

# Mustererkennung K-NN (K-nearest neighbour)

22.10.2018 Auf einem Tablet schreibt man mit einem Stift und erhält dadurch Informationen von Position und Zeit  $(x, y, t)$



Als Merkmal könnte man bei Zahlen z.B. die Koordinaten der Punkte als Vektoren darstellen.



Man könnte aber auch die Drehung bei der Stiftführung nehmen



Jetzt kann man einen Vektor von Features machen, um alle Merkmale zur Erkennung zu kombinieren.

$$\vec{x} = \begin{pmatrix} \theta \\ x \\ y \\ t \end{pmatrix}$$

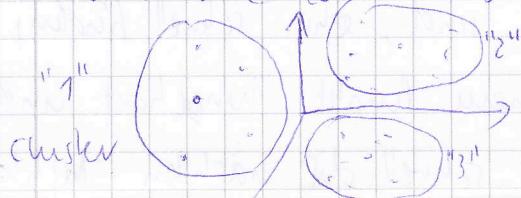
Ein anderes Merkmal ist die Öffnung, der Abstand vom Anfangs und Endpunkt geteilt durch die Länge des Zuges



$$\text{Öffnung} = \frac{d}{l}$$

Um die Länge  $L$  als eigenes Merkmal einen Linienzuges zu nehmen, muss man allerdings die Länge normieren, z.B. nur Werte zwischen 0 und 1.

Bei den Features spannt sich ein Raum  $\mathbb{R}^d$  auf, in dem die Zahlen dann durch Punkte identifizierbar sind,

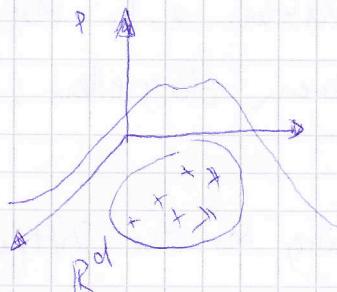


bzw. in bestimmten Bereichen landen.  
Trainingsmenge  $N \gg 1$ , die sich aufteilt.

Klassifizierung: 5 Merkmalsextraktion  $\Rightarrow \vec{x} = (\cdot)$

Dadurch kommt die Zahl irgendwo im Raum, und dann stimmen die nächsten Nachbarn ab, zu welcher Klasse die Zahl gehört. Wenn man 3 Nachbarn fragt, nennt man es 3-NN.

1) Bayes-Ansatz. Man betrachtet nur eine Klasse hier.



Hier kann man jetzt über die Wahrscheinlichkeitsverteilung im Raum Aussagen über die Klassifizierung treffen.

Man könnte auch unschägige Verteilungen nehmen oder Gaußverteilungen. Das gute an Gauß ist, dass durch den Mittelpunkt und die Kovarianz die Verteilung eindeutig bestimmt ist.

Bei KNN sind die Daten selber die Wahrscheinlichkeitsverteilung.

2) K-NN ist teuer. Man muss für jedes neue  $\vec{x}$  alle  $d(\vec{x}, \vec{x}_i)$  für alle N Testpunkte berechnen. Und unsere Trainingsmenge N ist sehr groß.

Wie kann man K-NN jetzt besser bzw. schneller machen?

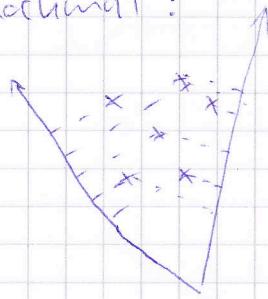
### Random projections

In unserem  $R^d$  erzeugen wir zufällig einen Vektor, der nur eine Richtung angibt. Dann werden alle Punkte auf den Richtungsvektor projiziert. Dann kann man wenn man



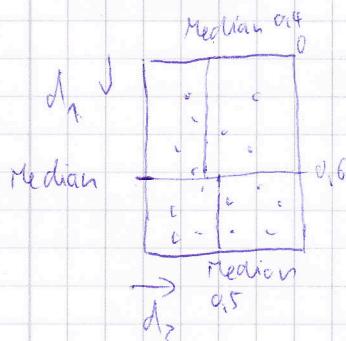
jetzt die Punkte sortiert, die auf der Linie sind, sehr schnell finden, wo der neue Punkt hingehört und dann sehr schnell die nächsten Nachbarn

auf der Linie finden. Das macht man jetzt nochmal:

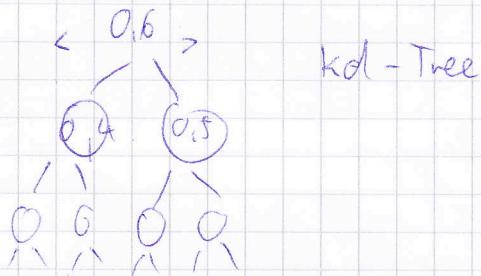


Diese zweite Projektion kann man wieder per Skalarprodukt machen

Ein anderes Verfahren, was besser zu machen wäre auch Kd-Trees.

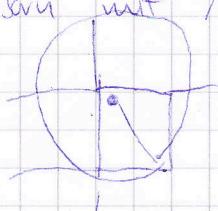


Man produziert alle Punkte in jeweils nur eine Dimension und sucht den Median, d.h. bei welchem Punkt liegen genauso viele Punkte oberhalb und unterhalb.



dann macht man so oft, bis in jedem Kästchen nur 1 Punkt übrigbleibt.

Wenn jetzt ein neuer Punkt kommt, dann kann der Baum einfach durchkämmen werden, und man weiß genau, wo er drin liegt. Dann weiß man immer nach  $\approx \log_2 N$  in welcher Region man ist. Jetzt kann man wieder den nächsten Nachbarn mit Abstand bestimmen.



Eventuell muss man auch benachbarte Regionen betrachten.

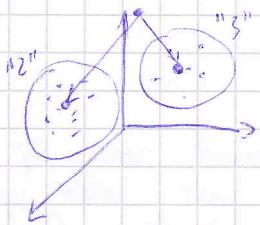
Dieses Verfahren ist aber auch nur gut, wenn die Dimension nicht so groß ist.

CLUSTERING

Vorlesung 3

Man will statt alle Datenpunkte zu behalten, für jede Klasse einen Vertreter haben.

29.10.2018



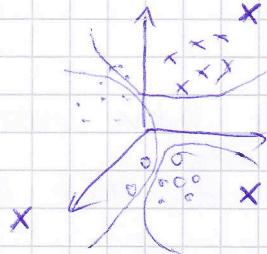
Ausserdem funktioniert es wie k-NN, aber dass man bei z.B. nur noch 10 Nachbarn betrachten muss.

Clustering ist nicht so mächtig wie k-NN, da man hier die Trennungslinien der einzelnen Cluster sehr einfach gestaltet. Dafür ist es schnell.

Unsupervised learning

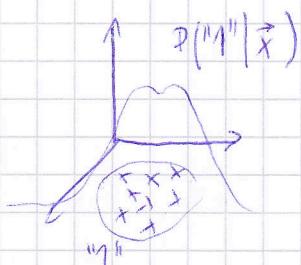
Wie können wir die Cluster bestimmen? Ein Algorithmus zum Clustering ist der LBG-Algorithmus. Wir gehen davon aus, dass wir  $k$ -Klassen im Datensatz haben.

A priori wählen wir  $k$  Vertreter  $\vec{e}_1, \vec{e}_2, \dots, \vec{e}_k$  zufällig.



Dann berechnen wir die Einteilung der Daten auf die Vertreter, indem man pro Punkt guckt, welcher Vertreter am nächsten liegt. Das ist der sogenannte Expectation-Schritt. Der Maximization-Schritt führt jetzt Korrekturen durch.  
 $\vec{e}_1$  := Schwerpunkt von Klasse 1  
 $\vdots$   
 $\vec{e}_k$  := Schwerpunkt von Klasse  $k$

Jetzt iterieren wir über diese zwei Schritte, bis es keine Korrekturen mehr gibt, bzw. nach bestimmter Anzahl von Schritten. Der Algorithmus muss nämlich nicht terminieren.



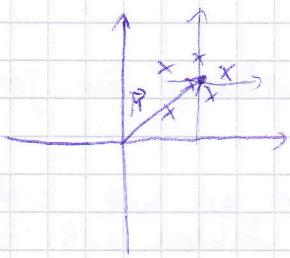
Für k-NN wissen wir, dass die Daten selber ihre Wahrscheinlichkeitsverteilung bestimmen.

Für Cluster nehmen wir die Normalverteilung an, da wir ja nur einen Vertreter haben.

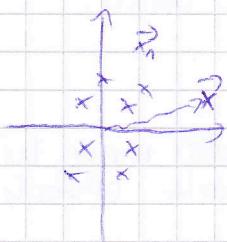


$$N(\vec{x}) = \frac{1}{2\pi} e^{-\frac{1}{2} (\vec{x} - \vec{M})^T \Sigma^{-1} (\vec{x} - \vec{M})}$$

Kovarianz



Als erster zeichnen wir die Daten. Das funktioniert durch  $\vec{x} - \vec{M}$ . Damit ergibt sich etwas:



Wir wollen einen Abstand von  $\vec{x}_1$  und  $\vec{x}$  messen, das machen wir aber nicht euklidisch, sondern mit dem Skalarprodukt.

$$\sum_{n=1}^N (\vec{x}_n^T \cdot \vec{x}_n) (\vec{x}_n^T \cdot \vec{x}) + (\vec{x}^T \cdot \vec{x}_1) (\vec{x}_1^T \cdot \vec{x}) + \dots$$

Die Formel sagt aus, wie gut der neue Vektor  $\vec{x}$  zu den vorangehenden Daten passt.

$$S = \vec{x}^T \underbrace{\frac{1}{N} (\vec{x}_1 \vec{x}_1^T + \vec{x}_2 \vec{x}_2^T + \dots + \vec{x}_N \vec{x}_N^T)}_{\Sigma \text{ Kovarianzmatrix}} \vec{x}$$

$$\Sigma = \frac{1}{N} \sum_{i=1}^N (\vec{x}_i \vec{x}_i^T)$$

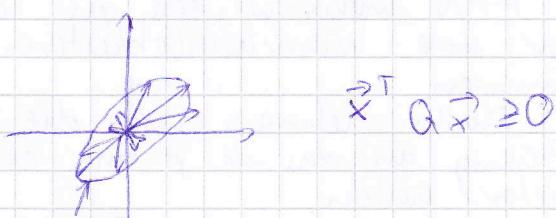
In der Kovarianzmatrix fassen wir alle Daten zusammen, und reduzieren dadurch die Komplexität.  
Bei  $N = 20^6$  und  $d = 20$  hat die Matrix die Größe  $20 \cdot 20 = 400$ .

