

Überwachtes Lernen:  $X = \{\langle \bar{x}^1, t^1 \rangle \dots \langle \bar{x}^N, t^N \rangle\}$

$\hookrightarrow NN$

SVM

Bayes Klassifikator

$N \dots \# \text{ Samples}$

$\bar{x}^n \dots \in \mathbb{N} \Rightarrow \text{Klassifikation}$

$t^n \dots \in \mathbb{R} \Rightarrow \text{Regression}$

Unüberwachtes Lernen:  $X = \{\bar{x}^1, \dots, \bar{x}^N\}$

- (Reinforcement Learning:)
- Explorative Datenanalyse: Hauptkomponententransformations (PCA)
  - Schätzung von  $w$ -Verteilung (ML Schätzungen)
  - ...

In der V3

$$P(X=x) = \lim_{N \rightarrow \infty} \frac{n_x}{N} \quad (= f_x(x) = P(X))$$

$$1 \geq P(X) \geq 0$$

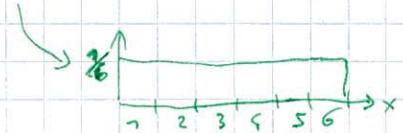
$N \dots \# \text{ Experimente}$

$n_x \dots \# X=x$

$X \dots \text{Zufallsvariable.}$

$$\sum_{x=1}^{|X|} P(X=x) = 1 \rightarrow \text{diskret}$$

$$\int P(X) dx = 1 \rightarrow \text{kontinuierlich}$$



Verbandswahrscheinlichkeit / Joint Prob.

$$\begin{aligned} P(X, Y) &= P(X|Y) \cdot P(Y) && \Leftarrow \text{Produktregel} \\ &= P(Y|X) \cdot P(X) \end{aligned}$$

Kettenregel  
↓

$$\text{Bsp: } P(X_1, \dots, X_N) = P(X_1) P(X_2|X_1) P(X_3|X_2, X_1) \dots P(X_N|X_{N-1})$$

Bayes:

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

$$\left. \begin{aligned} P(X) &= \sum_{Y=1}^{|Y|} P(X, Y=y) \\ P(Y) &= \sum_X P(X, Y) \end{aligned} \right\} \text{Summenregel}$$

$$P(X, Y) = P(X|Y) P(Y)$$

$$\left. \begin{aligned} &= P(X) P(Y) \\ P(X|Y) &= P(X) \end{aligned} \right\} \text{wenn } X \text{ statistisch unabhängig von } Y$$

# Schätzen von Wahrscheinlichkeitsverteilungen

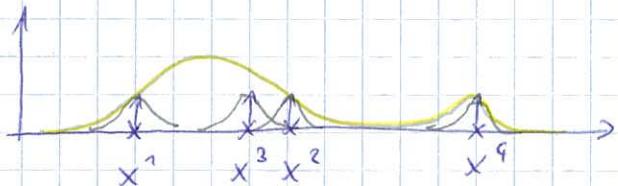
ges.:  $X = \{\bar{x}^1, \dots, \bar{x}^n\}$   $\bar{x}^n \in \mathbb{R}^d$

ges.: W!-Verteilung  $P(X)$

I)  $\rightarrow$  parametrisches Modell: Bsp: Gaußverteilung

II)  $\rightarrow$  nicht parametrisches Modell: ~~Modell~~

## II) Kernbasierter Schätzer:



Empirische Dichte fkt:  $P^e(x) = \frac{1}{N} \sum_{n=1}^N \delta(x - x^n)$

$\uparrow$   
Dirac

$$\int P^e(x) dx = 1$$

geglättete Dichte fkt:

$h(x) \dots$  Gaußkern  $\Rightarrow$  Glättungskern

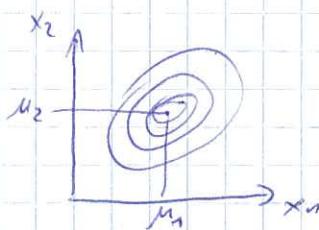
$$\begin{aligned} P^g(x) &= h(x) * P^e(x) \\ &= \int_{-\infty}^{\infty} h(x - \xi) P^e(\xi) d\xi \\ &= \int_{-\infty}^{\infty} h(x - \xi) \frac{1}{N} \sum_{n=1}^N \delta(\xi - x^n) d\xi \\ &= \underbrace{\int_{-\infty}^{\infty} h(x - \xi) d\xi}_{\text{Gaußkern}} \sum_{n=1}^N \int_{-\infty}^{\infty} \delta(\xi - x^n) d\xi \\ &= \frac{1}{N} \sum_{n=1}^N h(x - x^n) \end{aligned}$$

I) Multivariate Gaußverteilung  $\bar{x} \in \mathbb{R}^d$

$$N(\bar{x} | \Theta) = P(\bar{x} | \Theta) = \frac{1}{(2\pi)^{d/2} |\Sigma|^n} \cdot \exp(-\frac{1}{2} (\bar{x} - \bar{\mu})^\top \Sigma^{-1} (\bar{x} - \bar{\mu}))$$

$$\Theta \dots \{\bar{\mu}, \Sigma\}$$

$$\Theta = ?$$

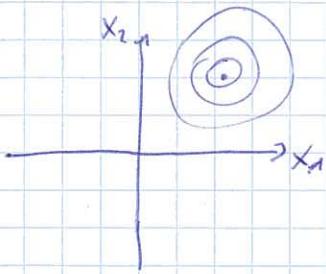


$$\bar{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

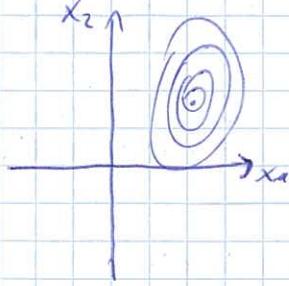
$\Leftrightarrow$  symmetrisch  $\Sigma = \Sigma^\top$   
 $\Leftrightarrow$  semi positiv definit

$$\Sigma = \begin{bmatrix} a & 0 \\ 0 & a \end{bmatrix}$$



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

$$b > a$$



### Maximum Likelihood (ML) Schätzer

ges:  $X = \{\bar{x}^1, \dots, \bar{x}^N\}$

ges:  $\Theta$

Likelihood:  $P(X | \Theta) = P(\bar{x}_1, \dots, \bar{x}_N | \Theta) =$

$$= P(\bar{x}_1 | \Theta) P(\bar{x}_2 | \bar{x}_1, \Theta) P(\bar{x}_3 | \bar{x}_1, \bar{x}_2, \Theta) \dots P(\bar{x}_N | \bar{x}_{N-1}, \dots, \bar{x}_1, \Theta)$$

i.i.d independent & identically distributed

$$= \prod_{n=1}^N P(\bar{x}_n | \Theta)$$

Likelihood:  $p(X | \Theta) \stackrel{\text{i.i.d}}{=} \prod_{n=1}^N p(\bar{x}_n | \Theta)$

Log Likelihood:  $L(X | \Theta) = \log \prod_{n=1}^N p(\bar{x}_n | \Theta) = \sum_{n=1}^N \log p(\bar{x}_n | \Theta)$

