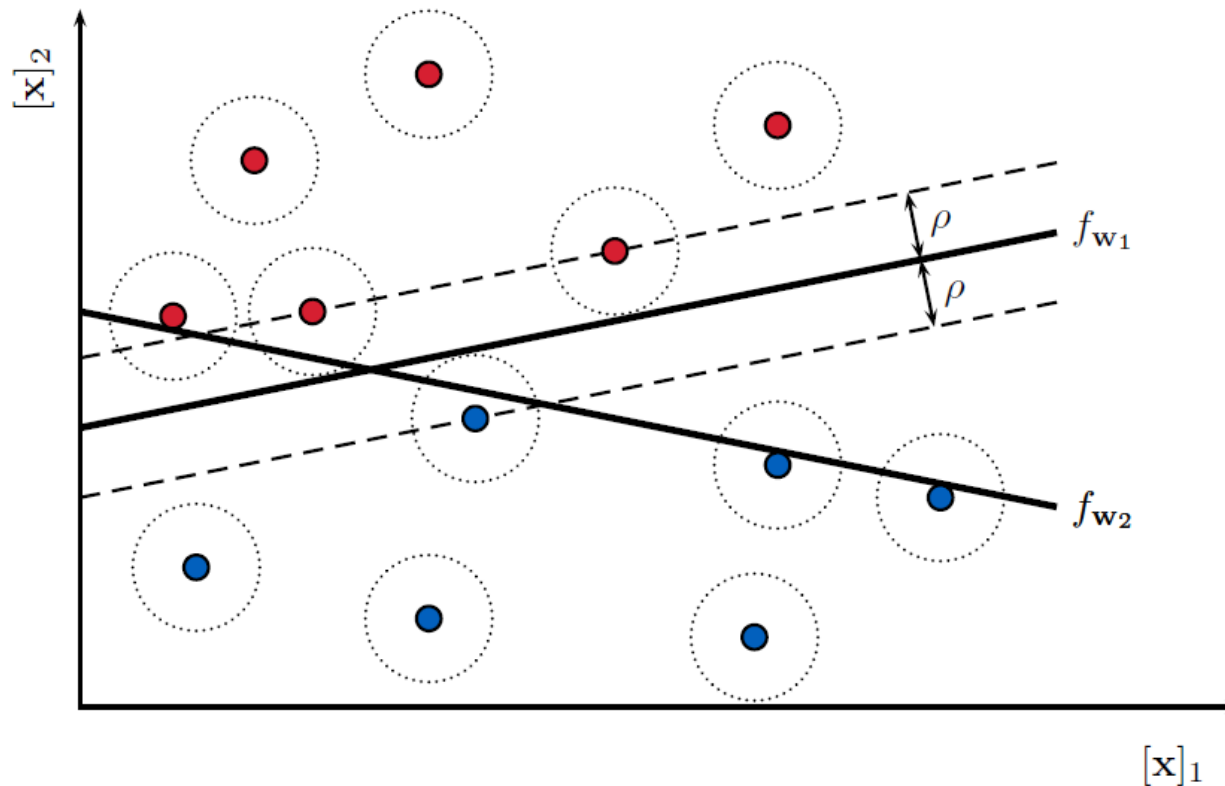
- Perzeptron findet immer eine Trennebene, wenn eine existiert (Optimierungskriterium ist konvex).
- Existiert immer eine Trennebene?
- Wenn keine Trennebene existiert, dann konvergiert der Perzeptron-Algorithmus auch nicht.

Probleme

1. Wenn linear-separierbar, dann mehrere Modelle:
Welches ist das beste?
 - ◆ Lösung: Maximal-Margin-Ansatz

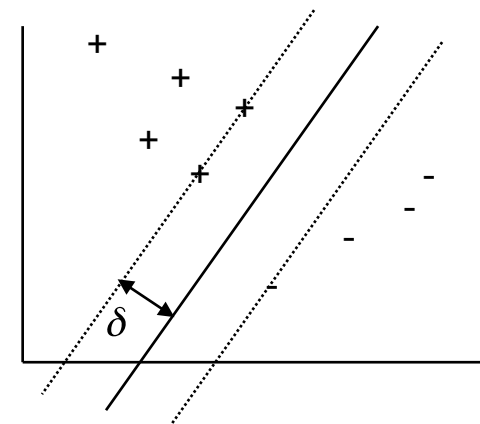
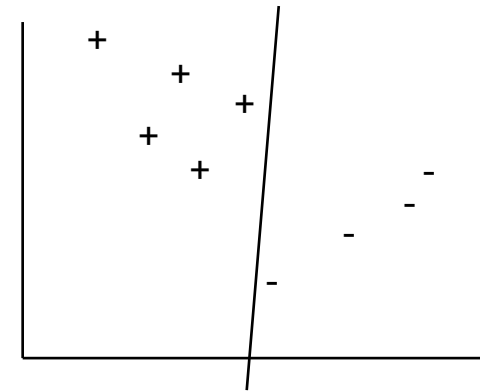
Probleme