Für einen Cluster betrachten wir also den Schwerpunkt  $\vec{M}$  und die Kovarianzmatrix  $\Sigma$ .

Zurück zur Formel der Normalverteilung:

$$N(\vec{x}) = \frac{1}{(2\pi)^d} \underbrace{\Sigma^{-\frac{1}{2}}}_{\text{Normierung auf Volumen der Gaußfläche auf 1.}} e^{-\frac{1}{2} (\vec{x} - \vec{M})^T \Sigma^{-1} (\vec{x} - \vec{M})}$$

Normierung auf Volumen der Gaußfläche auf 1.



$$P(E, \vec{x}) = P$$

Was sagt nun der Teil  $\frac{1}{2} (\vec{x} - \vec{M})^T \Sigma^{-1} (\vec{x} - \vec{M})$ ?

Dazu schauen wir uns mal einen Spezialfall an. Sei

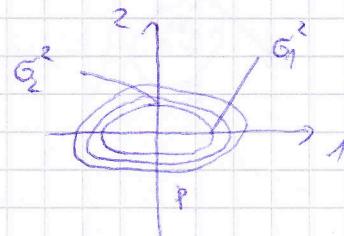
$$\Sigma = \begin{pmatrix} a & 0 \\ 0 & b \end{pmatrix}$$

$$= \begin{pmatrix} \Sigma(x_{11} - \vec{M}_1)^2 & 0 \\ 0 & \Sigma(x_{12} - \vec{M}_2)^2 \end{pmatrix}$$

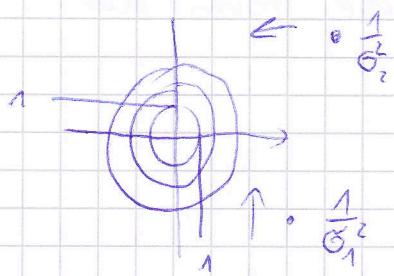
$$= \begin{pmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{pmatrix}$$

$$\Sigma^{-1} = \begin{pmatrix} \frac{1}{\sigma_1^2} & 0 \\ 0 & \frac{1}{\sigma_2^2} \end{pmatrix}$$

Standardabweichung  
in Dimension 2



Durch Skalierung der Standardabweichung erhält man dann so ein Bild:



In einer Formel:

$$(\vec{x} - \vec{M})^T \begin{pmatrix} \frac{1}{\sigma_1^2} & 0 \\ 0 & \frac{1}{\sigma_2^2} \end{pmatrix} (\vec{x} - \vec{M})$$

Das ist genau so, wie in der Formel.

Wenn die Kovarianzmatrix nicht diagonal ist, dann kann man sie durch

$$\Sigma = R^T D R$$

↑      ↑      ↗  
 Rotations- Diagona- Rotations-  
 matrix    matrix    matrix

$$\vec{x}^T \Sigma \vec{x} = (\vec{x}^T R^T) D (R \vec{x})$$

Dadurch löscht man die Korrelation der Daten durch Dehnung der Achsen.

## So. 11.2.2008 EM-Algorithmus

Vorlesung 4 Wir machen es überraschend, d.h. wir betrachten die Zettel mit:

$$\begin{array}{c} \vec{x}_1 \\ \vec{x}_2 \\ \vdots \\ \vec{x}_n \end{array}$$

Wir berechnen die Stellvertreter der Klassen und berechnen dann pro Klasse die Kovarianzmatrix und daraus kann man dann die Wahrscheinlichkeitsverteilung bekommen.

Ein neuer Punkt wird so klassifiziert, dass man für jedes Cluster die Wahrscheinlichkeit auf dem Punkt berechnet (durch die Formel von letzter Vorlesung) und dann die höchste nimmt und von dem Cluster das Label.

Spannender ist es, wenn wir unerwartet antreffen. Dann haben wir nur die Cluster ohne dessen Label zu kennzeichnen.

Wir machen jetzt folgendes: Wir wissen, dass es k Klassen sind und nehmen wie beim LBG zufällig k Stellvertreter  $\vec{x}_1$  bis  $\vec{x}_k$ . Als Kovarianzmatrizen wählen wir die Identitäten.

1) Teilung der Daten per Wahrscheinlichkeit (Expectation)

$$f_i(\vec{x}) = \frac{1}{2\pi} \sum_i e^{-\frac{1}{2} (\vec{x} - \vec{M}_i)^T \Sigma_i^{-1} (\vec{x} - \vec{M}_i)} \quad \text{mit } i=1 \dots k$$

Die Klasse von  $\vec{x}$  ist dann das  $i$  mit  $\max(f_i, x)$

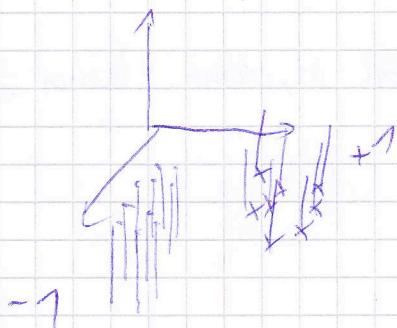
2) Berechne  $\vec{M}_1, \vec{M}_2, \dots, \vec{M}_k$  neu anhand der Erzeugten Klassifizierungen. Die Kovarianzmatrizen müssen natürlich auch angepasst werden. (Maximierung)

Das ganze iterieren wir nun wieder.

So ist das Clustering schöner, da auch die Form der Cluster berücksichtigt werden.

Jetzt was Neues - Lineare Regression.

Wir haben nun zwei Klassen



Wir hätten nun gerne eine Magicfunktion, die uns sagt, ob ein  $x$  -1 oder +1 hat.

Wir wollen, dass dieses  $f$  eine linear Funktion ist. Das nennt man Lineare Regression.

$$X = \begin{pmatrix} x_{11} & x_{12} & x_{13} & \dots & x_{1n} \\ x_{21} & x_{22} & x_{23} & \dots & x_{2n} \\ \vdots & & & & \\ x_{m1} & x_{m2} & x_{m3} & \dots & x_{mn} \end{pmatrix}$$

$$\vec{y} = X \cdot \vec{\beta}$$

$\vec{y}$  ist die Klassifizierung, die für jeden Datenpunkt  $x_i$  angibt, in welche Klasse der Punkt gehört.

$\vec{\beta}$  ist ein Vektor von Koeffizienten

Wie findet man nun  $\vec{\beta}$ ? Wir setzen mal  $\beta_0 = 0$ .

$$\vec{y} = X \cdot \vec{\beta}$$

$$x^{-1} \vec{y} = x^T X \cdot \vec{\beta} = \vec{\beta}$$

So leicht geht es leider nicht, da wir dazu eine quadratische Matrix brauchen, dass muss es aber nicht geben.

Daher optimieren wir diesen Ansatz.

Wir bauen uns nun also eine quadratische Matrix

$$x^T \vec{y} = x^T X \cdot \vec{\beta}$$

ist jetzt  
quadratisch

Das können wir nun invertieren; falls es existiert.

$$(x^T x)^{-1} x^T \vec{y} = (x^T x)^{-1} (x^T x) \vec{\beta} = \vec{\beta}$$

Falls unser  $\beta_0 \neq 0$  ist, müssen wir den  $\vec{y}$ -Vektor zentrieren, dann fällt  $\beta_0$  weg und es geht weiter wie oben.

Falls wir also nicht  $\vec{y} = x \cdot \vec{\beta}$  lösen können, suchen wir, ob wir eine Approximation minimal lösen können.

$$\text{dann } Q = (\vec{y} - X \vec{\beta})^T (\vec{y} - X \vec{\beta}) \text{ Quadratischer Fehler}$$

$$\frac{dQ}{d\beta} = -2\vec{x}^T(\vec{y} - \vec{x}\vec{\beta}) = 0$$

$$\vec{x}\vec{y} = \vec{x}^T\vec{x}\vec{\beta}$$

$$\vec{\beta} = (\vec{x}^T\vec{x})^{-1}\vec{x}\vec{y}$$

ist dann wieder das selbe wie vorne.

Wie lüftet man nun Vektoren ab?

$$\frac{d\alpha}{d\beta} = \begin{pmatrix} \frac{d\alpha}{dx_1} \\ \vdots \\ \frac{d\alpha}{dx_n} \end{pmatrix}$$

$$\text{wobei } \vec{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_m \end{pmatrix}$$

$$\frac{d\beta}{d\vec{x}} = \frac{d(\beta_1, \beta_2, \dots, \beta_n)}{d\begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_n \end{pmatrix}} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix} = I$$

12.11.2008

## VORVERARBEITUNG DURCH PRINCIPAL COMPONENT ANALYSIS (PCA)

(Hauptkomponentenanalyse)

Wir wollen, dass nach der Vorverarbeitung die Daten nicht mehr korreliert sind. Man kann auch Datenkompresion betreiben und es auch in der Computervision verwenden.



Wir haben hier also einen sichtbaren Cluster von Daten.

Die PCA soll uns nun gerade die Richtung aussuchen, in die das Cluster gerichtet ist.

Dadurch ist es mir z.B. möglich, durch eine Hauptkomponente ein Cluster zu beschreiben.

Besser ist natürlich, wenn man so viele wie Dimensionen findet, geordnet nach Aussagekraft in den Daten.

Algorithmus von Oja

- zufällig  $w_j, g$  wählen (vorher müssen die Daten zentriert sein)
- $\vec{w}_j^{(i+1)} = \vec{w}_j^{(i)} + \gamma \langle \vec{x}_j, \vec{w}_j^{(i)} \rangle (\vec{x}_j - \langle \vec{x}_j, \vec{w}_j^{(i)} \rangle \vec{w}_j^{(i)})$ ,  $0 < \gamma < 1$
- $s := \langle \vec{x}_j, \vec{w}_j^{(i)} \rangle$
- $\gamma^{(i+1)} = \gamma^{(i)} / 2$  oder anderes verkleinern

Das Iterieren wird nun länger. Entweder, bis wir keine Zeit oder Lust haben oder bis  $\Delta(w)$  hinreichend klein ist.