ML-Schätzer:  $\frac{\partial L(X | \Theta)}{\partial \Theta} \stackrel{!}{=} \emptyset$

$$\Theta_{ML} = \underset{\Theta}{\operatorname{argmax}} L(X | \Theta)$$

Bsp: Gauß:  $L(X | \Theta) = \sum_{n=1}^N \log N(\bar{x}_n | \Theta)$

$$\frac{\partial L(X | \Theta)}{\partial \bar{\mu}} \stackrel{!}{=} \emptyset$$

$$\frac{\partial L(X | \Theta)}{\partial \Sigma} \stackrel{!}{=} \emptyset$$

$$\bar{\mu} = \frac{1}{N} \sum_{n=1}^N \bar{x}_n$$

$$\Sigma = \frac{1}{N} \sum_{n=1}^N (\bar{x}_n - \bar{\mu})(\bar{x}_n - \bar{\mu})^\top$$

## Bayes'sche Schätzer:

ML  $\rightarrow \theta$  fix & unbekannt

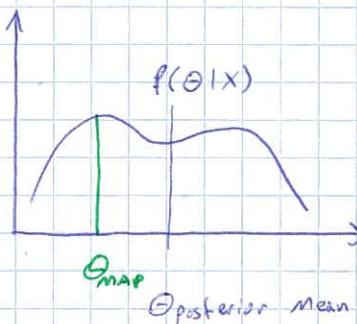
$\hookrightarrow \theta$  wird als RV (random variable) definiert.

$$P(\theta|x) = \frac{P(x|\theta) P(\theta)}{P(x)}$$

prior Wahrscheinlichkeit  
posterior Wahrscheinlichkeit

MAP-Schätzer (Maximum - a - posteriori)

$$\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmax}} P(\theta|x) = \underset{\theta}{\operatorname{argmax}} [P(x|\theta)P(\theta)]$$



Wenn  $P(\theta) \rightarrow$  uniform

$\hookrightarrow$  non-informative prior

$$P(\theta) \xrightarrow{\quad}$$

$$\hat{\theta}_{MAP} = \hat{\theta}_{ML}$$

Bayes Klassifikator

$$X = \{\langle \bar{x}^1, t^1 \rangle, \dots, \langle \bar{x}^n, t^n \rangle\}$$

$$\bar{x}^i \in \mathbb{R}^d$$

$$t \in \mathbb{N} \quad \# \text{... Samples}$$

$$t^n \in \{1, 2, \dots, C\} \quad C \# \text{Klassen}$$

**Bayes:**  $P(t|\bar{x}) = \frac{P(\bar{x}|t) P(t)}{\sum_{t'} P(\bar{x}|t') P(t')}$

Likelihood Prior prob

Posterior

$$= P(\bar{x})$$

Klassifikation

$$f: \bar{x} \rightarrow t^*$$

$P(\bar{x})$  wurde vernachlässigt, da es die Funktion nur Skaliert. Aber bei argmax kann man das vernachlässigen

$$t^* = \operatorname{argmax}_t P(t|\bar{x}) = \operatorname{argmax}_t [P(\bar{x}|t) P(t)] \quad \begin{array}{l} \text{-antw. parametr. } N(\bar{x}|\Theta_t) \\ \text{-nicht-parametr.} \end{array}$$

$$= \operatorname{argmax}_t [\ln P(\bar{x}|t) + \ln P(t)]$$

$$= \operatorname{argmax}_t g_t(\bar{x})$$

$g_t(\bar{x}) \dots$  Entscheidungsfkt.

Bsp: 2 Klassen:  $t \in \{1, 2\}$

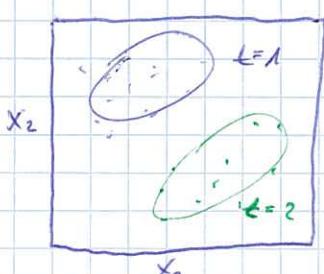
$$d=2 \quad \bar{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$p(\bar{x}|t) = N(\bar{x}|\Theta_t)$$

geg:  $\bar{x}$

$$\text{Klasse } t=1: X^1 = \{\bar{x}^n | t^n = 1\} \xrightarrow{\text{MC-schätzbar}} \Theta_1 = \{\mu_1, \Sigma_1\}$$

$$t=2: X^2 = \{\bar{x}^n | t^n = 2\} \Rightarrow \Theta_2 = \{\mu_2, \Sigma_2\}$$



$$P(t=i) = \frac{|X^i|}{|X|}$$

Gauß vert.

$$\mathcal{N}(\bar{x} | \Theta_t) = p(\bar{x} | t) = \frac{1}{(2\pi)^{d/2} |\Sigma_t|^{1/2}} \exp\left(-\frac{1}{2} (\bar{x} - \bar{\mu}_t)^T \Sigma_t^{-1} (\bar{x} - \bar{\mu}_t)\right)$$

g(t)

$$g_t(\bar{x}) = -\frac{1}{2} (\bar{x} - \bar{\mu}_t)^T \Sigma_t^{-1} (\bar{x} - \bar{\mu}_t) - \frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln(|\Sigma_t|) + \ln(p(t))$$

const

Euklidische Dst.  $(\bar{x} - \bar{\mu}_t)^T (\bar{x} - \bar{\mu}_t)$

$$g_t(\bar{x}) = -\frac{\|\bar{x} - \bar{\mu}_t\|^2}{2\sigma^2} + \ln(p(t))$$

$\Leftrightarrow$  varianz lösbar,  
wenn  $p(t)$  uniform.

$\Rightarrow$  lineare Entscheidungsgr.

3 Fälle:

$$1) \Sigma_t = \sigma^2 I \quad \text{Einheitsmatrix}$$

$$g_t(\bar{x}) = -\frac{\|\bar{x} - \bar{\mu}_t\|^2}{2\sigma^2} + \ln(p(t))$$

"Entscheidungsgr."

$$g_t(\bar{x}) = g_i(\bar{x}) \quad i \neq t$$

$$\Rightarrow P(t|\bar{x}) = P(i|\bar{x})$$

$$2) \Sigma_t = \Sigma \rightarrow \text{voll besetzt} \rightarrow \text{viele ohne Ellipse}$$

$$g_t(\bar{x}) = -\frac{1}{2} (\bar{x} - \bar{\mu}_t)^T \Sigma_t^{-1} (\bar{x} - \bar{\mu}_t) + \ln(p(t))$$

Mahalanobis Dist.

$\Rightarrow$  lineare Entscheidungsgr.

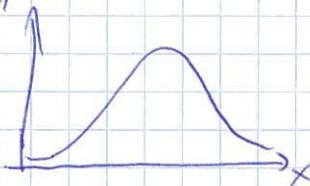
$$3) \Sigma_t \text{ beliebig}$$

$$g_t(\bar{x}) = -\frac{1}{2} (\bar{x} - \bar{\mu}_t)^T \Sigma_t^{-1} (\bar{x} - \bar{\mu}_t) - \frac{1}{2} \ln(|\Sigma_t|) + \ln(p(t))$$

$\Rightarrow$  hyperquadratische Entscheidungsgr.

Bis jetzt sind wir davon ausgegangen, dass die Daten Gauß verteilt sind

px



Doch was, wenn Daten so verteilt sind:

px

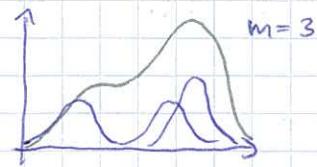
schlecht

besser - Gauß'sche Mischverteilung.

## Gaußsche Mischverteilung

$$p(\bar{x}_n | \Theta) = \sum_{m=1}^M \alpha_m N(\bar{x}_n | \bar{\mu}_m, \Sigma_m) \quad \Theta_m$$

$M \dots \# \dots$  Komponenten



Bed:

$$\sum_{m=1}^M \alpha_m = 1 = p(\bar{w}) \quad 0 \leq \alpha_m \leq 1$$

$$\int_{-\infty}^{\infty} p(\bar{x} | \Theta) d\bar{x} = 1$$

$$\Theta = \{\alpha_1, \dots, \alpha_M, \Theta_1, \dots, \Theta_M\}$$

?  $\Theta$ ?