1. Wenn linear-separierbar, dann mehrere Modelle:
Welches ist das beste?
 - ◆ Lösung: Maximal-Margin-Ansatz



Margin-Perzeptron

- Perzeptron-Klassifikation:
 - ◆ $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
- Perzeptron: für alle Beispiele muss gelten
 - ◆ $y_i \mathbf{w}^T \mathbf{x}_i > 0$
 - ◆ = Beispiel liegt auf der richtigen Seite der Ebene.
- Margin-Perzeptron:
 - ◆ $y_i \left(\frac{\mathbf{w}^T \mathbf{x}_i}{|\mathbf{w}|} \right) > \delta$
 - ◆ = Beispiel mindestens δ von Trennebene entfernt.

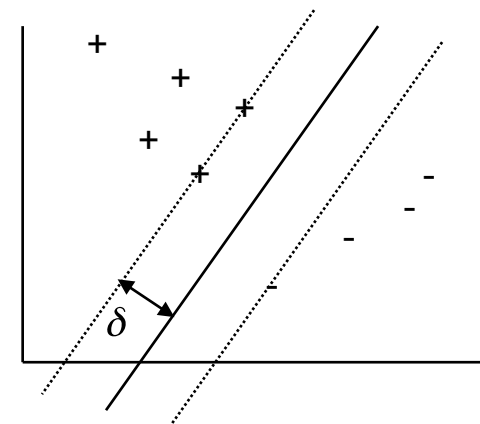
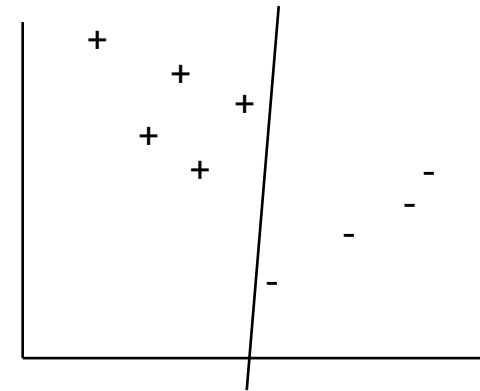


Margin-Perzeptron-Algorithmus

- Lineares Modell:
 - ◆ $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
- Perzeptron-Trainingsalgorithmus:
- Solange noch Beispiele (\mathbf{x}_i, y_i) mit der Hypothese inkonsistent sind, iteriere über alle Beispiele
 - ◆ Wenn $y_i \left(\frac{\mathbf{w}^T \mathbf{x}_i}{\|\mathbf{w}\|} \right) \leq \delta$ dann $\mathbf{w} \leftarrow \mathbf{w} + y_i \mathbf{x}_i$