Jetzt haben wir die 1. Hauptkomponente gefunden. Um die nächste zu finden, tössen wir nun den Einfluss der Hauptkomponente aus den Daten.

$$\forall j \text{ do: } \vec{x}_j^{(i+1)} = \vec{x}_j - \langle \vec{x}_j, \vec{w} \rangle \vec{w}$$

Usw. für alle Hauptkomponenten.

Hauptkomponenten aus der Kovarianzmatrix

Mittelpunkt der Daten  $\vec{\mu} = \frac{\sum_{i=1}^N \vec{x}_i}{N}$

$$\frac{\sum_{i=1}^N (\vec{w}^T \vec{x}_i - \vec{w}^T \vec{\mu})}{N} = \vec{w}^T \underbrace{\sum_i}_{\text{Kovarianzmatrix}} \vec{w}$$

Den Ausdruck wollen wir nun maximieren, mit der Maßgabe, dass  ~~$w^T w = 1$~~   $w^T w = 1$

$$z = \vec{w}^T \sum \vec{w} - \lambda (1 - \vec{w}^T \vec{w}) \quad (\text{la Grange Multiplikator})$$

$$0 = \vec{w}^T \sum - \lambda \vec{w} \Leftrightarrow \sum \vec{w} = \lambda \vec{w} \quad \text{Eigenwert}$$

Also maximiert  $\vec{w}$  die Kovarianzmatrix, wenn  $\vec{w}$  Eigenvektor von  $\sum$  ist.

$$\vec{w}^T \sum \vec{w} = \lambda$$

$\vec{w}$  muss Eigenvektor von  $\sum$  sein, damit die Varianz zu Maximieren.

$\lambda_i$  zum i-ten Eigenvektor muss maximal gewählt werden.

Wenn wir jetzt die  $w_k$  in Hauptkomponenten haben, können wir die Kovarianz entfernen.

$$\vec{x}^T \vec{w}_1 \vec{w}_1 = \vec{a}_1$$

$$\vec{x}^T \vec{w}_2 \vec{w}_2 = \vec{a}_2$$

$$\sum \vec{a}_i = \vec{x}$$

Wir drehen also das Koordinatensystem so, dass die Hauptkomponenten auf den Achsen liegen.

$$\vec{w} = \begin{pmatrix} w_1 & | & | & | \\ | & w_2 & w_3 & \dots & w_N \end{pmatrix}$$

$$y = L \vec{w}^T \vec{x} \quad \text{mit} \quad L = \begin{pmatrix} f_1 & & \\ & \ddots & \\ & & f_n \end{pmatrix} \quad (\text{eigene Vektoren})$$

# Merkverarbeitung FISHER - LINEAR - DISCRIMINANT

Vorlesung 6 Bisher haben wir mehrere lineare Funktionen

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

19.11.08

implizit ~~oder~~ oder explizit betrachtet, um  $x$  zu klassifizieren.

Die Klassifizierungregel ist

$$(1) \quad f(x) = \begin{cases} > 0 \Rightarrow x \in K_1 \\ < 0 \Rightarrow x \in K_2 \end{cases}$$

## Beispiele

- \* k-nn  $\mathbf{w} = p_1 - p_2$   
 $w_0 = \frac{\|p_2\|^2 - \|p_1\|^2}{2}$

$p_i$  ist der Prototyp der Klasse  $K_i$

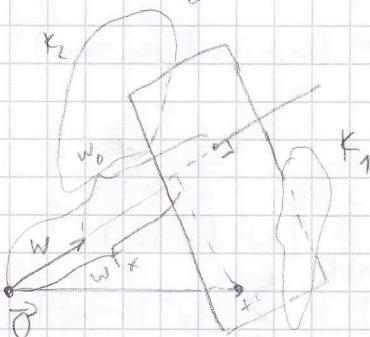
## \* Least squares

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$w_0 = \frac{1}{N} \sum_{i=1}^N y_i$$

Annahme:  $\|\mathbf{w}\|=1$

- Die Vektoren  $\mathbf{x}$ , die  $f(\mathbf{x})=0$  erfüllen, definieren eine lineare Hyperebene in  $\mathbb{R}^N$
- $\mathbf{w}$  ist orthogonal zur Hyperebene  $H$
- $w_0$  ist der orthogonale Abstand zwischen Ursprung und  $H$



Durch die Matrixmultiplikation kann die Dimension auf eine Ebene reduziert werden.

Die Regel (1) sagt uns, auf welcher Seite der Ebene  $H$  der Vektor  $x$  liegt.

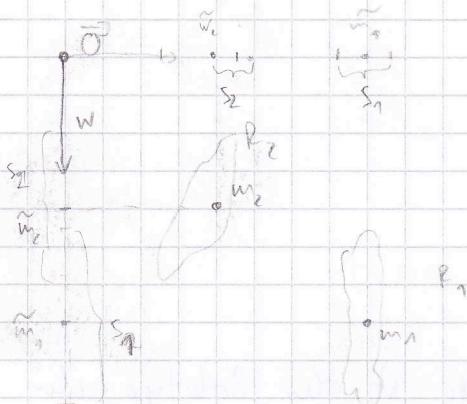
Dieser Ansatz versucht

- Probleme mit der Dimensionalität zu vermeiden, und zwar mit einer linearen Kombination (Transformation) der Merkmale + ein Qualitätskriterium

Fisher: Finde die Richtung  $w$ , wo die Projektion der Daten "am besten" zu trennen sind!

Wie ist das Kriterium zu definieren?

Zuerst geometrisch...



$$\max_w (\tilde{m}_1 - \tilde{m}_2)^2 S_i$$

streuung

Kriterium: Mittelwert und Varianz der Projektionen

groß  $(\tilde{m}_1 - \tilde{m}_2)$   
kleine Varianz

Jetzt algebraisch...

Mittelwert der Projektionen  $\tilde{m}_i = \frac{1}{n_i} \sum_{x \in K_i} \tilde{x}_i$  mit  $\tilde{x} = w^T x$  und  
 $n_i$  Anzahl der Elemente in  $K_i$ .

$$\begin{aligned} &= \frac{1}{n_i} \sum_{x \in K_i} w^T x \\ &= \frac{w^T}{n_i} \sum_{x \in K_i} x_i = w^T \left( \frac{1}{n_i} \sum_{x \in K_i} x \right) = w^T m_i \end{aligned}$$

$$\begin{aligned} (\tilde{m}_1 - \tilde{m}_2)^2 &= (w^T m_1 - w^T m_2)^2 = [w^T (m_1 - m_2)]^2 \\ &= w^T (m_1 - m_2) w^T (m_1 - m_2) = w^T (m_1 - m_2)(m_1 - m_2)^T w \\ &= w^T S_B w \quad \text{mit } S_B = (m_1 - m_2)(m_1 - m_2)^T \end{aligned}$$

Variante (oder etwas Ähnliches)  $S_i^2 = \sum_{x \in K_i} (\tilde{x} - \tilde{m}_i)^2$

$$\begin{aligned} &= \sum_{x \in K_i} (w^T x - w^T m_i)^2 = \sum_{x \in K_i} [w^T (x - m_i)]^2 \\ &= \sum_{x \in K_i} w^T (x - m_i) w^T (x - m_i) = \sum_{x \in K_i} w^T (x - m_i)(x - m_i)^T w \\ &= w^T \left( \sum_{x \in K_i} (x - m_i)(x - m_i)^T \right) w \end{aligned}$$

$$= w^T S_i w \quad \text{mit } S_i = \sum_{x \in K_i} (x - m_i)(x - m_i)^T \text{ der scat. Matrix. Das ist auch}$$

$$S_i = n_i \sum$$

$$(\tilde{m}_1 - \tilde{m}_2) = w^T S_B w \text{ maximieren}$$

$$S_1^2 + S_2^2 = w^T S_1 w + w^T S_2 w = w^T S_W w \text{ mit } S_W = S_1 + S_2 \text{ minimieren}$$

Numerisches Kriterium:

$$(0) \max_w J(w) = \frac{w^T S_B w}{w^T S_W w}$$

Beobachtung: Gegeben  $\alpha \neq 0$ , es gilt ein

$$(*) \quad J(\alpha \cdot w) = J(w)$$

Aus (\*), da das Optimierungsproblem (0) ist äquivalent zu

$$\max w^T S_B w \text{ unter } w^T S_W w = 1$$

$$\text{Lagrange-Multiplikation} \Rightarrow \max_{w, \lambda} L(w, \lambda) = w^T S_B w - \lambda (w^T S_W w - 1)$$

$$\text{Finden wir } w: \frac{\partial L(w, \lambda)}{\partial w} = 2 S_B w - 2 \lambda S_W w = 0$$

$$\Rightarrow S_B w = \lambda S_W w$$

Das ist ein allgemeineres Eigenwertproblem.

$$S_B w = \lambda S_W w$$

$$\begin{aligned} w &= \underbrace{\frac{1}{\lambda} S_W^{-1} S_B}_{\text{skalar}} w \\ &= \frac{1}{\lambda} S_W^{-1} (m_1 - m_2) (m_1 - m_2)^T w \\ &= \frac{1}{\lambda} S_W^{-1} (m_1 - m_2) w^T (m_1 - m_2) \\ &= \frac{1}{\lambda} S_W^{-1} (m_1 - m_2) (\tilde{m}_1 - \tilde{m}_2) \\ &= \frac{\tilde{m}_1 - \tilde{m}_2}{\lambda} S_W^{-1} (m_1 - m_2) \end{aligned}$$

$J(w)$  ist maximal, wenn  $w = \frac{(\tilde{m}_1 - \tilde{m}_2)}{\lambda} S_W^{-1} (m_1 - m_2)$  ist

$$J(w) = J\left(\frac{\tilde{m}_1 - \tilde{m}_2}{\lambda} S_W^{-1} (m_1 - m_2)\right)$$

$$\stackrel{(*)}{=} J(S_W^{-1} (m_1 - m_2))$$

$$\text{das heißt } w = S_W^{-1} (m_1 - m_2)$$

Das Wg behandeln wir im Tutorium...

