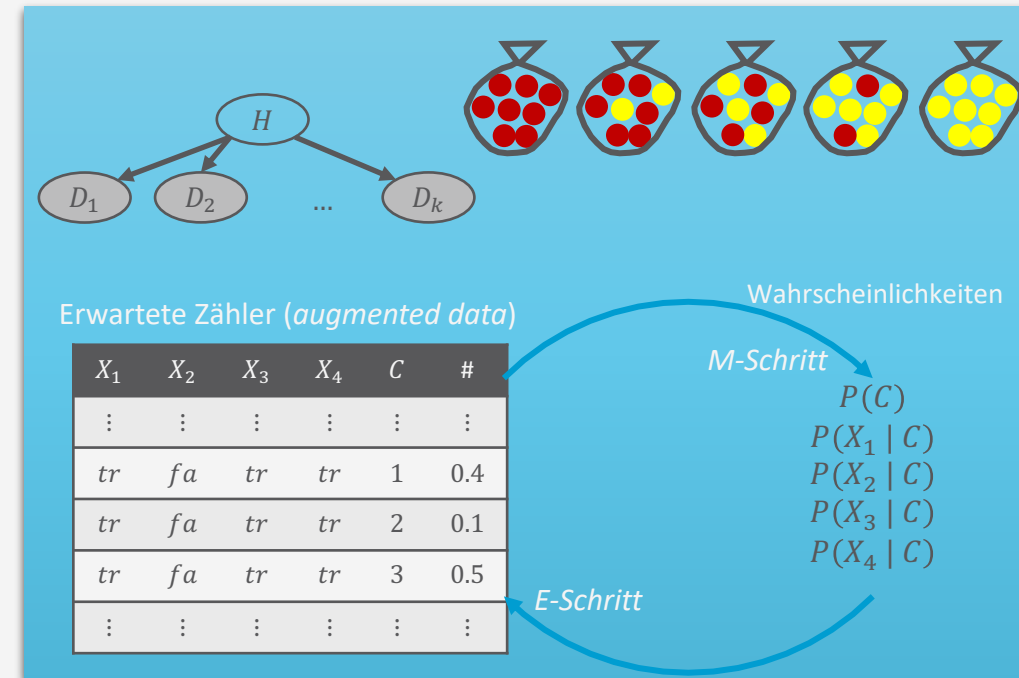


Lernalgorithmen für Episodische PGMs

Einführung in die
Künstliche Intelligenz



Inhalte

1. Künstliche Intelligenz & Agenten

- Agentenabstraktion, Rationalität
- Aufgabenumgebung

2. Episodische PGMs

- Gerichtetes Modell: Bayes Netze (BNs)
- Ungerichtete Modelle

3. Exakte Inferenz in episodischen PGMs

- Wahrscheinlichkeits- und Zustandsanfragen
- Direkt auf den Modellen, mittels Hilfsstrukturen

4. Approximative Inferenz in episodischen PGMs

- Wahrscheinlichkeitsanfragen
- Deterministische, stochastische Algorithmen

5. Lernalgorithmen für episodische PGMs

- Bei (nicht) vollständigen Daten, (un)bekannter Struktur

6. Sequentielle PGMs und Inferenz

- Dynamische BNs, Hidden-Markov-Modelle
- filtering / prediction / hindsight Anfragen, wahrscheinlichste Zustandssequenz
- Exakter, approximativer Algorithmus

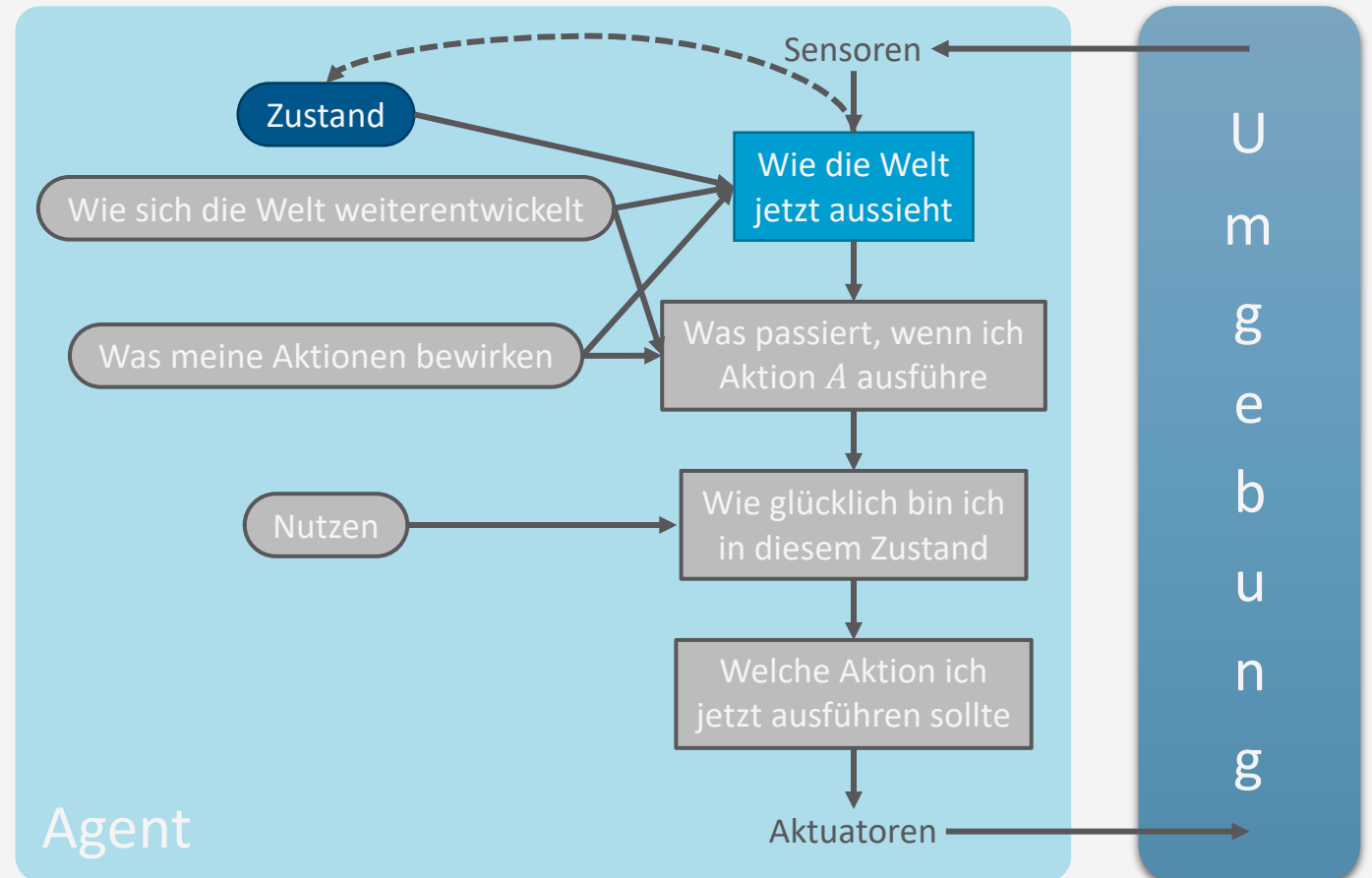
7. Entscheidungstheoretische PGMs

- Präferenzen, Nutzenprinzip
- PGMs mit Entscheidungs- und Nutzenknoten
- Berechnung der besten Aktion (Aktionssequenz)

8. Abschlussbetrachtungen

Einordnung der Vorlesung: *Modell- und nutzenbasierter Agent*

- Nachfolgende Themen der Vorlesung
 2. Episodische PGMs
 3. Exakte Inferenz in episodischen PGMs
 4. Approximative Inferenz in episodischen PGMs
 5. Lernalgorithmen für episodische PGMs
 6. Sequentielle PGMs und Inferenz
 7. Entscheidungstheoretische PGMs



Literaturhinweise

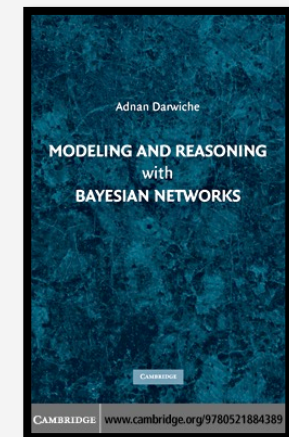
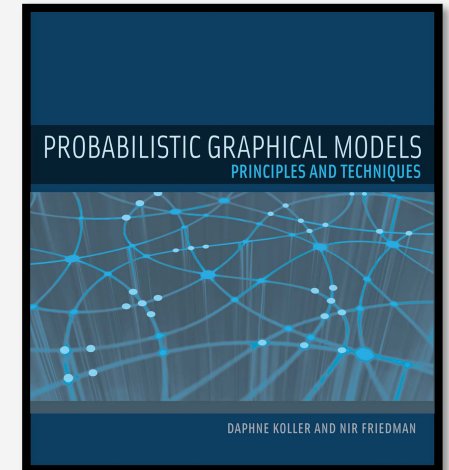
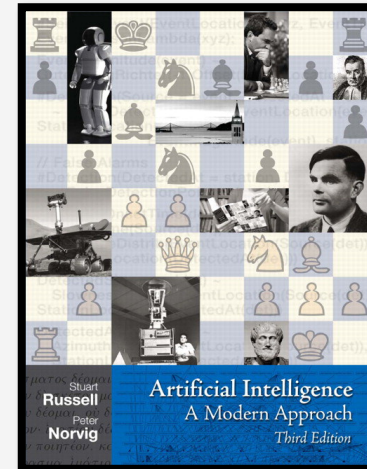
Inhalte dieses Themenblocks werden in den folgenden Kapiteln der Vorlesungsbücher behandelt

- AIMA(de)
 - Kap. 20: Lernen probabilistischer Modelle
- PGM
 - *Teil III: Lernen*

Dieser Thementeil basiert auf AIMA(de), Kap. 20, sowie den Papieren zu LDA und IFP wie auf den Folien zitiert, mit einigen Hinweisen rechts und links über den Tellerrand.

Hinweis: Das PGM Buch deckt sehr viel mehr ab, als die hier besprochenen Inhalte.

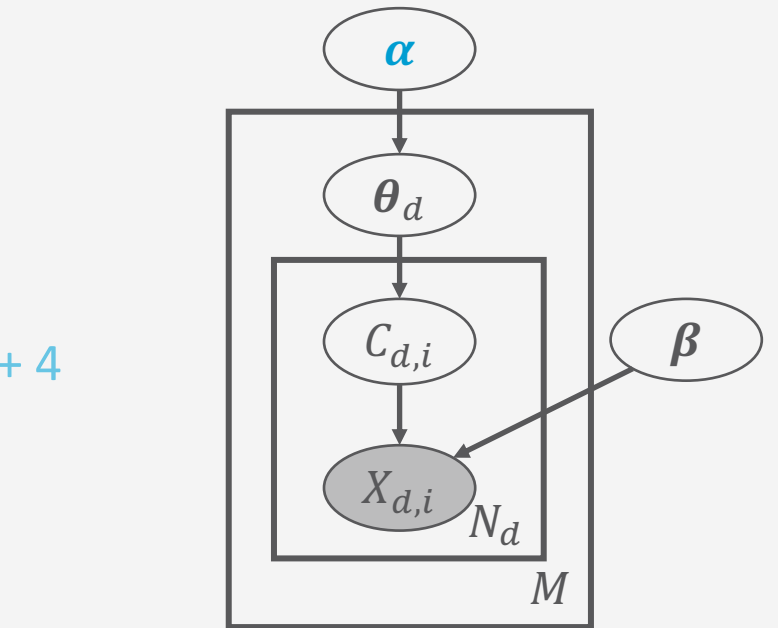
- Wer gerne ein anderes Buch ausprobieren möchte (Fokus auf BNs):
 - Adnan Darwiche, *Modelling and Reasoning with Bayesian Networks*, 2009.



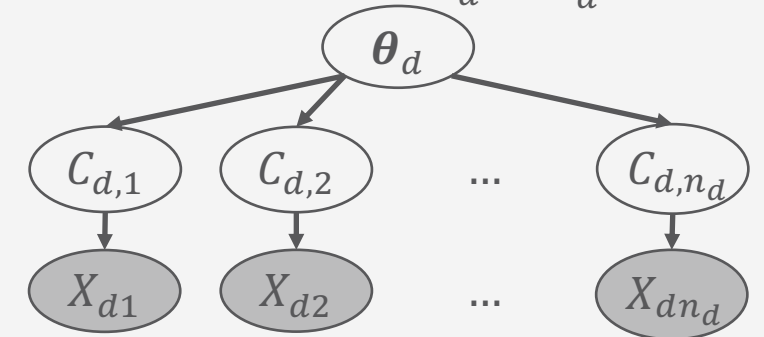
Inferenzaufgaben rund um Topic Modelle

- Gegeben ein Topic Modell für ein Korpus
 - Anfrage: Gegeben ein unbekanntes Dokument, welches Topic / welche Topic-Verteilung hat es?
 - Exakte + approximative Anfragenbeantwortung in PGMs: [Thema 3 + 4](#)
 - Generierung eines neuen Dokuments: Gegeben ein gewähltes Topic, generiere die Worte eines neuen Dokuments
 - Sampling aus PGMs: [Thema 4](#)
- Gegeben ein Korpus

- Lerne ein Topic Modell: [Thema 5](#)
 - Gegeben: Worte in den Dokumenten (nach Vorverarbeitung), Anzahl Topics K , Hyperparameter α
 - Lerne: Topicverteilungen pro Dokument, Wortverteilungen pro Topic (Basis-Verfahren: [Latent Dirichlet Allocation, LDA](#))



Pro Dokument d mit $\theta_d \neq \theta_{d'}$:



Überblick: 5. Lernalgorithmen für episodische PGMs

A. Überblick

- Full Bayesian Learning, MAP Learning, Maximum Likelihood (ML) Learning

B. ML-basiertes Parameter Lernen

- In BNs: ML (vollständige Daten), Expectation-Maximisation (EM; nicht-vollständige Daten)
- In Faktormodellen: Iterative Proportional Fitting (IPF), EM-IPF
- Anwendung: LDA

C. ML-basiertes Struktur Lernen

- Model Selection, Structural EM

D. Abschlussbemerkungen

- Alternativen, generative vs. diskriminative Modelle, Korrelation vs. Kausalität

Lernen

- Aufgabe: Lernen einer Verteilung
 - Gegeben Daten (als Evidenz)
 - Menge von zusammengesetzten Ereignissen
 - Annahme: aus Verteilung erzeugt
 - Parameter (Einträge in Verteilung) bestimmen, so dass die Daten am wahrscheinlichsten werden
 - **Hypothese**: Mögliche Parameterausprägung
 - Theorie darüber, wie Umgebung funktioniert
 - Hypothesenraum: alle möglichen Ausprägungen
- Lernproblem anhand Vorhersage eines nächsten Datums anschauen
 - Welche Verteilung generiert das Datum?
- Lernansätze
 - **Full Bayesian Learning**
 - Unter Berücksichtigung aller möglichen Hypothesen nächstes Datum vorhersagen
 - **Maximum A Posteriori (MAP) Learning**
 - Unter Berücksichtigung der MAP Hypothese nächstes Datum vorhersagen
 - **Maximum Likelihood (ML) Learning**
 - Unter Berücksichtigung der ML Hypothese nächstes Datum vorhersagen

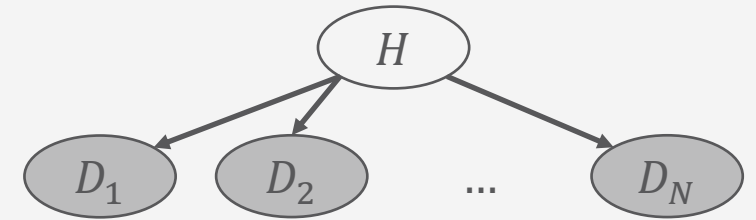
Quasi Sampling rückwärts:

- Sampling: Gegeben Verteilung, Daten erzeugen
- Lernen: Gegeben Daten, Verteilung erzeugen

Full Bayesian Learning

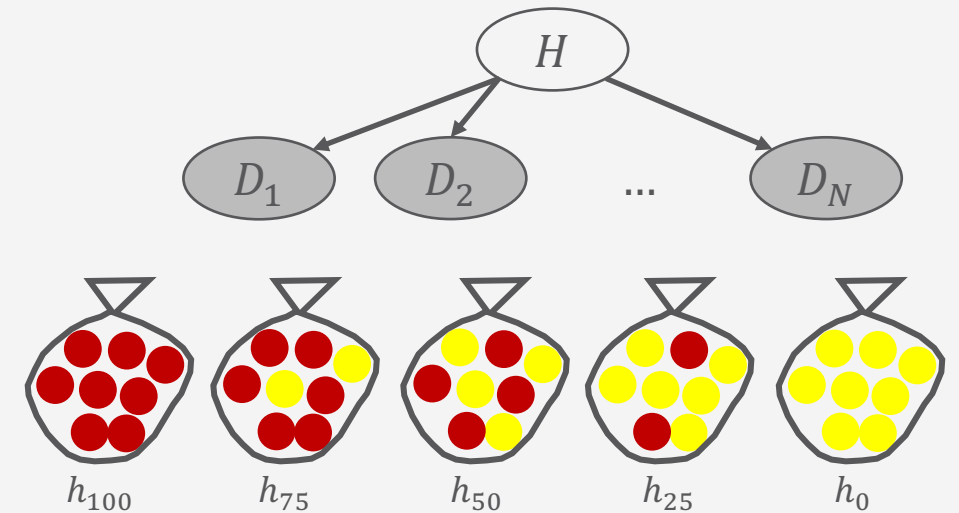
- Full Bayesian Learning
 - Aktualisierung einer Wahrscheinlichkeitsverteilung über den Hypothesenraum gegeben Daten
- H ist die Hypothesen-Zufallsvariable
 - Mögliche Hypothesen (Domäne von H) $\text{Val}(H) = \{h_1, \dots, h_n\}$
 - D.h. Zufallsvariable mit dem Hypothesenraum als Domäne
 - $P(H)$ = a priori Wahrscheinlichkeitsverteilung über Hypothesenraum
- Evidenzvariablen D_j für Beobachtungen über $j = 1, \dots, N$
 - Trainingsdaten $\mathbf{d} = d_1, \dots, d_N$
 - Bedingte Wahrscheinlichkeiten $P(\mathbf{d} | h_i)$: Gegeben h_i , Wahrscheinlichkeit der Daten \mathbf{d}
 - Bei unabhängig und identisch verteilten Daten (*independent and identically distributed, i.i.d.*):

$$P(\mathbf{d} | h_i) = P(d_1 | h_i) \cdot \dots \cdot P(d_N | h_i) = \prod_{j=1}^N P(d_j | h_i)$$



Beispiel

- Wir nehmen an, es gibt 5 Arten von Bonbon-Tüten
 - Bonbons gleich verpackt
 - 10% sind 100% Kirsch-Bonbons (h_{100})
 - 20% sind 75% Kirsch + 25% Zitronen-Bonbons (h_{75})
 - 40% sind 50% Kirsch + 50% Zitronen-Bonbons (h_{50})
 - 20% sind 25% Kirsch + 75% Zitronen-Bonbons (h_{25})
 - 10% sind 100% Zitronen-Bonbons (h_0)
- Bezeichnen wir mit θ den Parameter, der den Anteil an Kirsch-Bonbons spezifiziert und h_θ die dazugehörige Hypothese \rightarrow 5 Hypothesen $h_{100}, h_{75}, h_{50}, h_{25}, h_0$, $P(H) = (0.1, 0.2, 0.4, 0.2, 0.1)$
- Dann beobachten wir Bonbons aus einer Tüte: ●●●●●●●●●● ($d_1, \dots, d_{10} = 10$ mal Zitronen-Bonbons)
- Frage: Welche der 5 Arten von Tüten hat die 10 Beobachtungen generiert? $P(h_\theta | \mathbf{d})$
 - Welchen Geschmack wird der nächste Bonbon haben? Vorhersage $P(X | \mathbf{d}) = P(D_{N+1} | \mathbf{d})$

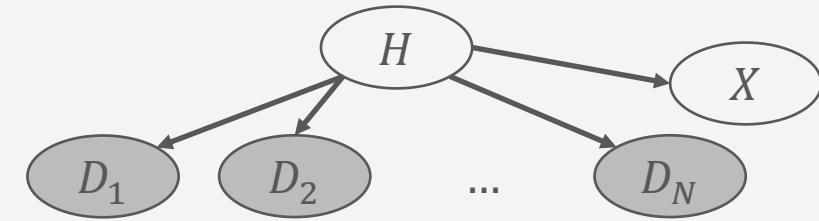


Full Bayesian Learning

- Gegeben die Daten soweit hat jede Hypothese h_i eine a posteriori Wahrscheinlichkeit:
 - $P(h_i | \mathbf{d}) \propto P(\mathbf{d} | h_i)P(h_i)$ (Bayes Theorem)
 - $P(\mathbf{d} | h_i)$ Likelihood der Daten unter Hypothese h_i genannt
- Vorhersagen über eine neue Entität X sind der gewichtete Durchschnitt über die Vorhersagen aller Hypothesen:

$$\begin{aligned}
 P(X | \mathbf{d}) &= \sum_{h_i \in \text{Val}(H)} P(X, h_i | \mathbf{d}) && \text{(Nicht-Anfragevariable aussummieren)} \\
 &= \sum_{h_i \in \text{Val}(H)} P(X | h_i, \mathbf{d})P(h_i | \mathbf{d}) && \text{(Kettenregel)} \\
 &= \sum_{h_i \in \text{Val}(H)} P(X | h_i)P(h_i | \mathbf{d}) && \text{(X unabhängig von } \mathbf{d} \text{ gegeben } h_i) \\
 &\propto \sum_{h_i \in \text{Val}(H)} P(X | h_i)P(\mathbf{d} | h_i)P(h_i) && \text{(Bayes Theorem – wie oben)}
 \end{aligned}$$

- Vorteil: Nicht eine Hypothese wählen müssen
- Nachteil: Iteration über alle Hypothesen

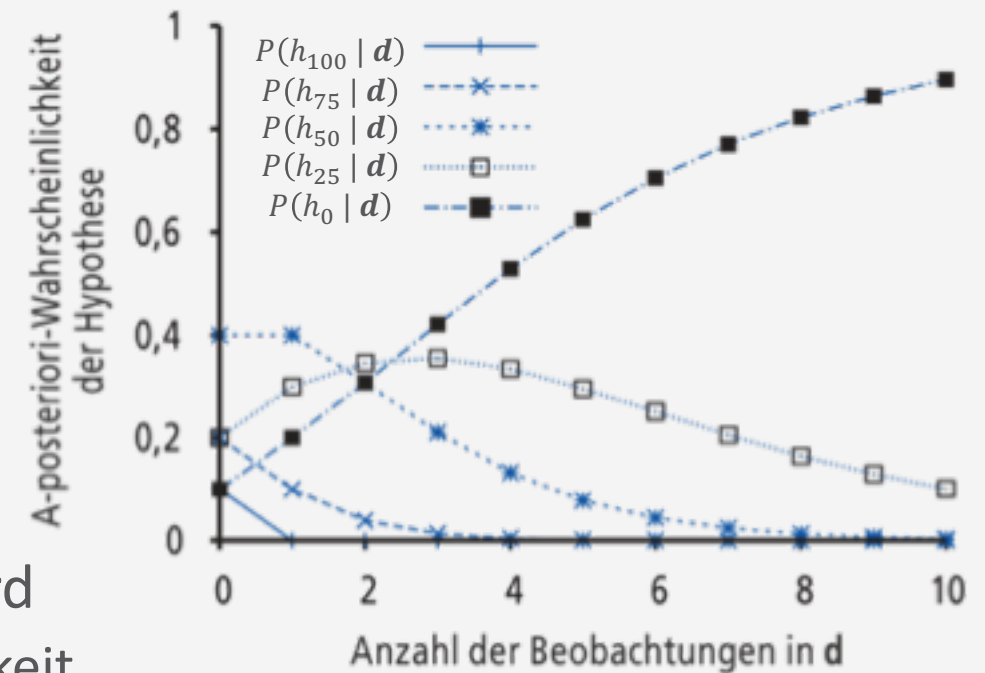


Beispiel: Likelihood

- Wenn unsere Tüten groß genug sind oder genauer, wir jeden Bonbon nach dem Check der Sorte wieder verpacken und in die Tüte zurücklegen, sind die zehn Beobachtungen i.i.d.:
 - $P(\mathbf{d} | h_\theta) = \prod_j P(d_j | h_\theta), j = 1, \dots, 10$
- Für eine gegebene Hypothese h_θ , der Wert von $P(d_j | h_\theta)$ ist
 - $P(d_j = \text{cherry} | h_\theta) = \theta, (d_j = \text{lime} | h_\theta) = 1 - \theta$
- Gegeben N Beobachtungen, von denen c Kirsche und $l = N - c$ Zitrone sind, gilt
 - $P(\mathbf{d} | h_\theta) = \prod_{j=1}^c \theta \prod_{j=1}^l (1 - \theta) = \theta^c (1 - \theta)^l$
 - Binomial-Verteilung: Wahrscheinlichkeit von # an Erfolgen in einer Reihe von N unabhängigen Versuchen mit binärem Ausgang, jeder Versuch ist erfolgreich mit Wahrscheinlichkeit θ
 - Beispiel: Nachdem man $c = 3$ mal Zitronen-Bonbons beobachtet hat:
 - $P([\text{lime}, \text{lime}, \text{lime}] | h_{50}) = 0.5^3$
 - Wahrscheinlichkeit Zitrone zu beobachten ist 0.5 unter der Hypothese h_{50}

