



Expectation Maximization



Expectation Maximization

Problem:

Lernen der Parameter einer Verteilung aus Trainingsdaten \mathcal{D} mit fehlenden oder schlechten Features. Die Daten werden eingeteilt in gute (\mathcal{D}_g) und schlechte (\mathcal{D}_b)

Beispiel: Punkte aus Gaußverteilung

$$\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4\} = \left\{ \begin{pmatrix} 0 \\ 2 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 \\ 2 \end{pmatrix}, \begin{pmatrix} * \\ 4 \end{pmatrix} \right\}, \mathcal{D}_b = x_{41}$$

Expectation Maximization

02. Juni 2017

Für eine Maximum-Likelihood-Schätzung der Verteilungsparameter θ suchen wir das θ , das die likelihood $p(\mathcal{D}_g | \theta)$ maximiert:

$$\theta = \arg \max_{\theta} p(\mathcal{D}_g | \theta) = \arg \max_{\theta} \sum_{\mathcal{D}_b} p(\mathcal{D}_g, \mathcal{D}_b | \theta)$$

Falls \mathcal{D}_b kontinuierlich ist, muss natürlich bei der Marginalisierung die Summe durch ein Integral ersetzt werden.

In der Praxis ist die Berechnung von $\sum_{\mathcal{D}_b} p(\mathcal{D}_g, \mathcal{D}_b | \theta)$ oft schwierig, oder nicht möglich, z.B. aufgrund der vielen Möglichkeiten für \mathcal{D}_b (exponentiell in der Anzahl der Dimensionen). Es kann aber gezeigt werden, dass $p(\mathcal{D}_g, \mathcal{D}_b | \theta)$ auch mit einer Hilfsfunktion $Q(\theta; \theta^i)$ optimiert werden kann.



Expectation Maximization

Wir stellen folgende Funktion auf:

$$\begin{aligned} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^i) &= E_{\mathcal{D}_b | \mathcal{D}_g, \boldsymbol{\theta}^i} \left[\ln p(\mathcal{D}_g, \mathcal{D}_b | \boldsymbol{\theta}) \right] \\ &= \sum_{\mathcal{D}_b} p(\mathcal{D}_b | \mathcal{D}_g, \boldsymbol{\theta}^i) \ln p(\mathcal{D}_g, \mathcal{D}_b | \boldsymbol{\theta}) \end{aligned}$$

Dabei ist $\boldsymbol{\theta}^i$ die zur Zeit beste Schätzung für die Verteilungsparameter.

$\boldsymbol{\theta}$ ist ein Kandidatenvektor für eine verbesserte Schätzung.

$Q(\boldsymbol{\theta}; \boldsymbol{\theta}^i)$ gibt uns den Erwartungswert für die log-Likelihood für die Daten an, wobei zu dessen Berechnung die schlechten Daten als mit der Verteilung $p(\mathbf{x} | \boldsymbol{\theta}^i)$ verteilt angenommen werden.



Expectation Maximization

Der Expectation Maximization Algorithmus läuft nun in zwei abwechselnd durchgeführten Schritten ab:

E-Step: Berechne $p(\mathcal{D}_b | \mathcal{D}_g, \theta^i)$ und damit

$$Q(\theta; \theta^i) = E_{\mathcal{D}_b | \mathcal{D}_g, \theta^i} \left[\ln p(\mathcal{D}_g, \mathcal{D}_b | \theta) \right]$$

M-Step: Berechne θ^{i+1} als $\arg \max_{\theta} Q(\theta; \theta^i)$

Diese Vorgehensweise, bei der für die fehlenden Daten eine Wahrscheinlichkeitsverteilung angenommen wird, wird als **Soft Expectation Maximization** bezeichnet.



Expectation Maximization

02. Juni 2017

Alternativ können statt einer Verteilung im E-Step auch die aktuell wahrscheinlichsten fehlenden Daten berechnet werden. Dann gilt:

E-Step: Berechne $\mathcal{D}_{b_i} = \arg \max_{\mathcal{D}_b} p\left(\mathcal{D}_g, \mathcal{D}_b \mid \boldsymbol{\theta}^i\right)$

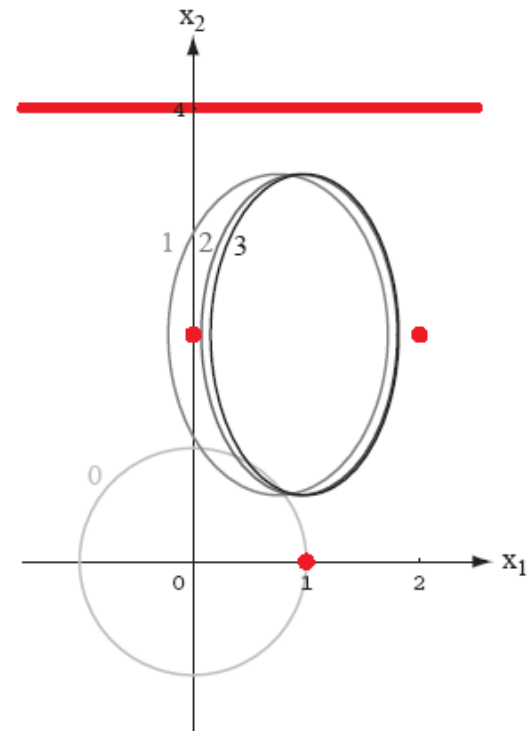
M-Step: Berechne $\boldsymbol{\theta}^{i+1}$ als $\arg \max_{\boldsymbol{\theta}} p\left(\mathcal{D}_g, \mathcal{D}_{b_i} \mid \boldsymbol{\theta}\right)$

Diese Vorgehensweise wird als **Hard Expectation Maximization** bezeichnet.

Expectation Maximization

Beispiel: Finden einer 2D-Normalverteilung
(achsenparallel) mit den angegebenen Daten:

$$\mathcal{D} = \left\{ \begin{pmatrix} 0 \\ 2 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 2 \\ 2 \end{pmatrix}, \begin{pmatrix} * \\ 4 \end{pmatrix} \right\}$$





Expectation Maximization

Generalized Expectation Maximization (GEM):

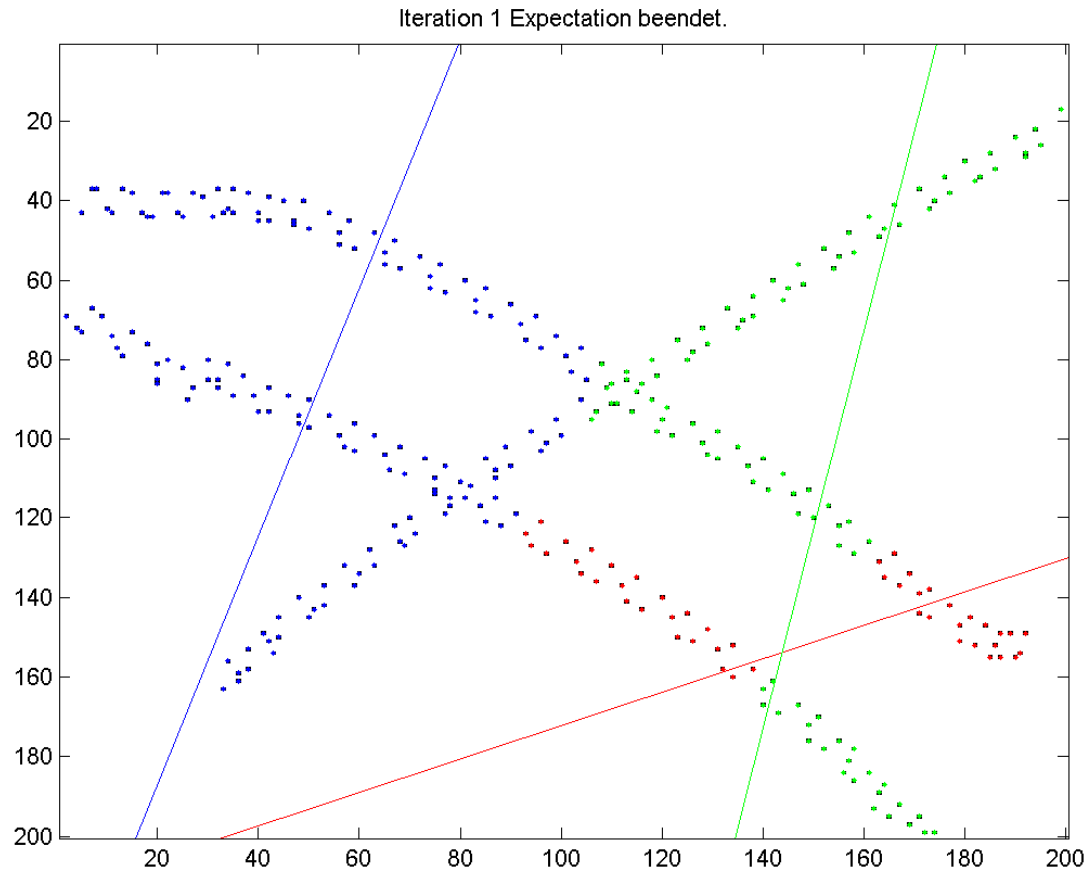
- Verbessere θ^i in jedem Schritt
- nicht unbedingt optimale Lösung
- Eine Version ist Hard EM, wo in jedem Schritt wahrscheinlichste Werte für die unbekannten Features geschätzt werden.
 - Beispiel [Egorova 02]:
 - Gesucht: Beste drei Geraden
 - Gegebene Daten: Punkte
 - Fehlende Daten: welcher Punkt gehört zu welcher Geraden

[Egorova 02] <http://www.cims.nyu.edu/~fischer/SemBayNet/>, Vortrag v. 5.12.02 (Link nicht mehr erreichbar)



3 Geraden – schönes Beispiel

02. Juni 2017

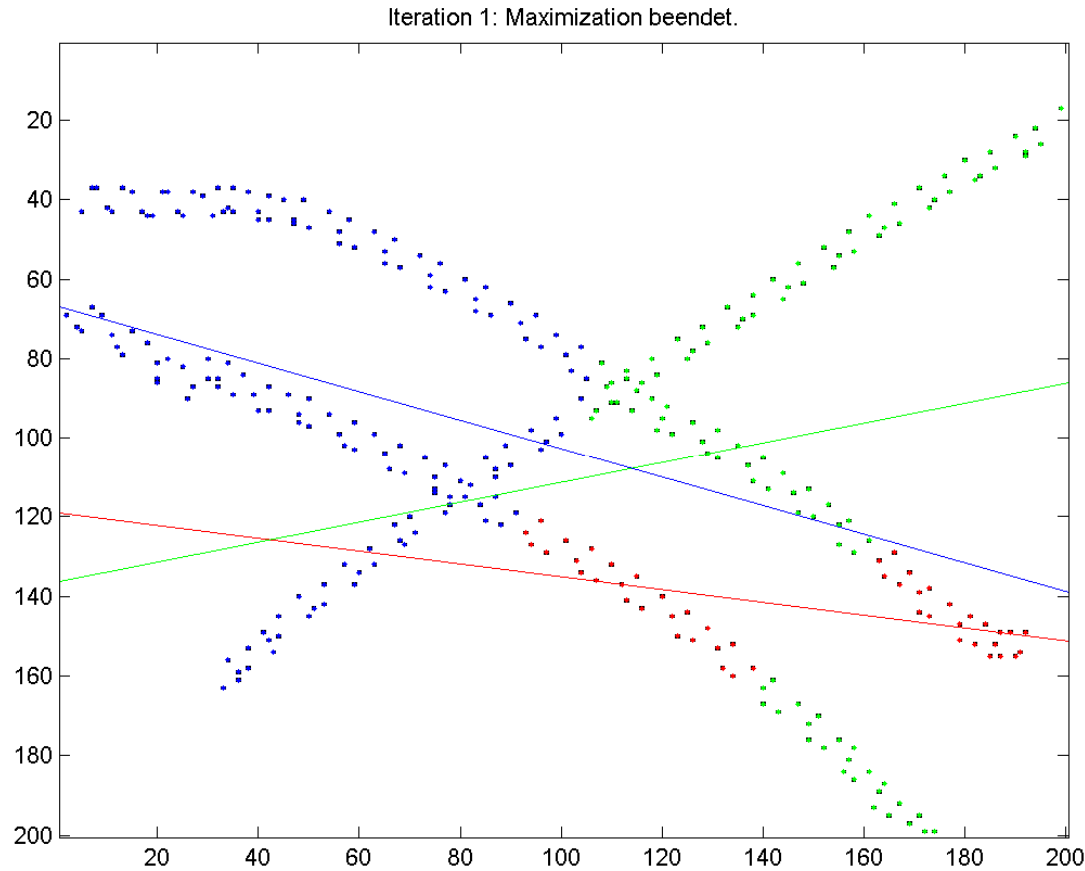


[Egorova 02]