$J(w)$  ist maximal, wenn

$$w = \frac{(\tilde{m}_1 - \tilde{m}_2)}{\|S_w^{-1}(m_1 - m_2)\|}$$

ist.

$$\begin{aligned} J(w) &= J\left(\frac{(\tilde{m}_1 - \tilde{m}_2)}{\|S_w^{-1}(m_1 - m_2)\|}\right) \quad (*) \\ &= J(S_w^{-1}(m_1 - m_2)) \end{aligned}$$

d.h.

$$w = S_w^{-1}(m_1 - m_2)$$

$$f(x) = w^T x + w_0$$

fehlt  
noch

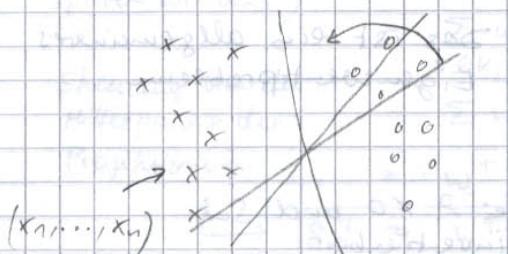
Hatten schon: Lineare Trennung

26.11.08

- Fisher-Distanzmaße

Ahnliches Konzept: Perzepton:

- iteratives Verfahren
- anfänglich ausgewählte Trennung wird immer weiter verbessert



[ $w_i x_i$  ist Multiplikation von Vektoren]

$$\begin{matrix} x_1 & \xrightarrow{w_1} \\ x_2 & \xrightarrow{w_2} \\ \vdots & \vdots \\ x_n & \xrightarrow{w_n} \end{matrix} \xrightarrow{+} \begin{cases} 1 & \\ 0 & \end{cases}$$

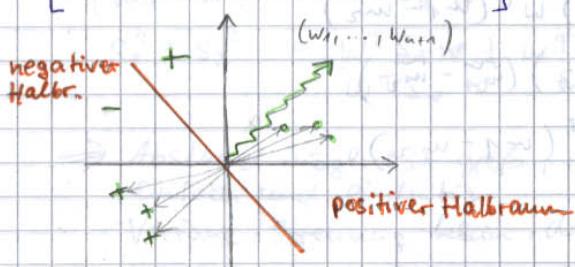
Lineare Kombination:

$$w_1 x_1 + w_2 x_2 + \dots + w_n x_n + w_{n+1} \geq 0$$

⇒ Ausgabe 1, sonst 0

↳ Äquivalente Schreibweise:

$$(x_1, x_2, \dots, x_n, 1) \cdot (w_1, \dots, w_n, w_{n+1}) \geq 0$$



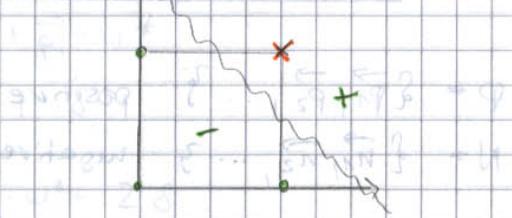
$x_1$	$x_2$	AND
0	0	0
0	1	0
1	0	0
1	1	1

Klassifizierung

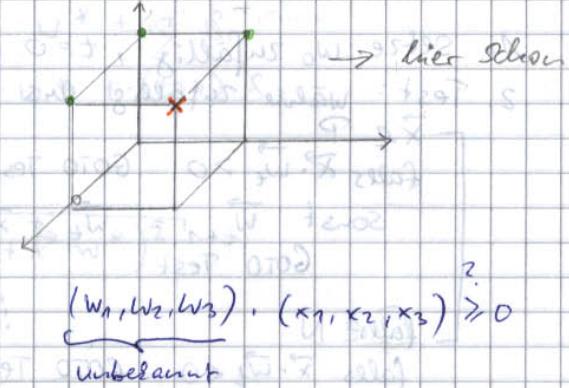
$x_1$	$x_2$	$x_3$	
0	0	1	0
0	1	1	0
1	0	1	0
1	1	1	1

zusätzl.  
Dimension  
einführen

→ keine Trennung durch Ursprung möglich



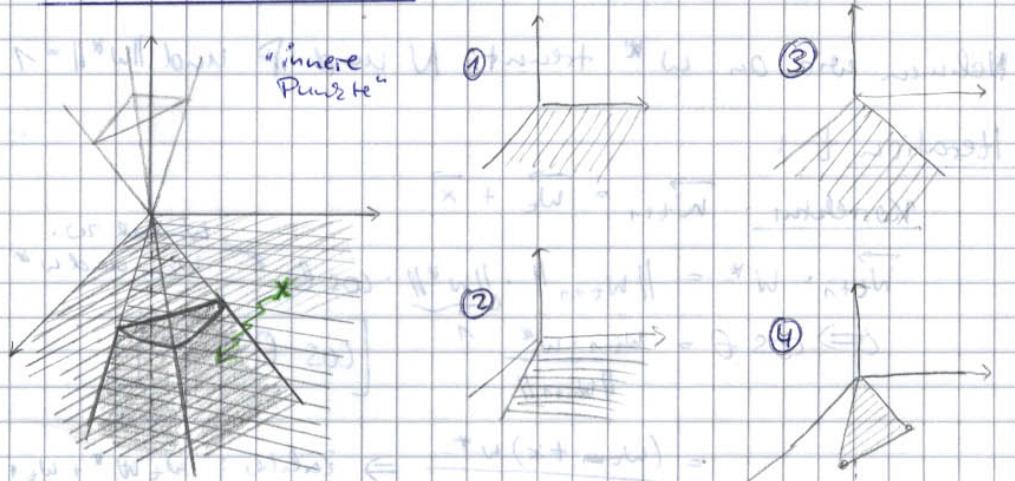
→ hier schon



- ①  $w_1 \cdot 0 + w_2 \cdot 0 + w_3 < 0$
- ②  $w_1 \cdot 0 + w_2 \cdot 1 + w_3 \cdot 1 < 0$
- ③  $w_1 \cdot 1 + w_2 \cdot 0 + w_3 \cdot 1 < 0$
- ④  $w_1 \cdot 1 + w_2 \cdot 1 + w_3 \cdot 1 \geq 0$

- o  $w_3$  muss negativ sein
- o  $w_2 + w_3$  muss negativ sein
- o  $w_1 + w_3$  muss negativ sein
- o alle zusammen positiv

Lösung geometrisch:



- o Lineare Programmierung zur Lösung bei vielen Variablen (schwierig)

oder

- o Approximativer Verfahren

## Approximativer Verfahrens ↳ Perceptron - Algorithmus

$P = \{\vec{p}_1, \vec{p}_2, \dots\}$  positive Vektoren

$N = \{\vec{n}_1, \vec{n}_2, \dots\}$  negative Vektoren



1. Setze  $w_0$  zufällig,  $t=0$

2. Test: wähle zufällig aus  $N \cup P$  ein  $\vec{x}$

$\vec{x} \in P$   
falls  $\vec{x} \cdot \vec{w}_t > 0$  GOTO Test

Sonst  $\vec{w}_{t+1} = \vec{w}_t + \vec{x}$ ,  $t = t+1$   
GOTO Test

$\vec{x} \in N$   
falls  $\vec{x} \cdot \vec{w}_t < 0$  GOTO Test

Sonst  $\vec{w}_{t+1} = \vec{w}_t - \vec{x}$ ,  $t = t+1$   
GOTO Test

- positives Skalarprodukt heißt,
- $\vec{x}$  liegt auf der richtigen Seite der Trennungsebene
- $\Rightarrow$  Winkel  $< 90^\circ$

Wenn es Lösung gibt, konvergiert der Algorithmus, dauert aber womöglich sehr lange.

Nehmen wir an  $w^*$  trennt  $N$  und  $P$  und  $\|w^*\| = 1$

Iteration  $t$ :

Korrektur:  $\vec{w}_{t+1} = \vec{w}_t + \vec{x}$

$$\vec{w}_{t+1} \cdot w^* = \|w_{t+1}\| \cdot \underbrace{\|w^*\|}_{\text{Winkel zw. } w_t \text{ und } w^*} \cdot \cos \theta$$

$$\Leftrightarrow \cos \theta = \frac{w_t \cdot w^*}{\|w_{t+1}\|} \quad [\cos \theta \leq 1]$$

$$= \frac{(w_t + x) \cdot w^*}{\|w_{t+1}\|} \Rightarrow \text{Zähler:}$$

$$w_t \cdot w^* + w^* \cdot x$$

$$\geq w_t \cdot w^* + \delta$$

$$\delta = \min \{p_1 \cdot w^*, \dots, p_n \cdot w^*\}$$

$$\delta > 0$$

Zähler:

$$\begin{aligned} w_t \cdot w^* &= (w_{t+n} \cdot \vec{p}_i) w^* \\ &= w_{t+n} \cdot w^* + \vec{p}_i \cdot w^* \\ &\geq w_{t+n} \cdot w^* + \delta \\ \Rightarrow w_{t+n} \cdot w^* &> w_{t+n} \cdot w^* + 2\delta \\ &\geq w_{t+n} \cdot w^* + 3\delta \\ &> w_{t+n} \cdot w^* + (t+1)\delta \end{aligned}$$

Neuer:

$$\|w_{t+1}\| = \sqrt{(w_t \cdot w^* + \vec{x} \cdot w^*)^2}$$

⋮

$$\sqrt{\dots (t+1) \dots}$$

$\Rightarrow t$  kann nicht  $\infty$  werden, irgendwo stoppt der Algorithmus

# Mustererkennung Heuk: Verallgemeinerung des Perzeption-Algorithmus.

Vorlesung 8

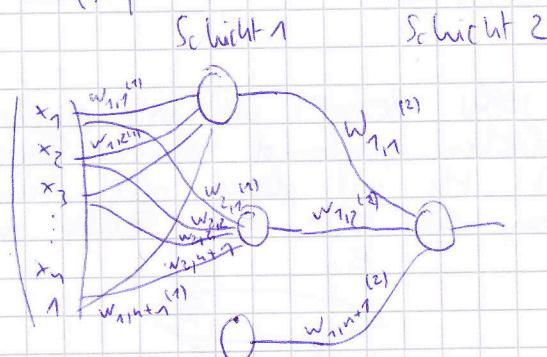
Hier wollen wir mehrere lineare Trennungen finden, die sich ergänzen. z.B. sinnvoll sei sowas:



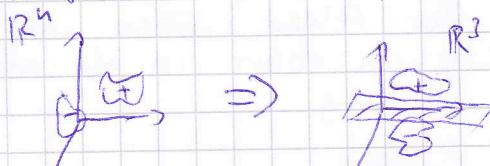
$$\vec{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \text{ wollen wir klassifizieren.}$$

Die  $w_{ij}^{(k)}$  sind Gewichte.