100% Zitrone: A Posteriori Wahrscheinlichkeitsverteilung von H

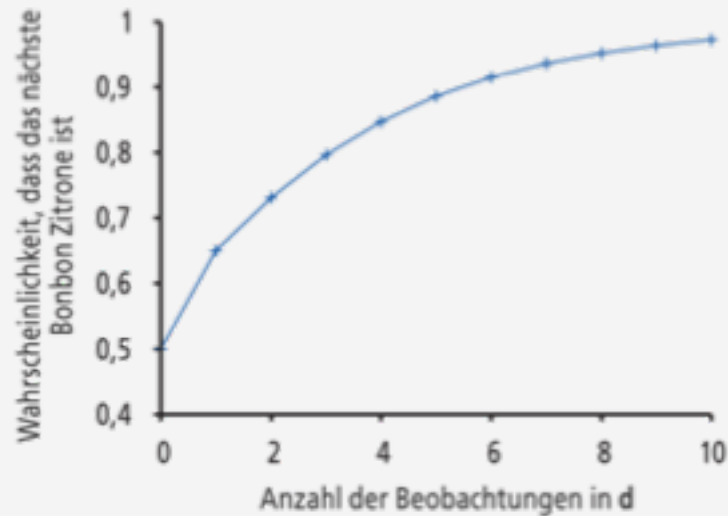
- Abbildung zeigt $P(h_i | \mathbf{d}) \propto P(\mathbf{d} | h_i)P(h_i)$ bei Beobachtungen rein von Zitrone
 - $\mathbf{d} = [\textit{lime}, \textit{lime}, \dots]$
- Initial dominieren die Hypothesen h_i mit größeren Apriori-Wahrscheinlichkeiten
 - h_{50} mit Apriori-Wahrscheinlichkeit 0.4
- Während Daten reinkommen, fängt die am Ende beste Hypothese (h_0) an zu dominieren, da die Wahrscheinlichkeit diese Daten zu beobachten gegeben die anderen Hypothesen immer kleiner wird
 - Nach drei Mal Zitrone in Reihe, nimmt Wahrscheinlichkeit extrem zu, dass wir eine 100% Zitrone Tüte haben



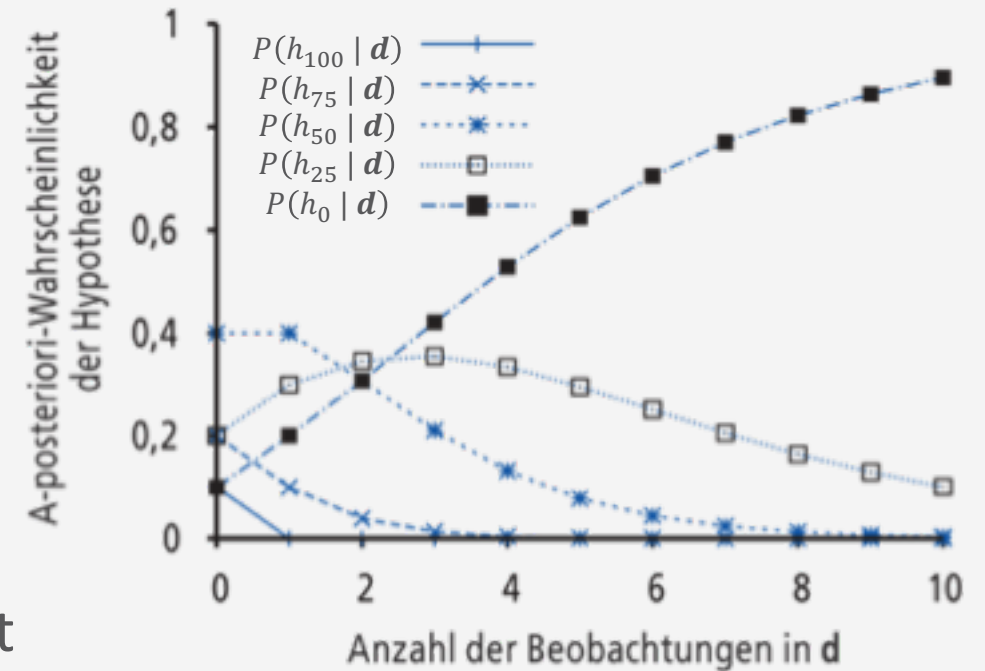
$$P(H) = (0.1, 0.2, 0.4, 0.2, 0.1)$$

Vorhersagewahrscheinlichkeiten

- Abbildung unten zeigt $P(D_{N+1} = \text{lime} \mid \mathbf{d})$
 $\propto \sum_{h_i \in \text{Val}(H)} P(D_{N+1} = \text{lime} \mid h_i) P(\mathbf{d} \mid h_i) P(h_i)$



- Wahrscheinlichkeit, dass nächster Bonbon Zitrone ist steigt mit der Wahrscheinlichkeit, dass die Tüte eine 100%-Zitrone-Tüte ist



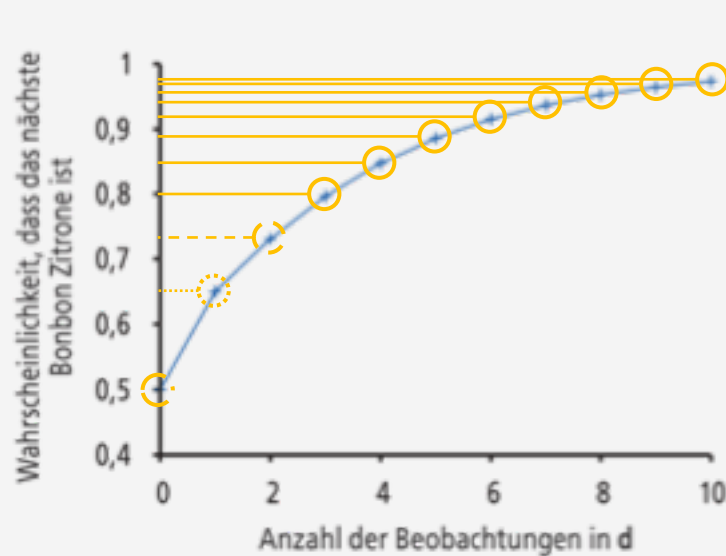
$$P(H) = (0.1, 0.2, 0.4, 0.2, 0.1)$$

MAP Approximation

- Full Bayesian Learning erscheint wie eine sehr sichere Wette, aber funktioniert nicht gut in der Praxis
 - Über den Hypothesen-Raum zu summieren (oder integrieren) ist meist intractabel
 - $2^{2^6} = 18,446,744,073,709,551,616$ Boolesche Funktionen für sechs Features
- Übliche Annäherung: **Maximum-a-posteriori (MAP) Learning**
 - Anstatt Vorhersagen zu treffen, indem man alle möglichen Hypothesen berücksichtigt, wie in
 - $P(X | \mathbf{d}) = \sum_{h_i \in \text{Val}(H)} P(X | h_i) P(h_i | \mathbf{d}) \propto \sum_{h_i \in \text{Val}(H)} P(X | h_i) P(\mathbf{d} | h_i) P(h_i)$
 - Vorhersagen treffen basierend auf h_{MAP} , welche $P(h_i | \mathbf{d})$ maximiert
 - I.e., maximiert $P(\mathbf{d} | h_i) P(h_i) \rightarrow h_{MAP} = \arg \max_{h_i \in \text{Val}(H)} P(\mathbf{d} | h_i) P(h_i)$
 - $P(X | \mathbf{d}) \approx P(X | h_{MAP})$

MAP Approximation

- MAP ist eine Schätzung, wenn $P(X | \mathbf{d}) \approx P(X | h_{MAP})$
 - Im Beispiel ist h_{MAP} die 100%-Zitrone-Tüte nach nur 3 Bonbons, mit der Vorhersage, dass der nächste Bonbon Zitrone ist mit Wahrscheinlichkeit 1
 - Full Bayesian Learner gab da noch die Vorhersage von Zitrone mit einer Wahrscheinlichkeit von 0.8 an



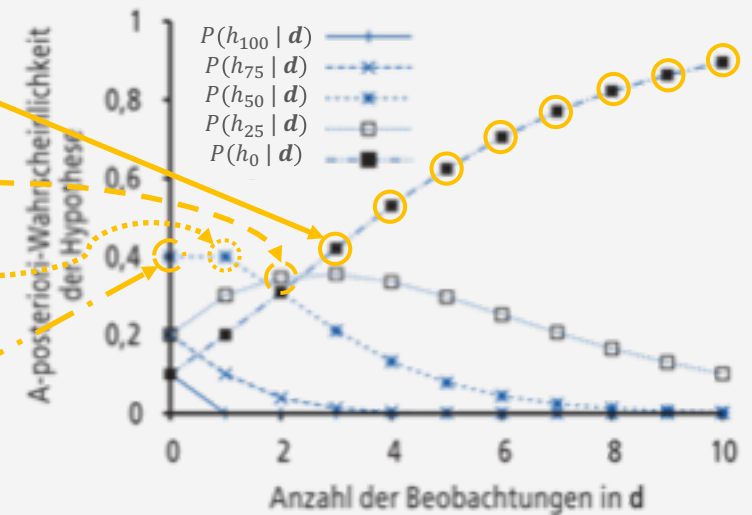
$$\vdots$$

$$h_{MAP}^3 = h_0 \rightarrow P(X | h_0) = 1.0$$

$$h_{MAP}^2 = h_{25} \rightarrow P(X | h_{25}) = 0.75$$

$$h_{MAP}^1 = h_{50} \rightarrow P(X | h_{50}) = 0.5$$

$$h_{MAP}^0 = h_{50} \rightarrow P(X | h_{50}) = 0.5$$

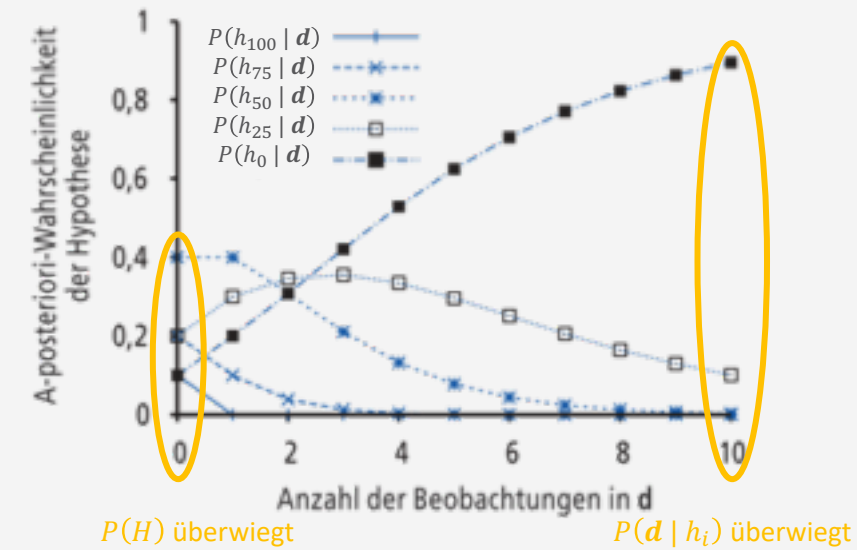


Bias

- Je mehr Daten eintreffen, desto näher beieinander liegen die MAP und Full Bayesian Vorhersagen, da die Nicht-MAP-Hypothesen weniger wahrscheinlich werden
- Meistens einfacher MAP Hypothesen zu finden, da Optimierungsproblem (statt Summierungsproblem)
- $P(H)$ spielt wichtige Rolle in MAP und Full Bayesian Learning
 - Gibt den Bias an: bestimmte Hypothesen werden a priori bevorzugt
 - Trade-off zwischen Modellkomplexität und Fähigkeit, Daten möglichst gut zu reproduzieren
 - Komplexere Modelle können die Daten meist besser erklären \rightarrow höhere $P(\mathbf{d} | h_i)$
 - Gefahr der **Überanpassung** (*overfitting*)
 - Komplexere Modelle weniger wahrscheinlich a priori, da es mehr komplexere Modelle gibt als einfachere \rightarrow geringere $P(h_i)$
 - I.e., üblicher Bias beim Lernen: **Komplexität bestrafen**

Maximum Likelihood (ML) Learning

- Weitere Vereinfachung gegenüber Full Bayesian und MAP Learning
 - A-priori-Verteilung über den Hypothesenraum uniform
 - MAP Learning vereinfacht sich zu ML Learning, da $P(h_i)$ immer gleich
 - MAP: Maximiere $P(\mathbf{d} | h_i)P(h_i) \rightarrow$ ML: Maximiere $P(\mathbf{d} | h_i)$
 - $h_{MAP} = \arg \max_{h_i \in \text{Val}(H)} P(\mathbf{d} | h_i)P(h_i) \rightarrow h_{ML} = \arg \max_{h_i \in \text{Val}(H)} P(\mathbf{d} | h_i)$
- Wann ist ML angemessen?
 - Wenn die Hypothesen tatsächlich gleich wahrscheinlich sind (z.B., gleiche Komplexität \rightarrow keine Hypothese zu bevorzugen)
 - In der Statistik die Standard (nicht-Bayes'sche) statistische Lernmethode, verwendet vor allem bei Misstrauen gegenüber subjektiven Natur der a priori Verteilung über die Hypothesen
 - Bei sehr großen Datenmengen neigt $P(\mathbf{d} | h_i)$ den Einfluss von $P(h_i)$ zu überwinden

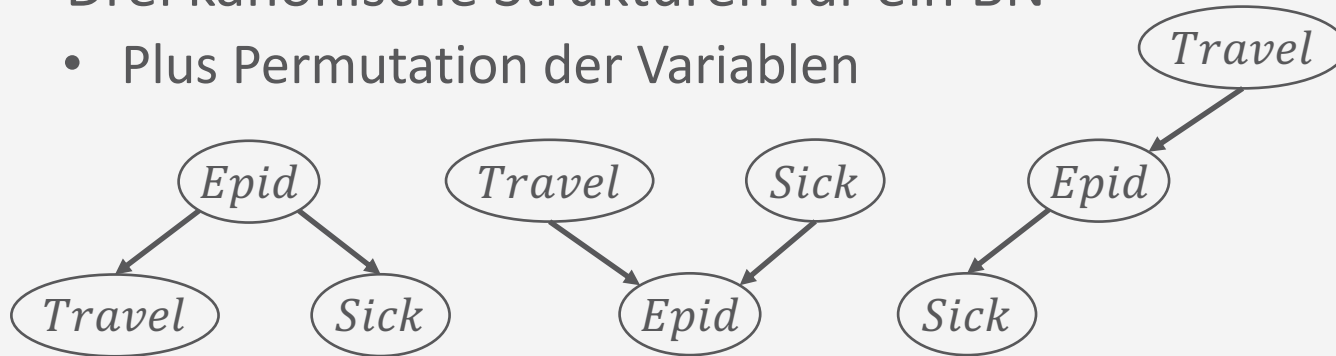


Lernen in PGMs: Dimensionen

- Annahme: Faktorisierung der vollständigen gemeinsamen Verteilung in irgendeiner Form gewünscht
- Hypothesenraum
 - Mögliche Strukturen
 - Welche Variablen?
 - Welche Kanten?
 - Mögliche Parameter
 - Welche Einträge in den CPDs / Faktoren?
 - Beispiel: Wie die folgenden Variablen kombinieren?
 - Travel*
 - Epid*
 - Sick*
- Daten
 - Vollständig vorhanden
 - Für jede existierende Zufallsvariable gibt es Beobachtungen
 - Dann Variablen in möglichen Strukturen bekannt
 - Nicht vollständig vorhanden
 - Zufallsvariable bekannt, aber keine Daten vorhanden
 - Nicht bekannt, ob es nur weitere Zufallsvariablen geben sollte (somit auch keine Daten dafür)
 - Beispiel: Ereignisse für alle drei Variablen oder nur für eine Untermenge von ihnen

Beispiel

- Drei kanonische Strukturen für ein BN
 - Plus Permutation der Variablen



- Struktur links: Verteilungen mit $\theta_0, \dots, \theta_4 \in [0,1]^5$ für $|\text{Val}(R)| = 2$

Epid	P
false	$1 - \theta_0$
true	θ_0

Epid	Travel	P
false	false	$1 - \theta_1$
false	true	θ_1
true	false	$1 - \theta_2$
true	true	θ_2

Epid	Sick	P
false	false	$1 - \theta_3$
false	true	θ_3
true	false	$1 - \theta_4$
true	true	θ_4

Travel	Epid	Sick
false	false	false
false	false	true
true	false	false
false	true	true
true	false	false
true	false	true
false	false	false
false	false	false

Menge von zusammengesetzten Ereignissen als Evidenz für Ausprägungen von $\theta_0, \dots, \theta_4 \in [0,1]^5$ bei vollständigen Daten

Bei unvollständigen Daten z.B. nur Beobachtungen zu *Travel* und *Sick*

In der Regel geht man davon aus, dass man eine Struktur vorgeben kann und nur noch die Parameter (bei vollständigen oder nicht vollständigen Daten) bestimmen muss

Lernen in PGMs: Lernansätze

- Im Unterschied zu vorher keine Datenvorhersage, sondern Daten im Vorhinein vorhanden
- Frage: Welche Modellausprägung ist die „richtige“?
 - Ausprägung der Modellstruktur und Parameter, welche die Daten am besten modelliert
 - Gegebene Daten am wahrscheinlichsten wieder reproduzieren kann (durch *Sampling*)
 - Full Bayesian Learning: alle Ausprägungen berücksichtigen
 - Pro Verteilung kann der Parameter bei Boole'schen Zufallsvariablen in $[0,1]$ liegen → kontinuierlich!
 - Erschwert zudem die Anfragebeantwortung, da nicht eine Anfrage an ein Modell gestellt wird, sondern die gleiche Anfrage an alle möglichen Modelle, deren Antworten gemittelt werden
 - MAP Learning: die MAP Ausprägung wählen
 - ML Learning: ML Ausprägung wählen
 - MAP + ML: ein Modell als Ausgabe → Anfragebeantwortung wie bisher

Zwischenzusammenfassung

- Lernproblem: Parameter einer Verteilung lernen
 - Parameter: Eintrag in CPD / Faktor
- Hypothesenraum
 - Hypothese: für jeden Parameter eine mögliche Ausprägung
 - Menge der Hypothesen: Menge der Kombinationen aller möglichen Ausprägungen
- Bestimmen der passenden Hypothese
 - Hypothese, die Daten am wahrscheinlichsten macht
- Methoden
 - Full Bayesian Learning: alle Hypothesen betrachten
 - MAP Learning: maximiert $P(\mathbf{d} | h_i)P(h_i)$
 - ML Learning: A priori Verteilung über Hypothesenraum uniform; maximiert $P(\mathbf{d} | h_i)$

Überblick: 5. Lernalgorithmen für episodische PGMs

A. Überblick

- Full Bayesian Learning, MAP Learning, Maximum Likelihood (ML) Learning

B. *ML-basiertes Parameter Lernen*

- In BNs: ML (vollständige Daten), Expectation-Maximisation (EM; nicht-vollständige Daten)
- In Faktormodellen: Iterative Proportional Fitting (IPF), EM-IPF
- Anwendung: LDA

C. *ML-basiertes Struktur Lernen*

- Model Selection, Structural EM

D. *Abschlussbemerkungen*

- Alternativen, generative vs. diskriminative Modelle, Korrelation vs. Kausalität

Gerichtete Modelle: BNs

ML-basiertes Parameter lernen

$P(F = \text{cherry})$
θ_0

Flavour

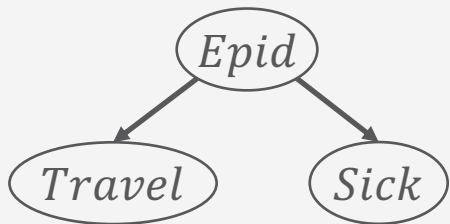
Wrapper

F	$P(W = \text{red} F)$
<i>cherry</i>	θ_1
<i>lime</i>	θ_2

Lernen in BNs: Data Science mit vollständigen Daten

- Wir starten mit dem einfachsten Problem in BN Lernen:
 - Bekannte Struktur
 - Daten beinhalten Beobachtungen für alle Variablen
 - Alle Variablen beobachtbar, keine fehlenden Daten
- Es müssen nur die BN Parameter in den CPDs gelernt werden

Modell



Daten

<i>Travel</i>	<i>Epid</i>	<i>Sick</i>
<i>false</i>	<i>false</i>	<i>false</i>
<i>false</i>	<i>false</i>	<i>true</i>
<i>true</i>	<i>false</i>	<i>false</i>
<i>false</i>	<i>true</i>	<i>true</i>
⋮	⋮	⋮

→ Wahrscheinlichkeitsverteilungen zu lernen

$P(\text{Epid}), P(\text{Travel} \mid \text{Epid}), P(\text{Sick} \mid \text{Epid})$

<i>Epid</i>	<i>P</i>
<i>false</i>	$1 - \theta_0$
<i>true</i>	θ_0

<i>Epid</i>	<i>Travel</i>	<i>P</i>
<i>false</i>	<i>false</i>	$1 - \theta_1$
<i>false</i>	<i>true</i>	θ_1
<i>true</i>	<i>false</i>	$1 - \theta_2$
<i>true</i>	<i>true</i>	θ_2

<i>Epid</i>	<i>Sick</i>	<i>P</i>
<i>false</i>	<i>false</i>	$1 - \theta_3$
<i>false</i>	<i>true</i>	θ_3
<i>true</i>	<i>false</i>	$1 - \theta_4$
<i>true</i>	<i>true</i>	θ_4

Maximum-Likelihood-Parameterschätzung

- Nehme an, die Struktur eines BNs sei bekannt
- Ziel: Schätze BN-Parameter θ
 - Einträge in CPDs $P(R \mid \text{Pa}(R))$
- Parametrisierung mit θ ist gut, falls hierdurch die beobachteten Daten wahrscheinlich generiert werden:

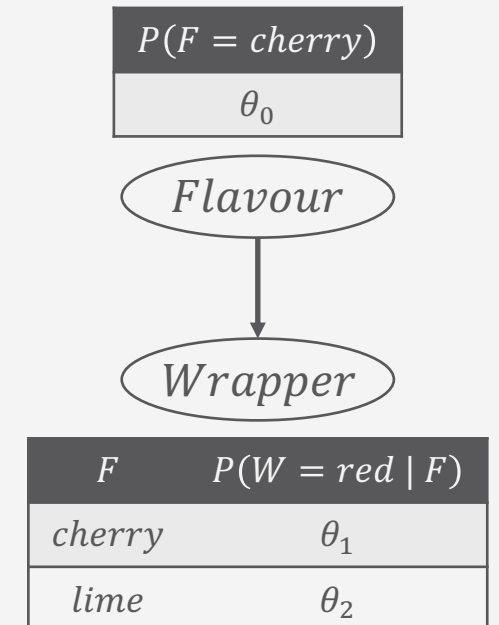
$$P(\mathbf{d} \mid \theta) = \prod_{j=1}^N P(d_j \mid \theta)$$

i.i.d. Datenpunkte

- **Maximum Likelihood Estimation** (MLE) Prinzip: Wähle θ so, dass $P(\mathbf{d} \mid \theta)$ maximiert wird

Beispiel

- Bonbonfabrik:
 - Farbe des Bonbonpapiers mit einer bestimmten Wahrscheinlichkeit je nach Geschmack, wobei die entsprechende Verteilung nicht bekannt sei
 - Wenn *Flavour* = *cherry*, wähle *Wrapper* = *red* mit Wahrscheinlichkeit θ_1
 - Wenn *Flavour* = *lime*, wähle *Wrapper* = *red* mit Wahrscheinlichkeit θ_2
 - Entsprechendes BN dazu enthält drei zu lernenden Parameter
 - $\theta_0, \theta_1, \theta_2$



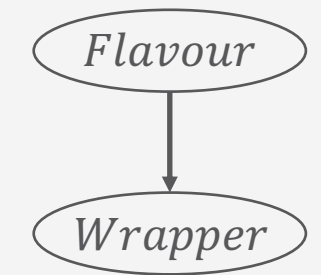
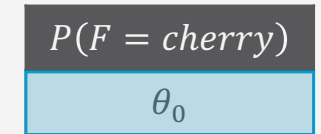
Beispiel

- Wahrscheinlichkeit für z.B. $W = yellow, F = cherry$

$$\begin{aligned}
 &P(W = yellow, F = cherry | h_{\theta_0\theta_1\theta_2}) \\
 &= P(W = yellow | F = cherry, h_{\theta_0\theta_1\theta_2})P(F = cherry | h_{\theta_0\theta_1\theta_2}) \\
 &= (1 - \theta_1)\theta_0
 \end{aligned}$$

- Wir packen N Bonbons aus
 - Jeder Versuch liefert eine Kombination aus Papier und Geschmack wie oben
 - c sind Kirsche, l sind Zitrone ($c + l = N$)
 - c_r sind Kirsche mit rotem Papier, c_y sind Kirsche mit gelbem Papier ($c_r + c_y = c$)
 - l_r sind Zitrone mit rotem Papier, l_y sind Zitrone mit gelbem Papier ($l_r + l_y = l$)
- Wahrscheinlichkeit der Daten:

$$P(\mathbf{d} | h_{\theta_0\theta_1\theta_2}) = \prod_{j=1}^N P(d_j | h_{\theta_0\theta_1\theta_2}) = \theta_0^c (1 - \theta_0)^l \theta_1^{c_r} (1 - \theta_1)^{c_y} \theta_2^{l_r} (1 - \theta_2)^{l_y}$$