3 Geraden – schönes Beispiel

02. Juni 2017

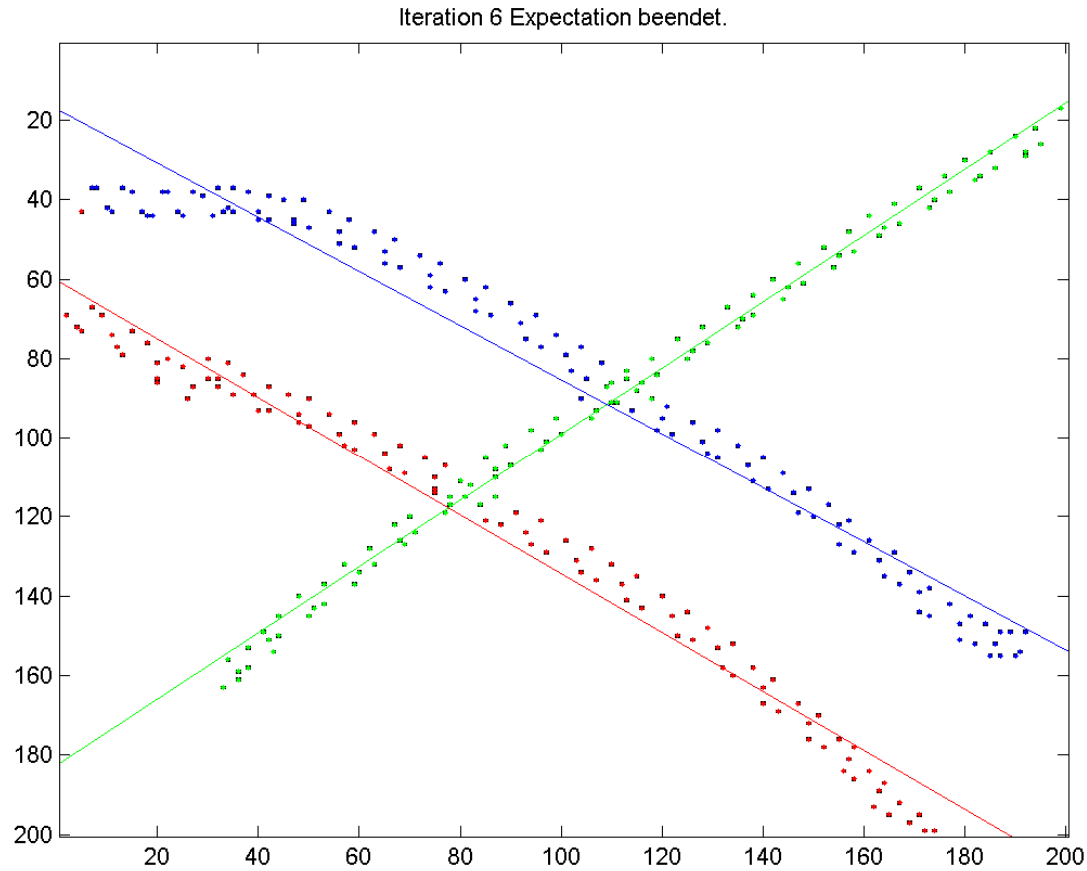


[Egorova 02]



3 Geraden – schönes Beispiel

02. Juni 2017

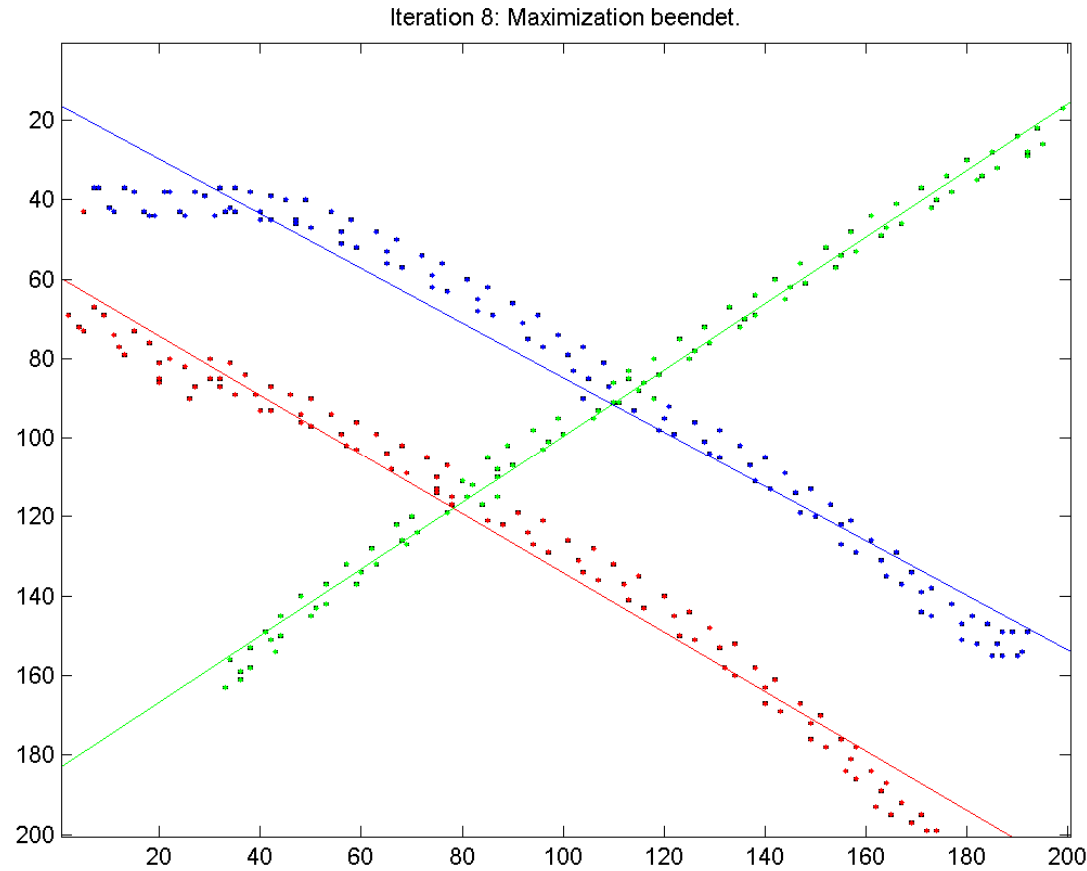


[Egorova 02]



3 Geraden – schönes Beispiel

02. Juni 2017

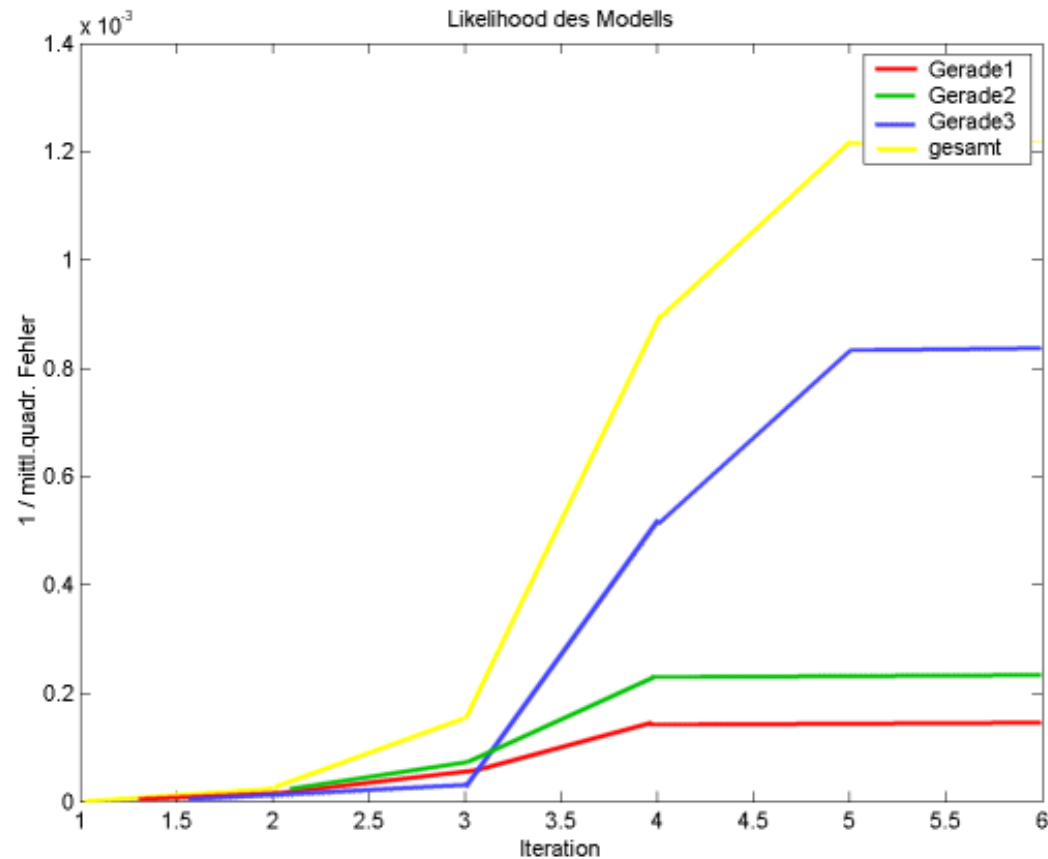


[Egorova 02]