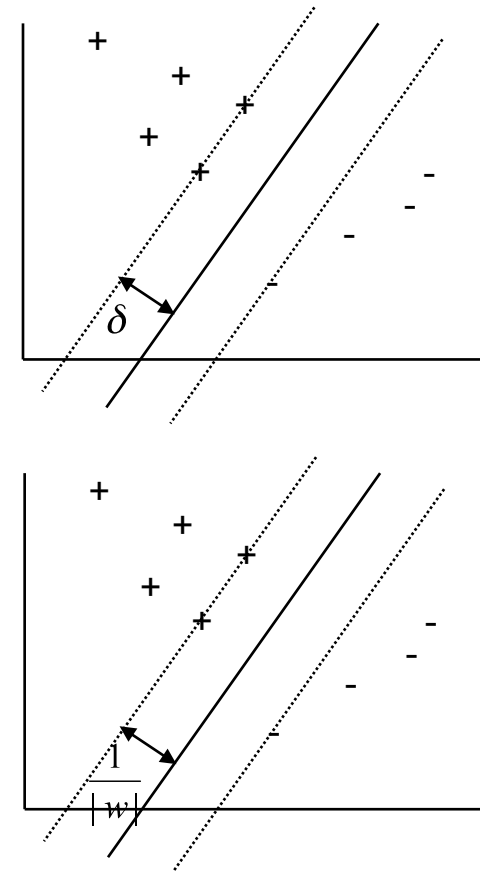
Margin-Maximierung

- Perzeptron: für alle Beispiele muss gelten $y_i \mathbf{w}^T \mathbf{x}_i > 0$
- Margin-Perzeptron: $y_i \left(\frac{\mathbf{w}^T \mathbf{x}_i}{|\mathbf{w}|} \right) > \delta$
 - ◆ Finde Ebene, die alle Beispiele mindestens δ von Ebene entfernt.
 - ◆ Fester, voreingestellter Wert δ .
- Margin-Maximierung:
 - ◆ Finde Ebene, die alle Beispiele mindestens δ von Ebene entfernt.
 - ◆ Für den größtmöglichen Wert δ .



Margin-Maximierung

- Margin-Maximierung: $y_i \left(\frac{\mathbf{w}^T}{|\mathbf{w}|} \mathbf{x}_i \right) \geq \delta$
 - ◆ Finde Ebene, die alle Beispiele mindestens δ von Ebene entfernt.
 - ◆ Für den größtmöglichen Wert δ .
- Maximiere δ unter der Nebenbedingung: für alle Beispiele (\mathbf{x}_i, y_i) gelte $y_i \left(\frac{\mathbf{w}^T}{|\mathbf{w}|} \mathbf{x}_i \right) \geq \delta$
- = Minimiere $|\mathbf{w}|$ unter der Nebenbedingung:
 - ◆ für alle Beispiele (\mathbf{x}_i, y_i) : $y_i \mathbf{w}^T \mathbf{x}_i \geq 1$

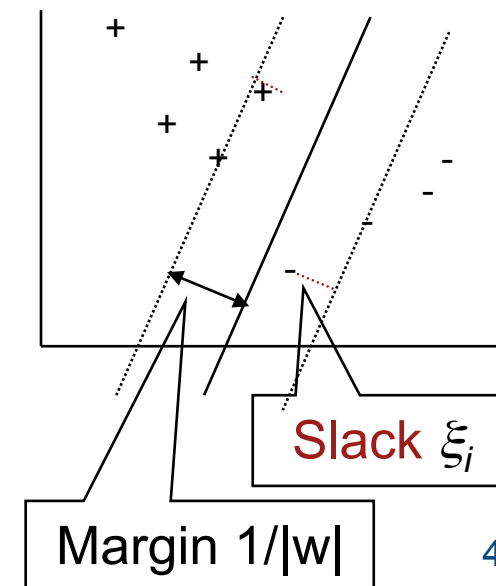
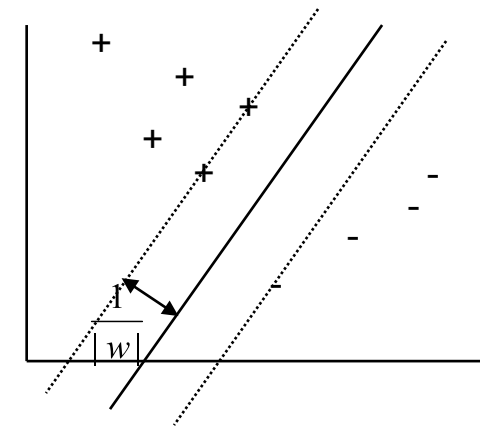


Probleme

1. Wenn linear-separierbar, dann mehrere Modelle:
Welches ist das beste?
 - ◆ Lösung: Maximal-Margin-Ansatz
2. Daten verrauscht / falsch gelabelte Daten
 - ◆ Lösung: Slack zulassen