F	$P(W = red F)$
cherry	θ_1
lime	θ_2

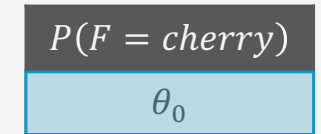
Beispiel

- ML-Hypothese: Werte für $\theta_0\theta_1\theta_2$, so dass $P(\mathbf{d} | h_{\theta_0\theta_1\theta_2}) = \prod_{j=1}^N P(d_j | h_{\theta_0\theta_1\theta_2})$ maximal
- Erhält man auch durch Maximierung der **Log-Wahrscheinlichkeit**:

$$L(\mathbf{d} | h_{\theta_0\theta_1\theta_2}) = \log P(\mathbf{d} | h_{\theta_0\theta_1\theta_2}) = \sum_{j=1}^N \log P(d_j | h_{\theta_0\theta_1\theta_2})$$

- Vorteil: Aus dem Produkt wird eine Summe
 - In der Regel leichter zu maximieren
 - Für das Beispiel heißt das:

$$\begin{aligned} &L(\mathbf{d} | h_{\theta_0\theta_1\theta_2}) \\ &= \log(\theta_0^c (1 - \theta_0)^l \theta_1^{c_r} (1 - \theta_1)^{c_y} \theta_2^{l_r} (1 - \theta_2)^{l_y}) \\ &= \log \theta_0^c + \log(1 - \theta_0)^l + \log \theta_1^{c_r} + \log(1 - \theta_1)^{c_y} + \log \theta_2^{l_r} + \log(1 - \theta_2)^{l_y} \\ &= c \log \theta_0 + l \log(1 - \theta_0) + c_r \log \theta_1 + c_y \log(1 - \theta_1) + l_r \log \theta_2 + l_y \log(1 - \theta_2) \end{aligned}$$



F	$P(W = \text{red} F)$
cherry	θ_1
lime	θ_2

Beispiel

- Bestimmung der Ableitungen bzgl. θ und Maximalwerte bei Ableitung gleich 0 finden
 - Im Beispiel
 - $L(d | h_{\theta_0\theta_1\theta_2}) = c \log \theta_0 + l \log(1 - \theta_0) + c_r \log \theta_1 + c_y \log(1 - \theta_1) + l_r \log \theta_2 + l_y \log(1 - \theta_2)$
 - Ableitung bzgl. θ_0

$$\frac{\partial L}{\partial \theta_0} = \frac{c}{\theta_0} - \frac{l}{1 - \theta_0} = 0 \quad \Rightarrow \quad \theta_0 = \frac{c}{c + l}$$

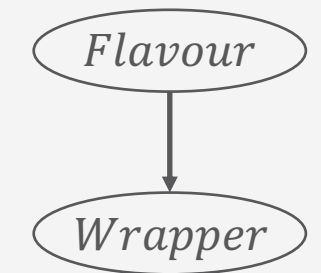
- Ableitung bzgl. θ_1

$$\frac{\partial L}{\partial \theta_1} = \frac{c_r}{\theta_1} - \frac{c_y}{1 - \theta_1} = 0 \quad \Rightarrow \quad \theta_1 = \frac{c_r}{c_r + c_y}$$

- Ableitung bzgl. θ_2

$$\frac{\partial L}{\partial \theta_2} = \frac{l_r}{\theta_2} - \frac{l_y}{1 - \theta_2} = 0 \quad \Rightarrow \quad \theta_2 = \frac{l_r}{l_r + l_y}$$

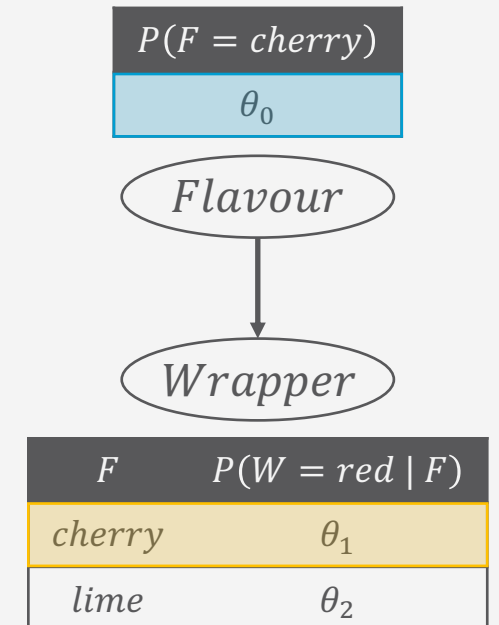
$P(F = \text{cherry})$
θ_0



F	$P(W = \text{red} F)$
cherry	θ_1
lime	θ_2

Maximum-Likelihood-Parameterschätzung

- Schätzung durch Bildung relativer Häufigkeiten
- Prozess auf jedes voll beobachtbare BN anwendbar
- Mit vollständigen Daten und MLE:
 - Parameter lernen zerfällt in separate Lernprobleme für jeden Parameter (CPD) durch Logarithmierung
 - Jeder Parameter wird durch die relative Häufigkeit eines Knotenwertes bei gegebenen Werten der Elternknoten bestimmt



Maximum-Likelihood-Parameterschätzung: Allgemeines Vorgehen

- Viel mathematischer Aufwand, um das Offensichtliche zu entdecken
- Skizziert aber ein generelles Vorgehen für das ML-Parameterlernen:
 - Aufgabe: Maximiere $P(\mathbf{d} | \boldsymbol{\theta})$, i.e., $h_{ML} = \arg \max_{\boldsymbol{\theta}} P(\mathbf{d} | \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \prod_{j=1}^N P(d_j | \boldsymbol{\theta})$
 - Vorgehen
 1. Aufschreiben eines Ausdrucks für die Wahrscheinlichkeit der Daten als Funktion der Parameter ($P(\mathbf{d} | \boldsymbol{\theta})$)
 2. Ableitung der Log-Wahrscheinlichkeit für jeden Parameter
 3. Ermittlung von Parameterwerten, für die die Ableitungen null sind
 - Danach in den Daten die benötigten relativen Häufigkeiten bestimmen (i.e., zählen)
- Nebenbemerkung: Sampling lässt sich auf diese Art interpretieren
 - Lernen der wahren Verteilung aus Samples mittels Bestimmen von relativen Häufigkeiten (ML)

Erweiterung auf MAP

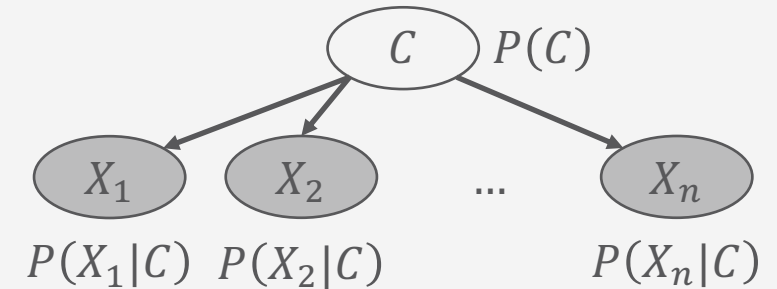
- Für die MAP-Hypothese erhalten wir

$$\begin{aligned}h_{MAP} &= \arg \max_{\boldsymbol{\theta}} P(\mathbf{d} | \boldsymbol{\theta})P(\boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \log P(d_j | \boldsymbol{\theta})P(\boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \log P(d_j | \boldsymbol{\theta}) + \log P(\boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \log \prod_{j=1}^N P(d_j | \boldsymbol{\theta}) + \log P(\boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \sum_{j=1}^N \log P(d_j | \boldsymbol{\theta}) + \log P(\boldsymbol{\theta})\end{aligned}$$

- Gleichungen identisch bis auf Einbindung von $P(\boldsymbol{\theta})$
 - Likelihood $P(d_j | \boldsymbol{\theta})$ nun gewichtet mit der Apriori-Wahrscheinlichkeit $P(\boldsymbol{\theta})$

Anwendung: Naives Bayes-Klassifikator

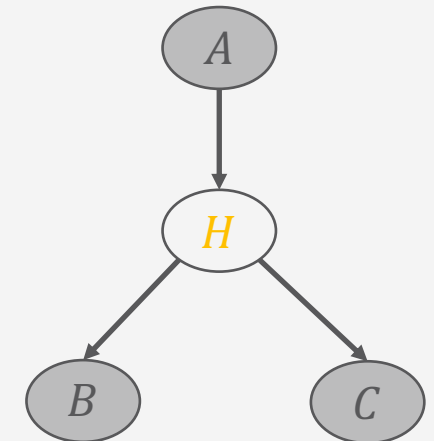
- Klassifikation:
 - Menge von Features X_1, \dots, X_n
 - Zuweisung eines Klassennamens $C = c$ auf Basis der Featurewerte x_1, \dots, x_n
 - Annahme: bedingte Unabhängigkeit gegeben C
- Vorteile: **Relativ wenige Parameter (Einträge in den CPDs) zu lernen**
 - **$2n + 1$ Parameter bei Booleschen Variablen benötigt (allgemein: $r(r - 1)n + r - 1$)**
 - $r = |\text{Val}(C)|$, $r - 1$ Einträge in $P(C)$, je $r(r - 1)$ Einträge in $P(X_i|C)$ (pro c , $(r - 1)$ Einträge lernen)
- Überwachtes Lernen (*supervised learning*; d.h., gelabelte Daten vorhanden):
 - Menge von Featurewerten mit dazugehörigem Klassennamen $\mathbf{D} = \left\{ (x_1^j, \dots, x_n^j, c^j) \right\}_{j=1}^N$
 - Zählen der relativen Häufigkeiten von jedem c in \mathbf{D} für $P(C)$, jedem x_i, c in \mathbf{D} für $P(X_i | C)$



Parameterlernen mit versteckten Variablen

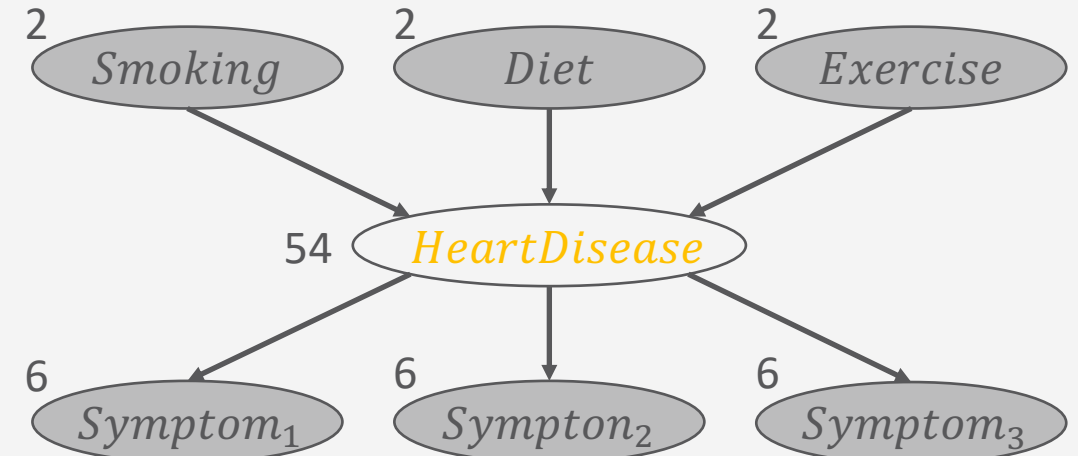
- Bisherige Annahme: Daten für jede Variable des BNs vorhanden
- Nächst schwierigeres Problem
 - Versteckte Variablen im Netz
 - Keine Beobachtungen vorhanden
 - Ansatz mit relativen Häufigkeiten so nicht übertragbar, da keine Zähler für diese Variablen vorhanden
 - Beispiel: Kein Zähler für H

A	B	C
<i>true</i>	<i>false</i>	<i>true</i>
<i>false</i>	<i>true</i>	<i>true</i>
<i>true</i>	<i>true</i>	<i>false</i>
⋮		



Vermeintliche Lösung

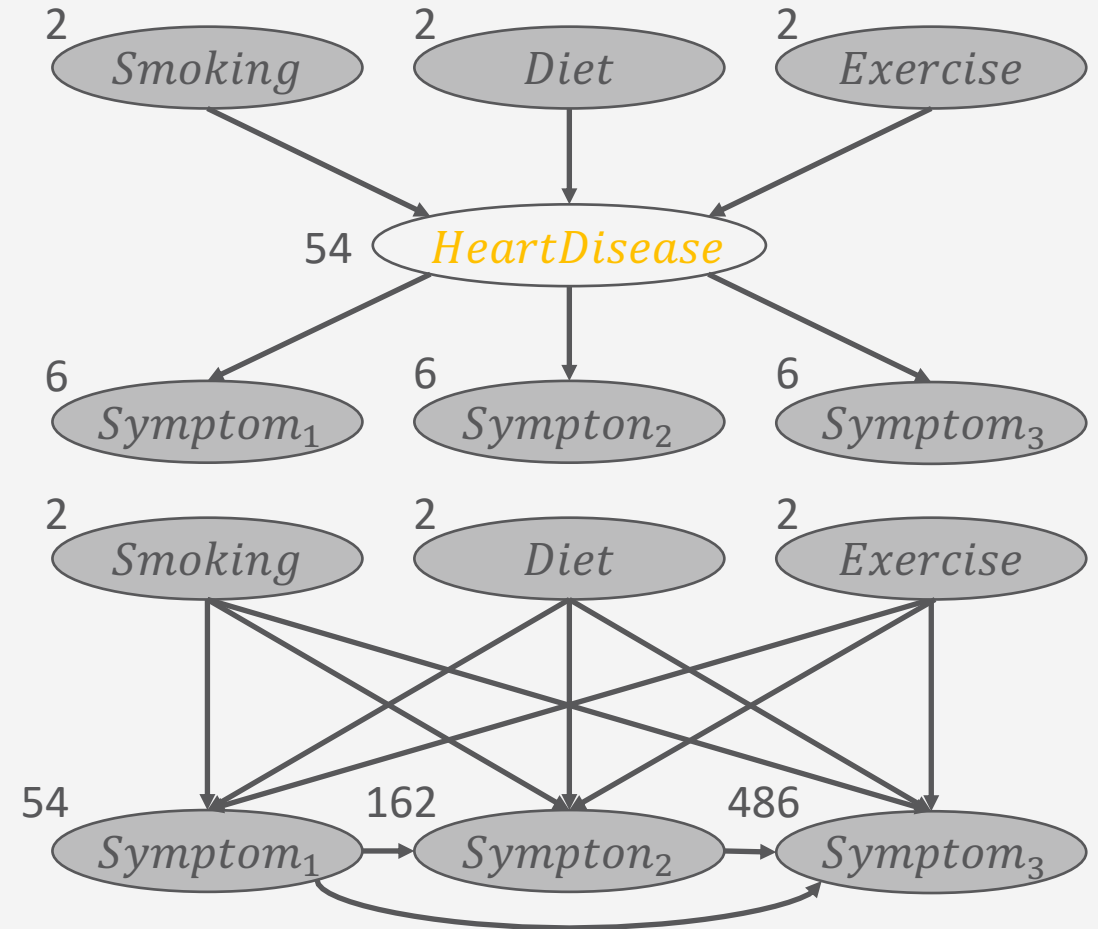
- Vermeide versteckte Variablen
- Könnte klappen bei kleinen Netzwerken
- Beispiel:
 - Jede Variable habe 3 Werte (*low, moderate, high*)
 - Zahlen an den Knoten repräsentieren, wie viele Parameter für die CPDs des jeweiligen Knotens bestimmt werden müssen
 - 78 Wahrscheinlichkeitswerte insgesamt



Was machen wir hier?

HeartDisease wegzulassen ist gar keine gute Lösung!

- Ohne versteckten Knoten müssen die indirekten Abhängigkeiten durch Kanten direkt abgebildet werden
- Symptome nicht länger bedingt unabhängig gegeben die Elternknotenwerte
- Sehr viel mehr Kanten, sehr viel mehr Wahrscheinlichkeiten zu spezifizieren: 708 insgesamt
- Brauchen sehr viel mehr Daten, um diese ganzen Werte vernünftig zu lernen



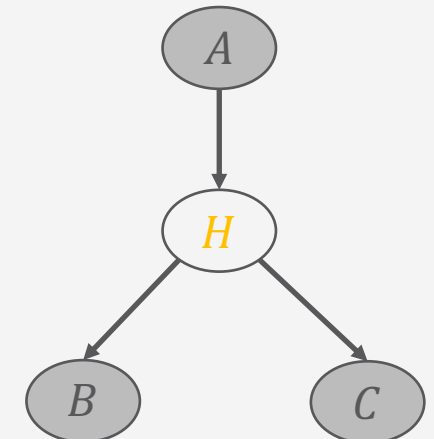
Expectation-Maximisation (EM)

- Wenn wir versteckte Variablen beibehalten und die Parameter aus Daten lernen wollen, haben wir eine Form des **unüberwachten Lernens** (*unsupervised learning*)
 - Daten geben nicht über alle Aspekte von Datenpunkten Auskunft
- **Expectation-Maximisation** (EM)
 - Genereller Algorithmus zum Lernen von Parameters bei unvollständigen Daten
 - Hier für diskrete Datenwerte im BN-Kontext gezeigt

EM: Generelle Idee

- Falls wir Daten für alle Variablen im BN hätten, könnten wir die Parameter mit ML- (oder MAP-) Ansätzen lernen
 - Relative Häufigkeiten bestimmen, wie vorher besprochen
- Falls wir die Parameter hätten, könnten wir die Wahrscheinlichkeiten jedes Ereignisses schätzen
 - $P(H \mid A, B, C)$
 - Könnten also Werte für H schätzen

<i>A</i>	<i>B</i>	<i>C</i>
<i>true</i>	<i>false</i>	<i>true</i>
<i>false</i>	<i>true</i>	<i>true</i>
<i>true</i>	<i>true</i>	<i>false</i>
⋮		



EM: Generelle Idee

- Algorithmus startet mit „erfundenen“ (zufällig generierten) Informationen, um das Lernproblem zu lösen
 - Erfundene Parameterwerte (virtuell notiert in den CPDs)
- Dann werden die Initialwerte in zwei Schritten verfeinert
 - **Expectation (E):**
Aktualisiere die Daten (virtuell) mit aus dem aktuellen Modell hergeleiteten Erwartungen
 - **Maximisation (M):**
Gegeben die aktualisierten Daten, aktualisiere die Parameter mittels ML Ansatz
 - Der gleiche Schritt wie im voll beobachtbaren Fall

EM: Naive Bayes-Modelle als Anwendung

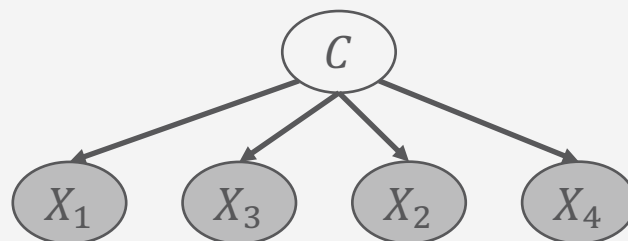
- Betrachte die folgenden Daten
 - N Beispiele mit Booleschen Attributen X_1, \dots, X_4
 - ... die wir kategorisieren wollen mit möglichen Werten einer Klasse C , $\text{Val}(C) = \{1,2,3\}$
- Naive Bayes Klassifikator mit versteckter Variable C

Daten

X_1	X_2	X_3	X_4
<i>true</i>	<i>false</i>	<i>true</i>	<i>true</i>
<i>false</i>	<i>true</i>	<i>true</i>	<i>false</i>
<i>true</i>	<i>false</i>	<i>true</i>	<i>true</i>
⋮			



X_1	X_2	X_3	X_4	Zähler
<i>false</i>	<i>false</i>	<i>false</i>	<i>false</i>	19
<i>false</i>	<i>false</i>	<i>false</i>	<i>true</i>	7
<i>false</i>	<i>false</i>	<i>true</i>	<i>false</i>	13
⋮				

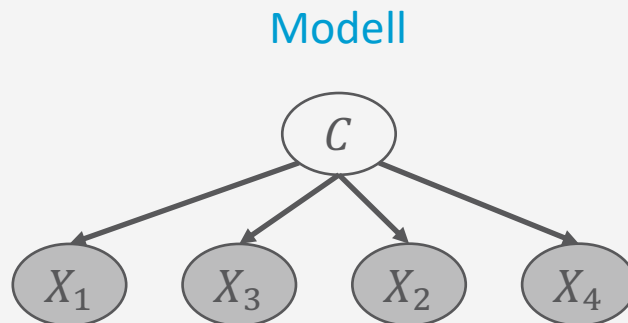
Modell


→ Wahrscheinlichkeitsverteilungen zu lernen

- $P(C)$?
- $P(X_1 | C)$?
- $P(X_2 | C)$?
- $P(X_3 | C)$?
- $P(X_4 | C)$?

EM: Initialisierung

- Algorithmus startet mit „erfundenen“ (zufällig generierten) Informationen, um das Lernproblem zu lösen
 - Erfundene Parameterwerte (virtuell notiert in den CPDs)

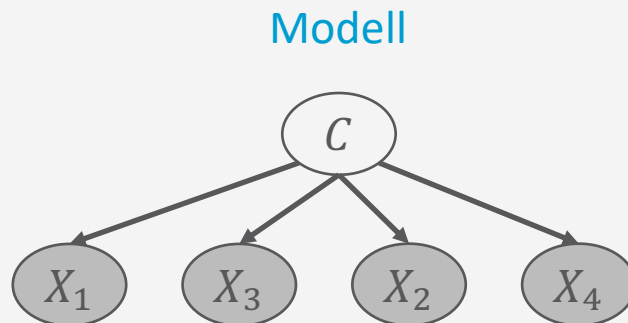


→ Wahrscheinlichkeitsverteilungen zu lernen

$P(C)$?	
$P(X_1 C)$?	Definiere zufällige Parameter
$P(X_2 C)$?	
$P(X_3 C)$?	
$P(X_4 C)$?	

EM: Expectation-Schritt (Bestimmung der erwarteter Zahlen)

- Was benötigen wir, um die Netzwerkparameter mit dem ML-Ansatz zu lernen?
 - $P(C = k) = \frac{\text{Anzahl}(\text{Datenpunkte mit } C=k)}{\text{Anzahl}(\text{alle Datenpunkte})}$ für jedes $k = 1,2,3$
 - $P(X_i = x \mid C = k) = \frac{\text{Anzahl}(\text{Datenpunkte mit } X_i=x \text{ und } C=k)}{\text{Anzahl}(\text{Datenpunkte mit } C=k)}$ für jedes $x \in \text{Val}(X_i), i = 1, \dots, 4, k = 1,2,3$
- Bisher haben wir nur: $N = \text{Anzahl}(\text{alle Datenpunkte})$
 - Wir approximieren alle weiteren Zähler mit erwarteten Zählerwerten, die aus dem Modell mit den „erfundenen“ Parametern abgeleitet werden



Daten

X_1	X_2	X_3	X_4	Zähler
false	false	false	false	19
false	false	false	true	7
false	false	true	false	13
⋮				

→ Wahrscheinlichkeitsverteilungen zu lernen

$$\begin{aligned}
 P(C) &? \\
 P(X_1 \mid C) &? \\
 P(X_2 \mid C) &? \\
 P(X_3 \mid C) &? \\
 P(X_4 \mid C) &?
 \end{aligned}$$

EM: Expectation-Schritt (Bestimmung der erwarteter Zahlen)

- Der erwartete Zähler $\hat{N}(C = k)$ ist die Summe, über alle N Beispieldaten, der Wahrscheinlichkeiten, dass jedes Beispiel in Kategorie k fällt

$$\hat{N}(C = k) = \sum_{j=1}^N P(C = k \mid \underbrace{x_1^j, x_2^j, x_3^j, x_4^j}_{\text{Featurewerte des } j\text{'ten Beispiels}})$$

- Schätzung dann

$$P(C = k) = \frac{\text{erwartete Anzahl(Datenpunkte mit } C = k)}{\text{Anzahl(alle Datenpunkte)}} = \frac{\hat{N}(C = k)}{N}$$

- Wie erhalten wir die Werte?