Schätzproblem:

$$\text{Beg: } X = \{\bar{x}_1, \dots, \bar{x}_n\} \quad \bar{x} \in \mathbb{R}^d$$

gs:  $\Theta$

ML-Schätzer

log := ln

$$\Theta = \underset{\Theta}{\operatorname{argmax}} \{ \log P(X | \Theta) \}$$

$$\left[ \frac{\partial \log P(X | \Theta)}{\partial \Theta} = 0 \right]$$

$$\log P(X | \Theta) \stackrel{i.i.d.}{=} \log \prod_{n=1}^N p(\bar{x}_n | \Theta) = \sum_{n=1}^N \log \left[ \sum_{m=1}^M \alpha_m N(\bar{x}_n | \bar{\mu}_m, \Sigma_m) \right]$$

Gaußsche Mischung.

Mittelwert:

$$\frac{\partial \log P(X | \Theta)}{\partial \bar{\mu}_m} = \sum_{n=1}^N \frac{1}{\sum_{m=1}^M \alpha_m N(\bar{x}_n | \bar{\mu}_m, \Sigma_m)} \cdot \frac{\partial \sum_{m=1}^M \alpha_m N(\bar{x}_n | \bar{\mu}_m, \Sigma_m)}{\partial \bar{\mu}_m}$$

$$= \sum_{n=1}^N \frac{\alpha_m N(\bar{x}_n | \bar{\mu}_m, \Sigma_m)}{\sum_{m=1}^M \alpha_m N(\bar{x}_n | \bar{\mu}_m, \Sigma_m)} \cdot \frac{\partial \log [\alpha_m N(\bar{x}_n | \bar{\mu}_m, \Sigma_m)]}{\partial \bar{\mu}_m} = \cancel{\times}$$

$=: r_m^n$

$$\frac{\partial \log N(\bar{x}_n | \bar{\mu}_m, \Sigma_m)}{\partial \bar{\mu}_m} =$$

$$\begin{aligned} \ln N(\bar{x}_n | \bar{\mu}_m, \Sigma_m) &= -\frac{1}{2} (\bar{x}_n - \bar{\mu}_m)^T \Sigma_m^{-1} (\bar{x}_n - \bar{\mu}_m) \\ &\quad - \frac{d}{2} \ln (2\pi) - \frac{1}{2} \ln |\Sigma_m| \end{aligned}$$

$$\frac{\partial (\bar{a} - \bar{x})^T C (\bar{a} - \bar{x})}{\partial \bar{x}} = (C + C^T) (\bar{a} - \bar{x})$$

$$= -\frac{1}{2} \left( \sum_m^{-1} + (\sum_m^{-1})^T \right) (\bar{x}_n - \bar{\mu}_m)$$

symmetrisch

$$= - \sum_m^{-1} \cdot (\bar{x}_n - \bar{\mu}_m)$$

$$\bigcirc = - \sum_{n=1}^N r_m^n \cdot \sum_m^{-1} \cdot (\bar{x}_n - \bar{\mu}_m)$$

$$\frac{\partial \dots}{\partial \bar{\mu}_m} \stackrel{!}{=} 0$$

$$- \sum_{n=1}^N r_m^n \cdot \sum_m^{-1} (\bar{x}_n - \bar{\mu}_m) = 0 \quad | \cdot \sum_m$$

$$- \sum_{n=1}^N \cdot r_m^n \cdot \bar{x}_n + \bar{\mu}_m \sum_{n=1}^N r_m^n = 0$$

$$\bar{\mu}_m = \frac{\sum_{n=1}^N r_m^n \bar{x}_n}{\sum_{n=1}^N r_m^n}$$

$\therefore N_m \dots$  Effektive # von Datenpunkten, die von der m-ten Komponente modelliert werden

$$r_m^n = P(m|x^n \theta) \Leftarrow \text{posterior Prob.}$$

$$\text{Bsp: } \bar{x} \in \mathbb{R}^2$$

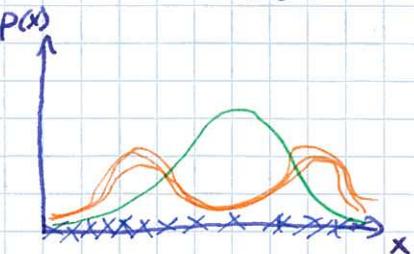
|  |     |                        |
|--|-----|------------------------|
|  |     | $P(m=1 x\theta) = 0,6$ |
|  | $x$ | $P(m=2 x\theta) = 0,4$ |
|  |     | $\theta$               |
|  |     |                        |

Henne-Ei-Problem:

für  $r_m^n$  brauche ich  $\theta$ .  
 Aber  $\theta$  will ich schätzen.  
 $\Rightarrow$  iterativ!

## Gauß'sche Mischverteilung

$$p(\bar{x}_n | \Theta) = \sum_{m=1}^M \alpha_m N(\bar{x}_n | \bar{\mu}_m, \Sigma_m)$$



$M \dots \# \text{ Komponenten}$

$$\bar{x}_n \in \mathbb{R}^d$$

$$\alpha_m = P(m)$$

$$0 \leq \alpha_m \leq 1$$

$$\sum_{m=1}^M \alpha_m = 1$$

$$\Theta = \{\alpha_m, \bar{\mu}_m, \Sigma_m\}_{m=1 \dots M}$$

↳ ML - Schätzer

### ML - Schätzer

$$X = \{\bar{x}_1, \dots, \bar{x}_N\} \quad N \dots \# \text{ Samples}$$

$$\Theta_{\text{ML}} = ?$$

$$\Theta_{\text{ML}} = \underset{\Theta}{\operatorname{argmax}} [\log P(X | \Theta)]$$

standard  
ML -  
Schätzer  
(Prüfungs-  
frage)

$$\text{Lösung: } \frac{\partial \log P(X | \Theta)}{\partial \Theta} = 0$$

$$\log P(X | \Theta) \stackrel{\text{i.i.d.}}{=} \log \prod_{n=1}^N p(\bar{x}_n | \Theta)$$

$$\log P(X | \Theta) = \sum_{n=1}^N \log \sum_{m=1}^M \alpha_m N(\bar{x}_n | \bar{\mu}_m, \Sigma_m)$$

$$\frac{\partial \log P(X | \Theta)}{\partial \bar{\mu}_m} = 0$$

$$\bar{\mu}_m = \frac{\sum_{n=1}^N r_m^n \bar{x}_n}{N_m}$$

$$N_m = \sum_{n=1}^N r_m^n \dots \text{Effektive } \# \text{ von Daten-} \\ \text{punkten die von Komponente } m \text{ modelliert werden.}$$

posterior

$$r_m^n = P(m | \bar{x}_n, \Theta) = \frac{\alpha_m N(\bar{x}_n | \bar{\mu}_m, \Sigma_m)}{\sum_{m=1}^M \alpha_m N(\bar{x}_n | \bar{\mu}_m, \Sigma_m)}$$

$$\bar{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} :$$



$$P(m=1 | \bar{x}_n, \Theta) = 0.6$$

$$P(m=2 | \bar{x}_n, \Theta) = 0.4$$

$$\rightarrow \frac{\partial \log P(x|\theta)}{\partial \Sigma_m} = 0$$

$$\Sigma_m = \frac{1}{N_m} \sum_{n=1}^N r_m^n \cdot (\bar{x}_n - \bar{\mu}_m) \cdot (\bar{x}_n - \bar{\mu}_m)^T$$

$$\rightarrow \frac{\partial \log P(x|\theta)}{\partial \alpha_m} = 0 \quad \sum_{m=1}^M \alpha_m = 1 \quad \leftarrow \text{Bedingung}$$