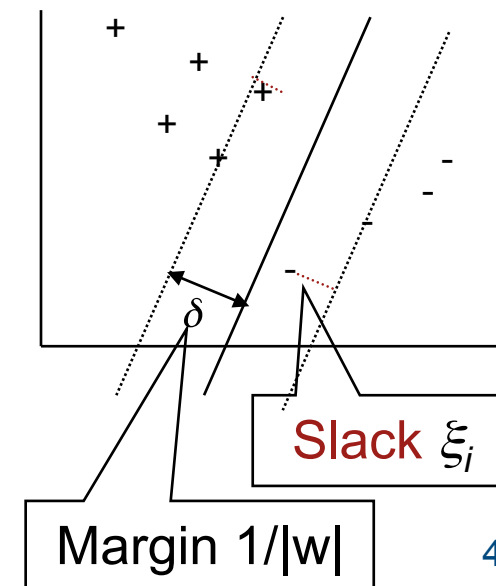
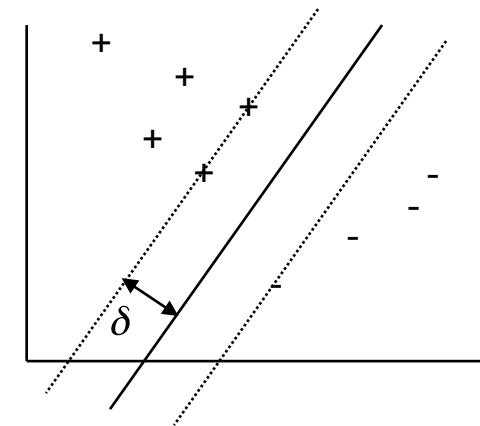
Soft-Margin-Maximierung

- Hard-Margin-Maximierung:
 - ◆ Minimiere $|\mathbf{w}|$ unter der Nebenbedingungen:
 - ◆ für alle Beispiele (\mathbf{x}_i, y_i) : $y_i \mathbf{w}^T \mathbf{x}_i \geq 1$
- Soft-Margin-Maximierung:
 - ◆ Minimiere $|\mathbf{w}| + C \sum_i \xi_i$ unter den Nebenbedingungen:
 - ◆ für alle Beispiele (\mathbf{x}_i, y_i) :
$$y_i \mathbf{w}^T \mathbf{x}_i \geq 1 - \xi_i$$
 - ◆ Alle $\xi_i \geq 0$.
- Soft-Margin-Ebene existiert immer, Hard-Margin-Ebene nicht!



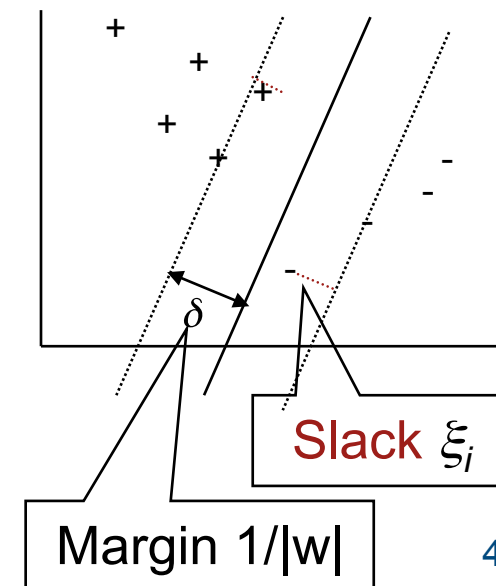
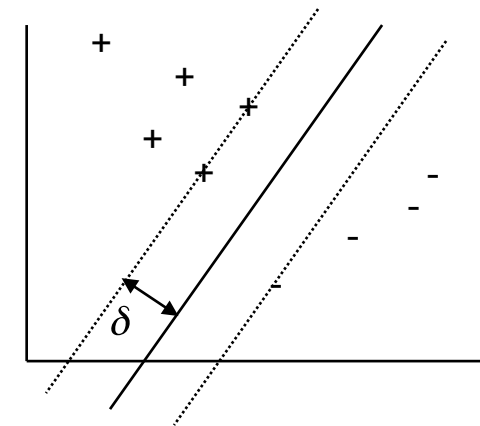
Soft-Margin-Maximierung

- Soft-Margin-Maximierung:
 - ◆ Minimiere $|\mathbf{w}| + C \sum_i \xi_i$ unter den Nebenbedingungen:
 - ◆ für alle Beispiele (\mathbf{x}_i, y_i) :
$$y_i \mathbf{w}^T \mathbf{x}_i \geq 1 - \xi_i$$
 - ◆ Alle $\xi_i \geq 0$.
- Einsetzen von ξ_i in Optimierungskriterium ergibt
 - ◆ Minimiere: $|\mathbf{w}| + C \sum_i \max\{0, 1 - y_i \mathbf{w}^T \mathbf{x}_i\}$



Soft-Margin-Maximierung

- Soft-Margin-Maximierung:
 - ◆ Minimiere: $|\mathbf{w}| + C \sum_i \max\{0, 1 - y_i \mathbf{w}^T \mathbf{x}_i\}$
- Minimierung mit Gradientenverfahren.
- Annäherung durch differenzierbare Variante (Huber- statt Hinge-Loss), dann Newton-Verfahren.
- Kriterium ist konvex, es gibt genau ein Minimum.
- **Primale Support Vector Machine.**
- $O(n^2)$, unter 100.000 Beispiele, unter 100.000 Attribute.