- Im Naive Bayes Modell einfach berechnen:

$$P(C = k \mid x_1^j, x_2^j, x_3^j, x_4^j) = \frac{P(C = k, x_1^j, x_2^j, x_3^j, x_4^j)}{P(x_1^j, x_2^j, x_3^j, x_4^j)} = \frac{\overbrace{P(x_1^j \mid C = k) \dots P(x_4^j \mid C = k) P(C = k)}^{\text{Erfundene Parameter für das Netzwerk nutzen}}}{P(x_1^j, x_2^j, x_3^j, x_4^j)}$$

EM: Expectation-Schritt (Bestimmung der erwarteter Zahlen)

- Durch einen ähnlichen Schritt erhalten wir die erwarteten Zähler für Beispiele mit Featurewert $X_i = x$ und Kategorie k

Auch alles über erfundene Parametern errechenbar

$$\hat{N}(X_i = x, C = k) = \sum_{j=1}^N P(C = k \mid x_i^j = x, x_{i'}^j, x_{i''}^j, x_{i'''}^j), i, i', i'', i''' \in \{1, 2, 3, 4\}, i \neq i' \neq i'' \neq i'''$$

- Zum Beispiel:

$$\hat{N}(X_1 = \text{true}, C = 1) = \sum_{e_j = (x_1, x_2, x_3, x_4), x_1 = \text{true}} P(C = 1 \mid x_1^j = \text{true}, x_2^j, x_3^j, x_4^j)$$

- Schätzung von $P(X_i \mid C)$

$$P(X_i = x \mid C = k) = \frac{\text{erwarteteAnzahl(Datenpunkte mit } X_i = x, C = k)}{\text{erwarteteAnzahl(Datenpunkte mit } C = k)} = \frac{\hat{N}(X_i = x, C = k)}{\hat{N}(C = k)}$$

- für jedes $x \in \text{Val}(X_i), i = 1, \dots, 4, k = 1, 2, 3$

EM: Generelle Idee

- Algorithmus startet mit „erfundenen“ (zufällig generierten) Informationen, um das Lernproblem zu lösen
 - Erfundene Parameterwerte (virtuell notiert in den CPDs)
- Dann werden die Initialwerte in zwei Schritten verfeinert
 - **Expectation (E):**
Aktualisiere die Daten (virtuell) mit aus dem aktuellen Modell hergeleiteten Erwartungen
 - **Maximisation (M):**
Gegeben die aktualisierten Daten, aktualisiere die Parameter mittels ML Ansatz
 - Der gleiche Schritt wie im voll beobachtbaren Fall

Maximisation-Schritt: (Verfeinerung der Parameter)

- Nun verfeinern wir die Parameter mittels ML-Lernen auf den erwarteten Zählern

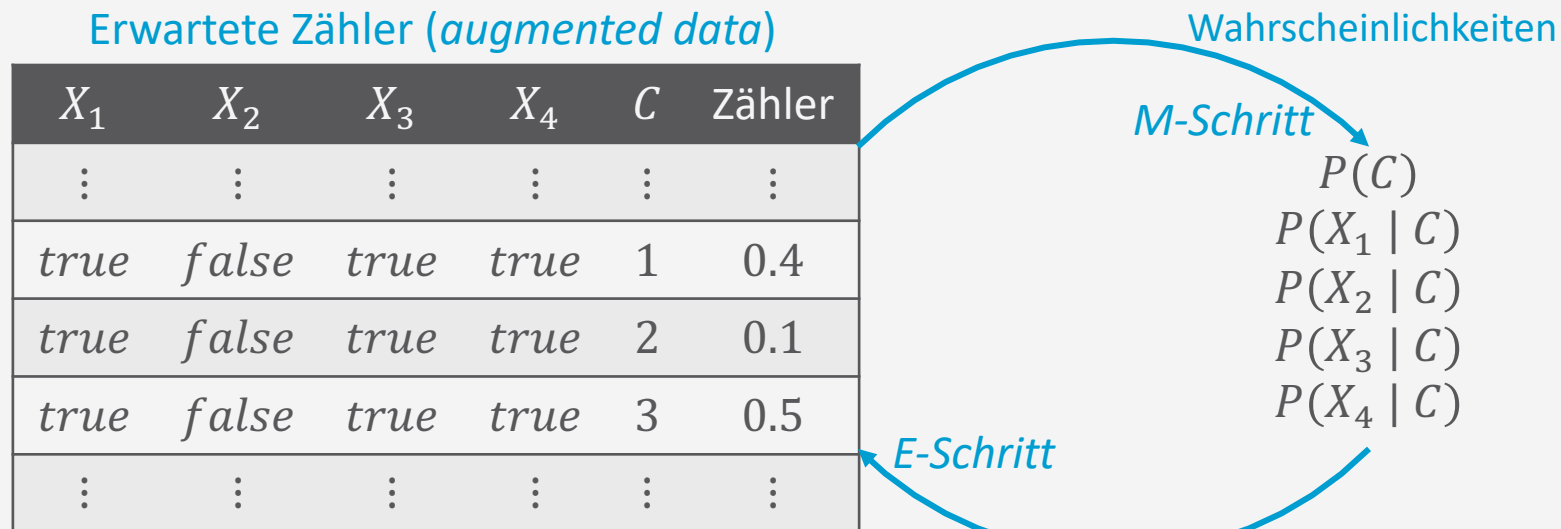
$$P(C = k) = \frac{\hat{N}(C = k)}{N}$$

$$P(X_i = x \mid C = k) = \frac{\hat{N}(X_i = x, C = k)}{\hat{N}(C = k)}$$

- für jedes $x \in \text{Val}(X_i), i = 1, \dots, 4, k = 1, 2, 3$

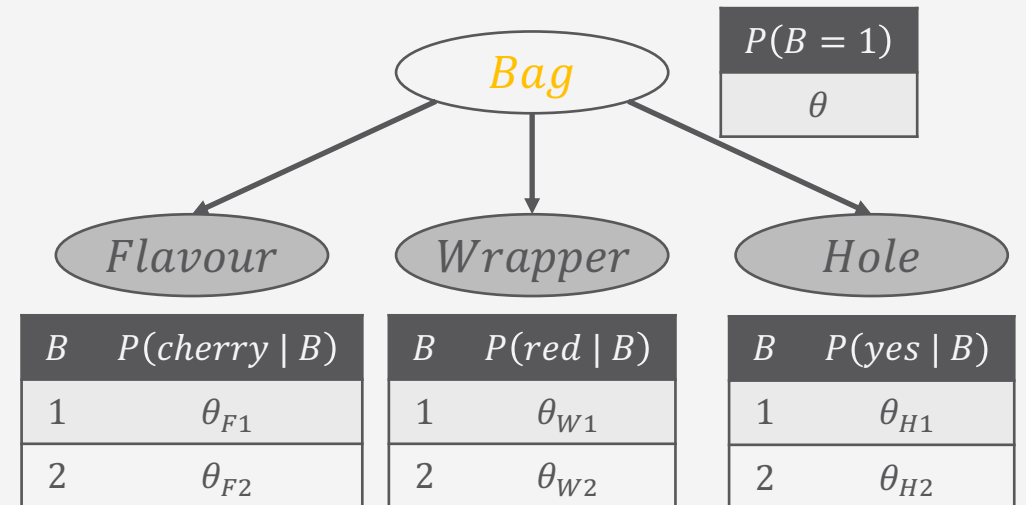
EM-Zyklus

- Nun kann der E-Schritt wiederholt werden
 - Neue erwartete Zähler berechnen gegeben der gerade aktualisierten Parameter
- Dann kann der M-Schritt wiederholt werden
 - Parameter aktualisieren gegebene der gerade neu berechneten erwarteten Zähler



Beispiel: Bonbonfabrik wieder einmal

- Zwei Bonbontüten (1 und 2) wurden vermischt
- Bonbons werden durch Eigenschaften beschrieben
 - *Flavour*
 - *Wrapper*
 - *Hole* (ob ein Loch in der Mitte ist)
 - Merkmale von Bonbon hängen mit bestimmten Wahrscheinlichkeiten von der Tüte ab, aus der sie kommen
- Wir wollen für jedes Bonbon vorhersagen, aus welcher Tüte es kam, je nach vorgefundenen Eigenschaft: Naive Bayes Klassifikator

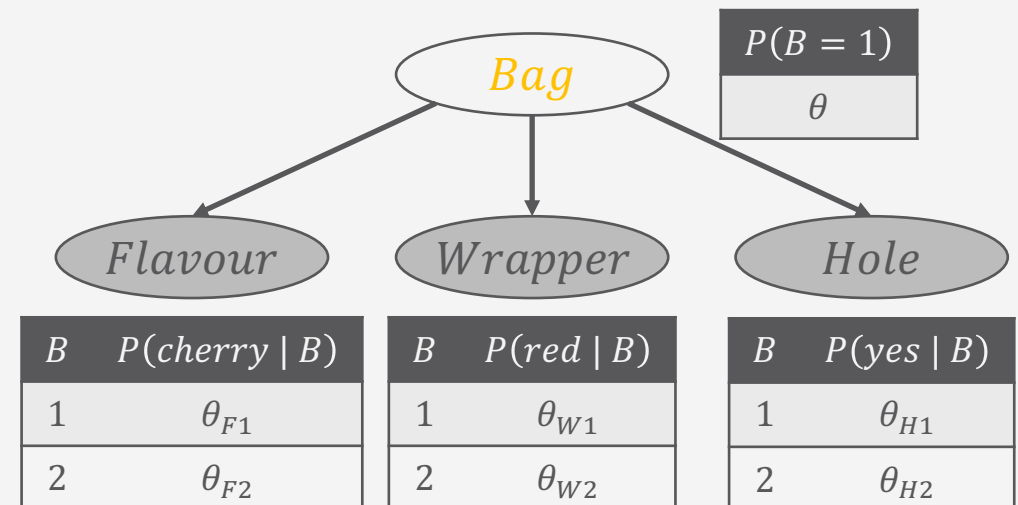


Daten

- Nehmen wir an, die wahren Parameter seien
 - $\theta = 0.5$
 - $\theta_{F1} = \theta_{W1} = \theta_{H1} = 0.8$
 - $\theta_{F2} = \theta_{W2} = \theta_{H2} = 0.3$
- Annahme: Die folgenden Zähler werden „gesampelt“ aus $P(B, F, W, H)$, $N = 1000$

Datensatz	<i>W = red</i>		<i>W = yellow</i>	
	<i>H = yes</i>	<i>H = no</i>	<i>H = yes</i>	<i>H = no</i>
<i>F = cherry</i>	273	93	104	90
<i>F = lime</i>	79	100	94	167

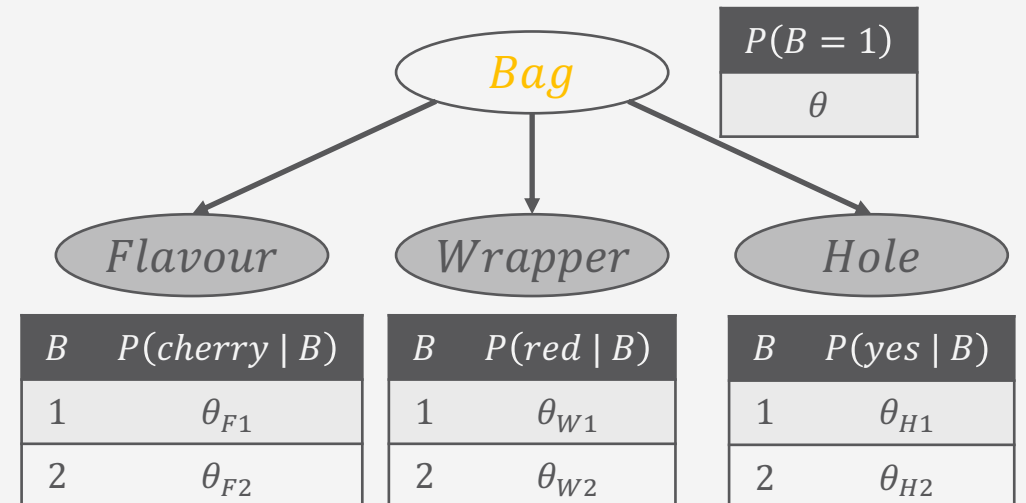
- Wir wollen nun die wahren Parameter aus diesen Daten mittels EM rekonstruieren



EM: Initialisierung

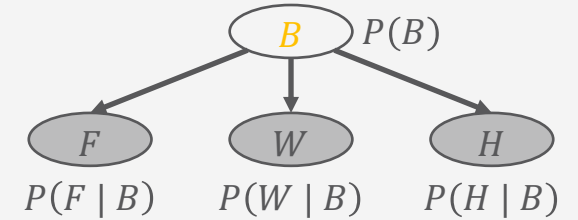
- Weise Parametern zufällige Initialwerte zu
 - Zur Vereinfachung der Darstellung verwenden wir:
 - $\theta^{(0)} = 0.6$
 - $\theta_{F1}^{(0)} = \theta_{W1}^{(0)} = \theta_{H1}^{(0)} = 0.6$
 - $\theta_{F2}^{(0)} = \theta_{W2}^{(0)} = \theta_{H2}^{(0)} = 0.4$
- Nun durchlaufen wir einen EM-Zyklus zur Berechnung von $\theta^{(1)}$
 - Beispielhaft werden wir hier uns zuerst E- und M-Schritt für $\theta^{(1)}$ anschauen, dann für $\theta_{F1}^{(1)}$
 - Konzeptionell würde man aber eher erst E für alle und dann M für alle machen

Datensatz	<i>W = red</i>		<i>W = yellow</i>	
	<i>H = yes</i>	<i>H = no</i>	<i>H = yes</i>	<i>H = no</i>
<i>F = cherry</i>	273	93	104	90
<i>F = lime</i>	79	100	94	167



E-Schritt

- Zuerst bauen wir die erwarteten Zähler für Bonbons von Tüte 1:
 - Addiere die Wahrscheinlichkeiten, dass jeder der N Datenpunkte aus Tüte 1 kommt
 - Seien f_j, w_j, h_j die Werte der entsprechenden Variablen für den j 'ten Datenpunkt



$$\begin{aligned}
 \hat{N}(Bag = 1) &= \sum_{j=1}^N P(Bag = 1 \mid f_j, w_j, h_j) \\
 &= \sum_{j=1}^N \frac{P(f_j, w_j, h_j \mid Bag = 1)P(Bag = 1)}{P(f_j, w_j, h_j)} \\
 &= \sum_{j=1}^N \frac{P(f_j \mid Bag = 1)P(w_j \mid Bag = 1)P(h_j \mid Bag = 1)P(Bag = 1)}{\sum_{k=1}^2 P(f_j \mid Bag = k)P(w_j \mid Bag = k)P(h_j \mid Bag = k)P(Bag = k)}
 \end{aligned}$$

E-Schritt

$$\theta^{(0)} = 0.6$$

$$\theta_{F1}^{(0)} = \theta_{W1}^{(0)} = \theta_{H1}^{(0)} = 0.6$$

$$\theta_{F2}^{(0)} = \theta_{W2}^{(0)} = \theta_{H2}^{(0)} = 0.4$$

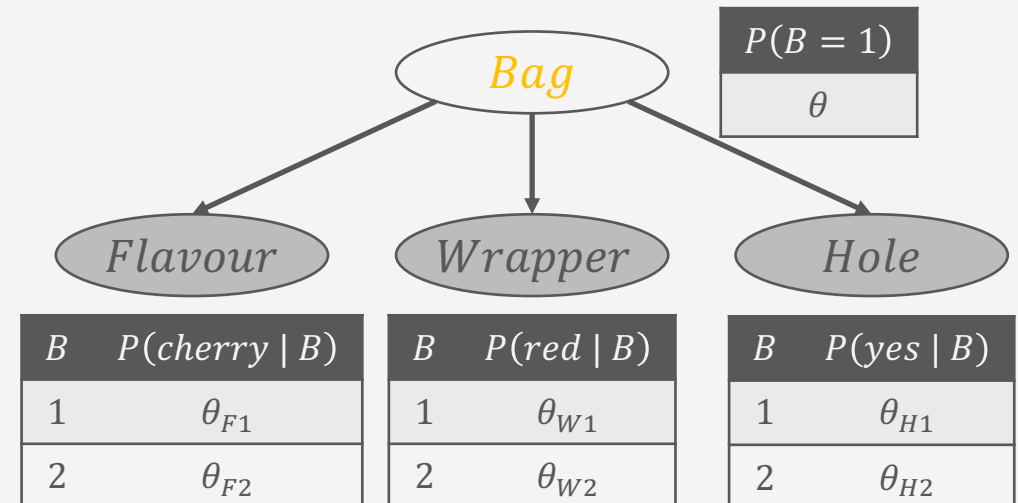
$$\hat{N}(Bag = 1) = \sum_{j=1}^N \frac{P(f_j | Bag = 1)P(w_j | Bag = 1)P(h_j | Bag = 1)P(Bag = 1)}{\sum_{k=1}^2 P(f_j | Bag = k)P(w_j | Bag = k)P(h_j | Bag = k)P(Bag = k)}$$

- Die Summation kann in die 8 Bonbongruppen aus der Tabelle aufgebrochen werden
- Zum Beispiel ergibt die Summe über **273** Kirsch-Bonbons mit rotem Papier und einem Loch (erster Eintrag in der Tabelle)

$$273 \frac{\theta_{F1}^{(0)} \theta_{W1}^{(0)} \theta_{H1}^{(0)} \theta^{(0)}}{\theta_{F1}^{(0)} \theta_{W1}^{(0)} \theta_{H1}^{(0)} \theta^{(0)} + \theta_{F1}^{(0)} \theta_{W1}^{(0)} \theta_{H1}^{(0)} (1 - \theta^{(0)})}$$

$$= 273 \frac{0.6^4}{0.6^4 + 0.4^4} = 273 \frac{0.1296}{0.1552} = 227.97$$

Datensatz	W = red		W = yellow	
	H = yes	H = no	H = yes	H = no
F = cherry	273	93	104	90
F = lime	79	100	94	167



M-Schritt

- Mit der Berechnung der anderen 7 Bonbongruppen erhalten wir

$$\hat{N}(Bag = 1) = 612.4$$

- Nun führen wir den M-Schritt aus, um θ zu verfeinern
- Hierzu nehmen wir den erwarteten Zähler der Datenpunkte aus Tüte 1 und teilen durch N

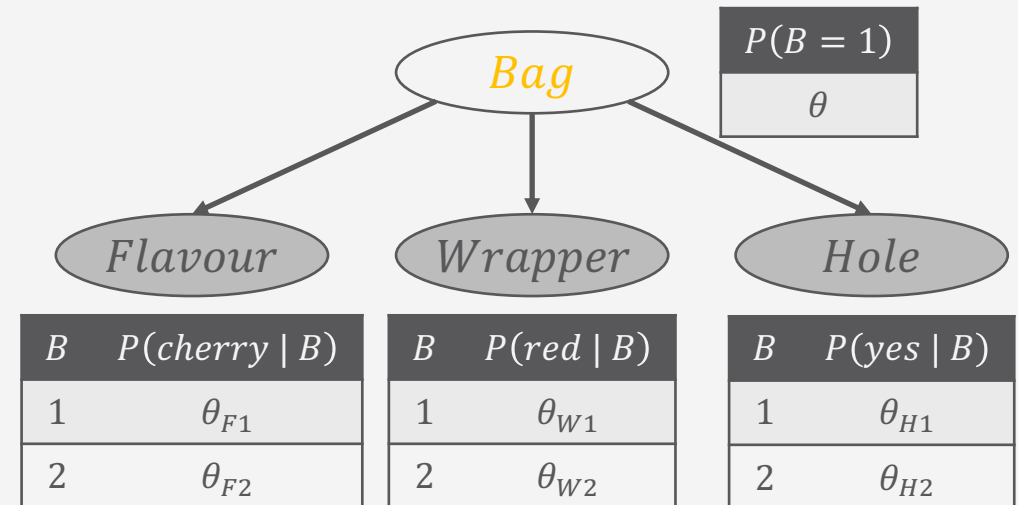
$$\theta^{(1)} = \frac{\hat{N}(Bag = 1)}{N} = \frac{612.4}{1000} = 0.6124$$

$$\theta^{(0)} = 0.6$$

$$\theta_{F1}^{(0)} = \theta_{W1}^{(0)} = \theta_{H1}^{(0)} = 0.6$$

$$\theta_{F2}^{(0)} = \theta_{W2}^{(0)} = \theta_{H2}^{(0)} = 0.4$$

Datensatz	$W = red$		$W = yellow$	
	$H = yes$	$H = no$	$H = yes$	$H = no$
$F = cherry$	273	93	104	90
$F = lime$	79	100	94	167



Noch ein Parameter: θ_{F1}

- E-Schritt: Bestimme erwarteten Zählerwert von Kirsch-Bonbons aus Tüte 1 basierend auf $\theta^{(0)}$

$$\hat{N}(B = 1, F = \text{cherry}) = \sum_{j:F=\text{cherry}} P(B = 1 | \text{cherry}_j, w_j, h_j)$$

- Bestimmbar wie vorher besprochen

$$\hat{N}(Bag = 1, F = \text{cherry})$$

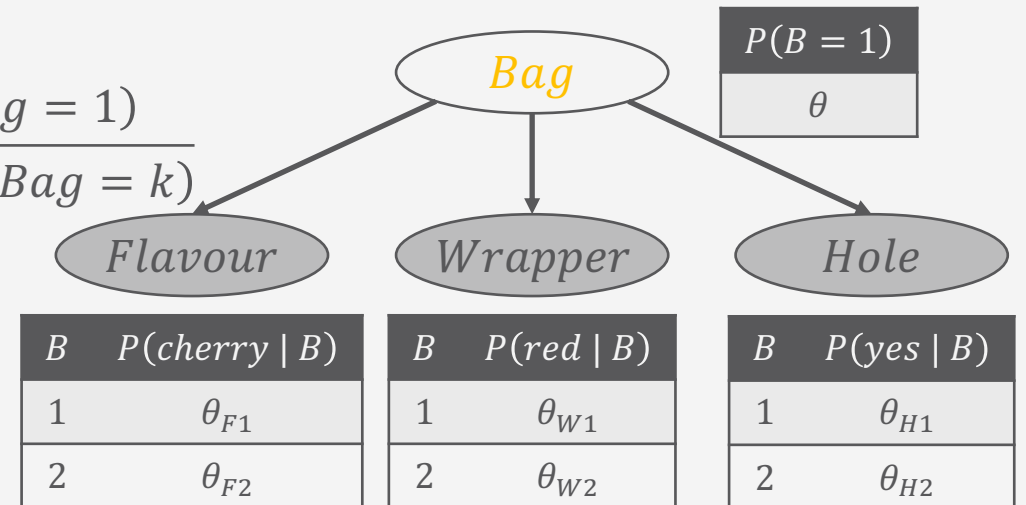
$$= \frac{\sum_{j=1}^N P(\text{cherry}_j | Bag = 1)P(w_j | Bag = 1)P(h_j | Bag = 1)P(Bag = 1)}{\sum_{k=1}^2 P(\text{cherry}_j | Bag = k)P(w_j | Bag = k)P(h_j | Bag = k)P(Bag = k)}$$

- M-Schritt: Verfeinere θ_{F1} durch Bestimmung der relativen Häufigkeiten

$$\theta_{F1}^{(1)} = \frac{\hat{N}(B = 1, F = \text{cherry})}{\hat{N}(B = 1)}$$

$$\begin{aligned} \theta^{(0)} &= 0.6 \\ \theta_{F1}^{(0)} &= \theta_{W1}^{(0)} = \theta_{H1}^{(0)} = 0.6 \\ \theta_{F2}^{(0)} &= \theta_{W2}^{(0)} = \theta_{H2}^{(0)} = 0.4 \end{aligned}$$

Datensatz	W = red		W = yellow	
	H = yes	H = no	H = yes	H = no
F = cherry	273	93	104	90
F = lime	79	100	94	167



Lernergebnis

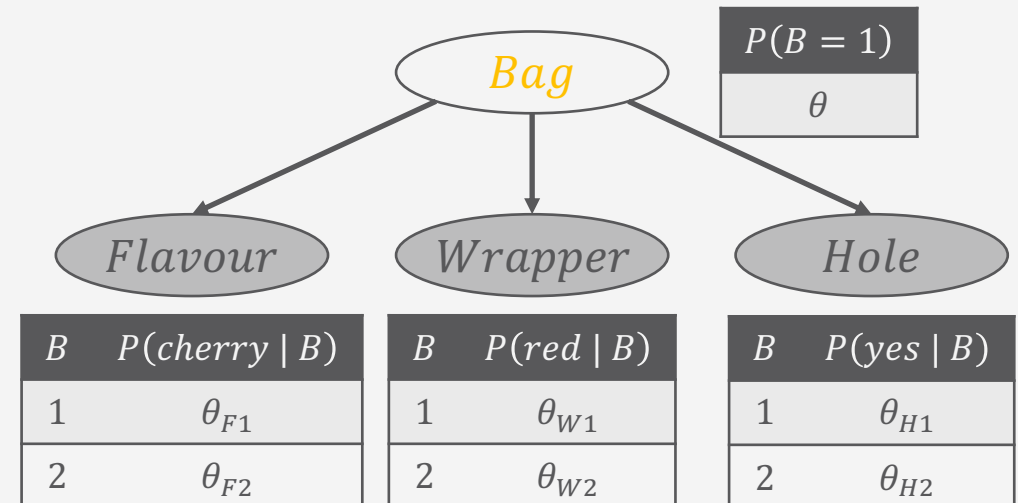
- Nach einem vollen Zyklus über alle Parameter
 - $\theta^{(1)} = 0.6124$
 - $\theta_{F1}^{(1)} = 0.6684, \theta_{W1}^{(1)} = 0.6483, \theta_{H1}^{(1)} = 0.658$
 - $\theta_{F2}^{(1)} = 0.3887, \theta_{W2}^{(1)} = 0.3817, \theta_{H2}^{(1)} = 0.3827$
- Für jeden Satz von Parametern Log-Wahrscheinlichkeiten bestimmbar, wie vorher

$$\theta^{(0)} = 0.6$$

$$\theta_{F1}^{(0)} = \theta_{W1}^{(0)} = \theta_{H1}^{(0)} = 0.6$$

$$\theta_{F2}^{(0)} = \theta_{W2}^{(0)} = \theta_{H2}^{(0)} = 0.4$$

Datensatz	$W = red$		$W = yellow$	
	$H = yes$	$H = no$	$H = yes$	$H = no$
$F = cherry$	273	93	104	90
$F = lime$	79	100	94	167