3 Geraden – schönes Beispiel

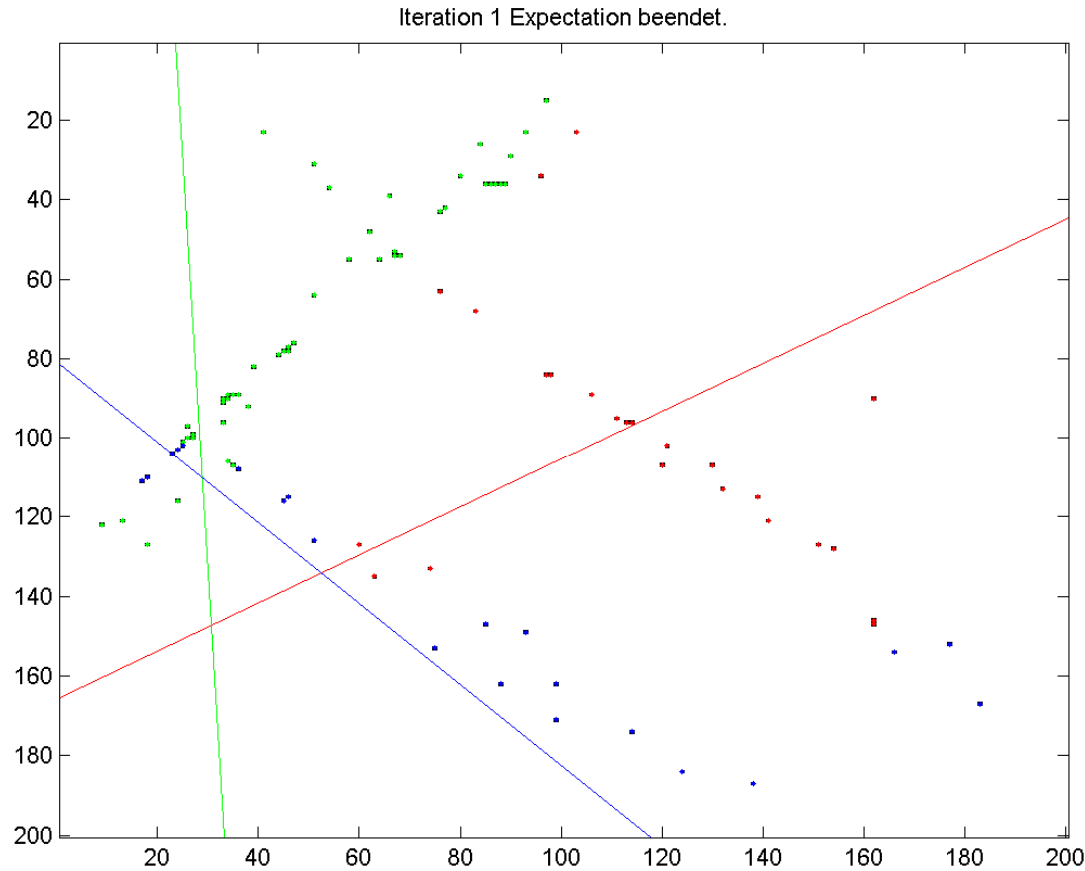
02. Juni 2017



[Egorova 02]

3 Geraden – noch ein Beispiel

02. Juni 2017

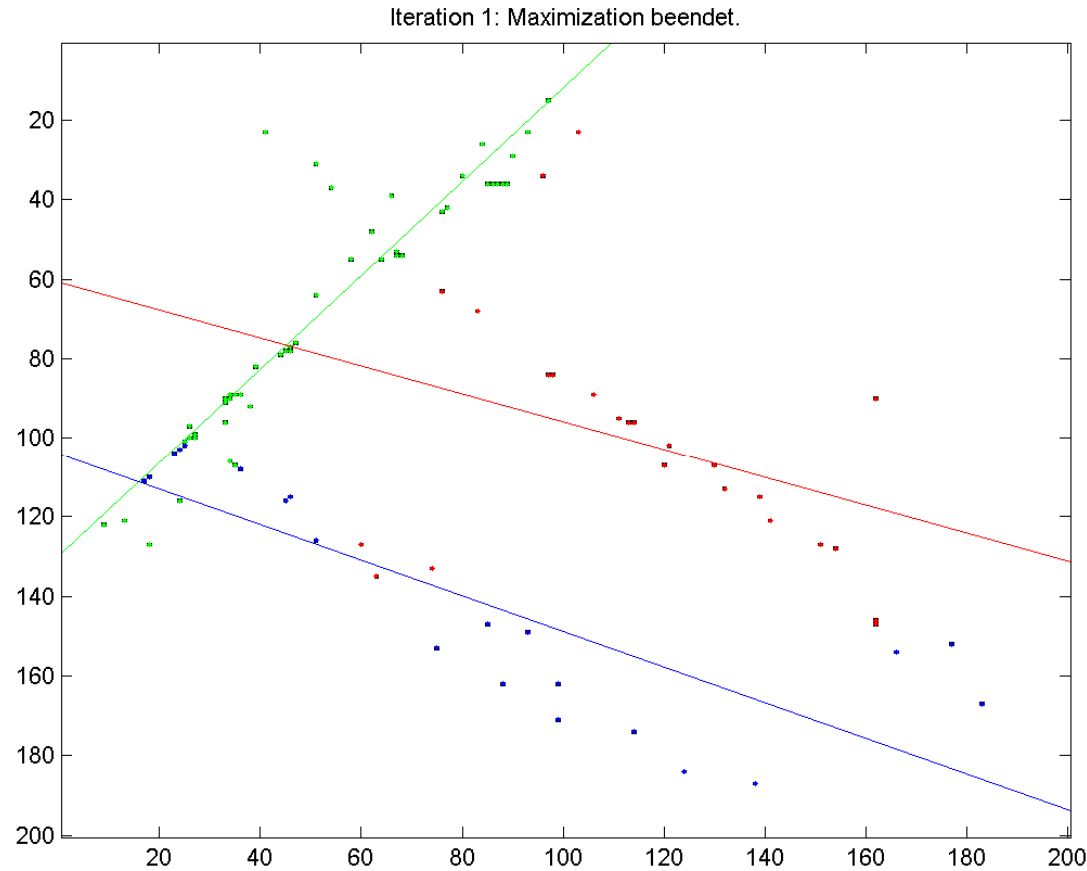


[Egorova 02]



3 Geraden – noch ein Beispiel

02. Juni 2017

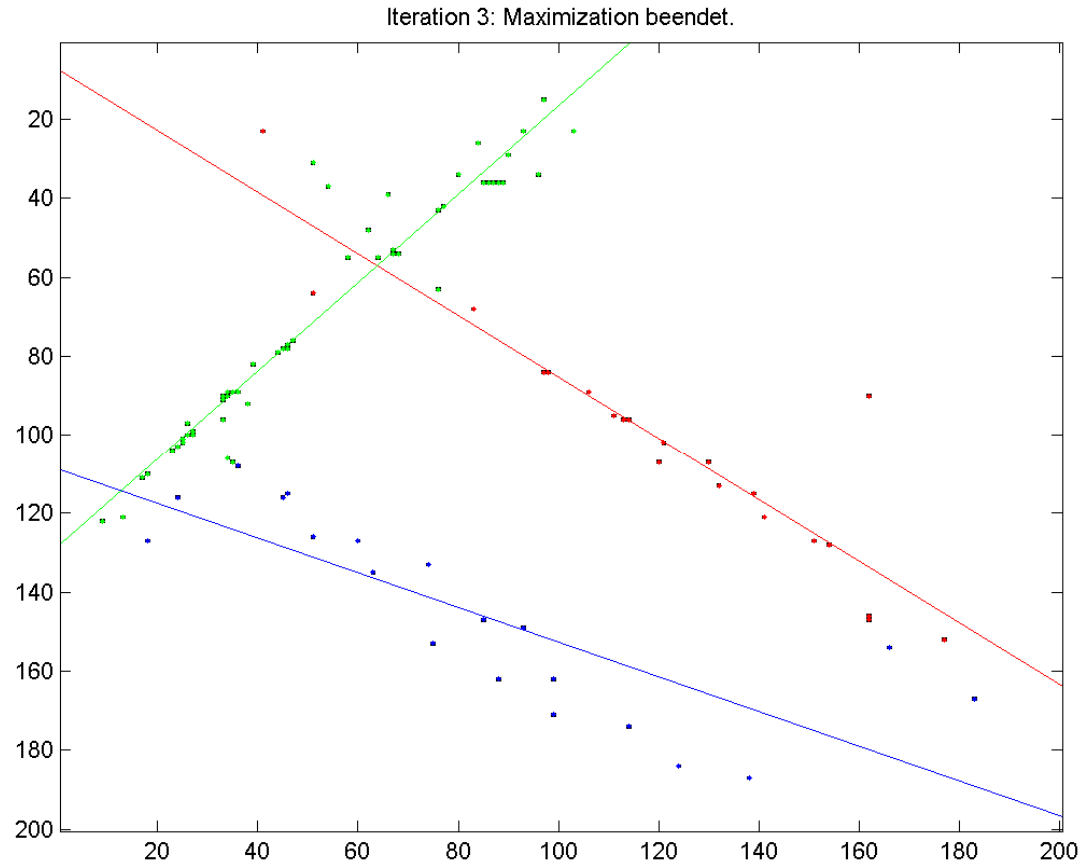


[Egorova 02]



3 Geraden – noch ein Beispiel

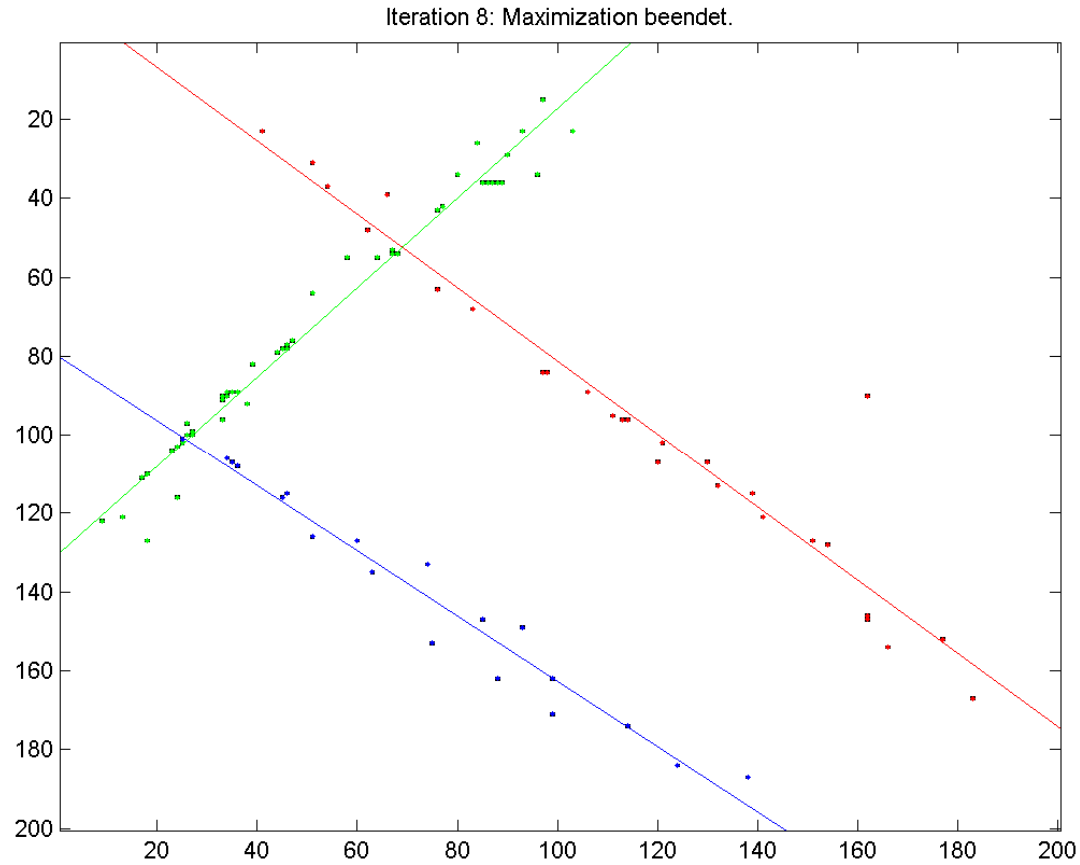
02. Juni 2017



[Egorova 02]