Optimierung mit Lagrange Multiplikator

$$J(m) = \log P(x|\theta) + \lambda \left( \sum_{m=1}^M \alpha_m - 1 \right)$$

$$\frac{\partial J(m)}{\partial \alpha_m} = 0$$

$$\frac{\partial \log P(x|\theta)}{\partial \alpha_m} = \sum_{n=1}^N \frac{1}{\sum_{m=1}^M \alpha_m N(\bar{x}_n | \bar{\mu}_m, \Sigma_m)} \cdot N(\bar{x}_n | \bar{\mu}_m, \Sigma_m)$$

$$\frac{\partial J(m)}{\partial \alpha_m} = \sum_{n=1}^N \left[ \frac{\alpha_m \cdot N(\bar{x}_n | \bar{\mu}_m, \Sigma_m)}{\sum_{m=1}^M \alpha_m N(\bar{x}_n | \bar{\mu}_m, \Sigma_m)} \right] + \lambda \cdot \alpha_m = 0 \quad | \cdot \alpha_m$$

$$\sum_{m=1}^M \underbrace{\sum_{n=1}^N r_m^n}_{N_m} + \lambda \alpha_m = 0 \quad | \sum_{m=1}^M$$

$$\underbrace{\sum_{m=1}^M N_m}_{N} + \lambda \underbrace{\sum_{m=1}^M \alpha_m}_{=1} = 0 \Rightarrow \boxed{\lambda = -N}$$

$$\underbrace{\sum_{n=1}^N r_m^n}_{N_m} - N \cdot \alpha_m = 0 \Rightarrow \boxed{\alpha_m = \frac{N_m}{N}}$$

Algorithmus für  $\Theta_{ML}$  von Gauß'schen Mischverteilungen

(EM - Algorithmus)  
Expectation  
Maximization

EM:

$$1) \text{ Init: } \Theta^{(0)} = \{\bar{\mu}_m^{(0)}, \Sigma_m^{(0)}, \alpha_m^{(0)}\}_{m=1}^M$$

$t=0$



$$2) \text{ E-Step: } r_m^n = p(m | \bar{x}_n, \Theta^{(t)})$$

$$\begin{array}{l} m=1..M \\ n=1..N \end{array}$$

$$3) \text{ M-Step: } \mu_m^{(t+1)} = \frac{1}{N_m} \sum_{n=1}^N r_m^n \bar{x}_n$$

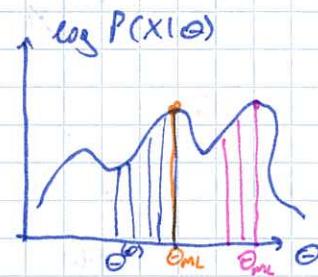
$$\Sigma_m^{(t+1)} = \frac{1}{N_m} \sum_{n=1}^N r_m^n (\bar{x}_n - \mu_m^{(t)}) \cdot (\bar{x}_n - \mu_m^{(t)})^T$$

$$\alpha_m^{(t+1)} = \frac{N_m}{N}$$

$$t = t + 1$$

4) Termination: if  $\log P(X|\Theta) \rightarrow$  konvergiert  $\Rightarrow \Theta_M = \Theta^{(t)}$

otherwise  $\rightarrow$  step 2



Eigenschaft:

- $\log P(X|\Theta)$  wird mit jedem Schritt monoton steiger.
- EM konvergiert gegen lokales Optimum.
- Lösung von  $\Theta^{(0)}$  abhängig

Init:  $\alpha_m \rightarrow$  uniform verteilen

$\Sigma_m = \Sigma \Rightarrow$  aus  $X$  berechnen

$\mu_m \rightarrow$  zufällig gewählte  $\bar{x}_n$

↳ oder: k-means Algorithmus.

Anwendung:

- Sprache

EM für GMM  $\rightarrow$  k-means

Annahme: a)  $\Sigma_m = G^2 I$

b)  $x_m \dots$  vernachlässigen

c)  $r_m^n \in \{0, 1\} \Leftarrow$  Klassifikation

E-step:

$$r_m^n = \frac{\alpha_m N(\bar{x}_n | \mu_m, \Sigma_m)}{\sum_{m'=1}^M \alpha_{m'} N(\bar{x}_n | \mu_{m'}, \Sigma_{m'})}$$

Bayes Klassifikator: CASE i  
(a)

$$\Rightarrow g_m(\bar{x}_n) = -\frac{\|\bar{x}_n - \mu_m\|^2}{2G^2} + c_m$$

$$(b) \Rightarrow g_m(\bar{x}_n) = -\|\bar{x}_n - \mu_m\|^2$$

Euklidische Distanz

$$(C) m^* = \arg\max_m (-\|\bar{x}_n - \bar{\mu}_m\|) = \arg\min_m \|\bar{x}_n - \bar{\mu}_m\|$$

k-means:

1) Init:  $\Theta^{(0)} = \{\bar{\mu}_k^{(0)}\}_{k=1}^K$

(k-means, daher  $m = k$   
 $M = K$ )

2) Step 1: Annahme (a) - (c)

$$y_k = \{\bar{x}_n \mid k = \arg\min_{k'} \|\bar{x}_n - \bar{\mu}_{k'}\|^2\} \quad n=1..N$$

$$k=1..K$$

3) Step 2:

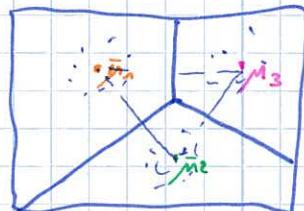
$$\bar{\mu}_k^{(t+1)} = \frac{1}{|y_k|} \sum_{x_n \in y_k} \bar{x}_n \quad t=t+1 \quad k=1..K$$

$|y_k| \dots \# \text{samples in } y_k$

4) Termination: if kumulative Distanz  $J = \sum_{k=1}^K \sum_{\bar{x} \in y_k} \|\bar{x} - \bar{\mu}_k^{(t)}\|$  konvergiert  
 otherwise Step 1  $\rightarrow$  Abbruch.

Eigenschaften:

- konvergiert gegen lokales Optimum
- abhängig von  $\Theta^{(0)}$
- kumulative Distanz wird minimiert.
- Clustergrenzen stückweise linear



Anwendung:

- Vektorquantisierung
- Init von  $\mu_m$  bei GMM.