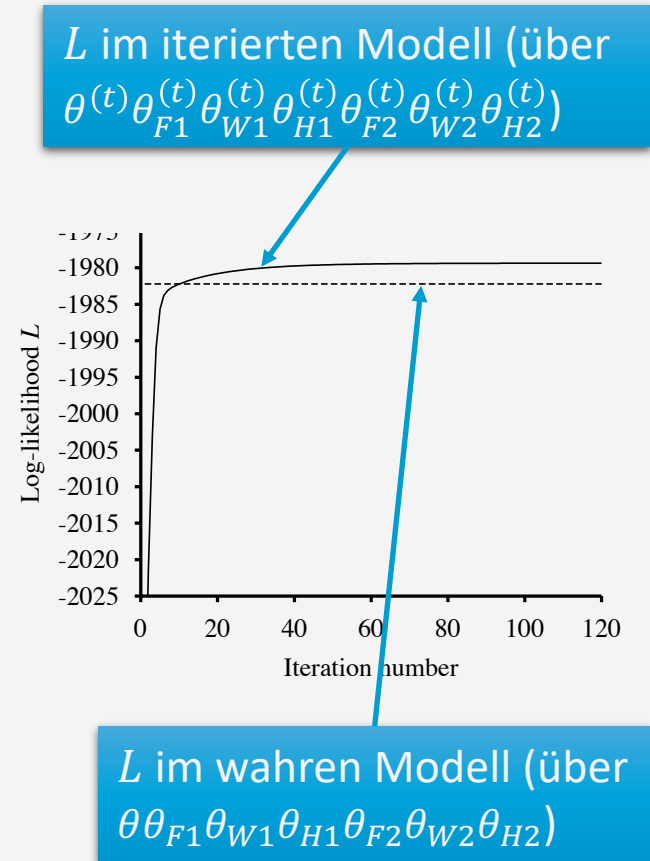
Lernergebnis

- Log-Wahrscheinlichkeiten der Parameter:

$$P(\mathbf{d} | h_{\theta^{(t)} \theta_{F1}^{(t)} \theta_{W1}^{(t)} \theta_{H1}^{(t)} \theta_{F2}^{(t)} \theta_{W2}^{(t)} \theta_{H2}^{(t)}}) = \prod_{j=1}^N P(d_j | h_{\theta^{(t)} \theta_{F1}^{(t)} \theta_{W1}^{(t)} \theta_{H1}^{(t)} \theta_{F2}^{(t)} \theta_{W2}^{(t)} \theta_{H2}^{(t)}})$$

$$\Rightarrow \log P(\mathbf{d}, h_{\theta}) = \sum_{j=1}^N \log P(d_j | h_{\theta^{(t)} \theta_{F1}^{(t)} \theta_{W1}^{(t)} \theta_{H1}^{(t)} \theta_{F2}^{(t)} \theta_{W2}^{(t)} \theta_{H2}^{(t)}})$$

- Man kann zeigen, dass dieser Wert mit jeder EM-Iteration steigt (**Konvergenz**)
 - Beispiel Bonbonfabrik: Log-Wahrscheinlichkeit $L = \log P(\mathbf{d}, h_{\theta})$ steigt in jeder Iteration an [Abbildung rechts]
- EM mit ML bleibt aber bei **lokalen Maxima** hängen
 - Mögliche Abhilfen: mit unterschiedlichen Startwerten anfangen, MAP Lernen nehmen



EM: Diskussion

- Für komplexere BNs ist der Algorithmus der gleich
 - Im Allgemeinen müssen wir die CPD-Einträge für jede Variable R_i gegeben die Elternknotenwerte $\text{Pa}(R_i)$ berechnen

$$\theta_{ijk} = P(R_i = r_{ij} \mid \text{Pa}(R_i) = \text{pa}_k)$$

$$\theta_{ijk} = \frac{\hat{N}(R_i = r_{ij}, \text{Pa}(R_i) = \text{pa}_k)}{\hat{N}(\text{Pa}(R_i) = \text{pa}_k)}$$

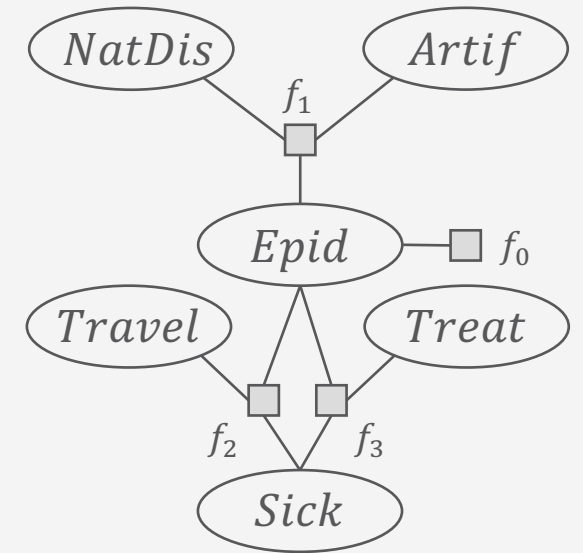
- Die erwarteten Zählerwerte werden durch Summierung über die Beispiele berechnet, nachdem alle notwendigen $P(R_i = r_{ij}, \text{Pa}(R_i) = \text{pa}_{ik})$ mittels BN-Algorithmen bestimmt sind
- Das Verfahren kann in eine **kombinatorische Explosion** laufen
 - Ggf. Approximationstechniken anwenden

Zwischenzusammenfassung

- ML-basierte Parameterschätzung in BNs bei vollständigen Daten
 - Maximierung der Log-Wahrscheinlichkeit der Daten
 - Zerfällt in Abschätzung der einzelnen Parameter
 - Vorkommen der Belegungen in den Daten zählen
- ML-basierte Parameterschätzung in BNs bei nicht vollständigen Daten: EM
 - Maximierung der Log-Wahrscheinlichkeit der Daten
 - Iterativer Algorithmus
 - Erwartete Zähler berechnen
 - Parameter optimieren
 - Läuft auf Fixpunkt (lokales Optimum)

Ungerichtete Modelle

ML-basiertes Parameter lernen



Ungerichtete Modelle

- BNs haben den Vorteil einer Normalisierungskonstante von $Z = 1$ zu haben
 - Parameter lernen reduziert sich auf Schätzung der Parameter der lokalen CPTs
 - Ungerichtete Modelle haben diesen Vorteil im Allgemeinen nicht:

$$P_F = \frac{1}{Z} \prod_{i=1}^n \phi_i(R_1, \dots, R_k)$$
$$Z = \sum_{r_1 \in \text{Val}(R_1)} \dots \sum_{r_n \in \text{Val}(R_n)} \prod_{i=1}^k \phi_i(r_1, \dots, r_k)$$

- Z verbindet alle Variablen im Modell in einer Funktion
- $Z \neq 1$ in den meisten Fällen

ML Vorgehen für ungerichtete Modelle

- Gegeben ein Faktormodel $F = \{f_i\}_{i=1}^n$ über Zufallsvariablen $\mathbf{R} = \text{rv}(F)$
- Lass \mathbf{r}_α die Projektion von \mathbf{r} auf die Zufallsvariablen einer Entität α bezeichnen
 - Beispiel:
 - $\mathbf{r}_f = \pi_{\text{rv}(f)}(\mathbf{r})$: wählt aus \mathbf{r} die Werte aus, die sich auf die Zufallsvariablen in f beziehen
- Lass ϕ_f sich auf die Potentialfunktion von f beziehen

- Gegeben ein Datensatz \mathbf{d} mit N zusammengesetzten Ereignissen für \mathbf{R}
 - I.e., vollständige Daten
 - Lass $\#(\mathbf{r})$ bezeichnen, wie häufig $\mathbf{r} \in \text{Val}(\mathbf{R})$ in \mathbf{d} vorkommt
 - \mathbf{d} als Multi-Menge: $\mathbf{d} = \{(\mathbf{r}, \#(\mathbf{r}))\}_{\mathbf{r} \in \text{Val}(\mathbf{R})}$
 - $N = \sum_{\mathbf{r} \in \text{Val}(\mathbf{R})} \#(\mathbf{r})$

ML Vorgehen für ungerichtete Modelle

1. Drücke die Wahrscheinlichkeit (Likelihood) $P(\mathbf{d}|\boldsymbol{\theta})$ der Daten \mathbf{d} als Funktion der Parameter $\boldsymbol{\theta}$ aus, die zu lernen sind

- $\boldsymbol{\theta}$ bezieht sich hier auf die Potentiale in den Faktoren von F

- $\mathbf{d} = \{r_j\}_{j=1}^N = \{(r, \#(r))\}_{r \in \text{Val}(\mathbf{R})}$

$$P(\mathbf{d}|\boldsymbol{\theta}) = \prod_{j=1}^N P(d_j|\boldsymbol{\theta}) = \prod_{r \in \text{Val}(\mathbf{R})} P(r|\boldsymbol{\theta})^{\#(r)}$$

- In der Zählerdarstellung ($\#(r)$) kann $P(\mathbf{d}|\boldsymbol{\theta})$ null werden, wenn ein $r \in \text{Val}(\mathbf{R})$ nicht beobachtet wurde
 - Abhilfe: initialisiere alle Zähler mit 1

und nehme den Logarithmus

$$\log P(\mathbf{d}|\boldsymbol{\theta}) = \log \prod_{r \in \text{Val}(\mathbf{R})} P(r|\boldsymbol{\theta})^{\#(r)}$$

$$\begin{aligned}
 \log P(D|\theta) &= \log \prod_{\mathbf{r} \in \text{Val}(\mathbf{R})} P(\mathbf{r}|\theta)^{\#(\mathbf{r})} && \phi_f(\mathbf{r}_f) \text{ ist ein Parameter } \in \theta \\
 &= \sum_{\mathbf{r} \in \text{Val}(\mathbf{R})} \#(\mathbf{r}) \log P(\mathbf{r}|\theta) = \sum_{\mathbf{r} \in \text{Val}(\mathbf{R})} \#(\mathbf{r}) \log \left(\frac{1}{Z} \prod_{f \in F} \phi_f(\mathbf{r}_f) \right) \\
 &= \sum_{\mathbf{r} \in \text{Val}(\mathbf{R})} \#(\mathbf{r}) \left(\log \left(\frac{1}{Z} \right) + \log \prod_{f \in F} \phi_f(\mathbf{r}_f) \right) \\
 &= \sum_{\mathbf{r} \in \text{Val}(\mathbf{R})} \#(\mathbf{r}) \left(-\log(Z) + \sum_{f \in F} \log \phi_f(\mathbf{r}_f) \right) \\
 &= \left(\sum_{\mathbf{r} \in \text{Val}(\mathbf{R})} \#(\mathbf{r}) \sum_{f \in F} \log \phi_f(\mathbf{r}_f) \right) - \sum_{\mathbf{r} \in \text{Val}(\mathbf{R})} \#(\mathbf{r}) \log(Z) \\
 &= \left(\sum_{\mathbf{r} \in \text{Val}(\mathbf{R})} \#(\mathbf{r}) \sum_{f \in F} \log \phi_f(\mathbf{r}_f) \right) - N \log(Z) \\
 &= \left(\sum_{f \in F} \sum_{\mathbf{r}_f \in \text{Val}(\text{rv}(f))} \#(\mathbf{r}_f) \log \phi_f(\mathbf{r}_f) \right) - N \log(Z)
 \end{aligned}$$

ML Vorgehen für ungerichtete Modelle

2. Leite die Log-Wahrscheinlichkeiten nach jedem Parameter ab, i.e., $\phi_f(\mathbf{r}_f)$

$$\log P(\mathbf{d}|\boldsymbol{\theta}) = \left(\sum_{f \in F} \sum_{\mathbf{r}_f \in \text{Val}(\text{rv}(f))} \underbrace{\#(\mathbf{r}_f) \log \phi_f(\mathbf{r}_f)}_{l_1} \right) - \underbrace{N \log(Z)}_{l_2}$$

- Ableitung von l_1

$$\frac{\partial l_1}{\partial \phi_f(\mathbf{r}_f)} = \#(\mathbf{r}_f) \cdot \frac{1}{\phi_f(\mathbf{r}_f)}$$

konstant bzgl. $\phi_f(\mathbf{r}_f)$

Ableiten von log:

$$\frac{\partial \log f(x)}{\partial x} = \frac{1}{f(x)} \cdot \frac{\partial f(x)}{\partial x}$$

Z.B., wenn $f(x) = x$:

$$\frac{\partial \log x}{\partial x} = \frac{1}{x} \cdot \frac{\partial x}{\partial x} = \frac{1}{x} \cdot 1 = \frac{1}{x}$$

- Ableitung von $l_2 = N \log(Z)$

$$\frac{\partial l_2}{\partial \phi_f(\mathbf{r}_f)} = N \cdot \frac{1}{Z} \cdot \frac{\partial Z}{\partial \phi_f(\mathbf{r}_f)} = N \cdot \frac{1}{Z} \cdot \underbrace{\frac{\partial}{\partial \phi_f(\mathbf{r}_f)} \cdot \sum_{\mathbf{s} \in \text{Val}(\text{rv}(F))} \prod_{e \in F} \phi_e(\mathbf{s}_e)}_{\text{...}}$$

Nur Summanden, die $\phi_f(\mathbf{r}_f)$ enthalten, bleiben; Ableitung der anderen ergibt 0; Indikatorfunktion stellt das dar:

$$\mathbf{1}(\mathbf{r}, \mathbf{s}) = \begin{cases} 1 & \text{if } \mathbf{r} = \mathbf{s} \\ 0 & \text{otherwise} \end{cases}$$

$$= N \cdot \frac{1}{Z} \cdot \sum_{\mathbf{s} \in \text{Val}(\text{rv}(F))} \mathbf{1}(\mathbf{r}_f, \mathbf{s}_f) \underbrace{\frac{\partial}{\partial \phi_f(\mathbf{r}_f)} \cdot \prod_{e \in F} \phi_e(\mathbf{s}_e)}_{\text{...}}$$

Eins der e 's ist f , i.e., $\phi_f(\mathbf{r}_f) = \phi_e(\mathbf{s}_e)$ für ein e ; der Rest ist

$$= \frac{\partial}{\partial \phi_f(\mathbf{r}_f)} \cdot \phi_f(\mathbf{s}_f) \cdot \prod_{\substack{e \in F \\ e \neq f}} \phi_e(\mathbf{s}_e) = \prod_{\substack{e \in F \\ e \neq f}} \phi_e(\mathbf{s}_e)$$

konstant bzgl. $\phi_f(\mathbf{r}_f)$

$$= \frac{\phi_f(\mathbf{r}_f)}{\phi_f(\mathbf{r}_f)} \cdot \prod_{\substack{e \in F \\ e \neq f}} \phi_e(\mathbf{s}_e) = \frac{1}{\phi_f(\mathbf{r}_f)} \cdot \prod_{e \in F} \phi_e(\mathbf{s}_e)$$

- Ableitung von $l_2 = k \log(Z)$

$$\begin{aligned}
 \frac{\partial l_2}{\partial \phi_f(\mathbf{r}_f)} &= N \cdot \frac{1}{Z} \cdot \frac{\partial Z}{\partial \phi_f(\mathbf{r}_f)} = N \cdot \frac{1}{Z} \cdot \frac{\partial}{\partial \phi_f(\mathbf{r}_f)} \cdot \sum_{\mathbf{s} \in \text{Val}(\text{rv}(F))} \prod_{e \in F} \phi_e(\mathbf{s}_e) \\
 &= N \cdot \frac{1}{Z} \cdot \sum_{\mathbf{s} \in \text{Val}(\text{rv}(F))} \mathbf{1}(\mathbf{r}_f, \mathbf{s}_f) \frac{\partial}{\partial \phi_f(\mathbf{r}_f)} \cdot \prod_{e \in F} \phi_e(\mathbf{s}_e) \\
 &= N \cdot \frac{1}{Z} \cdot \sum_{\mathbf{s} \in \text{Val}(\text{rv}(F))} \mathbf{1}(\mathbf{r}_f, \mathbf{s}_f) \frac{1}{\phi_f(\mathbf{r}_f)} \cdot \prod_{e \in F} \phi_e(\mathbf{s}_e) \\
 &= N \cdot \frac{1}{\phi_f(\mathbf{r}_f)} \cdot \sum_{\mathbf{s} \in \text{Val}(\text{rv}(F))} \mathbf{1}(\mathbf{r}_f, \mathbf{s}_f) \cdot \frac{1}{Z} \cdot \prod_{e \in F} \phi_e(\mathbf{s}_e) \\
 &= N \cdot \frac{1}{\phi_f(\mathbf{r}_f)} \cdot \sum_{\mathbf{s} \in \text{Val}(\text{rv}(F))} \mathbf{1}(\mathbf{r}_f, \mathbf{s}_f) \cdot P(\mathbf{s}) \\
 &= N \cdot \frac{1}{\phi_f(\mathbf{r}_f)} \cdot P(\mathbf{r}_f)
 \end{aligned}$$

ML Vorgehen für ungerichtete Modelle

2. Leite die Log-Wahrscheinlichkeiten nach jedem Parameter ab, i.e., $\phi_f(\mathbf{r}_f)$

$$\log P(\mathbf{d}|\boldsymbol{\theta}) = \left(\sum_{f \in F} \sum_{\mathbf{r}_f \in \text{Val}(\text{rv}(f))} \underbrace{\#(\mathbf{r}_f) \log \phi_f(\mathbf{r}_f)}_{l_1} \right) - \underbrace{N \log(Z)}_{l_2}$$

- Ableitung für jeden Parameter $\phi_f(\mathbf{r}_f)$

$$\frac{\partial l_1}{\partial \phi_f(\mathbf{r}_f)} = \frac{\#(\mathbf{r}_f)}{\phi_f(\mathbf{r}_f)}$$

$$\frac{\partial l_2}{\partial \phi_f(\mathbf{r}_f)} = N \cdot \frac{P(\mathbf{r}_f)}{\phi_f(\mathbf{r}_f)}$$

$$\frac{\partial \log P(\mathbf{d}|\boldsymbol{\theta})}{\partial \phi_f(\mathbf{r}_f)} = \frac{\#(\mathbf{r}_f)}{\phi_f(\mathbf{r}_f)} - N \cdot \frac{P(\mathbf{r}_f)}{\phi_f(\mathbf{r}_f)}$$

ML Vorgehen für ungerichtete Modelle

3. Ermittle die Parameterwerte, für die die Ableitungen null ist

$$\frac{\partial \log P(\mathbf{d}|\boldsymbol{\theta})}{\partial \phi_f(\mathbf{r}_f)} = \frac{\#(\mathbf{r}_f)}{\phi_f(\mathbf{r}_f)} - N \cdot \frac{P(\mathbf{r}_f)}{\phi_f(\mathbf{r}_f)} = 0 \rightarrow \frac{\#(\mathbf{r}_f)}{\phi_f(\mathbf{r}_f)} = N \cdot \frac{P(\mathbf{r}_f)}{\phi_f(\mathbf{r}_f)}$$
$$\#(\mathbf{r}_f) = N \cdot P(\mathbf{r}_f)$$
$$\frac{\#(\mathbf{r}_f)}{N} = P(\mathbf{r}_f)$$

- Gibt an, dass die ML Schätzungen so sein sollten, dass die Modellmarginalverteilungen $P(\mathbf{r}_f)$ gleich den normalisierten empirischen Zählern sind:

$$P^\#(\mathbf{r}_f) \stackrel{\text{def}}{=} \frac{\#(\mathbf{r}_f)}{N} \stackrel{!}{=} P(\mathbf{r}_f)$$

- Gibt nicht an, wie man die Schätzungen erhält

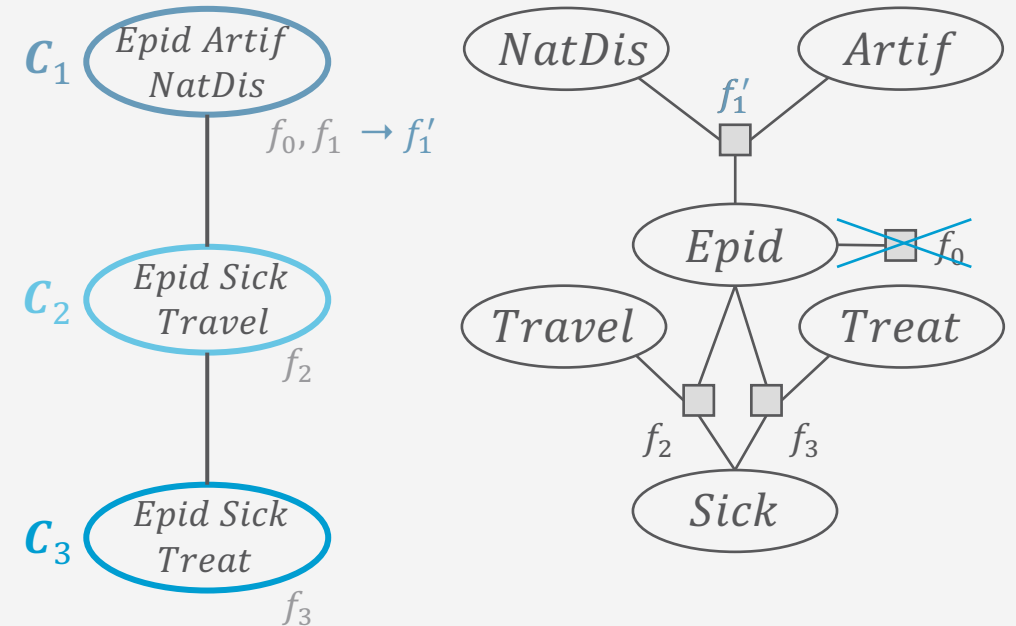
Lernen von Faktoren über maximale Cliques

- Gegeben ein Faktormodell mit Faktoren über maximale Cliques
 - Im entsprechenden Jtree (V, E) haben die lokalen Modelle einen Faktor f_C
- Vollständige gemeinsame Verteilung P_F :

$$P_F = \frac{1}{Z} \prod_{f \in F} f = \frac{1}{Z} \prod_{C \in V} f_C = \frac{\prod_{C \in V} P(C)}{\prod_{\{i,j\} \in E} P(S_{ij})}$$

- Für eine Welt $\mathbf{r} \in \text{Val}(\mathbf{R})$,

$$P(\mathbf{r}) = \frac{\prod_{C \in V} P(\mathbf{r}_C)}{\prod_{\{i,j\} \in E} P(\mathbf{r}_{S_{ij}})} = \frac{\prod_{f \in F} P(\mathbf{r}_f)}{\prod_{\{i,j\} \in E} P(\mathbf{r}_{S_{ij}})}$$



Wenn das Modell Faktoren f über nicht-maximale Cliques vorsieht, dann werden diese f pro maximale Clique in einen Faktor f_C zusammengefasst.