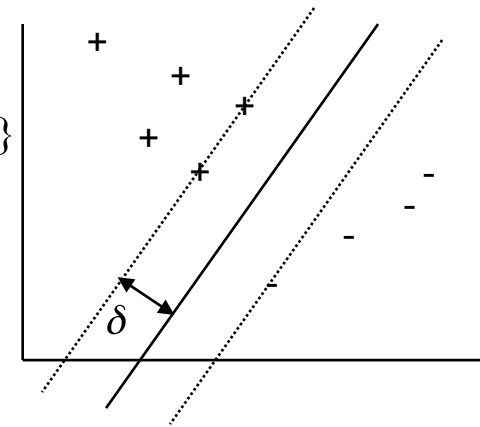


Soft-Margin-Maximierung

- Soft-Margin-Maximierung:
 - ◆ Minimiere: $E_H(\mathbf{w}) = |\mathbf{w}| + C \sum_i \max\{0, 1 - y_i \mathbf{w}^T \mathbf{x}_i\}$
- Minimierung mit Gradientenverfahren.
- Wiederhole:

- ◆ $\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial E_H(\mathbf{w})}{\partial \mathbf{w}}$

Enthält Summe über
alle Beispiele



Probleme

1. Wenn linear-separierbar, dann mehrere Modelle:
Welches ist das beste?
 - ◆ Lösung: Maximal-Margin-Ansatz
2. Daten verrauscht / falsch gelabelte Daten
 - ◆ Lösung: Slack zulassen
3. Keine Fehlerwahrscheinlichkeit der Klassifikation
 - ◆ Lösung: Logistic Regression

Logistische Regression

- SVM: großer Entscheidungsfunktionswert \sim hohe Sicherheit der Vorhersage.
- Aber: beim Lernen nicht auf korrekte Kalibrierung der Klassenwahrscheinlichkeiten optimiert.
- $f(\mathbf{x})=18.3 \rightarrow$ Risiko eines Fehlers?
- Logistische Regression: Vorhersage der Klassenwahrscheinlichkeit.

Logistische Regression

- Bayes' Regel:

- ◆
$$P(y = +1 | \mathbf{x}) = \frac{p(\mathbf{x} | y = +1)P(y = +1)}{p(\mathbf{x} | y = +1)P(y = +1) + p(\mathbf{x} | y = -1)P(y = -1)}$$
$$= \frac{1}{1 + \frac{p(\mathbf{x} | y = -1)P(y = -1)}{p(\mathbf{x} | y = +1)P(y = +1)}} = \frac{1}{1 + \exp(-a)} = \sigma(a)$$

- Log-odd ratio:

$$a = \ln \frac{p(\mathbf{x} | y = +1)P(y = +1)}{p(\mathbf{x} | y = -1)P(y = -1)}$$

- $\frac{1}{1 + \exp(-a)}$ nennt man auch Sigmoid-Funktion

Logistische Regression

- Likelihood jeder Klasse normalverteilt, gemeinsame Kovarianzmatrix für beide Klassen.

$$\diamond p(\mathbf{x} | y) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right]$$

- Logg-odds ratio:

$$\begin{aligned} a &= \ln \frac{\frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \mu_{+1})^T \Sigma^{-1}(\mathbf{x} - \mu_{+1})\right] P(y = +1)}{\frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \mu_{-1})^T \Sigma^{-1}(\mathbf{x} - \mu_{-1})\right] P(y = -1)} \\ &= \left[-\frac{1}{2}(\mathbf{x} - \mu_{+1})^T \Sigma^{-1}(\mathbf{x} - \mu_{+1})\right] - \left[-\frac{1}{2}(\mathbf{x} - \mu_{-1})^T \Sigma^{-1}(\mathbf{x} - \mu_{-1})\right] + \ln \frac{P(y = +1)}{P(y = -1)} \\ &= \underbrace{\left[\Sigma^{-1}(\mu_{+1} - \mu_{-1})\right]}_{\mathbf{w}} \cdot \mathbf{x} + \underbrace{\left[-\frac{1}{2}(\mu_{+1}^T \Sigma^{-1} \mu_{+1}) + \frac{1}{2}(\mu_{-1}^T \Sigma^{-1} \mu_{-1}) + \ln \frac{P(y = +1)}{P(y = -1)}\right]}_{w_0} \end{aligned}$$

Logistische Regression

- Likelihood jeder Klasse normalverteilt, gemeinsame Kovarianzmatrix für beide Klassen.