Markov Modell (MM)

Bsp.: Grammatikmodell für natürliche Sprache

$$q_n \in W^*$$

n... Zeitpunkt, Position in Sequenz

$$|W| = 1000$$

q... Zustand / State

W... Wörterbuch = { "id", "du", "ist", ... }

Satz: "Ich gehe einkaufen" =  $q_1 \dots q_n$

$$P(q_1 \dots q_n) = ?$$

$$P(q_1 \dots q_n) = P(q_1) P(q_2 | q_1) \dots P(q_n | q_{n-1} \dots q_1)$$

$$\# \text{ Parameter: } |W| \quad |W|^2 \quad \dots \quad |W|^n$$

Unigram - Modell:  $q_i \perp\!\!\! \perp q_j \quad i \neq j \quad \perp\!\!\! \perp \dots$  stat. unabhängig

$$P(q_1 \dots q_n) = \prod_{i=1}^n P(q_i) \quad \text{Problem: } P(\text{"Ich gehe einkaufen"}) = P(\text{"gehe einkaufen ich"})$$

Bigram - Modell

$$P(q_1 \dots q_n) = P(q_1) \cdot \prod_{i=2}^n P(q_i | q_{i-1})$$

$$\text{Bsp: } P(\dots) = P(\text{"id"}) \cdot P(\text{"gehe"} | \text{"id"}) P(\text{"einkaufen"} | \text{"gehe"})$$

↪ Markov Modell (MM) 1ster Ordnung

Trigram - Modell:

$$P(q_1 \dots q_n) = P(q_1) \cdot P(q_2 | q_1) \cdot \prod_{i=3}^n P(q_i | q_{i-1}, q_{i-2})$$

↪ MM - 2ter Ordnung

Bsp: Wettersequenzen

$$q_n \in \{ R, S, F \}$$

MM - 1ster Ordnung

$P(q_1)$  ... prior W!

$P(q_i | q_{i-1})$  ... Übergangsw! / Transition probability

# Hidden Markov Modell (HMM) (siehe Schiles auf Homepage)

$$\Theta = \{\pi, A, B\}$$

- HMM
- $S = \{s_1, \dots, s_n\}$  ... Menge der Zustände /  $N_s = |S|$
  - $P(q_1 = s_i) = \pi_i$  ... prior w!  $\pi = \{\pi_1, \dots, \pi_{N_s}\}$
  - $P(q_n = s_j | q_{n-1} = s_i) = a_{ij} = a_{q_{n-1}, q_n}$  ... Übergangs-w!

$$A = [a_{ij}]_{N_s \times N_s}$$

- Beobachtungsw!  $b_{i,x_n} = P(x_n | q_n = s_i)$   $\bar{x}_n \in \mathbb{N}^d$  (diskret)

$$B = [b_{i,x_n}]_{i=1 \dots N_s}$$

$\bar{x}_n \in \mathbb{R}^d$  (kontinuierlich)  
 $\Rightarrow$  bei Spracherkennung:  
 Gauß / Gauß'sche Mdlv.

$$P(x_n | q_1 \dots q_n, x_1 \dots x_{n-1}) = P(x_n | q_n) \Rightarrow x_n \perp\!\!\!\perp \{q_1, \dots, q_{n-1}, x_1, \dots, x_{n-1}\}$$

Annahme

Normalisierungsbedingung:

$$\sum_{i=1}^{N_s} \pi_i = 1$$

$$\sum_{j=1}^{N_s} a_{ij} = 1 \quad i = 1 \dots N_s$$

$$\int_{-\infty}^{\infty} b_{i,x_n} dx_n = 1 \quad i = 1 \dots N_s$$

Zustandssequenz / Pfad

$$Q = \{q_1 \dots q_n\}$$

Beobachtungssequenz

$$X = \{x_1, \dots, x_n\}$$

## 3 Probleme:

1) Evaluierungsproblem / Klassifikation

geg.:  $X, \Theta$   $\rightarrow$  Forward / Backward Algorithmus

ges.:  $P(X|\Theta)$  ... Likelihood / Produktionsw!

$$P(t|X) = \frac{P(X|t) P(t)}{P(X)}$$

$$t^* = \underset{t}{\operatorname{argmax}} P(t|X) = \underset{t}{\operatorname{argmax}} \frac{P(X|t) \cdot P(t)}{P(X|\Theta_t)}$$

## 2) Dekodierungsproblem

geg.:  $X, \Theta \rightarrow$  Viterbi - Algorithmus

ges.:  $Q^* = \arg\max_Q P(Q | X, \Theta)$

$\hookrightarrow$  Pfad der  $X$  bei  $\Theta$  am besten erklärt

## 3) Schätzproblem

geg.:  $X^{1:R} = \{X^1, \dots, X^R\}$   $R \dots \#$  Beobachtungssequenzen

ges.:  $\hat{\Theta}_{ML} = \arg\max_{\Theta} P(X^{1:R} | \Theta) = \arg\max_{\Theta} \prod_{i=1}^R P(X^i | \Theta)$   
siehe 1)

$\hookrightarrow$  EM - Algorithmus / Baum - Welch - Algorithmus

Zu 1)  $P(X | \Theta) = ?$

$P$  für state sequenz  $Q \Rightarrow MM$

$$P(Q | \Theta) = \prod_{i=1}^n a_{q_{i-1}, q_i} = P(q_1) \cdot \prod_{i=2}^n P(q_i | q_{i-1})$$

$P$  für  $X$  bei fsg.  $\Theta, Q$

$$P(X | Q, \Theta) = \prod_{i=1}^n b_{q_i, x_i}$$

$$P(X | Q | \Theta) = P(X | Q, \Theta) P(Q | \Theta)$$

$$P(X | Q) = \sum_{Q \in Q} P(X | Q | \Theta)$$

$Q \in Q \dots$  Menge von Pfaden

$$|Q| = (N_s)^n$$

# Rechenoperationen:  $\mathcal{O}(2 \cdot n (N_s)^n)$

Effizientere Methode: (Forward / Backward)

Ausnützen des Distributivgesetzes

$$\underbrace{a(b+c)}_{\text{Naive}} = \underbrace{a \cdot b + a \cdot c}_{\text{Effizient}}$$

$$\mathcal{O}(2 \cdot N_s^2 \cdot n)$$

zu 2) geg:  $\Theta, X$

ges:  $Q^*$

$$Q^* = \underset{Q}{\operatorname{argmax}} P(Q | X, \Theta) = \underset{Q}{\operatorname{argmax}} \frac{P(X, Q | \Theta)}{P(X | \Theta)} =$$

$$= \underset{Q \in Q}{\operatorname{argmax}} P(X, Q | \Theta)$$

$$P(Q^* | X | \Theta) = \max_Q P(X, Q | \Theta) \stackrel{!}{=} p^*(X | \Theta)$$

Viterbi:  $\rightarrow$  maximal erzielbare W! für Beobachtungen  $x_1 \dots x_n$  entlang eines

einzigen Pfades  $q_1 \dots q_{n-1}, q_n = s_i$ , der vom Zeitpunkt  $n$  im State  $q_n = s_i$  mündet.

$$S_n(i) = \max_{q_1 \dots q_{n-1}} P(q_1 \dots q_{n-1}, q_n = s_i, x_1 \dots x_n | \Theta)$$

$$S_n(i) = \max_{q_1} P(q_1 = s_i | x_n | \Theta) = \max_{q_1} \pi_{q_1} \cdot b_{q_1, x_n}$$

rekursiv

$$\downarrow S_2(i)$$

$$\psi_n(i) = \underset{q_1 \dots q_{n-1}}{\operatorname{argmax}} P(\text{---} | \text{---})$$

$\hookrightarrow$  Welches  $[S_{n-1}(i) \alpha_{ij}]$  hat zum Maximum geführt

Rekursionsformel:

$$S_n(j) = \max_{1 \leq i \leq N_s} [S_{n-1}(i) \alpha_{ij}] \cdot b_{j, x_n}$$

Algorithmus:

$$\text{Init: } S_1(i) = \pi_i \cdot b_{i, x_n} \quad i = 1 \dots N_s$$

$$\psi_1(i) = 0$$

$$\text{Rekurrenz: } S_n(j) = \max_{1 \leq i \leq N_s} [S_{n-1}(i) \alpha_{ij}] \cdot b_{j, x_n}$$

$$j = 1 \dots N_s$$

$$n = 2 \dots N \quad N \dots \text{Sequenzlänge}$$

$$\psi_n(j) = \underset{1 \leq i \leq N_s}{\operatorname{argmax}} [S_{n-1}(i) \alpha_{ij}]$$

$$\text{Termination: } p^*(X | \Theta) = \max_{1 \leq i \leq N_s} S_N(i)$$

$$q_N^* = \underset{1 \leq i \leq N_s}{\operatorname{argmax}} S_N(i)$$

Backtracking: Extraktion von  $Q^*$  aus  $\psi_n(i)$

$$\hookrightarrow q_n^* = \psi_{n+1}(q_{n+1}^*) \quad n = N-1 \dots 1$$

Beispiel von Viterbi ist auf den Slides (Homepage) zu finden

Graphical Models.