Lernen von Faktoren über maximale Cliques

$$P^\#(\mathbf{r}_f) \stackrel{\text{def}}{=} \frac{\#(\mathbf{r}_f)}{N} \stackrel{!}{=} P(\mathbf{r}_f)$$

- Mit $P(\mathbf{r}) = \frac{\prod_{C \in \mathcal{V}} P(\mathbf{r}_C)}{\prod_{\{i,j\} \in E} P(\mathbf{r}_{S_{ij}})} = \frac{\prod_{f \in F} P(\mathbf{r}_f)}{\prod_{\{i,j\} \in E} P(\mathbf{r}_{S_{ij}})}$, setze die Faktoren auf die normalisierten empirischen Zähler $P^\#(\mathbf{r}_f)$ und, für jeden Separator, wähle einen Nachbar und teile den Faktor durch den normalisierten empirischen Zähler des Separators $P^\#(\mathbf{r}_{S_{ij}})$, i.e.,

- Für jedes $f \in F$,
 - Setze $\phi_f(\mathbf{r}_f) \leftarrow P^\#(\mathbf{r}_f)$
- Für jedes $\{i, j\} \in E$, i.e., Separator S_{ij} ,
 - Wähle ein $h \in \{i, j\}$ zufällig
 - Setze $\phi_h(\mathbf{r}_h) \leftarrow \frac{\phi_h(\mathbf{r}_h)}{P^\#(\mathbf{r}_{S_{ij}})}$

ML-Lernen mit Max Cliques

- Sorgt für (bedingte) Wahrscheinlichkeitsverteilungen in den Faktoren \rightarrow Erzwingt $Z = 1$
 - Vergleich: Importance Sampling für Faktormodelle

Lernen von Faktoren über maximale Cliques: Beispiel

- Max-Cliquen-Faktoren auf normalisierte Zähler setzen

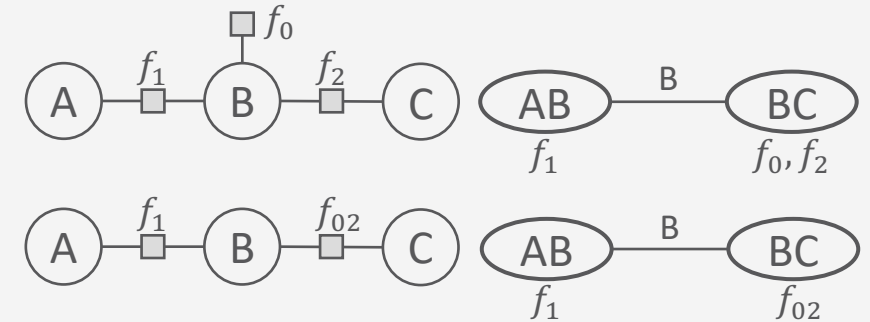
A	B	#	$P^\#(A,B) = f_1$
false	false	10	$10/70 = 0.14$
false	true	10	$10/70 = 0.14$
true	false	30	$30/70 = 0.43$
true	true	20	$20/70 = 0.29$

B	C	#	$P^\#(B,C)$
false	false	24	$24/70 = 0.34$
false	true	16	$16/70 = 0.23$
true	false	24	$24/70 = 0.34$
true	true	6	$6/70 = 0.09$

- Für jede Kante zwischen zwei maximalen Cliques, einen der beiden Faktoren durch normalisierten Zähler des Separators teilen

B	#	$P^\#(B)$
false	40	$40/70 = 0.57$
true	30	$30/70 = 0.43$

B	C	$P^\#(B,C)/P^\#(B) = f_{02}$
false	false	$0.34/0.57 = 0.60$
false	true	$0.23/0.57 = 0.40$
true	false	$0.34/0.43 = 0.79$
true	true	$0.09/0.43 = 0.21$



A	B	C	#
false	false	false	6
false	false	true	4
false	true	false	8
false	true	true	2
true	false	false	18
true	false	true	12
true	true	false	16
true	true	true	4

Lernen mit allgemeinen Faktoren

- Mit Faktoren über nicht-maximale Cliques, geschlossene Lösung nicht möglich

- Fixpunkt-Iteration:
$$\frac{\partial \log P(D|\theta)}{\partial \phi_f(\mathbf{r}_f)} = \frac{\#(\mathbf{r}_f)}{\phi_f(\mathbf{r}_f)} - N \cdot \frac{P(\mathbf{r}_f)}{\phi_f(\mathbf{r}_f)} = 0 \quad \rightarrow \quad \frac{\#(\mathbf{r}_f)}{\phi_f(\mathbf{r}_f)} = N \cdot \frac{P(\mathbf{r}_f)}{\phi_f(\mathbf{r}_f)}$$

$$\frac{\#(\mathbf{r}_f)}{N \cdot \phi_f(\mathbf{r}_f)} = \frac{P(\mathbf{r}_f)}{\phi_f(\mathbf{r}_f)}$$

$$\frac{P^\#(\mathbf{r}_f)}{\phi_f(\mathbf{r}_f)} = \frac{P(\mathbf{r}_f)}{\phi_f(\mathbf{r}_f)}$$

$$\frac{\phi_f(\mathbf{r}_f)}{P^\#(\mathbf{r}_f)} = \frac{\phi_f(\mathbf{r}_f)}{P(\mathbf{r}_f)}$$

$$\phi_f(\mathbf{r}_f) = \phi_f(\mathbf{r}_f) \frac{P^\#(\mathbf{r}_f)}{P(\mathbf{r}_f)}$$

- Update-Regel

$$\phi_f^{(t+1)}(\mathbf{r}_f) \leftarrow \phi_f^t(\mathbf{r}_f) \frac{P^\#(\mathbf{r}_f)}{P^t(\mathbf{r}_f)}$$

- $\phi_f^t(\mathbf{r}_f)$ momentane Potentiale
- $P^\#(\mathbf{r}_f)$ fest
- $P^t(\mathbf{r}_f)$ Anfrage zu berechnen auf Basis des momentanen $\phi_f^t(\mathbf{r}_f)$

Iterative Proportional Fitting: Iteratives Proportionales Anpassen

- Iterative proportional fitting (IPF) Prozedur
 - Initialisiere alle Faktoren mit $t = 0$ uniform, z.B., alle Potentiale auf 1
 - **for** $t = 1, 2, \dots$ **do**
 - **if** Konvergenzkriterium gilt nicht **then**
 - **for** alle $\mathbf{r}_f \in \text{Val}(\text{rv}(f)), f \in F$ **do**

$$\phi_f^{(t+1)}(\mathbf{r}_f) \leftarrow \phi_f^t(\mathbf{r}_f) \frac{P^\#(\mathbf{r}_f)}{P^t(\mathbf{r}_f)}$$
 - **else break**

IPF

- Konvergenzkriterium, z.B., Fehlergrenze ε
 - $\forall f \in F : \phi_f^{(t+1)}(\mathbf{r}_f) - \phi_f^t(\mathbf{r}_f) < \varepsilon$
- Viele $P^t(\mathbf{r}_f)$ Anfragen zu berechnen für alle Faktoren!
 - Effizientere Ausführung mittels Jtree

IPF mit Jtrees

- Jtree nutzen um $P^t(\mathbf{r}_f)$ effizient zu berechnen
 - Jtree-Konstruktion unabhängig von Parametern

- **IPF-JT: IPF mit Jtrees**

- Baue einen Jtree
- Initialisiere alle Faktoren und Nachrichten uniform für $t = 0$
- Wähle zufällig ein Cluster \mathbf{C}_i als momentanes Cluster
- **for** $t = 1, 2, \dots$ **do**
 - **if** Konvergenzkriterium gilt nicht **then**
 - **for** alle $\mathbf{r}_f \in \text{Val}(\text{rv}(f)), f \in F_i^t$ **do**

$$\phi_f^{(t+1)}(\mathbf{r}_f) \leftarrow \phi_f^t(\mathbf{r}_f) \frac{P^\#(\mathbf{r}_f)}{P^t(\mathbf{r}_f)}$$
 - **else break**
 - Wähle zufällig ein Nachbar $\mathbf{C}_j \in \text{Nb}(\mathbf{C}_i)$ als neues momentanes Cluster
 - Berechne und sende Nachricht m_{ij}^t zu \mathbf{C}_j

Berechne $P^t(\mathbf{r}_f)$ auf $\{m_{ji}^t\}_{\mathbf{C}_j \in \text{Nb}(\mathbf{C}_i)}$ und F_i^t

- Organisiere Berechnung der unterschiedlichen $P^t(\mathbf{r}_f)$ in sinnvoller Reihenfolge

Keine Reihenfolge vorgeschrieben;
implizite Anforderung, dass alle \mathbf{C}_i
häufig genug besucht werden

IPF-JT

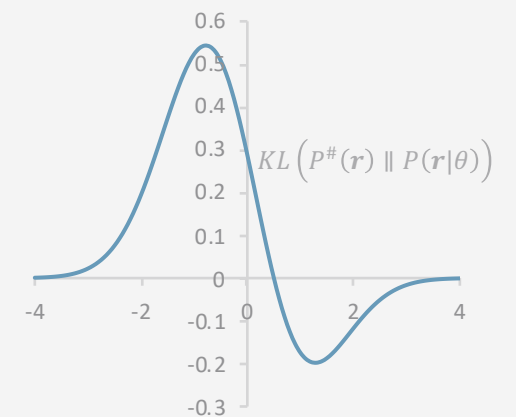
Eigenschaften von IPF Updates

- Ein Update pro Cluster bei Faktoren über maximale Cliques
- **Koordinatenanstieg**-Algorithmus
 - Koordinaten = Parameter θ in Clustern
 - In jedem Schritt erhöht das Update die Log-Wahrscheinlichkeit der Daten, $\log P(\mathbf{d}|\theta)$, und es konvergiert zu einem *globalen Maximum*
- Maximieren der Log-Wahrscheinlichkeit ist äquivalent zum **Minimieren der KL Divergenz** (Kreuzentropie, *cross entropy*)

$$\max \log P(\mathbf{d}|\theta) \Leftrightarrow \min KL \left(P^\#(\mathbf{r}) \parallel P(\mathbf{r}|\theta) \right)$$

$$KL \left(P^\#(\mathbf{r}) \parallel P(\mathbf{r}|\theta) \right) = \sum_{\mathbf{r}} P^\#(\mathbf{r}) \log \frac{P^\#(\mathbf{r})}{P(\mathbf{r}|\theta)}$$

- Max-Entropie-Prinzip für Parametrisierung: *duale Perspektive zu MLE*



IPF mit nicht vollständigen Daten

- Problem mit nicht vollständigen Daten:
 - Update-Regel

$$\phi_f^{(t+1)}(\mathbf{r}_f) \leftarrow \phi_f^t(\mathbf{r}_f) \frac{P^\#(\mathbf{r}_f)}{P^t(\mathbf{r}_f)}$$

- Versteckte Variablen haben kein $P^\#(\mathbf{r}_f)$
- **EM-IPF**: IPF Prozedur mit dem EM-Algorithmus verbunden
 - **E-Schritt**: Berechne *erwartete* Zähler für versteckte Variablen
 - **M-Schritt**: Aktualisiere Parameter $\phi_f^{(t+1)}(\mathbf{r}_f)$
 - Auch mit dem Jtree für Effizienz vereinbar

Zwischenzusammenfassung

- Parameter lernen in ungerichteten Modellen schwerer, weil $Z \neq 1$
- IPF
 - Fixpunkt-Iteration
 - Reduziert sich auf ein Update pro Faktor, wenn Faktoren über maximal Cliques gehen
 - Dann $Z = 1$ erzwingen
 - Jtrees für Effizienz
 - Update auf einem Cluster, dann aktualisierte Information an einen Nachbarn schicken
 - EM Variante für nicht vollständige Daten

Latent Dirichlet Allocation (LDA) zur Topic Modellierung

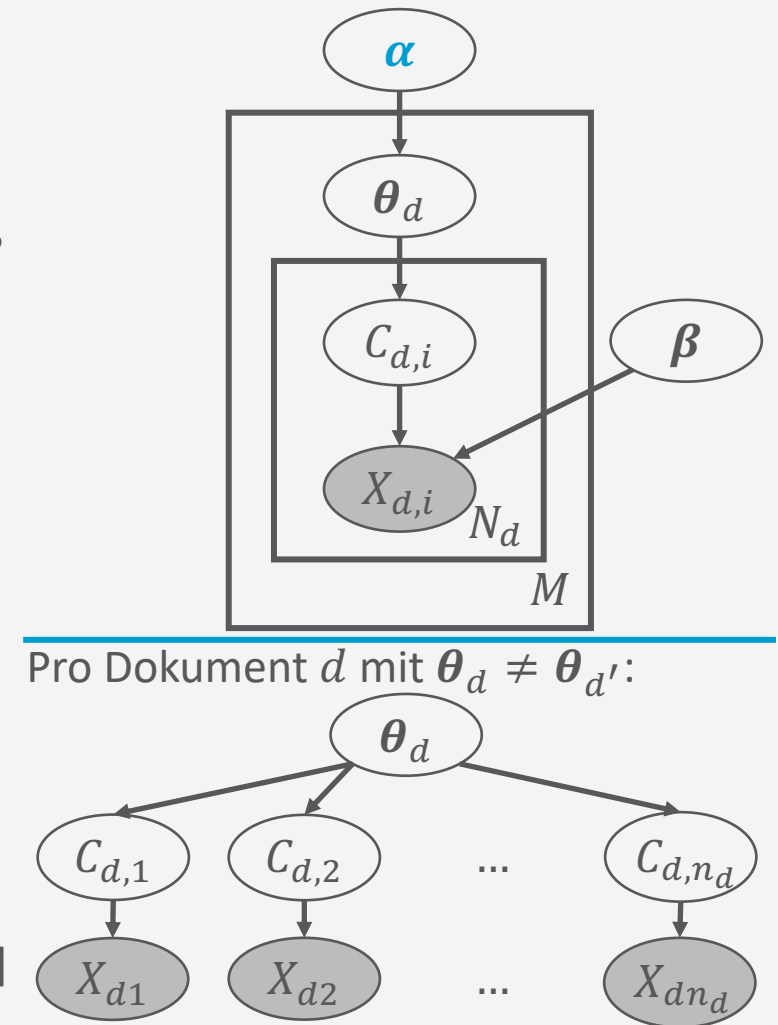
Umsetzung von Lernverfahren in komplexen Modellen

Topic Modellierung: Mixture of Topics

- Annahme: **Dokumente bestehen aus einem Mix von Topics**
- Pro Dokument $d \in \mathcal{C}$: Zufallsvariable X_i bezeichnet, welches Wort w aus einem Lexikon \mathcal{D} an Position i in d vorkommt
 - I.e., $\text{Val}(X_i) = \{w\}_{w \in \mathcal{D}}$
 - Bzw. $\text{Val}(X_i) = \mathcal{D}$ mit \mathcal{D} als Menge der Wörter im Lexikon
 - X_i unabhängig von X_j gegeben k :

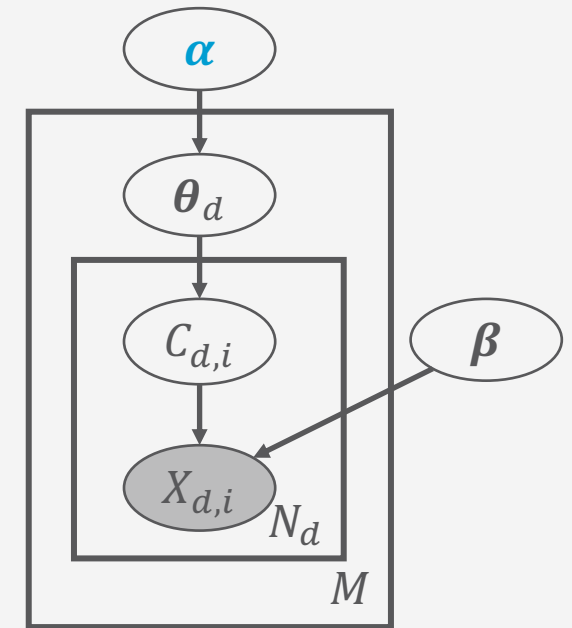
$$P(X_i, X_j | C) = P(X_i | C)P(X_j | C)$$
 - Annahme: Vorkommen von w unabhängig von Position:

$$P(X_i = w) = P(X_j = w)$$
- Verteilung über mögliche A-priori-Topic-Verteilung θ_d für jedes Dokument: **Hyperparameter α**
 - Vorwissen kodieren, welche Verteilungen wahrscheinlicher sind



Topic Modellierung: Mixture of Topics

- Dokumente bestehen aus einem Mix von Topics
 - M Dokumente, N_d Worte im Dokument d
 - K mögliche Topics, pro Dokument d Topic-Verteilung θ_d
- Annahme: Dokumente des Corpus sind Ergebnis des folgenden Generierungsprozesses
 - Sample eine Topic-Verteilung θ_d über K Topics aus $Dir(\alpha)$
 - Wähle N_d (z.B. sample aus Poisson-Verteilung mit Hyperparameter ξ)
 - Wiederhole N_d mal
 - Sample Topic $c_{d,i}$ aus $P(C_{d,i}) = \theta_d$
 - Sample Wort $x_{d,i}$ aus $P(X_{d,i}|c_{d,i}) = \beta[c_{d,i}]$
 - Häufige Schreibweise in der Forschungsgemeinde:
 $c_{d,i} \sim Mul(\cdot | \theta_d)$ und $x_{d,i} \sim Mul(\cdot | \beta_c)$



$$\theta_d = (0.32, 0.34, 0.25, 0.09)$$

$$N_d = 68$$

Topic Modellierung: Mixture of Topics

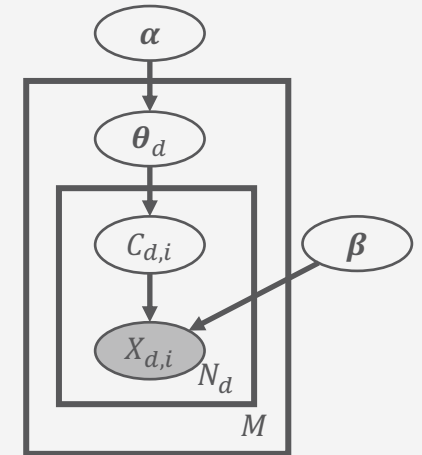
- LDA-Sicht auf ein Dokument
 - Worte des Dokuments kommen aus 4 Topics

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. “Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants and act every bit as important as our traditional areas of support in health, medical research, education and the social services;” Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center’s share will be \$200.000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400.000 each. The Juilliard School, where music and the performing arts are taught, will get \$250.000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100.000 donation, too.

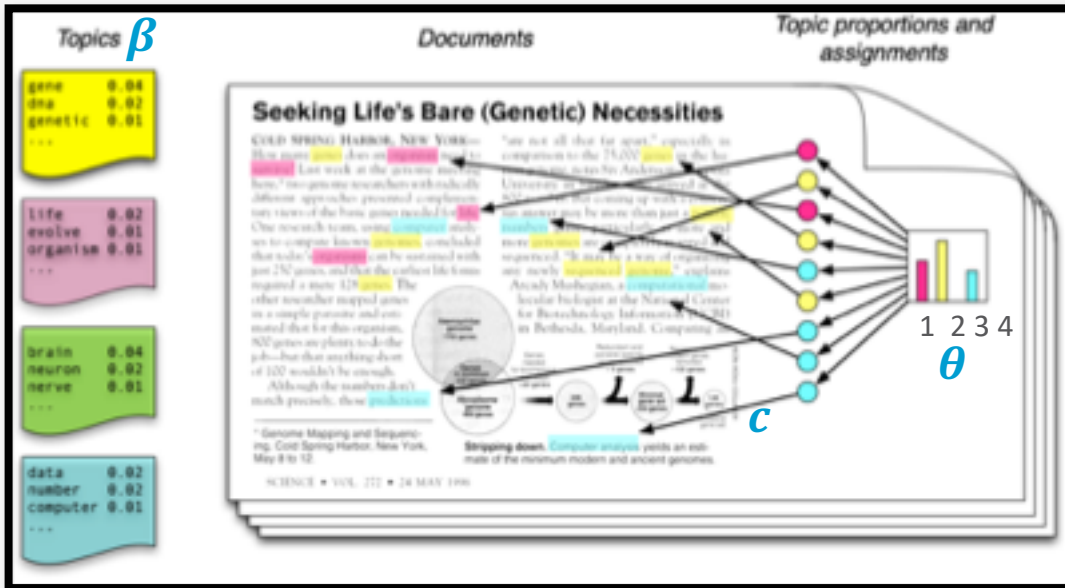
1 (Arts)	2 (Budgets)	3 (Children)	4 (Education)
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

Unbekannte Parameter in LDA

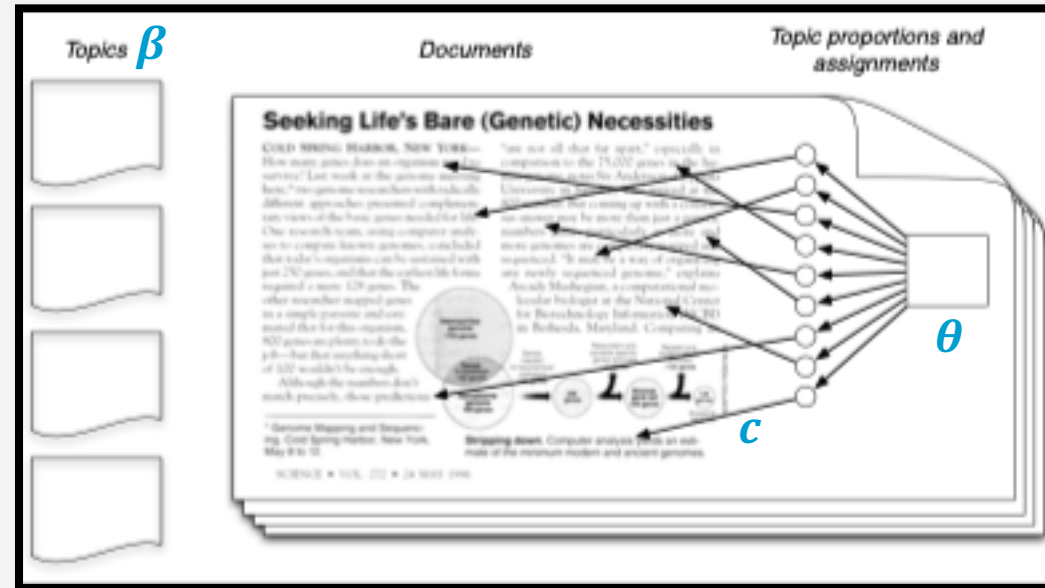
- Was sind die *Topics der Worte*, wie lauten die *Topicverteilungen der Dokumente* und wie lauten die *Wortverteilungen der Topics*?
- $P(\beta, \theta, c | x, \alpha)$



Ziel



Vorhanden

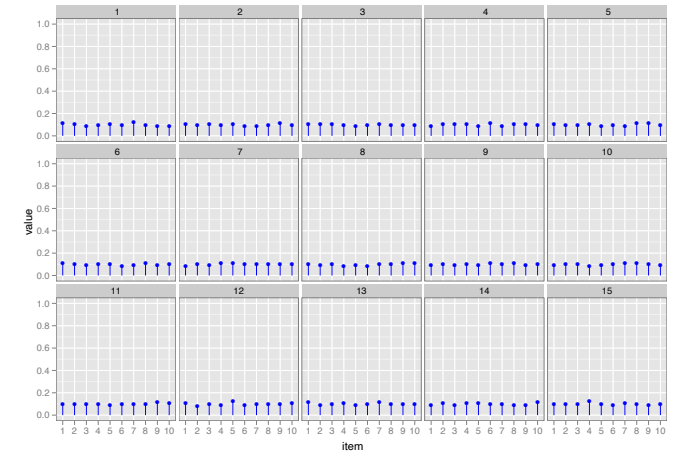
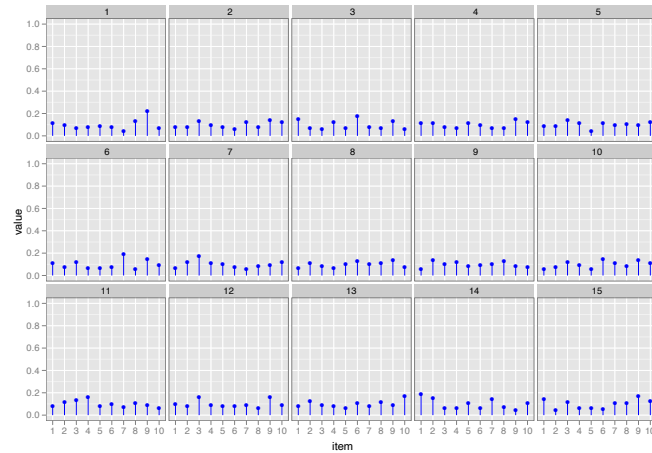
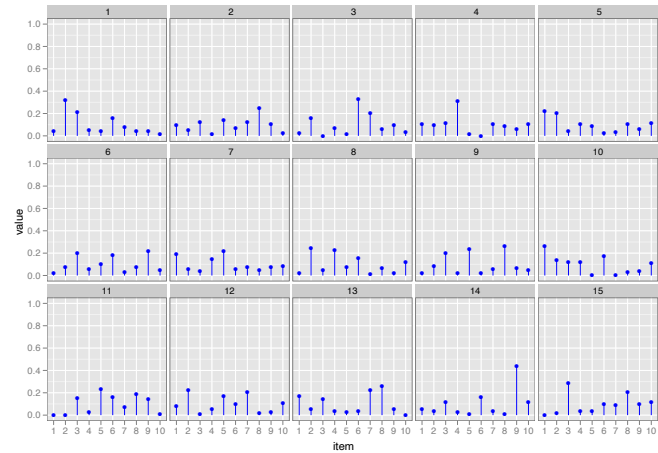


Hyperparameter α : Einfluss auf das Lernergebnis

$\alpha = 1$

$\alpha = 10$

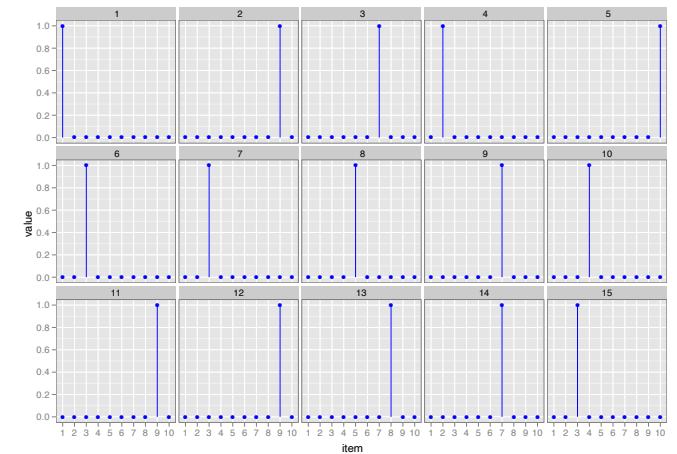
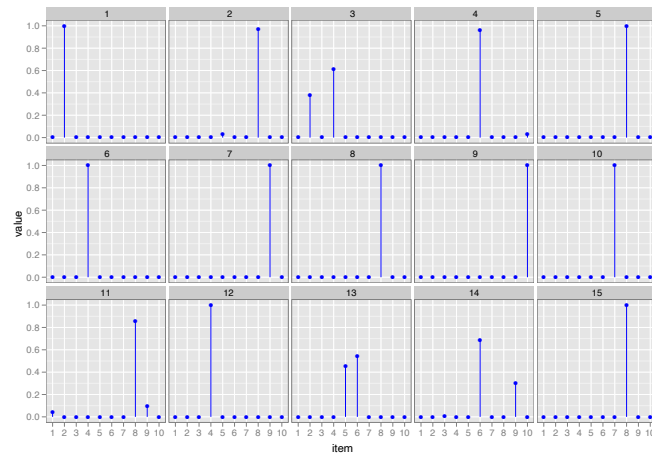
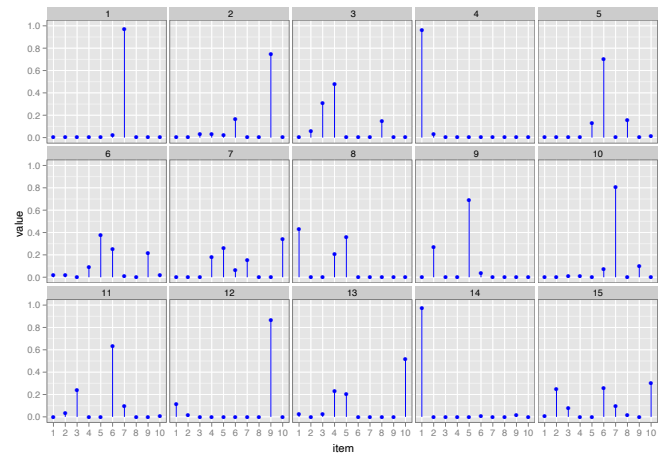
$\alpha = 100$