- ◆
$$p(\mathbf{x} | y) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \mu)^\top \Sigma^{-1}(\mathbf{x} - \mu)\right]$$

- Logg-odds ratio:

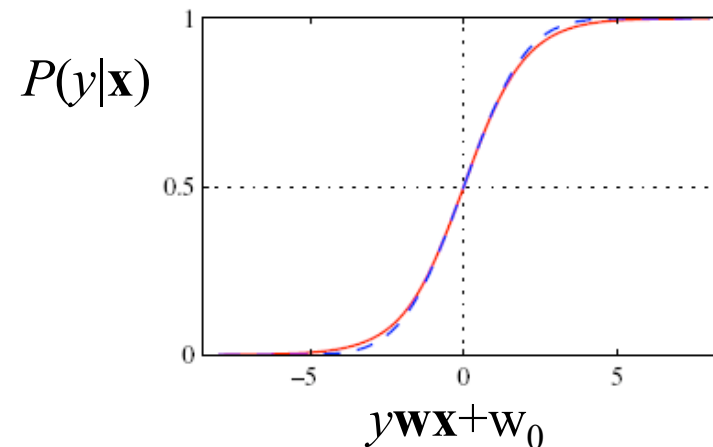
$$a = \underbrace{\left[\Sigma^{-1}(\mu_{+1} - \mu_{-1})\right]}_{\mathbf{w}} \mathbf{x} + \underbrace{\left[-\frac{1}{2}(\mu_{+1}^\top \Sigma^{-1} \mu_{+1}) + \frac{1}{2}(\mu_{-1}^\top \Sigma^{-1} \mu_{-1}) + \ln \frac{P(y = +1)}{P(y = -1)}\right]}_{w_0}$$

Logistische Regression

- Wenn zwei Klassen jeweils normalverteilte Likelihood mit derselben Kovarianzmatrix haben, dann nimmt $P(y|\mathbf{x})$ diese Form an:

$$\diamond P(y|\mathbf{x}) = \frac{1}{1 + \exp(-y\mathbf{w}\mathbf{x} + w_0)} = \sigma(\underbrace{y\mathbf{w}\mathbf{x} + w_0}_{\text{linearer Klassifikator}})$$

logistische Funktion



Logistische Regression

- Bisher: Motivation der Form des logistischen Klassifikationsmodells.
- Falls Klassenverteilungen bekannt wären, könnten wir \mathbf{w} und w_0 aus μ_{+1} , μ_{-1} und Σ herleiten.
- Sind aber nicht bekannt. Verteilungsannahme muss auch nicht stimmen.
- Jetzt: Wir finden wir tatsächlich Parameter \mathbf{w} und w_0 ?

Logistische Regression

- Prior über Parameter:

- ◆ Normalverteilung, $\mathbf{w} \sim N(0, \nu I)$

- ◆ Posterior:

- ◆
$$P(\mathbf{w} | L) = \prod_{i=1}^N p(y_i | \mathbf{x}_i, \mathbf{w}) p(\mathbf{w})$$

$$= \prod_{i=1}^N \sigma(\mathbf{w}^T \mathbf{x}_i)^{[[y_i = +1]]} (1 - \sigma(\mathbf{w}^T \mathbf{x}_i))^{[[y_i = -1]]} p(\mathbf{w})$$

- Verlustfunktion:

$$E(\mathbf{w}, L) = -\log p(\mathbf{w} | L)$$

$$= -\sum_{i=1}^N [[y_i = +1]] \log(\sigma(\mathbf{w}^T \mathbf{x}_i)) + [[y_i = -1]] \log(1 - \sigma(\mathbf{w}^T \mathbf{x}_i)) + \log\left(\frac{1}{\sqrt{(2\pi)^d |\nu I|}}\right) - \frac{\mathbf{w}^T I \mathbf{w}}{2\nu^2}$$

Logistische Regression

- Verlustfunktion ist konvex und differenzierbar.
- Gradientenabstieg führt zum Minimum.