3 Geraden – noch ein Beispiel

02. Juni 2017

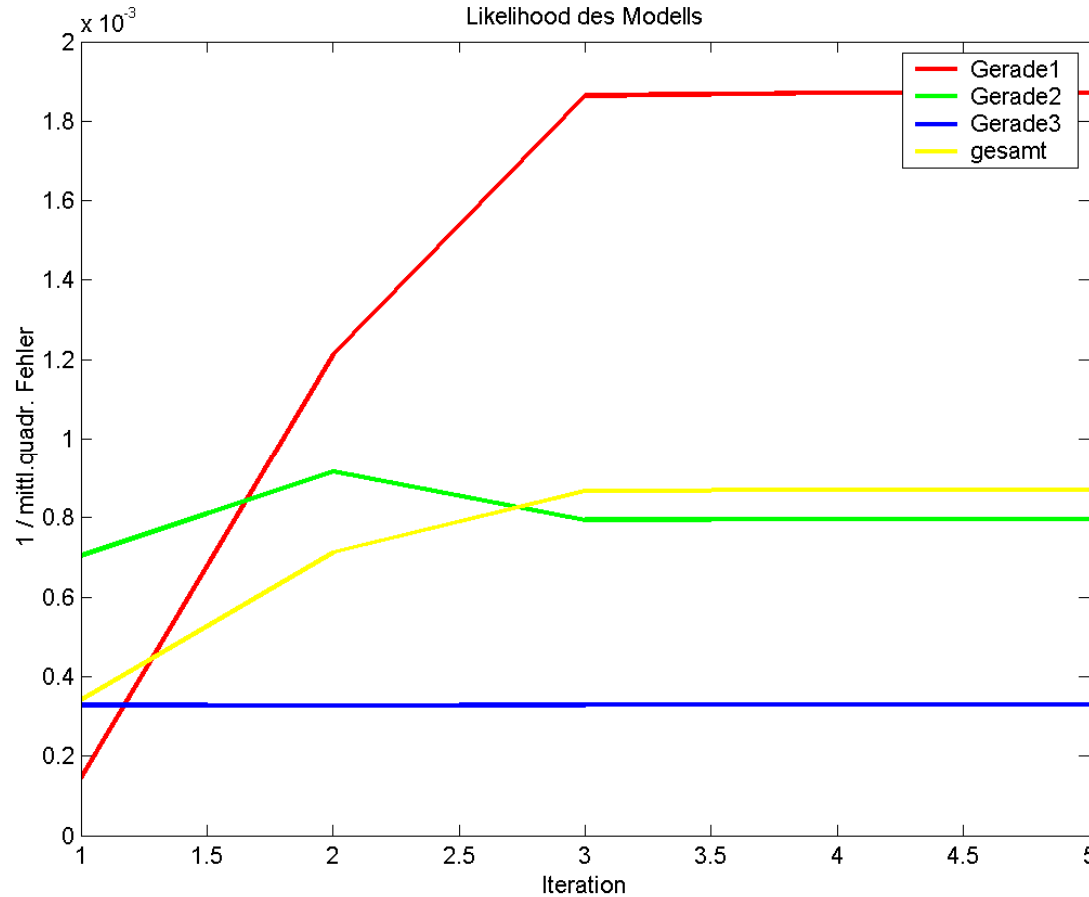


[Egorova 02]



3 Geraden – noch ein Beispiel

02. Juni 2017

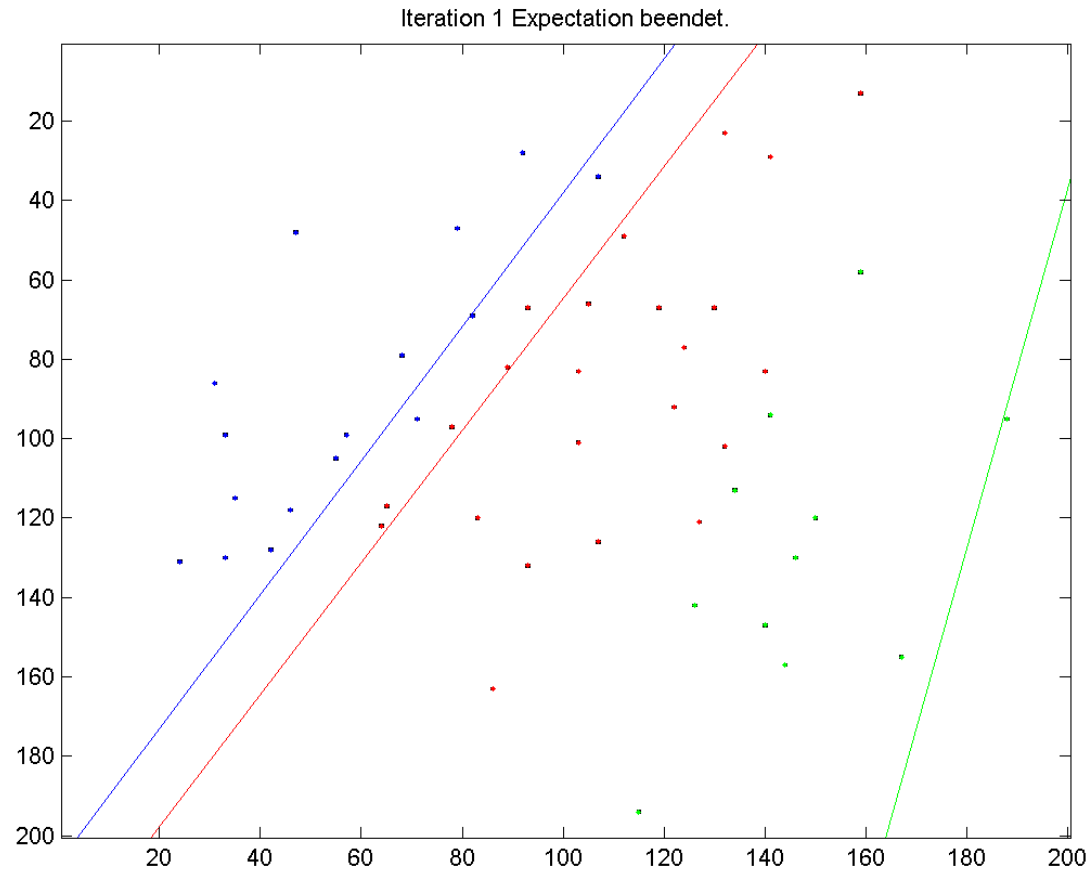


[Egorova 02]

drei Geraden – falsches Modell



02. Juni 2017

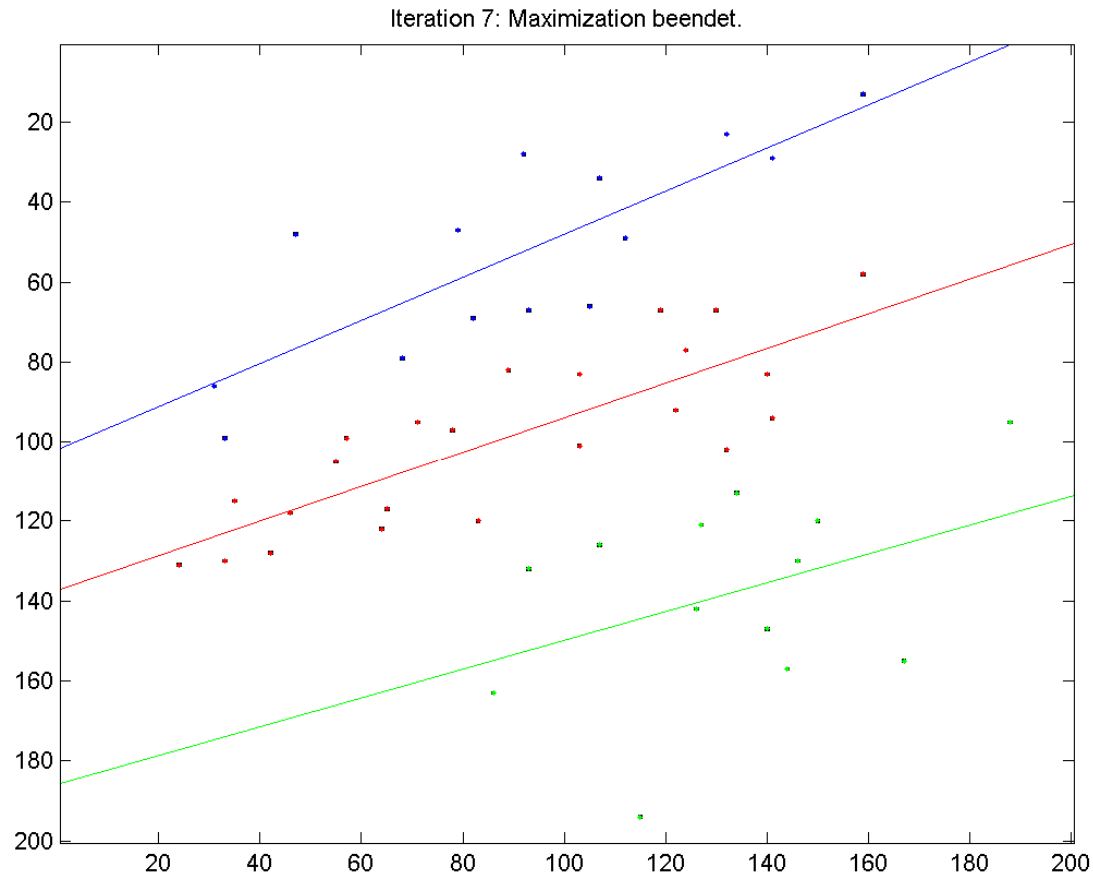


[Egorova 02]