$\alpha = 0.1$

$\alpha = 0.01$

$\alpha = 0.001$

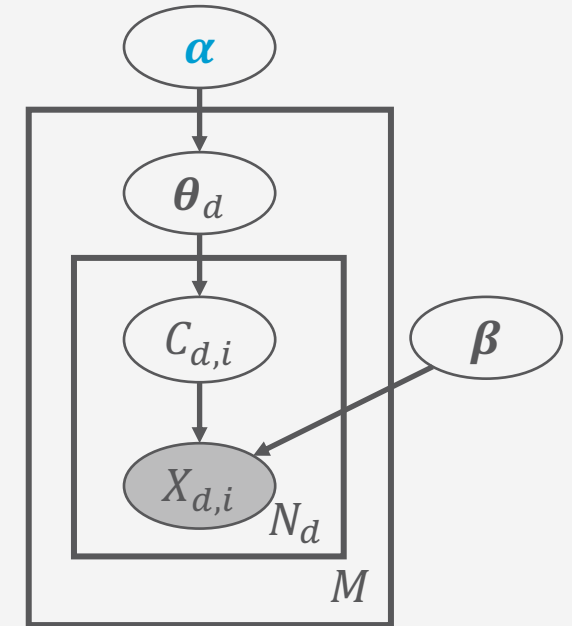
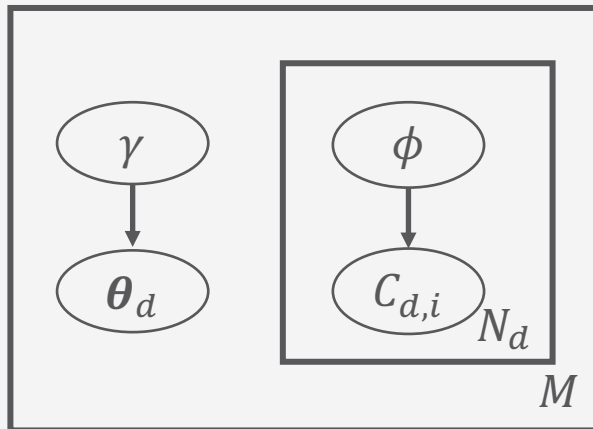


LDA: Parameter Lernen

- Zwei Methoden
 - Variational Inference / EM
 - Numerische Approximation über untere Schranken
 - Bias in der Lösung
 - Konvergenz mit numerischen Garantien
 - Gibbs Sampling
 - Stochastische Simulation
 - Kein Bias in der Lösung
 - Stochastische Konvergenz

LDA: Variational Inference – Skizze

- Ersetze LDA Modell P mit einem einfacheren Modell Q :



- Minimiere die KL Divergenz iterativ im EM-Stil:

$$(\gamma^*, \phi^*) = \arg \min_{(\gamma, \phi)} KL(Q(\theta, z | \gamma, \phi) \parallel P(\theta, z | w, \alpha, \beta))$$

LDA: Gibbs Sampling – Skizze

- MCMC Algorithmus
 - Halte alle momentanen Werte außer einem fest und sample den Wert, z.B., für $C_{i,d}$

$$P(C_{i,d} = k | \mathbf{c}_{-i,d}, \mathbf{w})$$
 - $\mathbf{c}_{-i,d}$ die Topic-Zuweisungen der Worte $X_{j \neq i}$ (ohne i)
- Notation
 - c entspricht z , x entspricht w
 - Alle Dokument gleiche Länge N
 - $n_{d,k}$: Häufigkeit von Topic k in Dokument d
→ über alle k normiert ergeben die Zähler θ_d
 - $n_{k,w}$: Häufigkeit von Wort w in Topic k → β_k
 - n_k : Häufigkeit von Topic k → α (angepasst)
- Konvergiert *irgendwann* zur wahren Verteilung

```

Input: words  $\mathbf{w} \in$  documents  $\mathbf{d}$ 
Output: topic assignments  $\mathbf{z}$  and counts  $n_{d,k}, n_{k,w}$ , and  $n_k$ 
begin
  randomly initialize  $\mathbf{z}$  and increment counters
  foreach iteration do
    for  $i = 0 \rightarrow N - 1$  do
       $word \leftarrow w[i]$ 
       $topic \leftarrow z[i]$ 
       $n_{d,topic} -= 1; n_{word,topic} -= 1; n_{topic} -= 1$ 
      for  $k = 0 \rightarrow K - 1$  do
         $p(z = k | \cdot) = (n_{d,k} + \alpha_k) \frac{n_{k,w} + \beta_w}{n_k + \beta \times W}$ 
      end
       $topic \leftarrow \text{sample from } p(z | \cdot)$ 
       $z[i] \leftarrow topic$ 
       $n_{d,topic} += 1; n_{word,topic} += 1; n_{topic} += 1$ 
    end
  end
  return  $\mathbf{z}, n_{d,k}, n_{k,w}, n_k$ 
end
  
```

Algorithm 1: LDA Gibbs Sampling

LDA: Parameter lernen

- Zwei Methoden
 - Variational Inference / EM
 - Numerische Approximation über untere Schranken
 - Bias in der Lösung
 - Konvergenz mit numerischen Garantien
 - Gibbs Sampling
 - Stochastische Simulation
 - Kein Bias in der Lösung
 - Stochastische Konvergenz
- Vergleich
 - Gibbs Sampling ist langsamer
 - Läuft Tage für einen Datensatz mittlerer Größe
 - Variational Inference braucht nur ein paar Stunden dafür
 - Gibbs Sampling ist genauer
 - Gibbs Sampling Konvergenz schwer zu testen

Warum LDA funktioniert?

- Trade-off zwischen zwei Zielen:
 1. Für jedes Dokument: Worte so wenig Topics wie möglich zuschreiben
 2. Für jedes Topic: Hohe Wahrscheinlichkeiten so wenigen Worten wie möglich zuweisen
- Diese Ziele stehen im Konflikt
 - Worte eines Dokuments d einem einzigen Topic k zuschreiben macht Ziel 2 schwer erfüllbar: Alle Worte in d müssen eine nicht vernachlässigbare Wahrscheinlichkeit in k haben
 - Wenige Worte einem Topic k zuschreiben, so dass (nur) diese eine hohe Wahrscheinlichkeit in k haben, macht Ziel 1 schwer erfüllbar: Um die Worte eines Dokuments d mit hoher Wahrscheinlichkeit zu erzeugen, müssen die Worte vielen Topics zugeschrieben sein
- Diese Ziele gegeneinander abwägen ermöglicht Gruppen von häufig miteinander auftretenden Worten zu finden

LDA: Anwendung zur Analyse von sozialen Netzwerken

- “follow”-Beziehung zwischen Nutzer*innen sieht erstmal chaotisch aus
 - Netzwerk willkürlich entstanden und nicht zentral kontrolliert
- Strukturen in dieser Art von Beziehung finden, indem man
 - Nutzer*innen basierend auf ihren „Topic“-Interessen gruppiert und
 - Jede Beziehung mit der identifizierten „Topic“-Gruppe labelt
- Nutzt Gibbs Sampling
- Getestet auf Twitter Daten

LDA: Anwendung zur Analyse von sozialen Netzwerken

- Spezifische Topic-Gruppe
 - g : Nutzer*in, $|e'(g)|$: Follower-Zahl, Bio: Twitter Bio

g	$ e'(g) $	Bio
lancearmstrong	2370	7-time Tour de France winner, full time cancer fighter - LIVESTRONG!
levileipheimer	285	Pro Cyclist
ghincapie	234	<i>[No Bio] Professional Cyclist</i>
johanbruyneel	203	9X winning Tour De France Sports Director
barackobama	7410	44th President of the United States
teamradioshack	141	<i>[No Bio] Pro Cycle Team RadioShack</i>
dzabriskie	137	Professional Bike Racer
philliggett	126	Pro-cycling commentator.
christianvdv	114	Professional Cyclist
stephenfry	3384	British Actor, Writer, Lord of Dance, Prince of Swimwear & Blogger

- Topic-Gruppe „mobile gadget blog“ aus dem „non-popular“ Datensatz
 - Nur Twitter Nutzer*innen mit $|e'(g)| \leq 100$

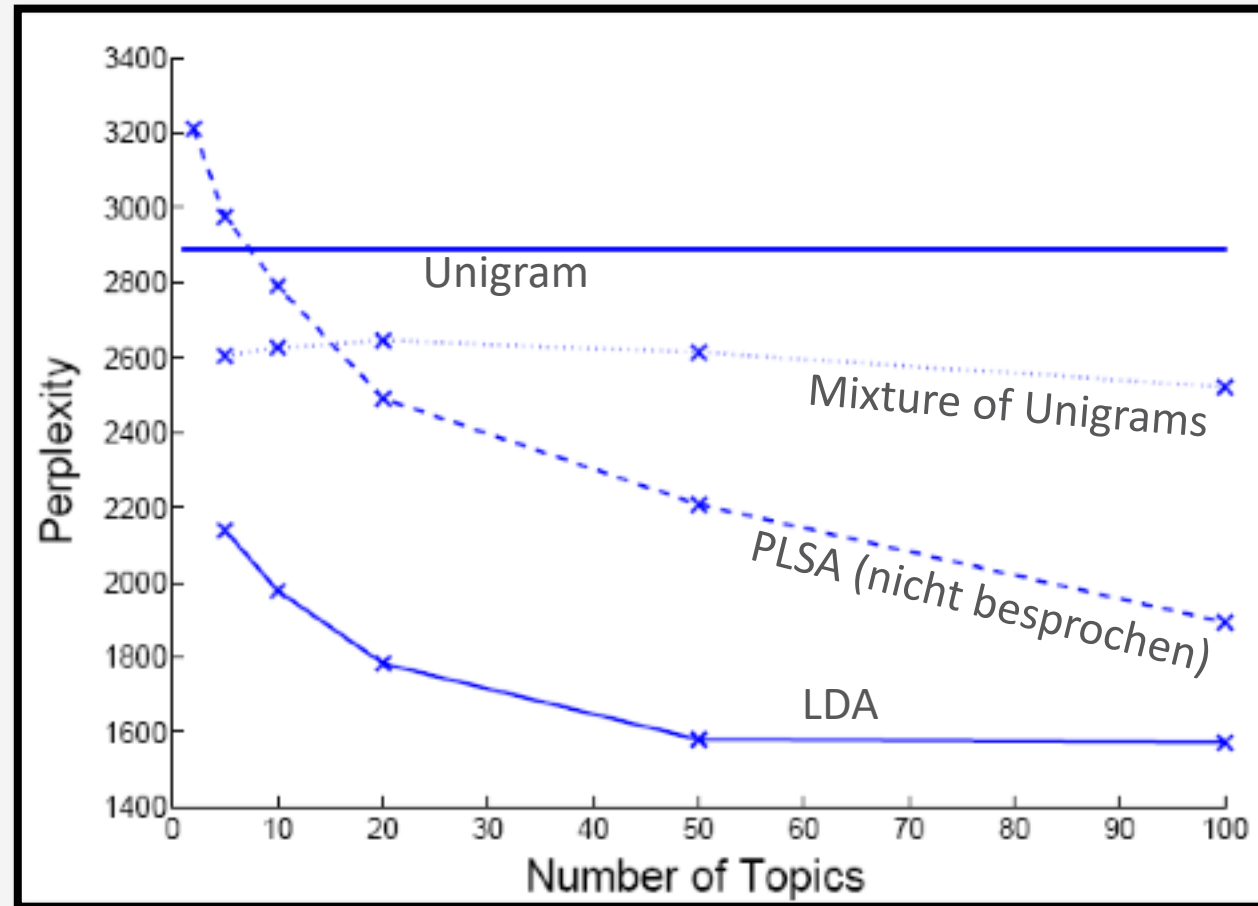
g	$ e'(g) $	Bio
googleimages	94	News, tips, and tricks direct from the Google Images team
androidandme	84	Meet your new Android friend.
androidguys	73	Your trusted source for Android news and opinion since November 5, 2007.
prethinking	66	PreThinking.com - your source for all things Palm Pre
engadgetmobile	97	Official Twitter account of Engadget Mobile!
skyfire	88	Skyfire 3.0 for Android and Skyfire 2.0 for iPhone are changing the way people browse on thier mobile devices! www.get.skyfire.com
precentral	75	All about Palm webOS devices
boygenius	95	President, General Manager, Editor-in-chief of BGR.com. Professional smack talker.
theandroidsite	71	News and reviews of Google's Android OS
nokconv	58	A team of Nokia tweeters. Follow us for various updates, info, and the occasional banal banter.

LDA: Performanz bewerten mittels Perplexität

- Wie gut sagt eine Wahrscheinlichkeitsverteilung oder probabilistisches Modell ein Sample vorher?
 - Q : Modell einer unbekanntes Wahrscheinlichkeitsverteilung P gelernt aus einem Datensatz bestehend aus Datenpunkten gesampelt aus P
 - Werte Q aus, indem man fragt, wie gut Q ein gesondertes Test-Sample x_1, \dots, x_n (auch aus P gesampelt) vorhersagt
 - **Perplexität** (*perplexity*) von Q definiert als
$$b^{-\frac{1}{n} \sum_{i=1}^n \log_b Q(x_i)}$$
 - b üblicherweise 2
 - Ein besseres Q tendiert dazu höhere Wahrscheinlichkeiten $Q(x_i)$ zuzuweisen \rightarrow geringere Perplexität („weniger überrascht – perplex – vom Sample“)

Perplexität unterschiedlicher Modelle

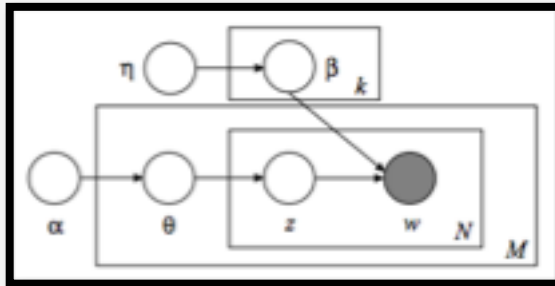
Weniger ist besser →



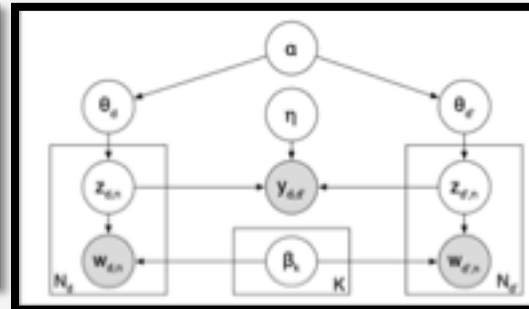
LDA Erweiterungen

Smoothed LDA

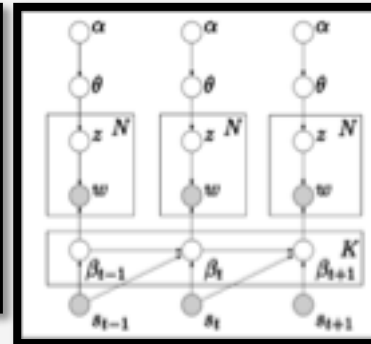
(Hyperparameter wie α zur Beeinflussung der Wortverteilungen pro Topic)



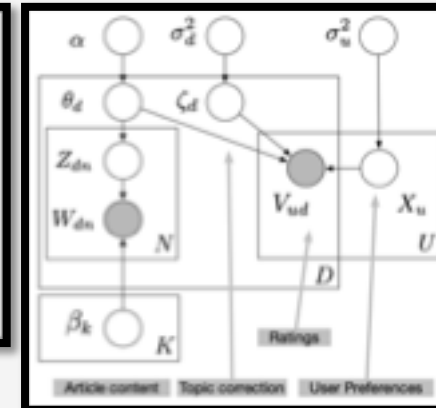
Relational Topic Model
(Links zwischen Papieren)



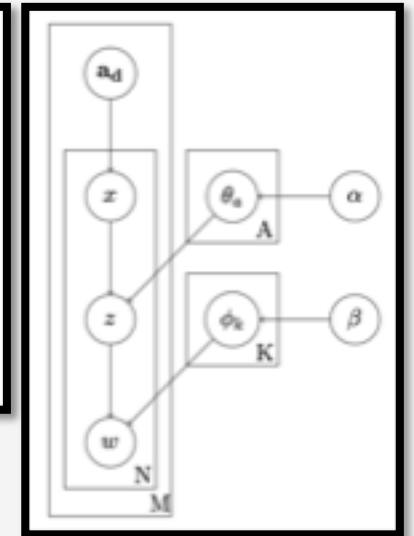
Dynamic Topic Model
(über die Zeit)



Topic Model für Empfehlungen



Author-Topic Modell (Übung 2)



- Viele Erweiterungen sind im Laufe der Zeit vorgeschlagen worden
- Implementierungen in Python vorhanden (mindestens LDA, aber auch andere Varianten)

Zwischenzusammenfassung

- Komplexes Modell, in dem viele versteckte Parameter zu lernen sind
- Annahme: Dokumente sind durch einen generativen Prozess entstanden
 - Worte als Evidenz dieses Prozesses
- Parameter lernen mittels
 - Variational Inference + EM
 - Einfacheres Modell nehmen, KL Divergenz minimieren mit einem EM-Vorgehen
 - Gibbs Sampling
 - Relative Häufigkeiten durch Samples abschätzen

Überblick: 5. Lernalgorithmen für episodische PGMs

A. Überblick

- Full Bayesian Learning, MAP Learning, Maximum Likelihood (ML) Learning

B. ML-basiertes Parameter Lernen

- In BNs: ML (vollständige Daten), Expectation-Maximisation (EM; nicht-vollständige Daten)
- In Faktormodellen: Iterative Proportional Fitting (IPF), EM-IPF
- Anwendung: LDA

C. ML-basiertes Struktur Lernen

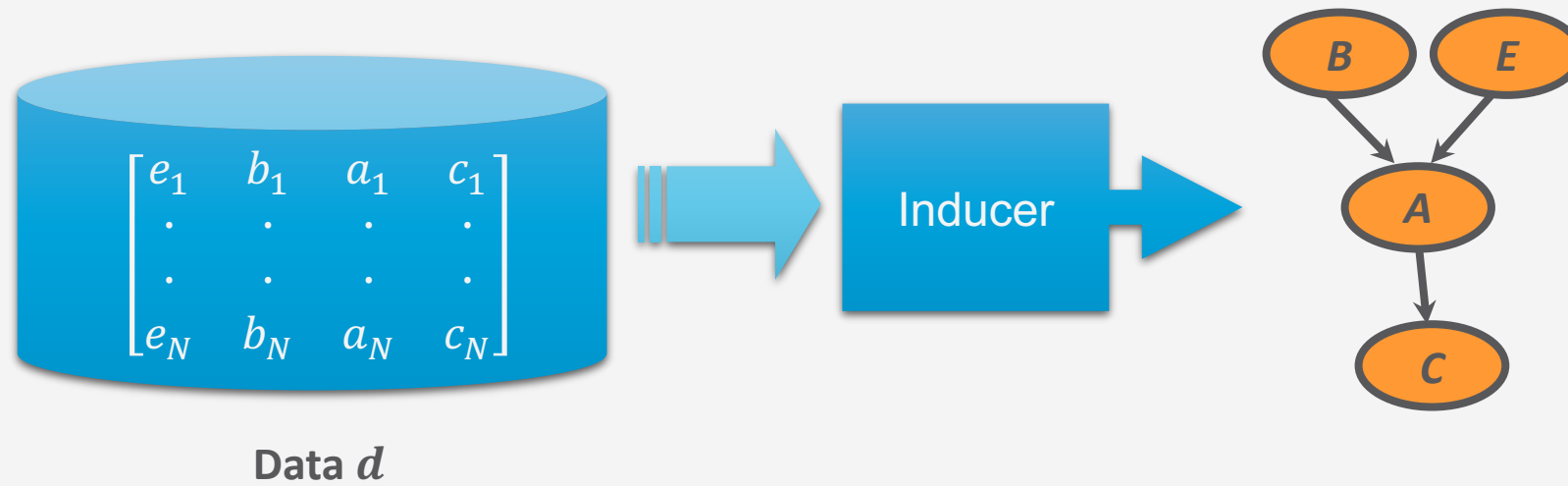
- Model Selection, Structural EM

D. Abschlussbemerkungen

- Alternativen, generative vs. diskriminative Modelle, Korrelation vs. Kausalität

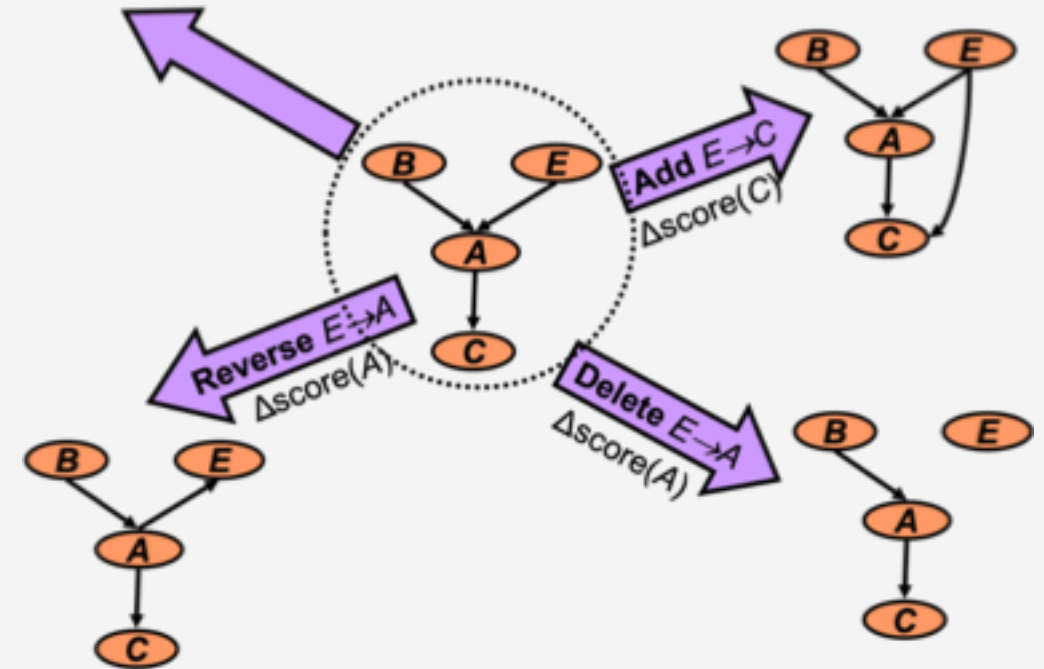
Struktur lernen

- Unbekannte Graphstruktur
- Gegeben Datensatz d
- Finde Modell, das am besten zu d passt; beinhaltet:
 - *Model Selection*: Modellwahl
 - Parameter Schätzung (wie auf den vorherigen Folien)



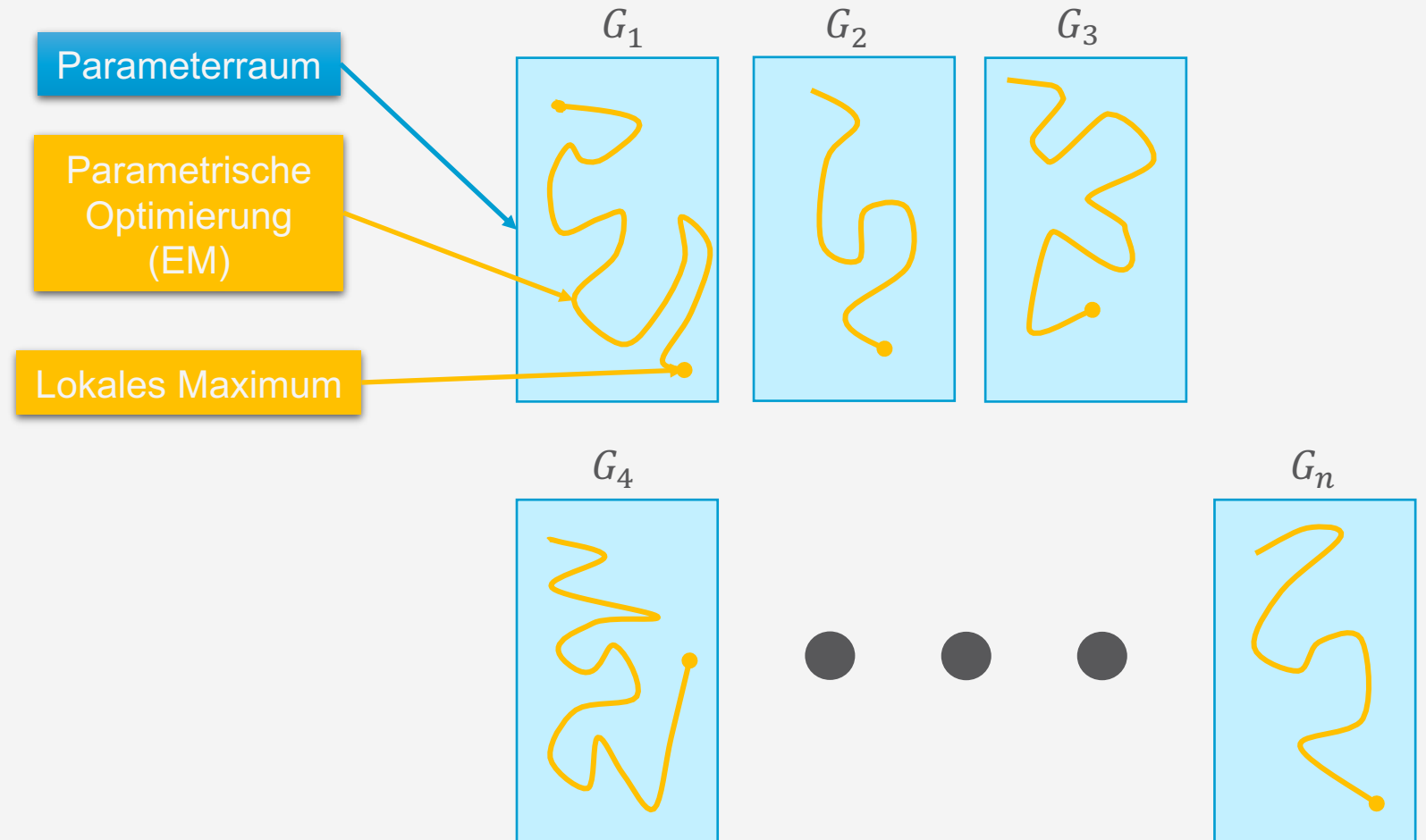
Model Selection: Modellauswahl

- Ziel: Beste Graphstruktur wählen
 - Eingabe: Daten, Bewertungsfunktion (*scoring function*)
 - Ausgabe: Graphstruktur mit maximaler Bewertung
1. Führe eine heuristische Suche nach Modellkandidaten durch
 2. Schätze in jedem Modell die Parameter mit EM
 - Bei vollständigen Daten \rightarrow z.B. MLE
 3. Bewerte jedes Modell
 4. Wähle das Modell mit der höchsten Bewertung