Das nennt man dann "hierarchisches Netz"



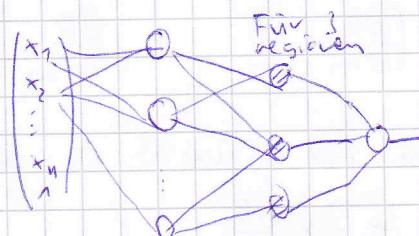
Was wir hier machen ist eine Projektion von  $\mathbb{R}^n$  zu  $\mathbb{R}^3$



Bei sowas



macht man das so:

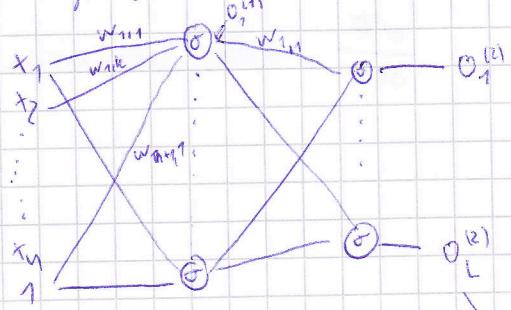


Wir machen das ganze stetig, so dass die Ausgabe einer Klassifikation nur 1 oder -1 ist, statt der normalen Signum-Funktion beobachten wir aber

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

So sind die Ausgaben der Elemente "grauwerthaft", so dass man sie ein wenig mit Wahrscheinlichkeiten vergleichen kann.

Im Allgemeinen oft:

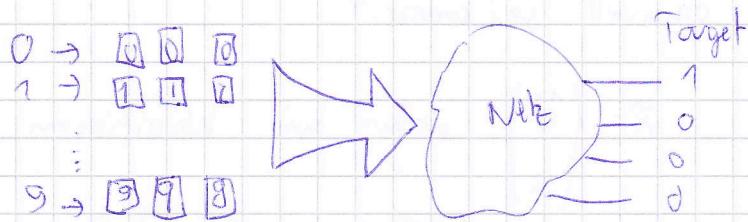


hElement | Element Mehrp  
klasse.

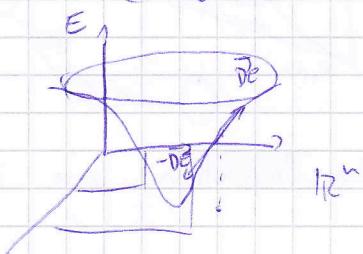
Bei Ziffern könnte es so aussehen:

$$\vec{x} \rightarrow \begin{cases} 1 & \text{if } x \in [0, 1] \\ 0 & \text{else} \end{cases}$$

Die Aufgabe ist jetzt, die Trainingsmenge so durch das Netz zu jagen, dass das Target erreicht wird.



Für jedes Beispiel, was wir so durchjagen, berechnen wir den euklidischen Abstand zum Target, der uns den Fehler verrät. Ich will dann die Kombination aus Gewichten finden, für die der Fehler minimal ist.



Wir fangen mit zufälligen Gewichten an und schauen uns den Fehler an und suchen dann, in welche Richtung der Fehler höher wird, und gehen dann in die andere Richtung.

Fehlerfunktion ist Funktion aller Gewichte.

$$E(w_1, w_2, \dots, w_m)$$

$$\vec{\nabla} E = \left( \frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \dots, \frac{\partial E}{\partial w_m} \right)$$

Die Korrektur für Gewicht  $w_i$  sei  $\Delta E$  ist

$$\Delta w_i := -\gamma \frac{\partial E}{\partial w_i}$$

Die Frage ist, wie man die partiellen Ableitungen bekommt.

$x \rightarrow \{w_i\} \rightarrow E$  wir verändern ein wenig  $w_i$  aus  $w_i + \epsilon$  und suchen was dann am, wie sich der Fehler ändert:  $E \rightarrow E + \epsilon'$

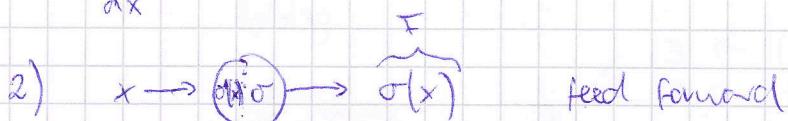
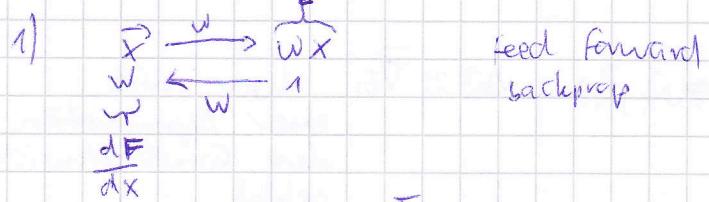
$$\text{Dann gilt } \frac{\partial E}{\partial w_i} \approx \frac{(E + \epsilon') - E}{(w_i + \epsilon) - w_i} = \frac{\epsilon'}{\epsilon}$$

Hilfsproblem:



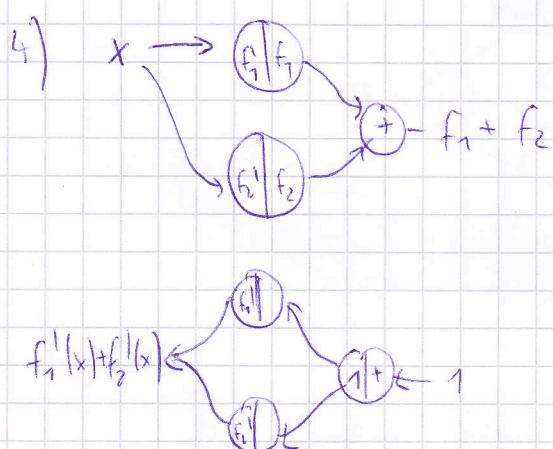
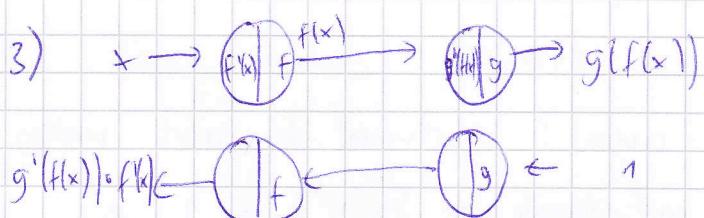
$$\frac{\partial E}{\partial x}$$
 wollen wir wissen.

Wir unterscheiden mehrere Fälle: wir rechnen mit Backprop

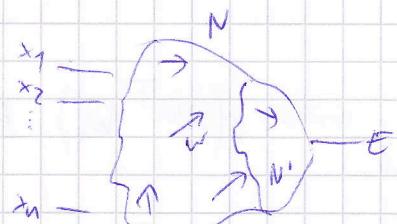


$\sigma'(x) = \sigma(x)(1-\sigma(x))$

backprop. Die Regel ist hier: linke Seite des Neurons verwenden und mit der Einheit multiplizieren.



Jetzt können wir wieder zurück zum Hauptproblem gehen.



Zunächst machen wir feed forward. Dann schreien wir es wieder rückwärts, indem wir eine 1 reingetragen.

Da wir keine Kreise haben, gibt es ein Subgraph  $N^i$ , das nur von der Veränderung beeinflusst wird. Man macht also beim Backprop nur  $N^i$  betrachten und man ist fertig.

Das macht man für alle Gewichte dann. Und für alle Beispiele.

$$\frac{\partial E}{\partial w} = \frac{\partial (E_1 + E_2 + \dots)}{\partial w}$$

# Zwei Varianten von Backprop

Vorlesung 9

10.12.2008

1. On-line

$$(\vec{x}_1, \vec{t}_1)$$

$$\rightarrow (\vec{x}_1, \vec{t}_1) \quad \Delta \vec{w} \approx \vec{\nabla} E_i$$

hier wird

eine Approximation  
der Gradientenrichtung  
geleistet

2. Off-line  
(batch)

$$(\vec{x}_1, \vec{t}_1) \rightarrow E_1$$

$$(\vec{x}_2, \vec{t}_2) \rightarrow E_2$$

$$\vdots$$

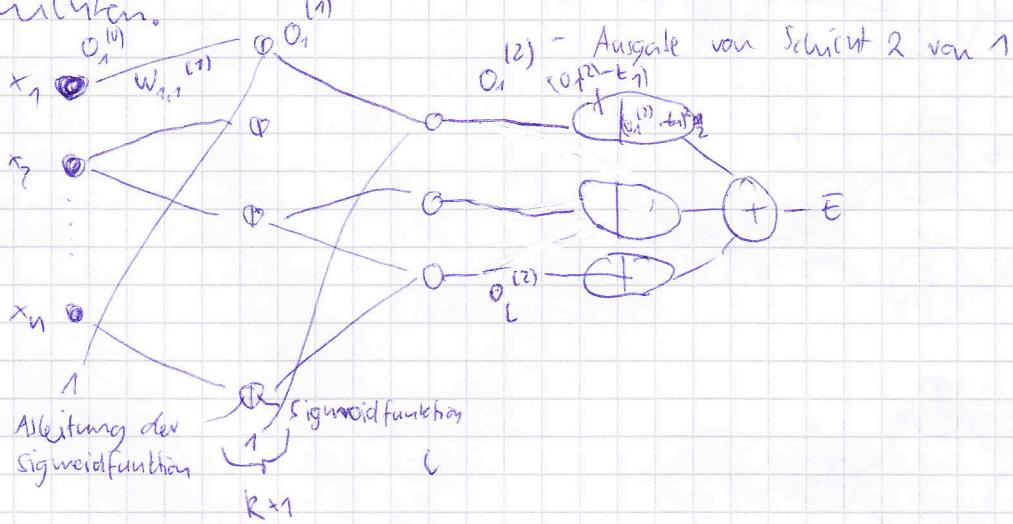
$$\Delta \vec{w} \approx -\vec{\nabla} E$$

$$(\vec{x}_N, \vec{t}_N) \rightarrow E_N$$

Gradientenabschlag

Viele benutzen on-line Verfahren, da man hier nicht so schnell  
in ein lokales Minimum fällt, da man nicht direkt den  
Gradientenrichtung folgt.

Wie implementiert man das nun? Erstmal benutzt man  
Schichten.