drei Geraden – falsches Modell



02. Juni 2017

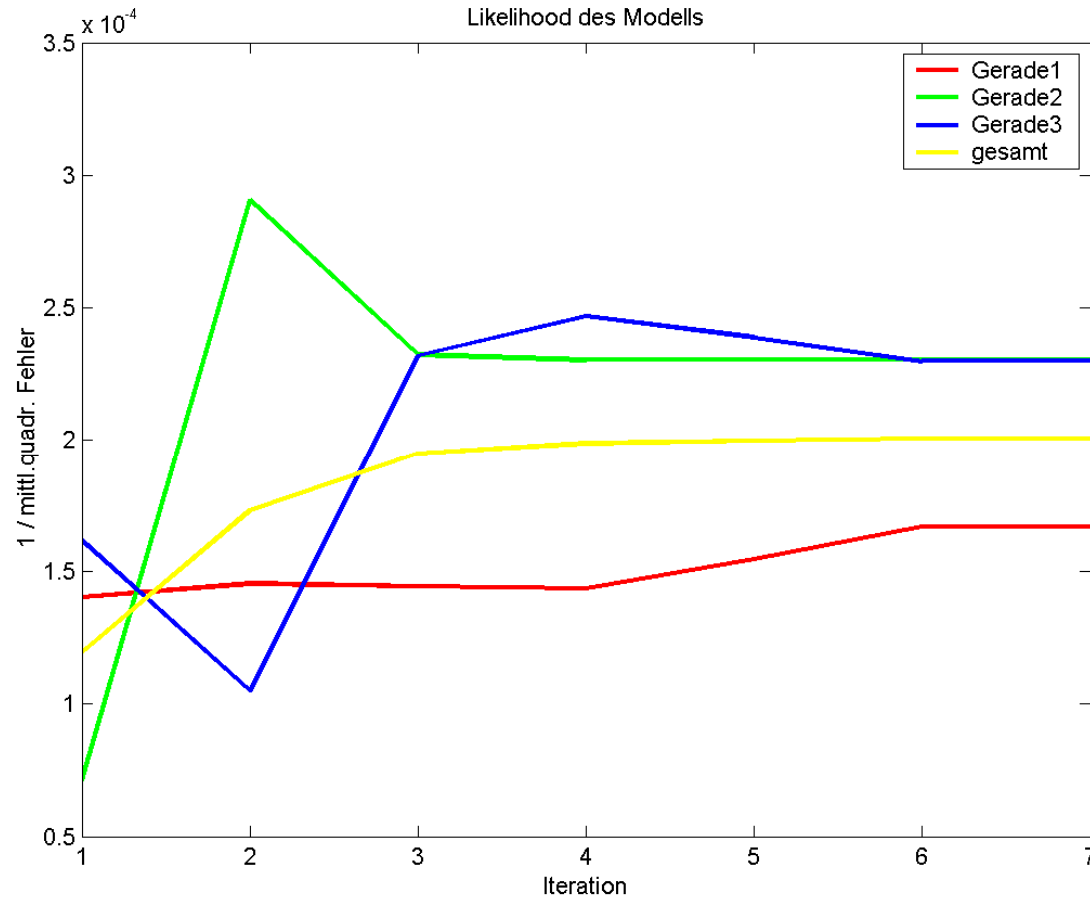


[Egorova 02]



drei Geraden – falsches Modell

02. Juni 2017

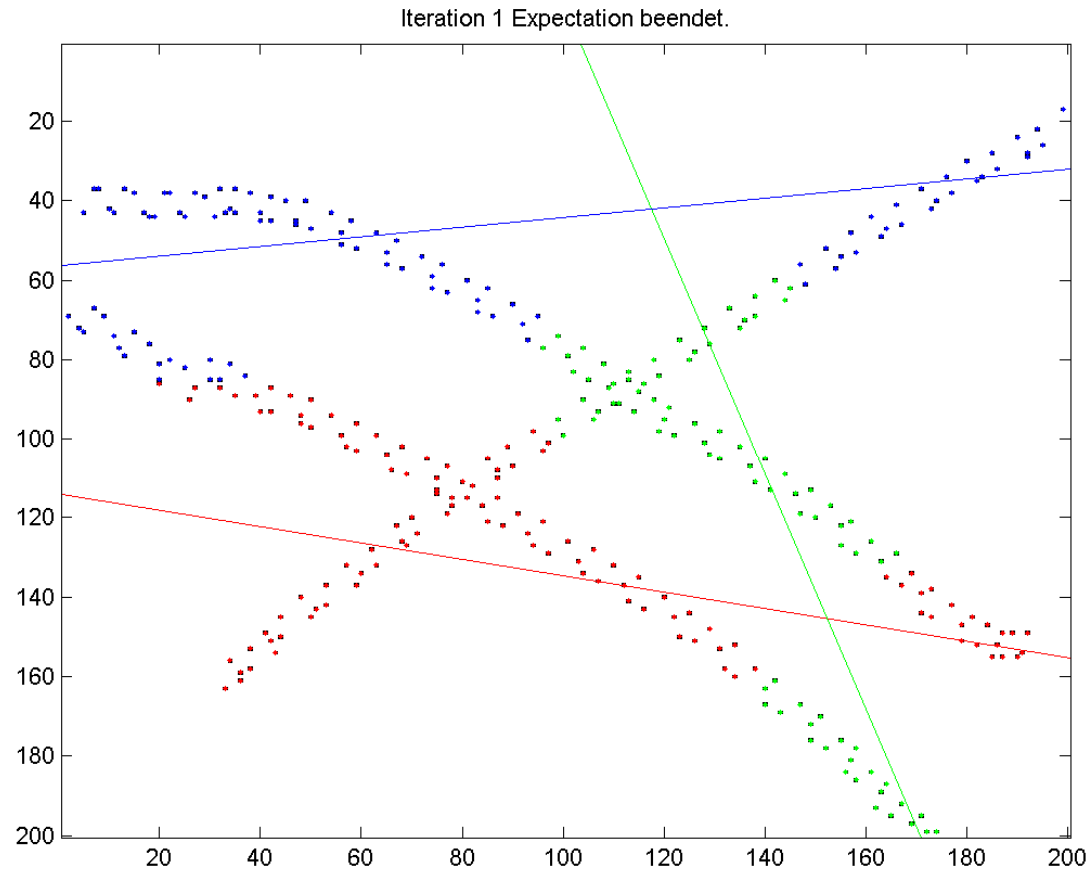


[Egorova 02]

drei Geraden – falsches Maximum



02. Juni 2017

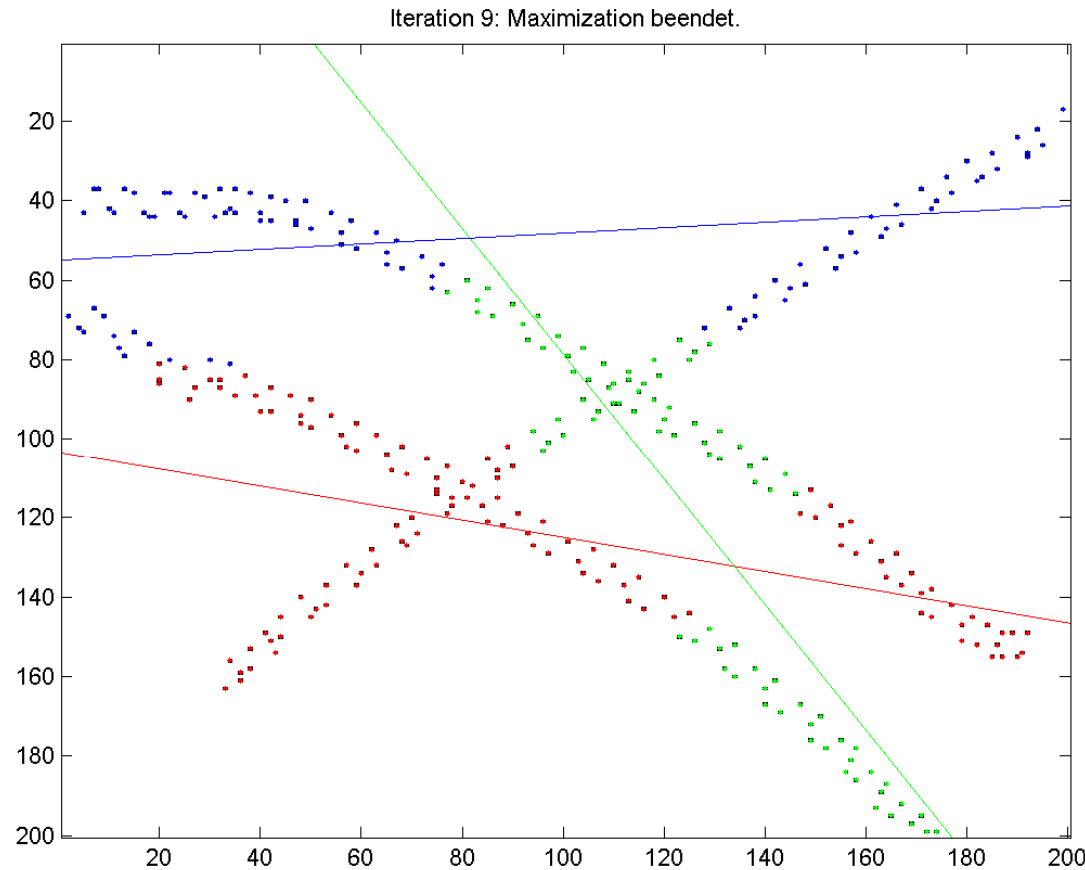


[Egorova 02]

drei Geraden – falsches Maximum



02. Juni 2017

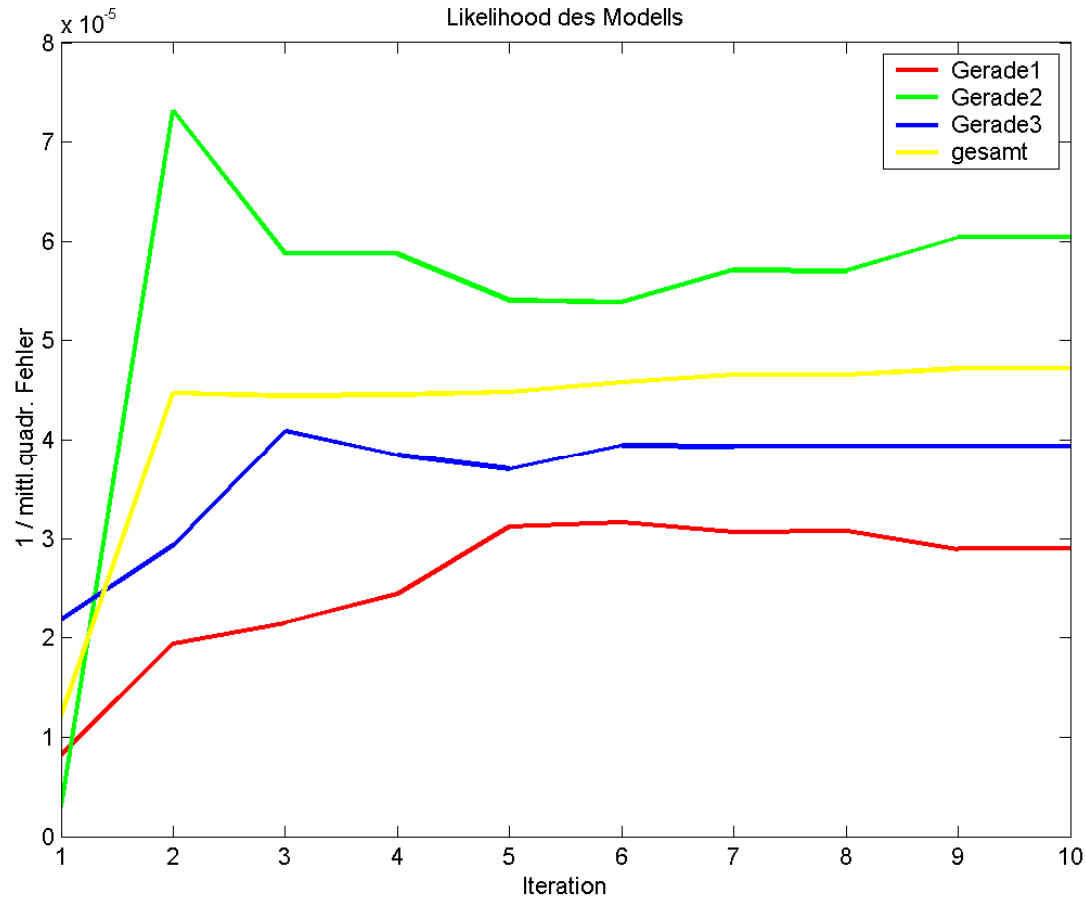


[Egorova 02]



drei Geraden – falsches Maximum

02. Juni 2017



[Egorova 02]



02. Juni 2017

Referenzen:

Wie immer:

R.O. Duda, P.E. Hart, and D.G. Stork, Pattern Classification, New York: John Wiley & Sons, 2001

Zusätzlich:

Wikipedia contributors, "Expectation–maximization algorithm," Wikipedia, The Free Encyclopedia, http://en.wikipedia.org/w/index.php?title=Expectation%E2%80%93maximization_algorithm (accessed June 25, 2013).