Lokale Suche in der Praxis

- Schätze in jedem Modell die Parameter mit EM
- Berechnungstechnisch teuer
 - Parameter optimieren mit EM – nicht trivial
 - EM für alle Kandidaten ausführen
 - Verbringt Zeit auch auf schwachen Kandidaten
 - Daten können unvollständig sein (nicht repräsentativ)

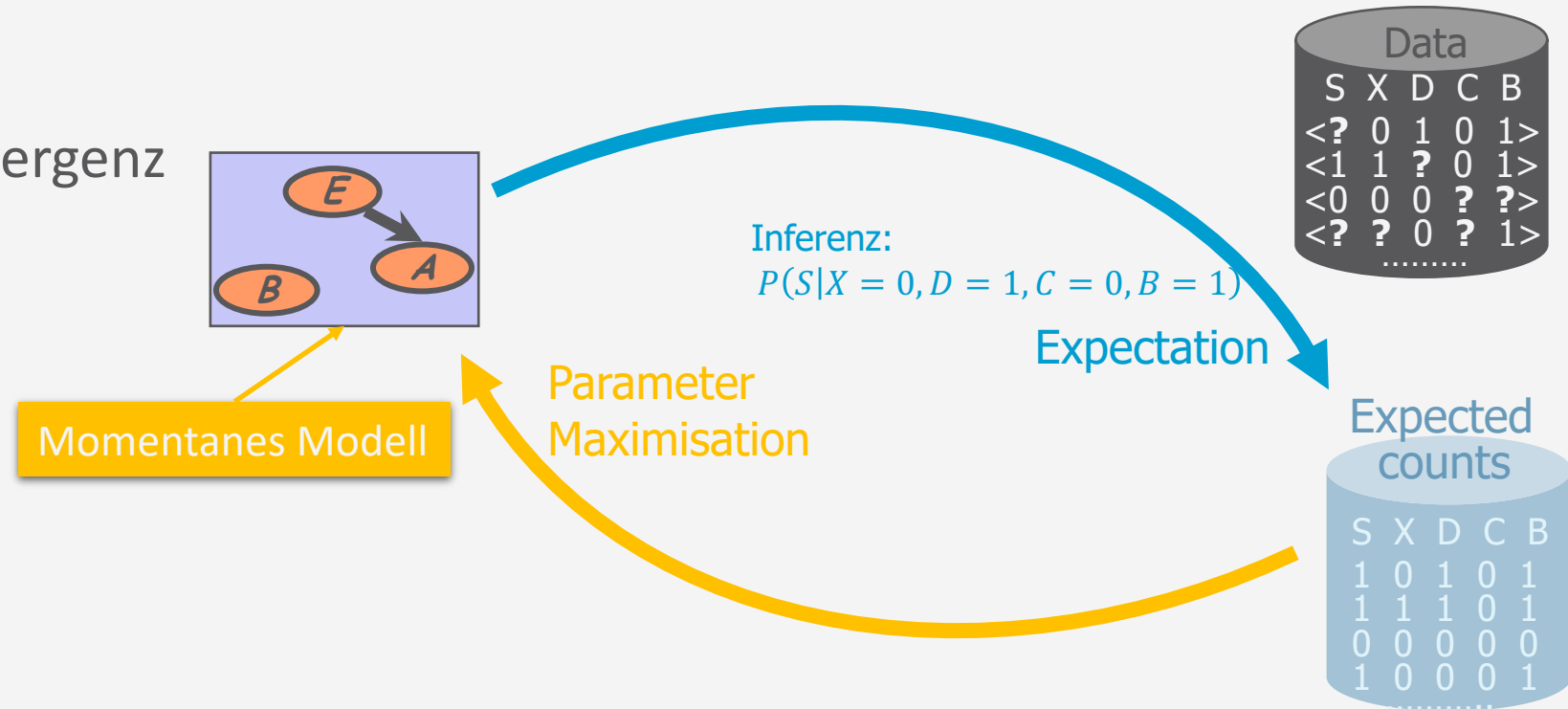


Struktur lernen: Nicht vollständige Daten

- Problem: versteckte, unbekannte Variablen
 - Keine Daten für die versteckten Variablen
 - Suchraum sehr groß
- Idee:
 - Nutze das momentane Modell um neue Strukturen zu bewerten
- Übersicht:
 - Suche im (Struktur, Parameter)-Raum durchführen
 - In jeder Iteration, das momentane Modell nutzen zum Finden von:
 - Besser bewertete Parameter: „parametrischer“ EM-Schritt
oder
 - Besser bewertete Struktur: „struktureller“ EM-Schritt

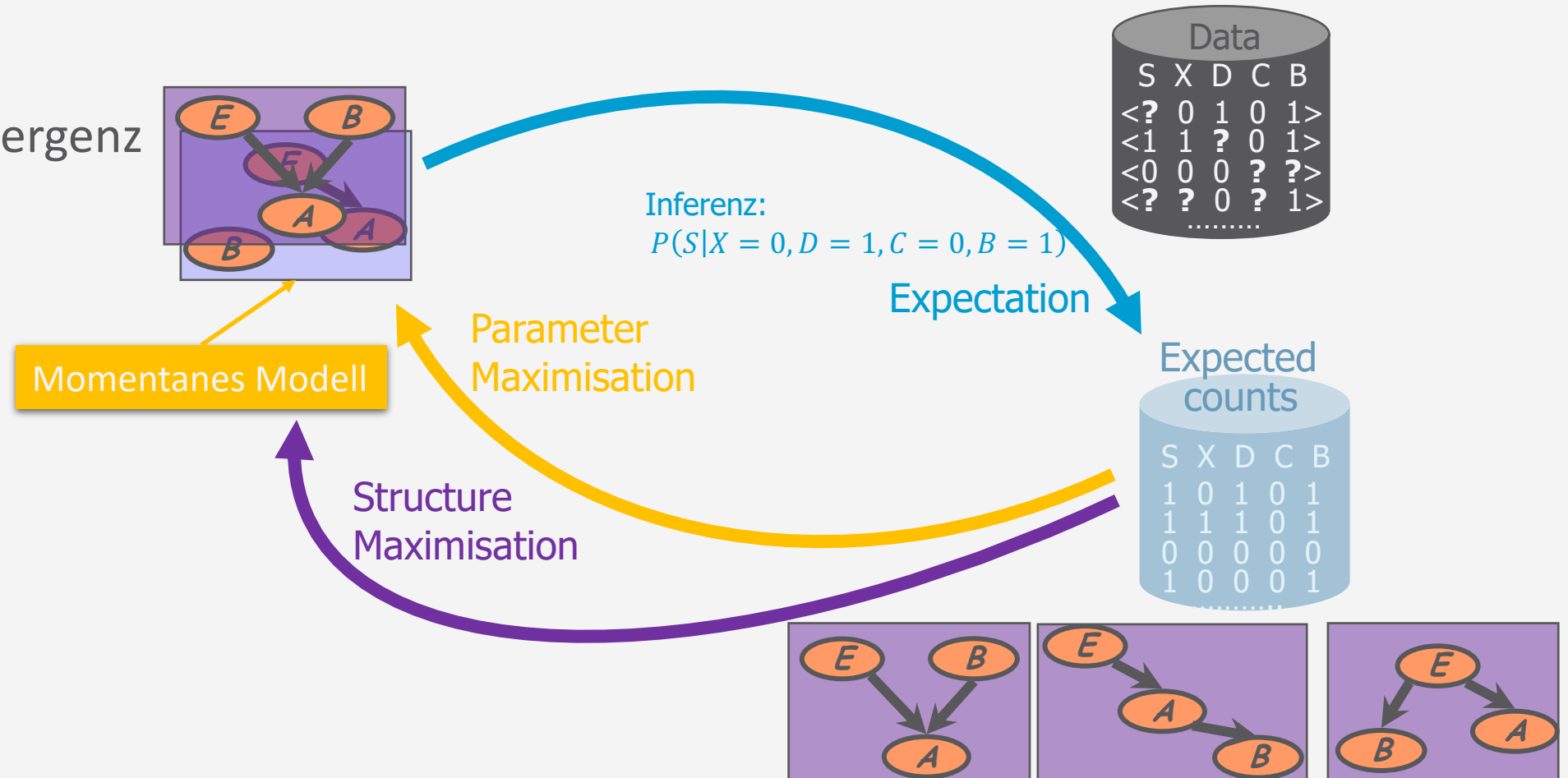
EM

EM-Algorithmus:
Iteriere bis zur Konvergenz



Structural EM

SEM-Algorithmus:
Iteriere bis zur Konvergenz



Zwischenzusammenfassung

- Heuristische Suche nach geeigneten Graphstrukturen
- Model Selection
 - Suche nach Kandidaten
 - EM für jeden Kandidaten, um beste Parameter zu bestimmen
 - Den Kandidaten mit höchster Bewertung nehmen
- Structural EM
 - Immer nur ein momentanes Modell, welches man entweder bezogen auf die Parameter aktualisiert oder bezogen auf die Struktur anpasst
 - Abwechselnd Parameter und Struktur optimieren

Überblick: 5. Lernalgorithmen für episodische PGMs

A. Überblick

- Full Bayesian Learning, MAP Learning, Maximum Likelihood (ML) Learning

B. ML-basiertes Parameter Lernen

- In BNs: ML (vollständige Daten), Expectation-Maximisation (EM; nicht-vollständige Daten)
- In Faktormodellen: Iterative Proportional Fitting (IPF), EM-IPF
- Anwendung: LDA

C. ML-basiertes Struktur Lernen

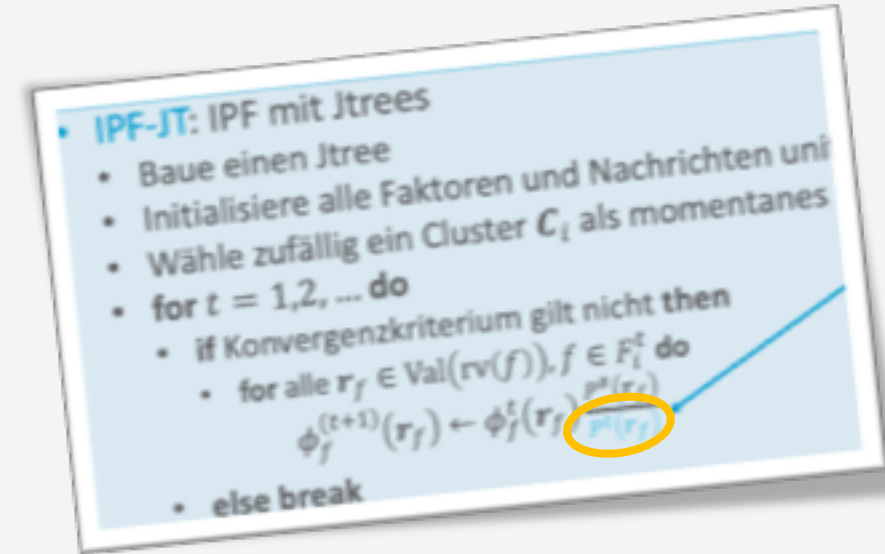
- Model Selection, Structural EM

D. Abschlussbemerkungen

- Alternativen, generative vs. diskriminative Modelle, Korrelation vs. Kausalität

Weitere Möglichkeiten

- Approximative Inferenz um Anfragen in Lern-Prozeduren zu beantworten
 - Z.B. Sampling: Beispiel Gibbs-LDA
- Andere Optimierungskriterien nutzen
 - Pseudo-Likelihood
 - Kontrastive Divergenz
- Struktur lernen
 - Ähnliche Ideen für ungerichtete Modelle verfolgen
 - Model Selection + Parameter Estimation verschränken wie in structural EM
 - Unabhängigkeitstests nutzen
 - (Bedingt) unabhängige Variablen finden, Kanten löschen
 - In Vorlesungsteil 2 Episodische PGMs skizziert
- Mehr Informationen in PGM, Kapitel 20



Generative & Diskriminative Modelle

- Methoden hier sind dafür da, Modelle zu lernen, die eine vollständige gemeinsame Wahrscheinlichkeitsverteilung $P(\mathbf{R})$ kodieren
 - Jede Art von Inferenzaufgabe lösen: Vorhersage, Klassifizierung, Wahrscheinlichkeiten von Ereignissen, wahrscheinlichste Zustände, neue Daten *generieren* (Sampling)

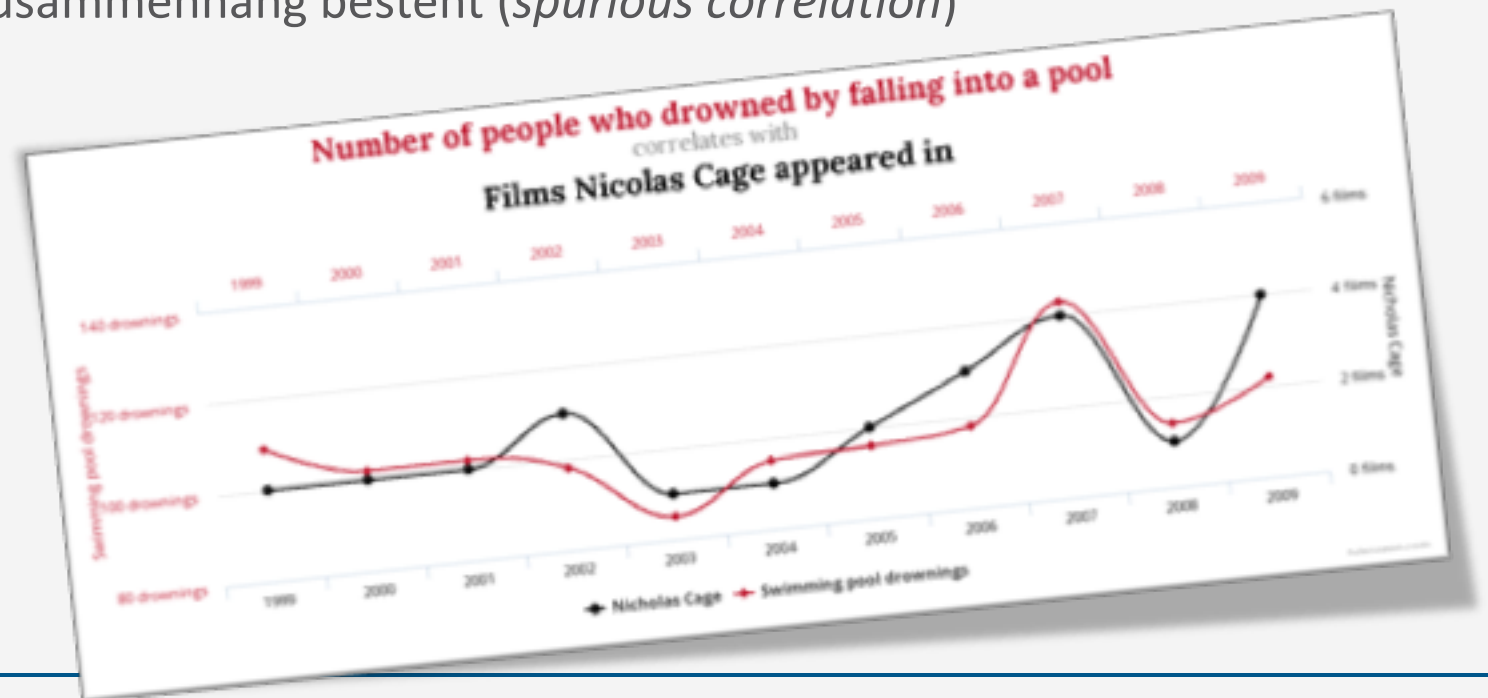
→ Generative Modelle

- Im Gegensatz dazu stehen: **Diskriminative** Modelle (wie neuronale Netze)
 - Spezifisch dafür gebaut und trainiert um Klassifikationsleistung zu maximieren: $P(C|\mathbf{X})$
 - Keine allgemeinen Inferenzaufgaben damit lösbar
 - Bei Klassifikationsanfragen im Allgemeinen bessere Performanz dieser Modelle im Vergleich zu generativen Modellen, wenn eine angemessene Menge an Trainingsdaten zur Verfügung steht
 - Fokus auf Abbildung einer bedingten Verteilung (muss nicht $P(C, \mathbf{X})$ darstellen, nur $P(C|\mathbf{X})$)

In beiden Fällen gilt: Modelle sind eine Abstraktion / Generalisierung der Daten, welche die Daten nicht mit 100% Genauigkeit repräsentieren kann.
Nur die Daten selbst können das.

Vorsicht: Korrelation vs. Kausalität

- Korrelative Zusammenhänge beim Lernen ausgenutzt
 - Aussage: Wenn A vorkommt, kommt auch B häufig vor
 - Keine Aussage darüber, ob A B verursacht oder B A oder ob es eine dritte Variable C gibt, die Verursacher ist, oder ob gar kein Zusammenhang besteht (*spurious correlation*)
 - ML-basiertes Parameter lernen: Wir haben beim Zählen für bedingte Wahrscheinlichkeiten nur gezählt, wie häufig $R = r$ vorkam gegeben jeweils die Werte von $P_a(R)$
 - Reine Korrelation, keine Aussage über Kausalität
- **Kausaler Zusammenhang unbekannt!**



Vorsicht: Korrelation vs. Kausalität

- **Kausaler Zusammenhang unbekannt!**
 - Kann vorhanden sein und bei der Interpretation der Daten helfen
 - Kann einen aber auch zu falschen Schlüssen verleiten
 - Beispiele
 - Zwei Variablen: Cholesterinwerte und Herzerkrankung
 - Positiv korreliert: Cholesterinwerte höher, häufiger Herzerkrankung vorhanden
 - Verbindung zwischen Essgewohnheiten und Herzerkrankungen komplex, aber wahrscheinlich
 - Zwei Variablen: Verkäufe von Sonnenbrillen und Hautkrebsvorkommen
 - Positiv korreliert, aber erscheint albern, dass Sonnenbrillenkäufe Hautkrebs verursachen
 - Dritter Einflussfaktor: Klima, welches dafür sorgt, dass Menschen sich der Sonne aussetzen
 - Expertenwissen kann helfen, kausale Zusammenhänge zu bestimmen



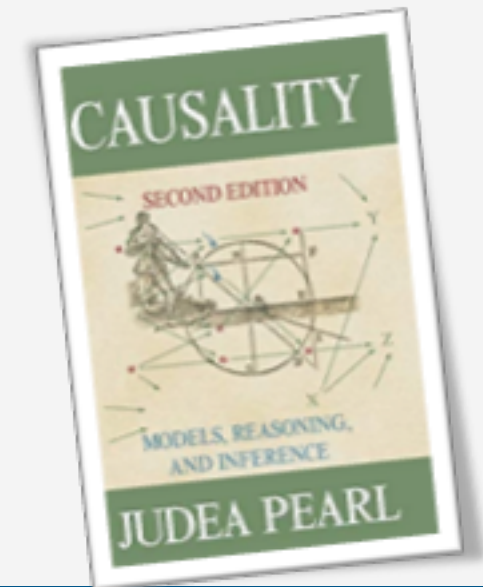
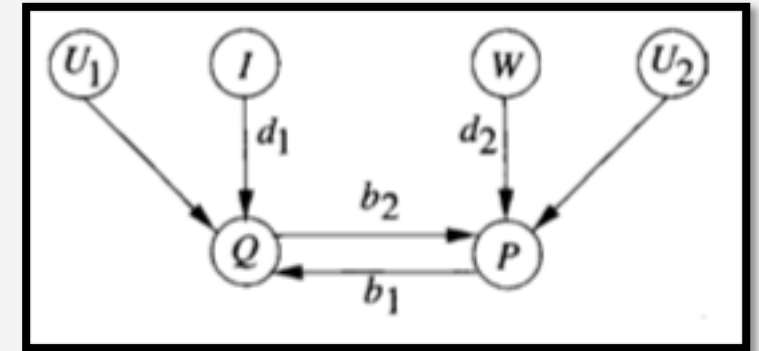
Vorsicht: Korrelation vs. Kausalität

- Problem: Auf Kausalität schwer zu prüfen
 - Wenn A B verursacht, kann B nicht ohne A auftreten
 - Krater (B) treten nur auf, wenn ein Asteroid (A) ihn verursacht hat durch eine Crash-Landung an der Stelle
 - Wenn wir nur A und B zusammen beobachten, kann das ein Indiz dafür sein
 - Es ist aber ungewiss, ob das Auftreten Ausdruck des kausalen Zusammenhangs ist oder ob B (ohne A) einfach noch nicht beobachtet wurde
 - Bzw. ob B durch etwas anderes, nicht beobachtetes hervorgerufen werden kann
- Möglichkeit: Hypothesentests einsetzen, um kausale Zusammenhänge zu prüfen
 - Quantifizieren, ob es genug Evidenz zur Unterstützung einer Hypothese gibt
- Negativbeispiele würden helfen, sind aber selten vorhanden



Vorsicht: Korrelation vs. Kausalität

- Forschungsgebiet: *Kausale Inferenz*
 - Untersuchung von Ursache-Effekt-Beziehungen
 - Modellierung wahrscheinlichkeitsbasiert
 - Structural equation models als mathematische Darstellung
 - Graph-basierte Darstellungen an BNs angelehnt
- Fragestellungen
 - Was ist Ursache, was ist Wirkung? → Richtung der Pfeile
 - Gibt es einen spürbaren Effekt bei der Durchführung einer Handlung im Vergleich zur Nicht-Durchführung der Handlung? → *do-Kalkül, Interventionen*
 - Was wäre gewesen, wenn...? → *Counterfactuals*
- Mehr zum Thema zum Beispiel von Judea Pearl
 - Tech Report: *Causal Inference in Statistics*
https://ftp.cs.ucla.edu/pub/stat_ser/r350.pdf



Zwischenzusammenfassung

- ML-basierte Methoden nicht der einzige Weg zum Ziel
- Generative Modelle erlauben beliebige Anfragen, diskriminative Modelle Klassifikationsanfragen
 - Diskriminative Modelle spezialisiert auf Klassifikation, deshalb meist besser bei der speziellen Aufgabe als generative Modelle, die alle Anfragetypen mit annehmbarer Genauigkeit beantworten können sollen
- Korrelation vs. Kausalität
 - Kausalität häufig Intuition hinter bedingten Wahrscheinlichkeiten, vor allem in BNs
 - Gemeinsamer Effekt, gemeinsame Ursache, etc.
 - Beim Lernen aus Daten nur Korrelation genutzt, die eine reine Scheinkorrelation sein kann
 - Zusätzlicher Aufwand oder Expertenwissen notwendig, um kausale Zusammenhänge zu bestätigen

Überblick: 5. Lernalgorithmen für episodische PGMs

A. Überblick

- Full Bayesian Learning, MAP Learning, Maximum Likelihood (ML) Learning

B. ML-basiertes Parameter Lernen

- In BNs: ML (vollständige Daten), Expectation-Maximisation (EM; nicht-vollständige Daten)
- In Faktormodellen: Iterative Proportional Fitting (IPF), EM-IPF
- Anwendung: LDA

C. ML-basiertes Struktur Lernen

- Model Selection, Structural EM

D. Abschlussbemerkungen

- Alternativen, generative vs. diskriminative Modelle, Korrelation vs. Kausalität

→ Sequentielle PGMs und Inferenz

Einordnung der Vorlesung: *Modell- und nutzenbasierter Agent*

- Nachfolgende Themen der Vorlesung
 2. Episodische PGMs
 3. Exakte Inferenz in episodischen PGMs
 4. Approximative Inferenz in episodischen PGMs
 5. Lernalgorithmen für episodische PGMs
 6. **Sequentielle PGMs und Inferenz**
 7. Entscheidungstheoretische PGMs