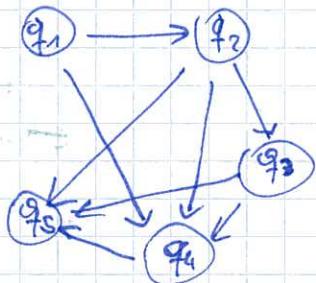
→ Tutorial <sup>Slides</sup> auf der Homepage!

Definitionen zu Graphen.

directed GM

HMM kann als Bayes Network geschrieben werden.

Bsp: Bayes Network

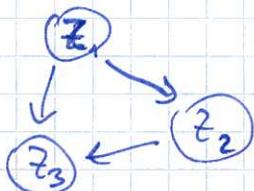


Beim Tutorial sind  
Section 1, 2 (2.1 & 2.2),  
Section 3 (3.1 - 3.4), und  
Section 5 + 5.1  
relevant ☐

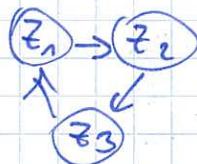
geschrieben werden.

→ nicht voll verbunden:  
statistische Unabhängigkeit  
zv. einzelnen Samples.

Directed Acyclic Graph (DAG):



$$P(z_1 z_2 z_3) = P(z_1) \cdot P(z_2 | z_1) \cdot P(z_3 | z_2 z_1)$$



$$P(z_1 z_2 z_3) = P(z_1 | z_3) \cdot P(z_2 | z_1) \cdot P(z_3 | z_2)$$

↳ geht nicht.

stat. independence:  $x_i \perp\!\!\!\perp x_j$

$$P(x_i x_j) = P(x_i) P(x_j) \Rightarrow P(x_i | x_j) = P(x_i)$$

$$P(x_j | x_i) = P(x_j)$$

cond. independence:  $x_i \perp\!\!\!\perp x_j | x_z$

$$P(x_i x_j | x_z) = P(x_i | x_z) \cdot P(x_j | x_z) \Rightarrow P(x_i | x_j x_z) = P(x_i | x_z)$$

$$P(x_j | x_i x_z) = P(x_j | x_z)$$

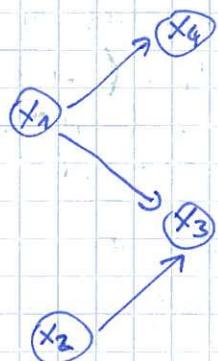
Canonical examples → siehe Folien.

ad converging connection:  $z_i \perp\!\!\!\perp z_j$  aber  $z_i \not\perp\!\!\!\perp z_j | z_k$



Bsp:

Mit Hilfe der canonical examples ablesen:



$$\rightarrow x_1 \perp\!\!\!\perp x_2$$

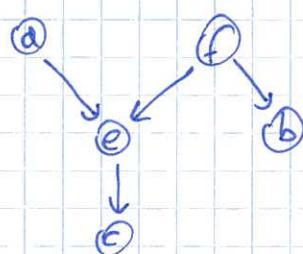
$$\rightarrow x_4 \perp\!\!\!\perp x_3 | x_1$$

$$\rightarrow x_2 \perp\!\!\!\perp x_4$$

$$\rightarrow x_2 \perp\!\!\!\perp x_4 | x_1$$

:

Bsp:



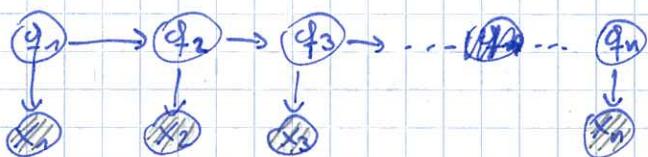
$$? a \perp\!\!\!\perp b | f \quad \checkmark \leftarrow \text{stimmmt}$$

$$? a \perp\!\!\!\perp b | e \quad \text{b} \rightarrow a \perp\!\!\!\perp b$$

$$? a \perp\!\!\!\perp b | c \quad \text{b} \leftarrow \text{Statement stimmt nicht} \quad \square$$

Bsp: HMM - revisited

$$P(X|Q) = \underbrace{P(q_1)}_{\prod q_1} \cdot \underbrace{P(x_1|q_1)}_{b_{q_1 x_1}} \cdot \prod_{n=2}^N \underbrace{P(q_n|q_{n-1})}_{a_{q_{n-1} q_n}} \underbrace{P(x_n|q_n)}_{b_{q_n x_n}}$$



$\text{X}_n \dots$  schraffiert heißt beobachtet.

$$q_1 \perp\!\!\!\perp q_3 | q_2$$

$$q_{n+1:N} \perp\!\!\!\perp q_{1:n-1} | q_n$$

$$x_2 \perp\!\!\!\perp x_1 | q_2$$

Produktionsw!:  $p(x) = \sum_{Q \in Q} P(X|Q)$

$\downarrow$

inferenz.

Bsp: GMM - revisited

$$P(x) = \sum_{m=1}^M \underbrace{P(m)}_{\alpha(m)} \cdot \underbrace{\frac{N(x|\theta_m)}{P(x|m)}}_{P(x|m)}$$

$$P(x) = \sum_{m=1}^M \underbrace{P(m) P(x|m)}_{P(x,m)}$$



## Inferenz:

$P(m|X) = ? \hat{=} \underline{\text{E-Step im EM-Algo.}}$

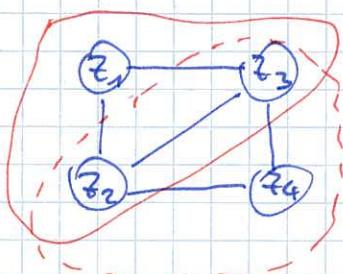
$$P(m|X) = \frac{P(m, X)}{P(X)} = \frac{P(m) P(X|m)}{\sum_{m'} P(m') \underbrace{P(X|m')}_{\propto^{(m')} N(x|\theta_m)}}$$

HMM und GMM  $\rightarrow$  Prüfungsrelevant!