David McAllester, Expectation Maximization (EM), lecture note, TTIC 103 (CMSC 35420): Statistical Methods for Artificial Intelligence, Autumn 2007
<http://ttic.uchicago.edu/~dmcallester/ttic101-07/lectures/em/em.pdf>



Bayes'sche Netze



02. Juni 2017

Wissen wird durch Wahrscheinlichkeitsverteilungen repräsentiert. Manchmal sind kausale Abhängigkeiten bekannt. Dann bieten sich **Bayesian belief nets** an:

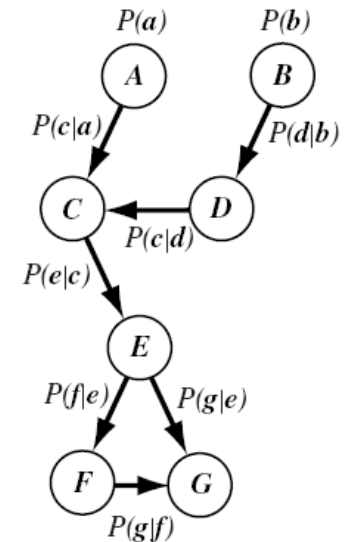


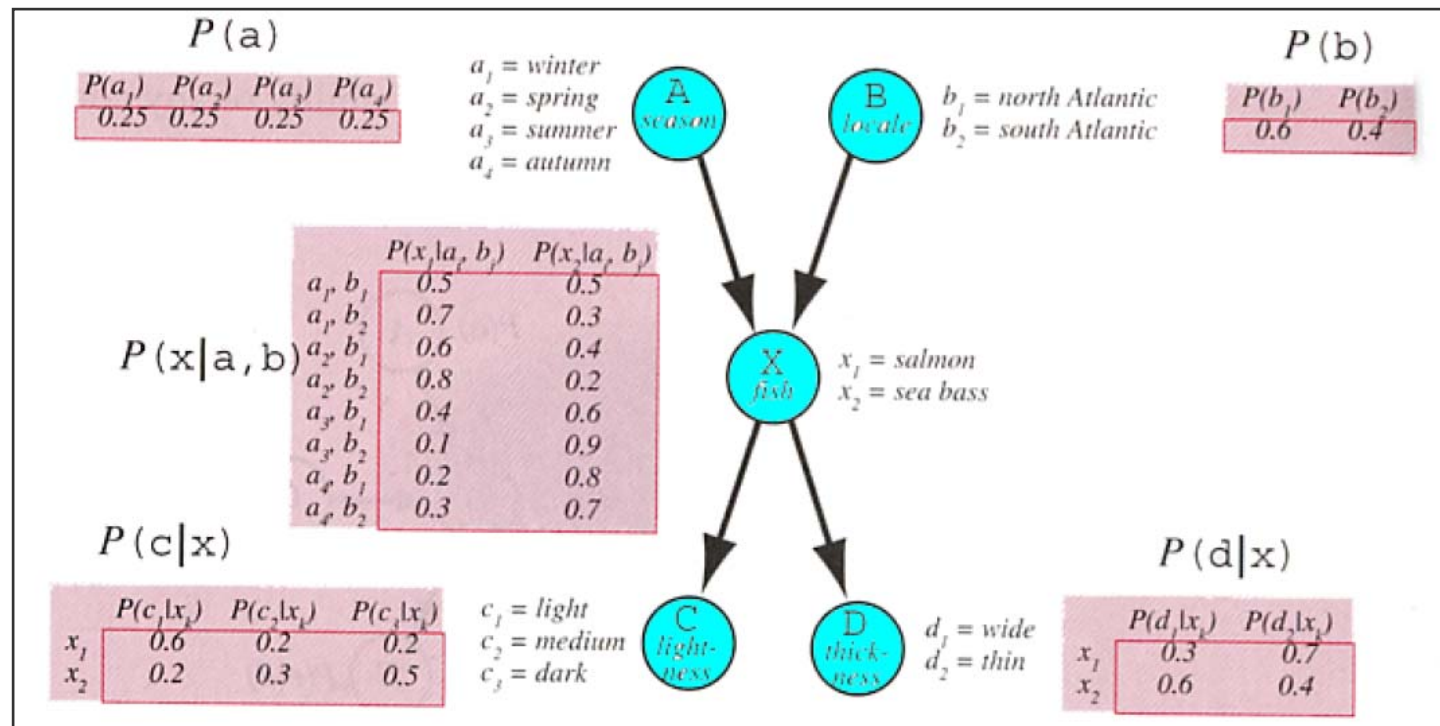
FIGURE 2.24. A belief network consists of nodes (labeled with uppercase bold letters) and their associated discrete states (in lowercase). Thus node **A** has states a_1, a_2, \dots , denoted simply **a**; node **B** has states b_1, b_2, \dots , denoted **b**, and so forth. The links between nodes represent conditional probabilities. For example, $P(c|a)$ can be described by a matrix whose entries are $P(c_i|a_j)$. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.



Kausalnetze

02. Juni 2017

Außerdem brauchen wir natürlich die bedingten Wahrscheinlichkeiten, die im diskreten Fall als Tabellen gegeben sind:





Wenn Abhängigkeit nicht bekannt, wird oft
Unabhängigkeit angenommen
(naive Bayes' rule, idiot Bayes' rule):

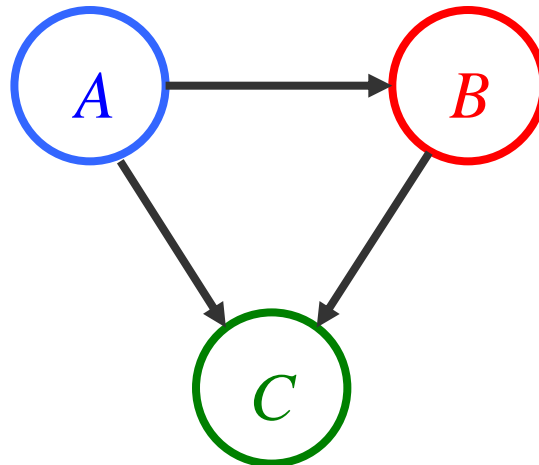
$$P(x | \mathbf{a}, \mathbf{b}) = P(x | \mathbf{a}) P(x | \mathbf{b})$$



Nocheinmal grundsätzlich:

Produktregel

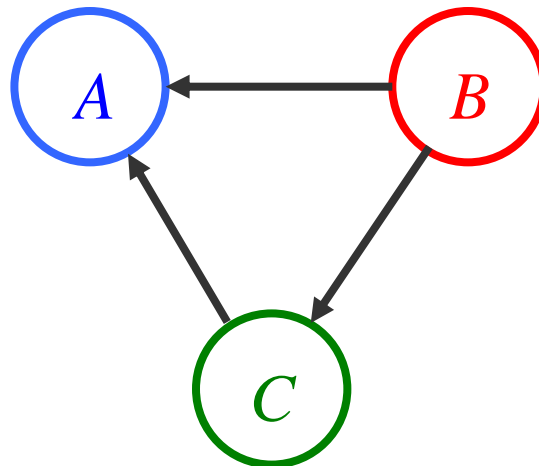
$$P(A, B, C) = P(A) P(B | A) P(C | A, B)$$





Andere, ebenso mögliche Faktorisierungen:

$$P(A, B, C) = P(B)P(C|B)P(A|B, C)$$



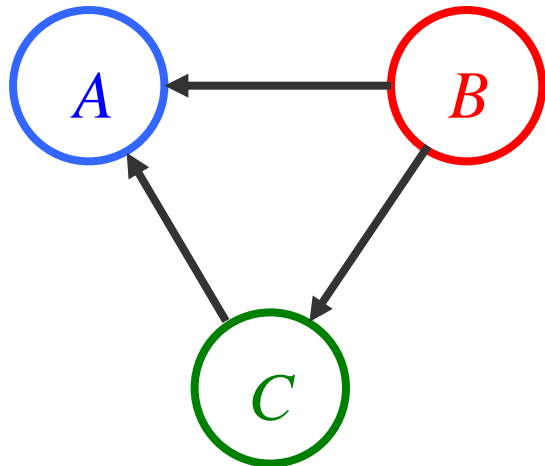
Kausalnetze

02. Juni 2017

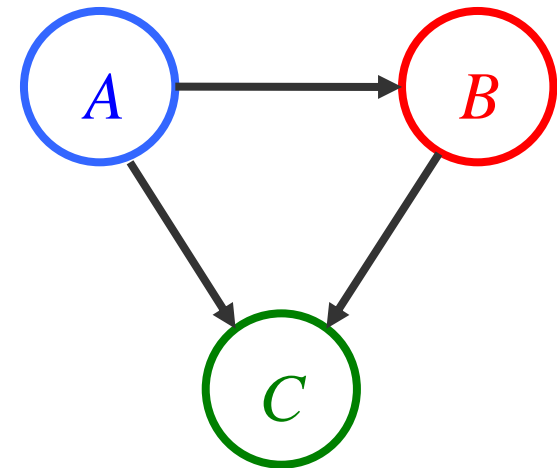
A – Fanggebiet, B – Fisch heute, C – Fisch gestern

Welche Faktorisierung ist besser,

$$P(A, B, C) = P(B)P(C|B)P(A|B, C)$$



$$P(A, B, C) = P(A)P(B|A)P(C|A, B)$$



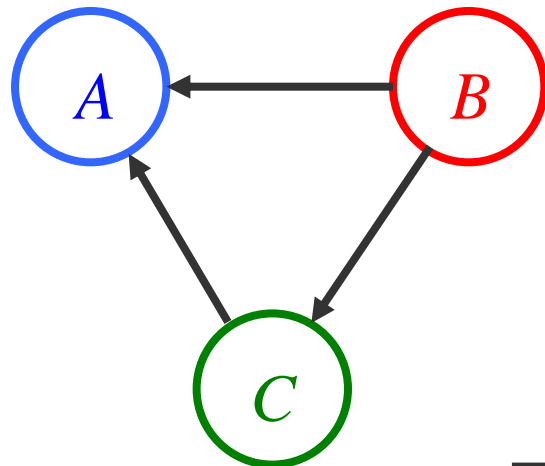
oder

?