Matrix der  
Gewichte mit  
 $(n+1) \times k$

Der Feedforward Schritt besteht daran, dass man  $s(\vec{o}^{(1)} \top W_1)$  berechnet.  
Das gibt alle Ausgänge für die Hidden Layer.

$$o_1^{(1)\top} = s(\vec{o}^{(1)\top} W_1)$$

Wenn man noch die konstante 1 hinzufügt zum weiteren  
nennt man das jetzt  $\vec{o}^{(1)\top}$ . Für die letzte Schicht ist es dann

$$\vec{o}^{(2)\top} = s(\vec{o}^{(2)\top} W_2)$$

Wir müssen uns natürlich noch die Ausleitungen mitreichen. Das  
machen wir auch in Matrizen.  $D_1$  für die 1. Schicht,  $D_2$  für die  
2.

$$D_1 = \begin{pmatrix} O_1^{(1)} (1-O_1^{(1)}) & & & \\ & O_2^{(1)} (1-O_2^{(1)}) & & \\ & & \ddots & \\ & & & O_k^{(1)} (1-O_k^{(1)}) \end{pmatrix}$$

$$D_2 = \begin{pmatrix} O_1^{(2)} (1-O_1^{(2)}) & & & \\ & O_2^{(2)} (1-O_2^{(2)}) & & \\ & & \ddots & \\ & & & O_k^{(2)} (1-O_k^{(2)}) \end{pmatrix}$$

Für den Fehler bzw. die Ableitung benutzen wir einen Vektor

$$\vec{E} = \begin{pmatrix} O_1^{(1)} - t_1 \\ O_2^{(1)} - t_2 \\ \vdots \\ O_k^{(1)} - t_k \end{pmatrix}$$

$\vec{f}^{(2)} = D_2 \vec{E}$ , jetzt können wir schon Korrekturen berechnen.

$$\Delta w_{1,1}^{(2)} = -\gamma \vec{f}_1^{(2)} O_1^{(1)}$$

Jetzt für alle Gewichte.

$$\Delta w_2^T = -\gamma \vec{f}^{(2)} \vec{O}^{(1)T}, \quad \Delta w_1^T = -\gamma \vec{O}^{(1)} \vec{f}^{(2)T}$$

$$\vec{f}^{(1)} = D \vec{w}_2 \vec{f}^{(2)}$$

$w_2$  ohne letzte Zeile!

$$\Delta w_1^T = -\gamma \vec{f}^{(1)} \vec{O}^{(0)T}$$

Man kann also mit nur ein paar Zeilen das Programmieren.

Jetzt kümmern wir uns um das  $\gamma$ . Für jedes  $w_{ij}$  haben wir ein  $r_{ij}$ .

$$r_i^{(k+1)} = \begin{cases} r_i^{(k)} u & : \nabla_i E^{(k)}, \frac{\partial E^{(k+1)}}{\partial w_{ij}} \geq 0 \\ r_i^{(k)} d & " \quad < 0 \end{cases}$$

$$\Delta w_i^{(k)} = r_i^{(k)} \nabla_i E^{(k)}, \quad u=1,5 \quad d=0,8 \text{ oder so}$$

(Algorithmus von Silva & Almeida)

Das heißt, wenn die partielle Ableitung  $\frac{\partial E^{(k)}}{\partial w_i}$  in der  $k$ -ten und  $k+1$ -ten Iteration das gleiche Vorzeichen hat, dann wird  $\gamma$  größer gemacht.

Eine andere Heuristik wäre "Momentum".

$$\Delta w_i^{(k)} = -\eta \frac{\partial E}{\partial w_i} + \alpha \Delta w_i^{(k-1)} \quad \alpha < 1$$

## RPROP - Algorithmus

Es gibt eine Verbesserung für Silva & Almeida. Da ist das Problem, dass die Schritte die man macht, sehr klein um das Minimum zu finden. Zudem kann es passieren, dass man immer das Minimum überspringt.

$$g_i^{(k+1)} = \begin{cases} \min(g_i^{(k)} u, g_{\max}) : & \nabla_i E_i^{(k)} \nabla_i^{\text{next}} E_i^{(k+1)} > 0 \\ \max(g_i^{(k)} d, g_{\min}) : & " \quad \quad \quad \leq 0 \\ g_i^{(k)} & " \quad \quad \quad = 0 \end{cases}$$

$$\Delta w_i^{(k)} = -g_i^{(k)} \operatorname{sgn}(\nabla_i E^{(k)})$$

## Blind source separation (ICA)

Vorlesung 1  
17.12.2008

Quelle



:

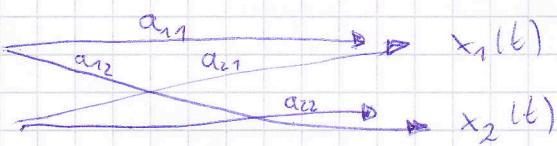
$$s_1(t)$$

$$s_2(t)$$

Mischungen



$$H M$$



Mischmatrix A

Die Mischungen  $x_i(t)$  sind eine lineare Funktion von dem Quellsignal  $s_i(t)$

$$x_1(t) = a_{11}s_1(t) + a_{12}s_2(t) \quad (*)$$

$$x_2(t) = a_{21}s_1(t) + a_{22}s_2(t)$$

$$A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \text{ Mischmatrix. } \begin{pmatrix} x_1(t) \\ x_2(t) \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \cdot \begin{pmatrix} s_1(t) \\ s_2(t) \end{pmatrix}$$

$$x(t) = A^T s(t)$$

Was ist jetzt das blind-source-separation Problem?

(Gegeben  $x(t)$  mit  $t=1, \dots, T$ . Finden den Quellsignalen und die Mischmatrix.

Wir gucken es uns einmal geometrisch an:

Für jeden Zeitpunkt  $t = 1, \dots, T$  betrachten wir den Vektor  $x(t)$  als ein Datum, d.h. wir haben eine Trainingsmenge

$$\{x(1), x(2), \dots, x(T)\}$$

Nehmen wir an  $s_i$  ist dünnsetzt, d.h.  $s_i(t) \neq 0$  für wenige  $t$ .

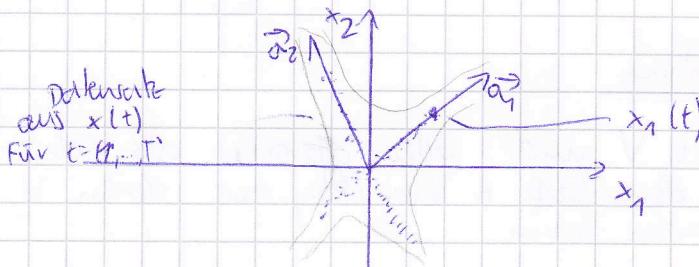
Das bedeutet,

$s_1(t) \neq 0 \Rightarrow s_2(t) = 0$  mit hoher Wahrscheinlichkeit

Aus (\*) folgt dann

$$\begin{aligned} x_1(t) &= a_{11} s_1(t) + 0 & \vec{a}_1 &= \begin{pmatrix} a_{11} \\ a_{12} \end{pmatrix} \\ x_2(t) &= a_{21} s_1(t) + 0 \end{aligned}$$

d.h.  $x(t)$  ist ein Vielfaches von  $\vec{a}_1$



Die Daten akkumulieren sich um  $\vec{a}_1$  herum.

Die  $\vec{a}_i$  sind Zeilen von  $A$ .

$\vec{a}_i$  sind die Richtungen mit viel "Informationsgehalt" (Wir denken  $\vec{a}_i$  an PCA, Faktoren usw.)

Eine Lösung des BSS Problems ist eine unmixing Matrix  $W$  und die Signale  $y(t)$ , die die Quellen approximieren.

$$y(t) = W^\top x(t)$$

$$y(t) = \begin{pmatrix} a_1 s_1(t) \\ a_2 s_2(t) \end{pmatrix}$$

(Wir denken nur  $W$  so ein wenig als  $A^{-1}$ )

Der Ansatz zum das BSS Problem zu lösen, ist folgendes:

Wir nehmen an, dass die Quellsignale eine Eigenschaft  $P$  erfüllen. Diese Eigenschaft lässt sich numerisch mit einer Qualitätsfunktion  $q$  messen.

Dann maximieren wir  $q(y(t)) = q(W^\top x(t))$ .

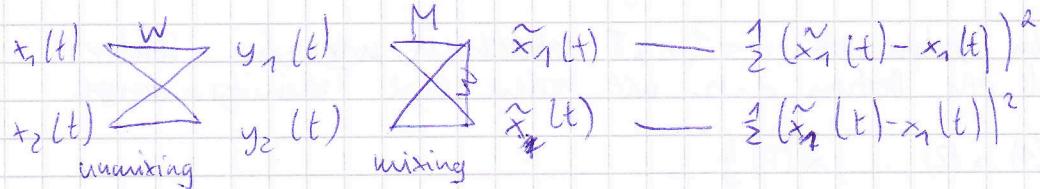
$$W^* = \arg \max_W q(W^\top x(t))$$

$$\tilde{y}_i(t) = \theta_{y_i}(t_1) + \dots + \theta_{y_i}(t_k) \quad \text{die } t_i \text{ sind dabei Zeitpunkte in der Vergangenheit}$$

$$(\tilde{y}_i(t) - y_i(t))^2 \cdot \frac{1}{2} = \text{Fehler}$$

je kleiner der Fehler, desto linearer ist  $y_i(t)$

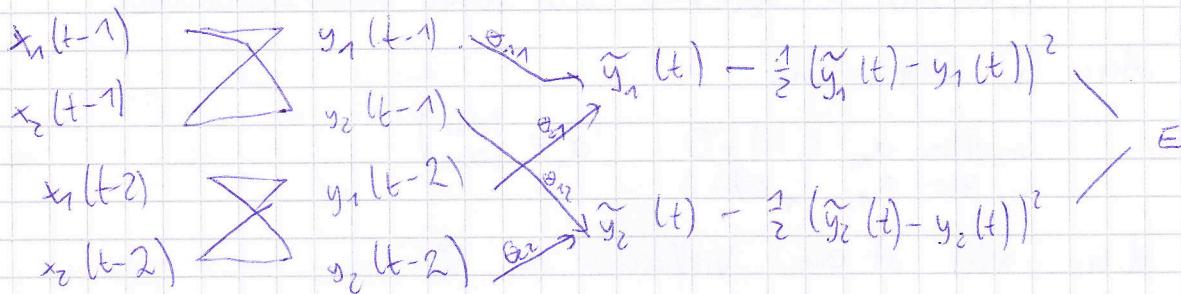
## Nekroholle Nebe



$$y(t) = W^T x(t)$$

$$W^{-1} = M$$

$$\tilde{x}(t) = M^T y(t)$$



## Non-negative matrix factorization

Vorlesung 1  
7.1.2009

- Wir haben unsere Datenbank  $V$  und wollen daraus alle Szenen Merkmale (Basis)  $W$  und den Code  $H$  ermitteln, aus denen  $V$  am besten rekonstruiert werden kann

$$V \approx V' = W \cdot H$$

- Das Verfahren heißt nicht negativ, da man negative Komponenten verhindern. Faktorisierung heißt gerade das Produkt  $W \cdot H$  aus  $V$  zu finden

Der Algorithmus funktioniert wieder mit Expectation-Maximization und zwar Ändern wir  $W$  und  $H$  immer abwechselnd, bis wir gut gelung sind.