## undirected GM

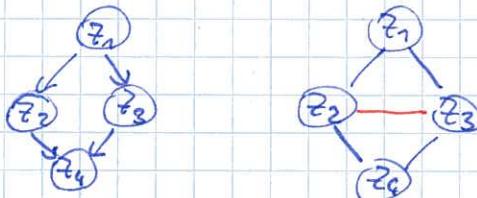
Markov Netzwerk. (siehe Folien)

Bsp: Clique



Clique:  $(z_1 - z_3), (z_1 - z_2), (z_3 - z_4), (z_2 - z_4),$   
 $(z_1 - z_2 - z_3), (z_2 - z_3 - z_4)$   
 max. Clique

Convert DGM to UGM (siehe Folien)



$$P(z_1 \dots z_4) = \underbrace{P(z_1)}_1 \underbrace{P(z_2|z_1)}_1 \underbrace{P(z_3|z_1)}_1 \underbrace{P(z_4|z_2 z_3)}_1$$

$$P(z_1 \dots z_4) = \frac{1}{w} \psi_{123}(z_1 z_2 z_3) \cdot \psi_{234}(z_2 z_3 z_4)$$

## Factor Graph

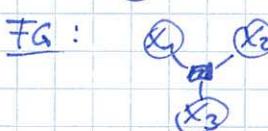
$$\text{DGM: } P(z_1 \dots z_4) = \underbrace{P(z_1)}_1 \underbrace{P(z_2|z_1)}_1 \underbrace{P(z_3|z_1)}_1 \underbrace{P(z_4|z_3)}_1$$

$$\text{UGM: } P(z_1 \dots z_4) = \frac{1}{w} \psi_{12}(z_1 z_2) \underbrace{\psi_{23}(z_2 z_3)}_1 \underbrace{\psi_{34}(z_3 z_4)}_1$$

Bsp:

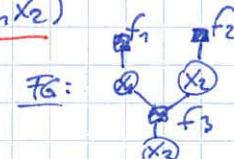


$$P(x_1 \dots x_3) = \frac{1}{w} \underbrace{f_1(x_1)}_1 \underbrace{f_2(x_2)}_1 \underbrace{f_3(x_1 x_2 x_3)}_1$$



$$P(x_1 x_2 x_3) = \frac{1}{w} f(x_1 x_2 x_3)$$

zwei mögliche Factor Graphen!



$$P(x_1 x_2 x_3) = \frac{1}{w} \underbrace{f_1(x_1)}_1 \underbrace{f_2(x_2)}_1 \underbrace{f_3(x_3)}_1$$

Lernen:  $\Rightarrow$  DGM / BN

$\hookrightarrow$  Was:  $\Rightarrow$  Parameter  $\Rightarrow P(z_i | z_{\pi_i}) \Rightarrow$  ML  
Bayes Schätzer

$\Rightarrow$  Graph-Struktur  $\Rightarrow$  Faktorisierung

$\hookrightarrow$  Greedy - Heuristiken  $\Rightarrow$  Score

Inference (siehe Tutorial)

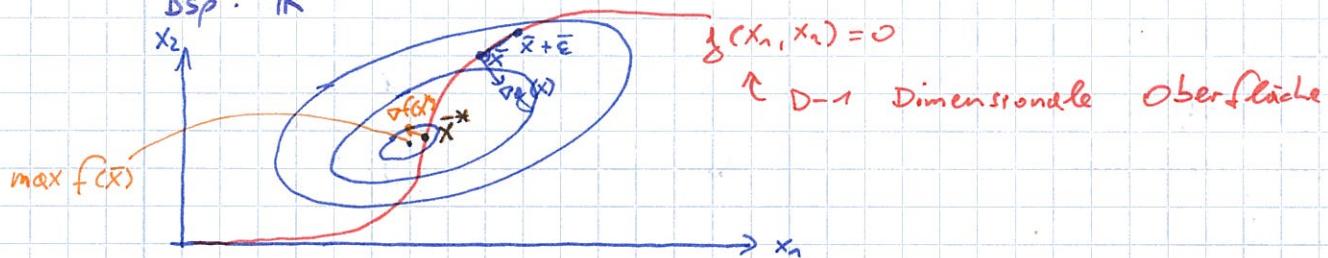
## Lineare Transformationen - Dimensionsreduktion

Exkurs: Lagrange-Multiplikatoren

→ maximieren einer Funktion, mit Bedingungen (Constraints)

Problem:  $\max_{\mathbb{R}^2} f(\bar{x})$  (Objective)     $\bar{x} = [x_1, \dots, x_D]^T$   
 s.t.  $g(\bar{x}) = 0$  (Constraints)    s.t. ... subject +

Bsp:  $\mathbb{R}^2$



Taylor von  $g(\bar{x})$

$$g(\bar{x} + \bar{\epsilon}) \approx g(\bar{x}) + \bar{\epsilon}^T \nabla g(\bar{x})$$

$$g(\bar{x}) = g(\bar{x} + \bar{\epsilon}) = 0$$

$$\Rightarrow \bar{\epsilon}^T \nabla g(\bar{x}) \approx 0 \quad (\lim_{\|\bar{\epsilon}\| \rightarrow 0} \bar{\epsilon}^T \nabla g(\bar{x}) = 0)$$

→  $\bar{\epsilon}$  parallel zu  $\nabla g(\bar{x})$

$\nabla g(\bar{x})$  normal auf  $\nabla g(\bar{x})$

$\nabla f(\bar{x}^*)$  normal auf  $\nabla g(\bar{x})$

⇒  $\nabla f$  und  $\nabla g$  parallel an  $\bar{x} = \bar{x}^*$

$$\Rightarrow \nabla f(\bar{x}) + \lambda \nabla g(\bar{x}) = 0 \quad (\lambda \neq 0)$$

→ Lagrange-Funktion

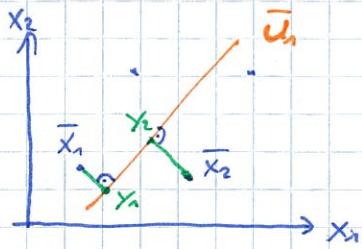
$$\mathcal{L}(\bar{x}, \lambda) = f(\bar{x}) + \lambda \underset{\text{Lagrange-Multipl.}}{\uparrow} g(\bar{x})$$

$$\nabla_{\bar{x}} \mathcal{L} = \nabla_{\bar{x}} f(\bar{x}) + \lambda \cdot \nabla_{\bar{x}} g(\bar{x}) \stackrel{!}{=} 0$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = g(\bar{x}) \stackrel{!}{=} 0$$

Hintergrundinfo,  
kein Prüfungsstoff

$$\text{geg.: } X = \{\bar{x}_n\}_{n=1}^N, \quad \bar{x}_n \in \mathbb{R}^D$$



→ Projektion auf:

$M = 1$ -dim. Raum

$$y_n = \bar{u}_n^T \bar{x}_n$$

(1x1) (1xD) (Dx1)

$$\rightarrow \bar{u}_n = ?$$

$M \leq D$ -dim. Raum

$$\bar{Y}_n = \bar{U} \bar{X}_n = \begin{bmatrix} \bar{u}_1^T \\ \vdots \\ \bar{u}_M^T \end{bmatrix} \bar{X}_n$$

Statistische Eigenschaften der transf. Daten ( $M=1$ )

$$\text{Mittelwert: } m_y = \frac{1}{N} \sum_{n=1}^N \bar{u}_n^T \bar{x}_n = \underline{\bar{u}_n^T \bar{m}_x}$$

$$\begin{aligned} \text{Varianz: } \sigma_y^2 &= \frac{1}{N} \sum_{n=1}^N (\bar{u}_n^T \bar{x}_n - \bar{u}_n^T \bar{m}_x)^2 \\ &= \frac{1}{N} \sum_{n=1}^N (\bar{u}_n^T (\bar{x}_n - \bar{m}_x))^2 \\ &= \bar{u}_n^T \underbrace{\frac{1}{N} \sum_{n=1}^N (\bar{x}_n - \bar{m}_x)(\bar{x}_n - \bar{m}_x)^T}_{= S_x} \bar{u}_n \\ &= \underline{\bar{u}_n^T S_x \bar{u}_n} \end{aligned}$$