Repräsentation der Verbundwahrscheinlichkeit durch Tabellen. Faktorisierung 1:

$$P(A, B, C) = P(B)P(C|B)P(A|B, C)$$



Lachs	Barsch
44,5%	55,5%

B,C	Nord	Süd
L,L	54,6%	45,4%
L,B	59,5%	40,5%
B,L	59,5%	40,5%
B,B	64,3%	35,7%

B	Lachs	Barsch
Lachs	44,6%	55,4%
Barsch	44,4%	55,6%

Achtung: $P(B) \neq P(C|B)$

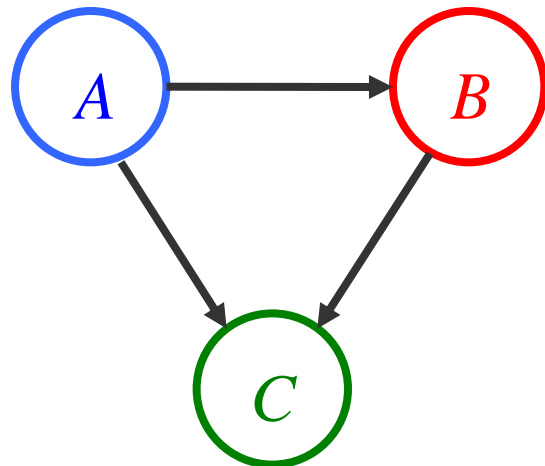


Kausalnetze

02. Juni 2017

Repräsentation der Verbundwahrscheinlichkeit durch Tabellen. Faktorisierung 2:

$$P(A, B, C) = P(A)P(B|A)P(C|A, B)$$



A,B	Lachs	Barsch
Nord,L	42,5%	57,5%
Nord,B	42,5%	57,5%
Süd,L	47,5%	52,5%
Süd,B	47,5%	52,5%

Nord	Süd
60%	40%

A	Lachs	Barsch
Nord	42,5%	57,5%
Süd	47,5%	52,5%

Kausalnetze

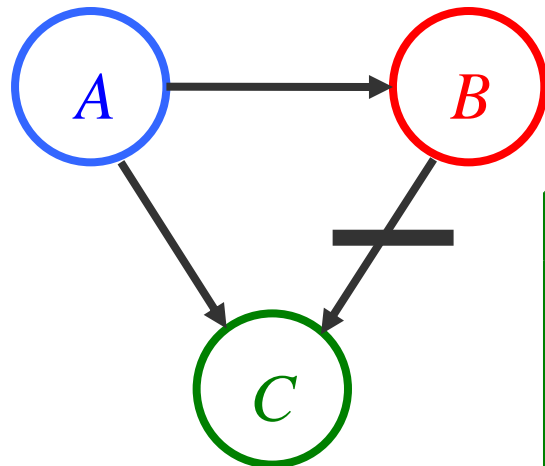
02. Juni 2017

Repräsentation der Verbundwahrscheinlichkeit durch Tabellen. Faktorisierung 2:

$$P(A, B, C) = P(A)P(B|A)P(C|A, B)$$

$$\rightarrow P(A, B, C) = P(A)P(B|A)P(C|A)$$

Nord	Süd
60%	40%



A,B	Lachs	Barsch
Nord,L	42,5%	57,5%
Nord,B	42,5%	57,5%
Süd,L	47,5%	52,5%
Süd,B	47,5%	52,5%

A	Lachs	Barsch
Nord	42,5%	57,5%
Süd	47,5%	52,5%

A	Lachs	Barsch
Nord	42,5%	57,5%
Süd	47,5%	52,5%

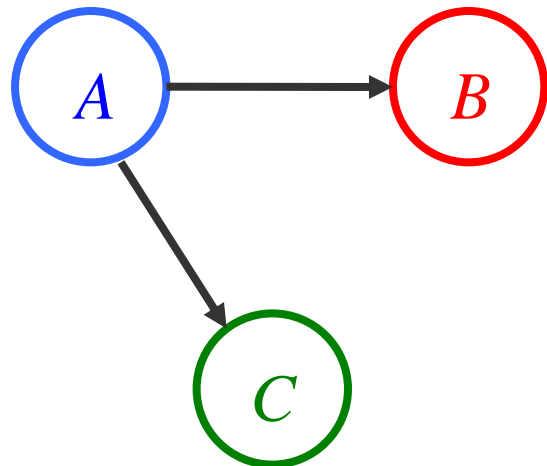


Kausalnetze

02. Juni 2017

Repräsentation der Verbundwahrscheinlichkeit durch Tabellen. Faktorisierung 2:

$$P(A, B, C) = P(A)P(B|A)P(C|A)$$



Faktorisierung 2
 ist besser, nur
 so kann Pfeil
 weggelassen
 werden, und
 außerdem gilt:

$$P(B|A) = P(C|A)$$

Nord	Süd
60%	40%

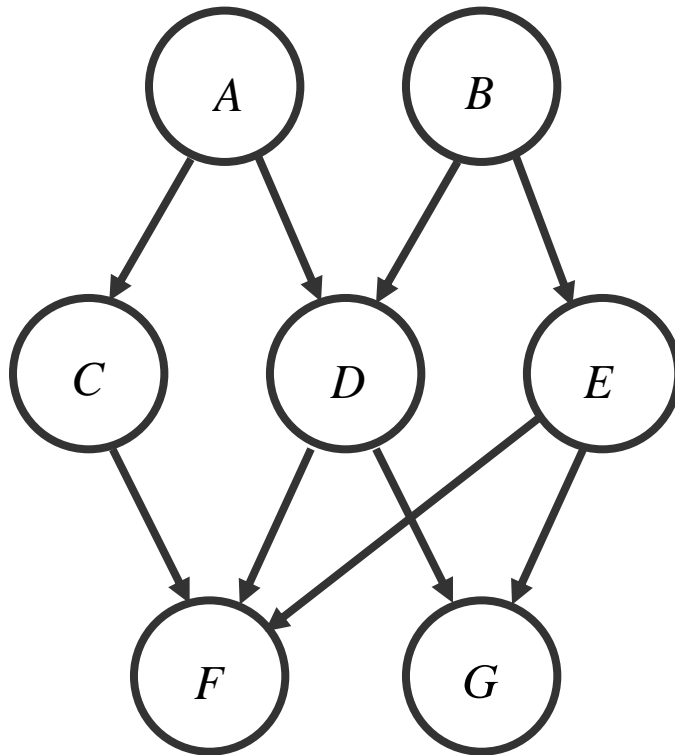
A	Lachs	Barsch
Nord	42,5%	57,5%
Süd	47,5%	52,5%

=

A	Lachs	Barsch
Nord	42,5%	57,5%
Süd	47,5%	52,5%



Umwandlung von Graph in Faktorisierung



$$\begin{aligned} &P(A, B, C, D, E, F, G) \\ &= \prod_{X \in \{A, B, C, D, E, F, G\}} P(X \mid \text{Eltern}(X)) \\ &= P(A)P(B) \\ &\quad * P(C \mid A)P(D \mid A, B)P(E \mid B) \\ &\quad * P(F \mid C, D, E)P(G \mid D, E) \end{aligned}$$



Angenommen, ich kenne die Verbundwahrscheinlichkeit:

- Wie berechne ich Wahrscheinlichkeit für einzelne Zufallsvariable?
- Was lerne ich, wenn ich den Zustand einer Zufallsvariablen erfahre, wie verwerte ich dieses Wissen?



Wie berechne ich Wahrscheinlichkeit für einzelne Zufallsvariable, wenn $\text{prob}(A, B, C)$ bekannt ist?

- Durch **Marginalisierung** über alle anderen unbekannten Zufallsvariablen:

$$\text{prob}(A) = \sum_{b_i: \text{ alle möglichen Werte von } B} \left(\sum_{c_k: \text{ alle möglichen Werte von } C} \text{prob}(A, b_i, c_k) \right)$$



Was lerne ich, wenn ich den Zustand einer Zufallsvariablen erfahre ($\text{prob}(A, B, C)$ bekannt)?

Antwort: Conditioning bzw. Inferenz

Conditioning:

- Es interessieren nur entsprechende Zeilen/Ebenen... aus der Tabelle. Beispiel: es wird bekannt, dass $A = a_r$. Dann wird Verteilung zu:

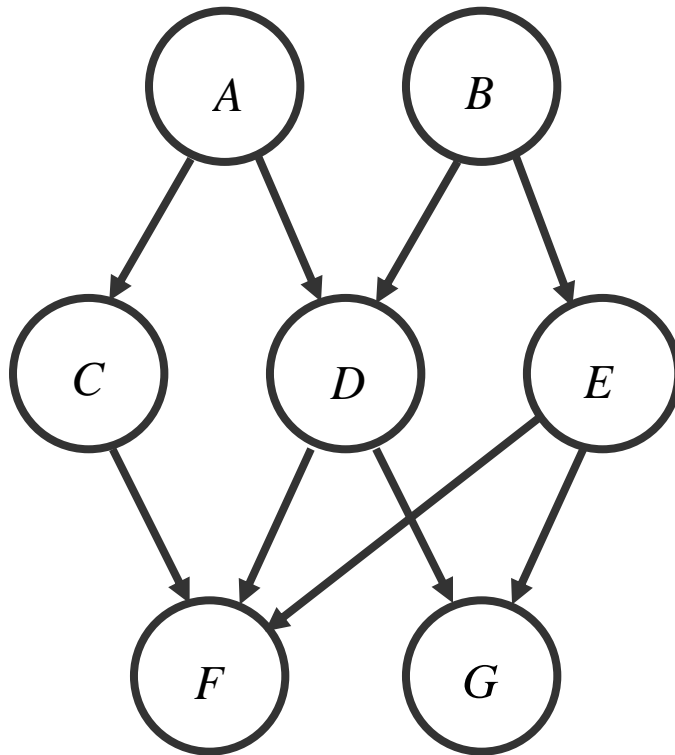
$$\text{prob}(B, C \mid A = a_r) = \frac{\text{prob}(B, C, A = a_r)}{\text{prob}(A = a_r)}$$



Conditioning, Marginalisierung, Inferenz

02. Juni 2017

Wenn Verteilung als Bayesnetz gegeben, ist Conditioning einfach, wenn Zustand von Wurzelknoten bekannt wird, da bedingte Wahrscheinlichkeiten gegeben sind:



$$\begin{aligned}
 &P(A, B, C, D, E, F, G) \\
 &= \prod_{X \in \{A, B, C, D, E, F, G\}} P(X \mid \text{Eltern}(X)) \\
 &= P(A)P(B) \\
 &\quad * P(C \mid A)P(D \mid A, B)P(E \mid B) \\
 &\quad * P(F \mid C, D, E)P(G \mid D, E)
 \end{aligned}$$



02. Juni 2017

Conditioning und Marginalisierung sieht dann z.B. so aus
(wenn $\text{prob}(B | a_r)$ gesucht ist:

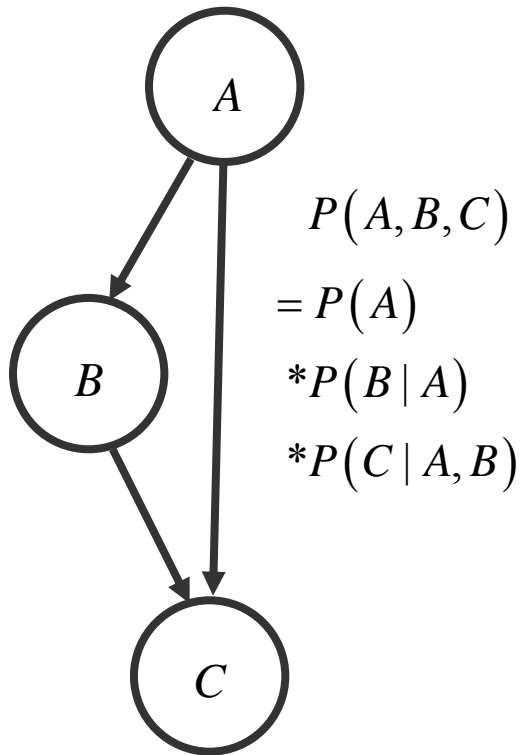
$$\text{prob}(B | a_r) = \sum_{c_k: \text{alle möglichen Werte von } C} \text{prob}(B, c_k | a_r)$$



Inferenz

02. Juni 2017

Schließen vom Zustand der Kindknoten auf Zustand bzw. Wahrscheinlichkeiten für Elternknoten wird **Inferenz** genannt.
Beispiel: Zustand von C bekannt ($C = c_r$).