Den Unterschied zweier nacheinander folgender Iterationen kann man entweder über den normalen euklidischen Abstand oder über Faktoren.

$$1. W_{ia} \leftarrow w_{ia} \sum_{j=1}^{V_{im}} \frac{v_{ij}}{\sum_{i=1}^{V_{im}} v_{ij}} h_{aj}$$

Korrektur: Der 1. Index ist die Zeile im Merkmal, der 2. gibt den Wert und der 3. gibt das Merkmal an

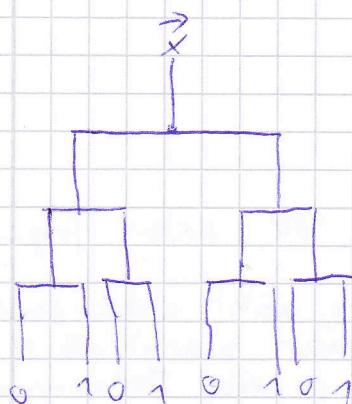
$$2. w_{ia} \leftarrow \frac{w_{ia}}{\sum_j w_{ja}}$$

Normieren

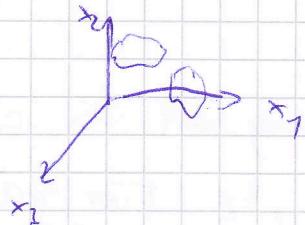
$$3. h_{aj} \leftarrow h_{aj} \sum_i \frac{w_{ia}}{\sum_{i=1}^{V_{im}} w_{ia}} v_{im}$$

Code korrigieren

## ENTSCHEIDUNGSBÄUME



So was liefet sich an, wenn man mehrere Dimensionen hat



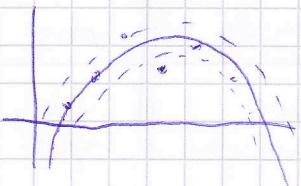
Man kann hier dann z.B. so fragen, ob im 1. Zweig  $x_1 > 0,5$  ist und wann je nach dem einen ist nimmt. Dann kann man in der nächsten Ebene über  $x_2$  fragen usw.

Jetzt ist die Frage, wie man die Reihenfolge der Frage geschickt wählt.

### Bagging

### Bootstrap aggregation:

Die Idee ist folgende:



Wir wollen die beste Anpassung an die Punkte haben, am besten ohne Ausreißer. Dann wollen wir ein Konfidenzintervall, in das z.B. 95% der Punkte reinfallen.

$$\text{Trainingsmenge } T = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_N\}$$

Die Bootstrap Trainingsmenge funktioniert so, dass man  $N$  mal zufällig ein Element selektiert und Ersetzung wieder zurücklegt.

Dann wird das Polygon ausgeglichen

Das gibt uns in Funktionen, die uns den Durchschnitt von diesen  $m$  liefern, die wir dann nehmen.

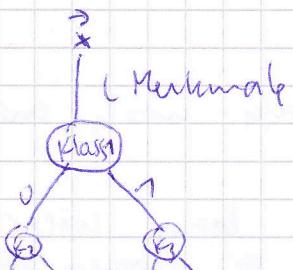
Das praktische ist, dass sich das Konfidenzintervall automatisch aus allen Funktionen ergibt.

Der Entscheidungsbaum sieht so aus



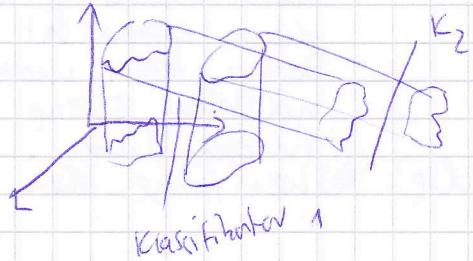
### Random Forests

Hier hat man  $N$  Punkte,  $K$  Merkmale, und die Lernkurve hat Fehler. D.h. wir können damit umgehen, dass Merkmal bei Beispielen fehlen.



Nur 2 klassen

Die Idee ist so ein wenig wie bei Fisher. Nur haben folgendes:



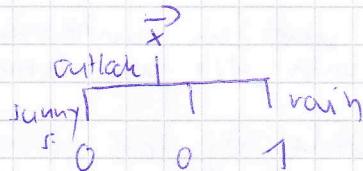
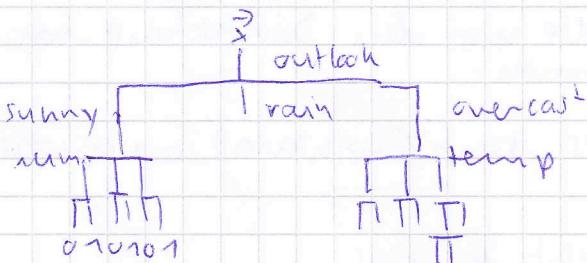
Für jede Ebene wählt man verschiedene Merkmale. Man macht das so lange bis man mal so eine Trennung erhält, wo auf einer Seite des Klassifikators nur noch ein Punkt liegt. Dann treffen wir dann tatsächlich die Entscheidung.

Die Höhe des Baums hängt von der Trainingsmenge ab, wie gut sie sich trennen lässt.

### ID3 (Inductive - Decision - Tree)

outlook	temperature	humidity	windy	play (Klasse)
sunny	cold	85	yes	no
overcast	hot	70	no	yes
rain	temp	50	:	:
	.	.	:	:
	.	.	:	:

Aus dieser Tabelle von Daten über Baseballspiele kann man nun einen Baum machen, der Prognosen abgibt, ob es bei bestimmten Bedingungen gespielt wird oder nicht.



Welches Merkmal sollte man nun als erstes im Baum benutzen. Da ist es geschickt das zu nehmen, welches die Trainingsmenge am besten trennt. Das Mat, was ID3 verwendet ist daten-orientiert Entropie der Trainingsmenge.

Die Entropie misst, wie viel Information man bekommen will.

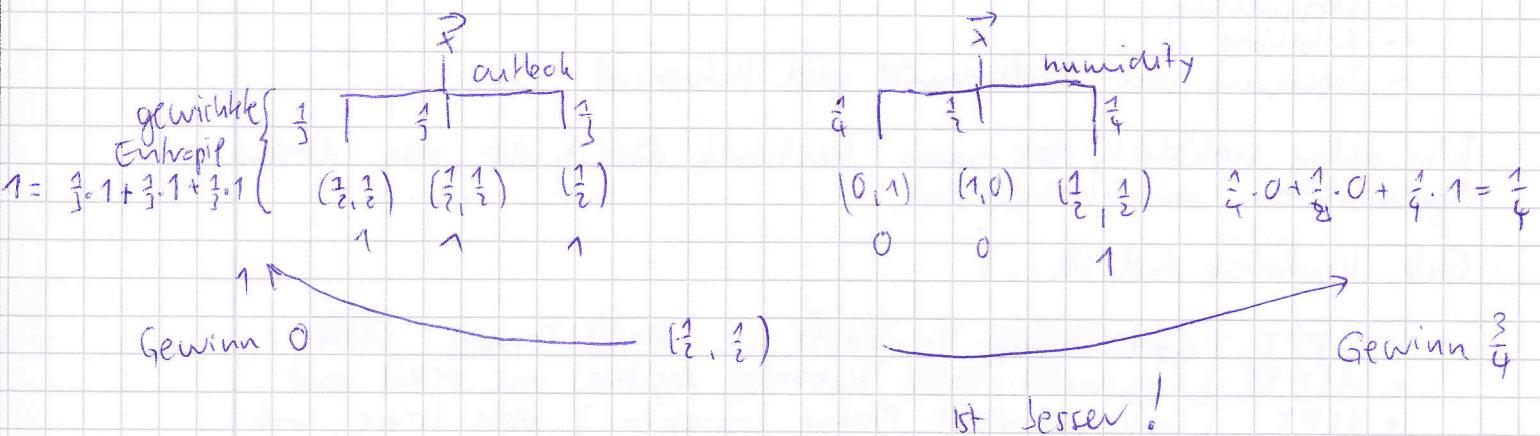
Ja  $(1,0)$  Null information  
nein  $(1,1)$

A.  $(0,5,0,5)$  1 Bit

B.  $(0,25,0,25,0,25,0,25)$  2 Bit

$T = (p_1, p_2, \dots, p_N)$  (Trainingsmenge nach Verteilung)

$$E(T) = - \sum_{i=1}^n p_i \log_2 p_i \quad (\text{Entropie})$$



Man muss natürlich alle Paare von Merkmalen gegeneinander abheften lassen.

Pseudocode für ID3:

$R$  = Menge der Merkmale

$C$  = Attribute

$S$  = Trainingsmenge

if ( $\forall \vec{x} \in S, \vec{x}$  von einer Klasse)  
return Knoten mit Klasse A

if ( $R = \emptyset$ )  
return Knoten mit häufigster Klasse in  $S$

Nehme Merkmal  $D$  mit maximalem Gewinn

return Knoten mit Markierung  $D$  und  $d$  Ausgängen (für jedes Attribut von  $D$  einen) und  $d$  Kindern (ID3 Baum)

$ID3(R - \{D\}, C, S_1), \dots, ID3(R - \{D\}, C, S_m)$

Element,  
für alle  $D = d_i$

$D = d_m$

Der Maximale Gewinn ergibt sich so:

$$\text{Gewinn} = \text{Info}(T) - \text{Info}_{\text{Merkm}}(x, T)$$

$$\text{Info}(x, T) = \sum_{i=1}^l \frac{|T_i|}{|T|} \text{Info}(T_i) \quad (\text{Anzahl der Attribute des Merkmals } x)$$

Die Info ist die Entropie.