Bestimmen von  $\bar{u}_n^T$

Annahme: Varianz (Leistung)  $\hat{=}$  Information

$$\begin{aligned} \bar{u}_n &= \underset{\bar{u}_n}{\arg \max} \bar{u}_n^T S_x \bar{u}_n && (\text{würde dazu führen, dass } \|\bar{u}_n\| \rightarrow \infty) \\ \text{s.t. } \bar{u}_n^T \bar{u}_n &= 1 && \text{Nur Richtung von } \bar{u}_n \text{ wichtig} \\ \rightarrow \bar{u}_n^T \bar{u}_n - 1 &= 0 && \rightarrow \|\bar{u}_n\| = 1 \rightarrow \bar{u}_n^T \bar{u}_n = 1 \end{aligned}$$

$$\text{Lagrange: } \mathcal{L}(\bar{u}_n, \lambda) = \bar{u}_n^T S_x \bar{u}_n + \lambda(1 - \bar{u}_n^T \bar{u})$$

$$\frac{\partial \mathcal{L}(\bar{u}_n, \lambda)}{\partial \bar{u}_n} = \cancel{\lambda S_x \bar{u}_n} + \cancel{\lambda \bar{u}_n} \stackrel{!}{=} 0$$

$$S_x \bar{u}_n = \tilde{\lambda} \bar{u}_n \quad (\bar{u}_n \text{ ist EV von } S_x \text{ zu EW } \tilde{\lambda})$$

$$\text{oder: } \bar{u}_n^T S_x \bar{u}_n = \tilde{\lambda} \quad (\rightarrow \text{größter EW } \tilde{\lambda})$$

→ Principal Component Analysis (PCA)

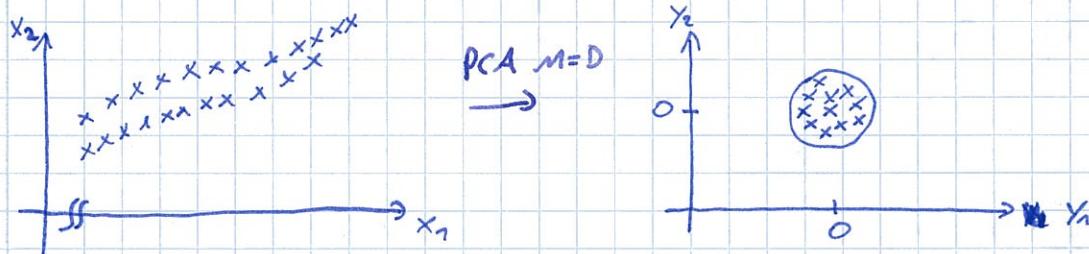
## PCA

$$y_n = \bar{u}_n^T \bar{x}_n \quad (1. \text{ Hauptkomponente})$$

→ M Hauptkomponenten

$$\bar{y}_n = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_D \end{bmatrix} = \begin{bmatrix} \bar{u}_1^T \\ \bar{u}_2^T \\ \vdots \\ \bar{u}_D^T \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}$$

PCA als Vorverarbeitung



$$L = \begin{bmatrix} \lambda_1 & \phi \\ \phi & \lambda_D \end{bmatrix} \quad U = \begin{bmatrix} \bar{u}_1^T \\ \vdots \\ \bar{u}_D^T \end{bmatrix}$$

$$S \cdot U = U \cdot L$$

PCA (modifiziert)

$$\bar{y}_n = L^{-\frac{1}{2}} \cdot U^T (\bar{x}_n - \bar{m}_x)$$

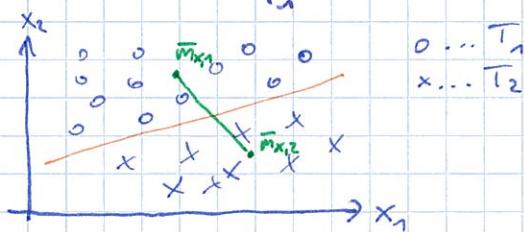
$$1) \bar{m}_y = \bar{0}$$

$$\begin{aligned} S_y &= \frac{1}{N} \sum_{n=1}^N \bar{y}_n \bar{y}_n^T \\ &= \frac{1}{N} \sum_{n=1}^N L^{-\frac{1}{2}} U^T (\bar{x}_n - \bar{m}_x) \underbrace{(\bar{x}_n - \bar{m}_x)^T}_{\text{mit Summe } = S_x} U L^{-\frac{1}{2}} \\ &= L^{-\frac{1}{2}} \underbrace{U^T S_x U}_{=L} L^{-\frac{1}{2}} = I \end{aligned}$$

⇒  $y_i, y_j$  paarweise dekorreliert (heißt i.A. nicht statistisch unabhängig!)

Linear Discriminant Analysis (LDA)

$$X = \left\{ \underbrace{\bar{x}_1, \dots, \bar{x}_{N_1}}_{T_1}, \underbrace{\bar{x}_{N_1+1}, \dots, \bar{x}_N}_{T_2} \right\}$$



$$\bar{m}_{x,1} = \frac{1}{N_1} \sum_{\bar{x} \in T_1} \bar{x}_n ; \quad \bar{m}_{x,2} = \dots$$

Projektion:  $y = \bar{u}^T \bar{x}$

$$\begin{aligned} \text{Idee: } \bar{u} &= \underset{\bar{u}}{\operatorname{argmax}} (m_{y,2} - m_{y,1}) \\ &= \underset{\bar{u}}{\operatorname{argmax}} \bar{u}^T (\bar{m}_{x,2} - \bar{m}_{x,1}) \\ \text{s.t. } \bar{u}^T \bar{u} &= 1 \end{aligned}$$

$$d(\bar{u}, \lambda) = \bar{u}^T (\bar{m}_{x,2} - \bar{m}_{x,1}) + \lambda(1 - \bar{u}^T \bar{u})$$

:

$\rightarrow$  nicht optimal!

$$\boxed{\bar{u} \propto \bar{m}_{x,2} - \bar{m}_{x,1}}$$

→ Varianzen einbeziehen!

$$G_{y,1}^2 = \sum_{n \in T_1} (y_n - m_{y,1})^2$$

Summen-Variante:  $G_y^2 = G_{y,1}^2 + G_{y,2}^2 \rightarrow$  soll "klein" sein.

→ Als  $f(\bar{u}, \bar{x})$

$$S_w = \sum_{n \in T_1} (\bar{x}_n - \bar{m}_{x,1})(\bar{x}_n - \bar{m}_{x,1})^T + \sum_{n \in T_2} (\bar{x}_n - \bar{m}_{x,2})(\bar{x}_n - \bar{m}_{x,2})^T \text{ "within-class cov"}$$

$$\rightarrow G_{y,1}^2 + G_{y,2}^2 = \bar{u}^T S_w \bar{u}$$

maximieren: "Between-class cov" 

$$(m_{y,2} - m_{y,1})^2 = (\bar{u}^T (\bar{m}_{x,2} - \bar{m}_{x,1}))^2 = \dots = \bar{u}^T S_B \bar{u}$$

$$\bar{u} = \underset{\bar{u}}{\operatorname{argmax}} \frac{\bar{u}^T S_B \bar{u}}{\bar{u}^T S_w \bar{u}}$$

s.t.  $\bar{u}^T \bar{u} = 1$

:

$$\boxed{\bar{u} \propto S_w^{-1} (\bar{m}_{x,2} - \bar{m}_{x,1})}$$

↳ LDA

$$\bar{u} \propto \bar{m}_{x,2} - \bar{m}_{x,1} \quad \text{wenn } S_w \propto I$$