Inferenz

02. Juni 2017

Wir wollen z.B. $P(A | C = c_r)$ berechnen.

Das geschieht durch Einsetzen von c und Marginalisierung:

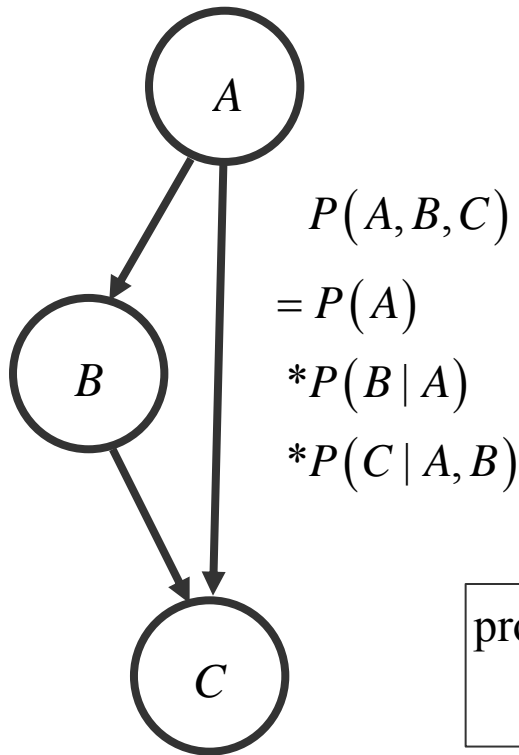
$$\text{prob}(A, c_r) = \sum_{b_k: \text{alle möglichen Werte von B}} \text{prob}(A, B = b_k, C = c_r)$$

Anschließend muss normiert werden:

$$\text{prob}(A | c_r) = \frac{\text{prob}(A, c_r)}{\text{prob}(c_r)}$$

Dabei ist

$$\text{prob}(c_r) = \sum_{a_j: \text{alle möglichen Werte von A}} \sum_{b_k: \text{alle möglichen Werte von B}} \text{prob}(A = a_j, B = b_k, C = c_r)$$

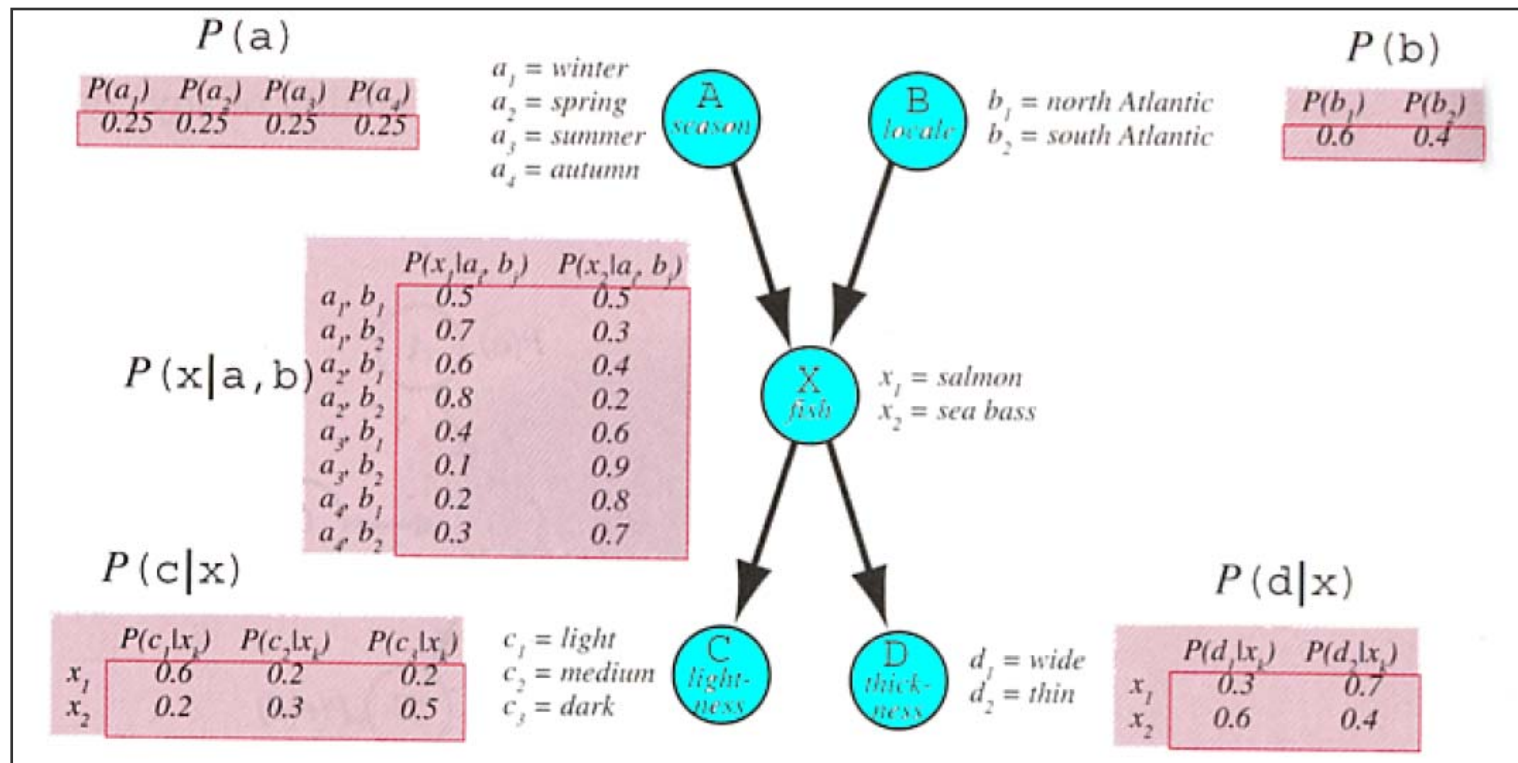




Kausalnetze

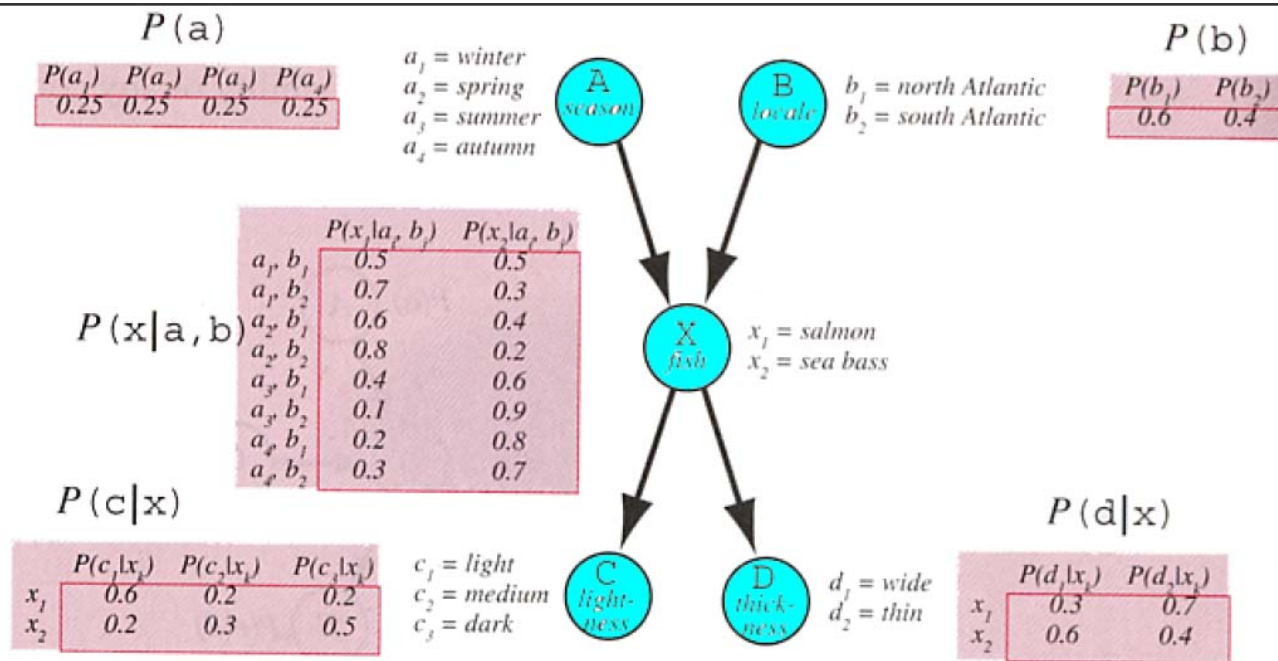
02. Juni 2017

Beispiel: Fisch ist hell und im Südatlantik gefangen, mehr wissen wir nicht.





02. Juni 2017



$$\begin{aligned}
 P(x_1|c_1, b_2) &= \frac{P(x_1, c_1, b_2)}{P(c_1, b_2)} \\
 &= \alpha \sum_{a, d} P(x_1, a, b_2, c_1, d) \\
 &= \alpha \sum_{a, d} P(a)P(b_2)P(x_1|a, b_2)P(c_1|x_1)P(d|x_1) \\
 &= \alpha P(b_2)P(c_1|x_1) \\
 &\quad \times \left[\sum_a P(a)P(x_1|a, b_2) \right] \left[\sum_d P(d|x_1) \right] \\
 &= \alpha P(b_2)P(c_1|x_1) \\
 &\quad \times [P(a_1)P(x_1|a_1, b_2) + P(a_2)P(x_1|a_2, b_2) + P(a_3)P(x_1|a_3, b_2) + P(a_4)P(x_1|a_4, b_2)] \\
 &\quad \times \underbrace{[P(d_1|x_1) + P(d_2|x_1)]}_{=1} \\
 &= \alpha(0.4)(0.6) [(0.25)(0.7) + (0.25)(0.8) + (0.25)(0.1) + (0.25)(0.3)] 1.0 \\
 &= \alpha 0.114.
 \end{aligned}
 \tag{100}$$



Fortsetzung Beispiel

Entsprechende Rechnung ergibt:

$$P(x_2|c_1, b_2) = \alpha 0.066$$

Normierung auf 1:

$$P(x_1|c_1, b_2) = 0.63 \text{ and } P(x_2|c_1, b_2) = 0.37.$$

Also ist Fisch wahrscheinlich Lachs.