Eine Erweiterung von ID3 ist C4.5.

# Merkmale in Bildern extrahieren

Vorlesung 1

21.1.2009

- Helligkeitsunterschiede
- Farben unterschiedlich
- Translation
- Skalierung
- Rotation
- Trennung von Vordergrund und Hintergrund

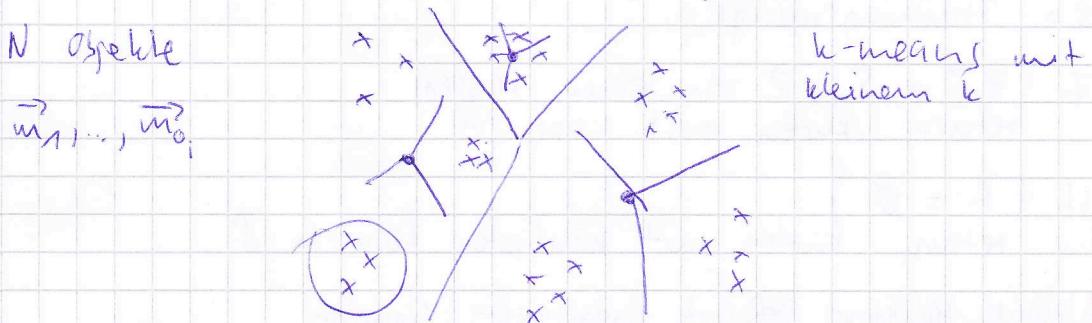
Wir wollen natürlich, dass man Merkmale findet, die von diesen Eigenschaften unabhängig sind.

Gute Merkmale sind z.B.

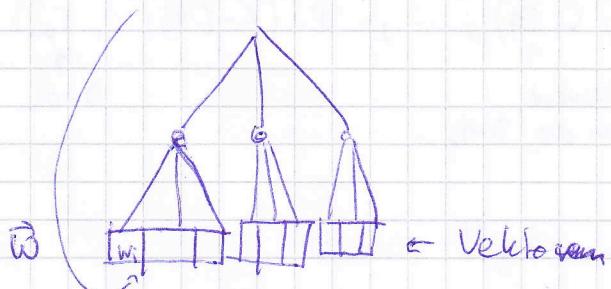
- GFTT (Good Features To Track) von Shi & Tomasi, 1994
- MSERs (Maximally Stable Extremal Regions) von Matas, 2002
- SIFT (Scale Invariant Feature Transform) von Lowe, 2004
- PCA-SIFT von Ke, 2004
- SURF (Speeded-up Robust Feature) von Bay, 2006

Angenommen, man findet nun also solche guten Merkmale, hofft aber für mehrere Objekte unterschiedliche Merkmalsdimensionen. Alles, was wir bisher hatten, funktioniert nicht.

Daniel Nister hat das 2006 wie folgt gelöst:



Daraus bekommt man rekursiv einen Baum



Die einzelnen Zellen kann man mit Gewichten  $w_i = \ln(\frac{N}{n_i})$  beschriften. Das ist das Gütekriterium für Merkmale.

$d_i = w_i n_i$  mit  $n_i$  ist die Anzahl der Merkmale bei Index i

Das ist der "Fingerabdruck".

Für jedes Objekt i der Datenbank wird ein Fingerprint  $\vec{d}_i$  erstellt.

Für ein Anfragobjekt wird ebenfalls solch ein Fingerprint  $\vec{q}$  erstellt

$$\vec{q} := w_1 \vec{x}_1 + \dots + w_N \vec{x}_N$$

Als Ähnlichkeitsmaß nehmen wir ein  $s$ .

$$s(d_i, \vec{q}) = \left\| \frac{\vec{q}}{\|\vec{q}\|} - \frac{d_i}{\|d_i\|} \right\|$$

Das macht man für alle Objekte der Datenbank, d.h. wir wollen

$$\underset{i}{\operatorname{argmin}} s(d_i, \vec{q})$$

Als Optimierung kann man auf  $\operatorname{argmin}_i s(d_i, \vec{q})$  berechnen.

Vorlesung 14  
28.1.2009

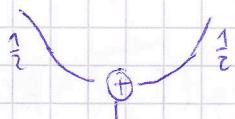
## Boosting

Weak classifiers. Wir haben ein "Komitee" von Klassifikatoren,  $f_1, f_2, \dots, f_m$ , wobei jeder  $f_i(x) \in \{-1, 1\}$  ist.  
Das Komitee ist dann  $f(x) = \sum_{i=1}^m \alpha_i f_i(x)$  mit  $\alpha_i = \frac{1}{m}$

In der Praxis kann man z.B. das so machen:



Klassifikator ist Neuron  
Net



Komitee

Der Auskunft für das Boosting ist nun, dass wir einen Gesamt-Klassifikator haben möchten.

$$G(x) = \operatorname{Sign} \left( \sum_{m=1}^M \alpha_m G_m(x) \right)$$

stark

schwacher Klassifikator (binär)  
mit Trefferquote > 50%

Der Algorithmus, den wir uns anschauen, heißt Adaboost.

Dabei haben wir eine Trainingsmenge

$\{(x_1, y_1), (x_2, y_2), \dots\}$ , die Gewichtet ist durch  $\{w_1, w_2, \dots\}$

1. Gewichte setzen auf  $w_i = \frac{1}{N}$  für  $i = 1, \dots, N$

2. For  $m=1$  to  $M$  (das sind die schwachen Klassifikatoren)

a) Trainiere schwachen Klassifikator  $G_m$  mit Trainingsmenge  $T$  und dem Gewichten  $w_i, i = 1, \dots, N$

$$b) \text{ evr}_m = \frac{\sum_{i=1}^N w_i I(y_i \neq G_m(x_i))}{\sum_{i=1}^N w_i}$$

$I$  ist dabei eine Indikatorvariable, die 1 ist, wenn die Bedingung zutrifft, 0 sonst.

$$c) \alpha_m = \log \frac{1 - \text{evr}_m}{\text{evr}_m}$$

$$d) w_i^{(m+1)} = w_i^{(m)} \cdot \exp(\alpha_m I(y_i \neq G_m(x_i))) \quad \text{für } i = 1, \dots, N$$

$$3) G(x) = \text{sign} \left( \sum_{m=1}^M \alpha_m G_m(x) \right)$$

Wir wollen natürlich zeigen, dass das auch gut geht und stimmt.

$$\text{Fehlerfunktion: } L(y, f(x)) = \exp \left( \begin{array}{c} -y \cdot f(x) \\ \hline -1,1 \end{array} \right) \quad \begin{array}{c} \\ \hline -1,1 \end{array}$$

die Fehlerfunktion gilt also  $e^{-1}$  bei Übereinstimmung und  $e^1$  sonst aus.

Wir wollen folgendes lösen:

Klassifikatoren  
bis jetzt

$$(\beta_m, G_m) = \underset{(\beta_m, G_m)}{\operatorname{arg\,min}} \left( \sum_{i=1}^N \exp(-y_i (\underbrace{f_m(x_i) + \beta_m G_m(x_i)}_{\text{neuer Klassifikator}})) \right)$$

$$\Rightarrow (\beta_m, G_m) = \underset{(\beta_m, G_m)}{\operatorname{arg\,min}} \left( \sum_{i=1}^N \exp(-y_i f_m(x_i)) \right) \quad \text{mit } w_i^{(m)} = \exp(-y_i f_m(x_i))$$

Zunächst suchen wir nun das beste  $G_m$ . Wir behaupten, dass

$$G_m = \underset{G_m}{\operatorname{arg\,min}} \sum_{i=1}^N w_i^{(m)} I(y_i \neq G_m(x_i))$$

$$(\beta_m, G_m) = \underset{(\beta_m, G_m)}{\operatorname{arg\,min}} \left( e^{-\beta_m} \sum_{y_i = G_m(x_i)} w_i^{(m)} + e^{\beta_m} \sum_{y_i \neq G_m(x_i)} w_i^{(m)} \right)$$

$$= \underset{(\beta_m, G_m)}{\operatorname{arg\,min}} \left( (e^{\beta_m} - e^{-\beta_m}) \sum_{i=1}^N w_i^{(m)} + e^{-\beta_m} \sum_{i=1}^N w_i^{(m)} I(y_i \neq G_m(x_i)) \right)$$

Wir schlussfolgern, dass wir  $\sum w_i^{(m)} I(y_i \neq G_m(x_i))$  minimieren. Dann finden wir das optimale  $\beta$  für A

$$(e^\beta - e^{-\beta}) \left( \sum_{i=1}^N w_i^{(m)} I(y_i \neq G_m(x_i)) \right) + e^{-\beta} \left( \sum_{i=1}^N w_i^{(m)} \right)$$

$$\text{Dann setzen wir } (e^{\beta} + e^{-\beta}) A - e^{-\beta} \beta = 0$$

$$(e^{2\beta} + 1) A - \beta I = 0$$

$$e^{2\beta} \cdot A = \beta I - A$$

$$e^{2\beta} = \frac{\beta - A}{A}$$

$$\beta = \frac{1}{2} \log \frac{\beta I - A}{A}$$

$$= \frac{1}{2} \log \left( \frac{1 - \text{err}_m}{\text{err}_m} \right)$$

$$\frac{A}{\beta} = \frac{\sum w_i^{(m)} I(y_i \neq G(x_i))}{\sum w_i^{(m)}} \\ = \text{err}_m$$

$$w_i^{(m+1)} = w_i^{(m)} \cdot e^{-\beta_m y_i G_m(x_i)}$$

$$- y_i G_m(x_i) = 2(I(y_i \neq G_m(x_i))) - 1$$

$$w_i^{(m+1)} = w_i^{(m)} \cdot e^{\alpha_m I(y_i \neq G_m(x_i))} \cdot e^{-\beta_m}$$

$\underbrace{\quad}_{\text{Multiplikative Konstante für die ganze Trainingsmenge}}$   
Konstante für die ganze Trainingsmenge

In AdaBoost lässt man diese Konstante einfach weg.