

# **Syntaktische und Statistische Mustererkennung**

VO 1.0 840.040  
(UE 1.0 840.041)

Bernhard Jung

bernhard@jung.name  
<http://bernhard.jung.name/VUSSME/>

# Rückblick

- Nicht lineare Entscheidungsfunktionen
- SVM, Kernel Trick
- Sequentielle Klassifikation
- Boosting
- Clustering

# K-means

- Vorgabe:
  - Anzahl an Klassen  $k$
  - Funktion zur Bestimmung eines Clustermittelpunkts
  - Algorithmus:
    1. Wahl von  $k$  (beliebigen) Zentren
    2. Jedes Objekt wird dem ihm am nächsten liegenden Zentrum zugeordnet
    3. Clusterzentren werden neu berechnet
    4. Wiederholung der Schritte 2 und 3 bis sich die Zuordnung nicht mehr ändert

# K-means

- Probleme:
  - Konvergiert nicht notwendigerweise
  - Ein Cluster kann leer bleiben und somit kann sein Zentrum nicht berechnet werden
  - Lösungsansatz: Abbruch und Neustart mit anderen Clusterzentren
- Vorteil: Praktisch einfach zu berechnen und gute Resultate

# Expectation-Maximization (EM)

- Basierend auf k-means mit parametrischen Wahrscheinlichkeitsverteilungen
  - Initialisierung der Parameter der Verteilungen
  - Wiederhole bis Verfahren konvergiert:
    - Bestimme die Klassenzugehörigkeit aufgrund der Likelihood
    - Schätze die Verteilungsparameter aus der Stichprobe (Maximum Likelihood Estimation)

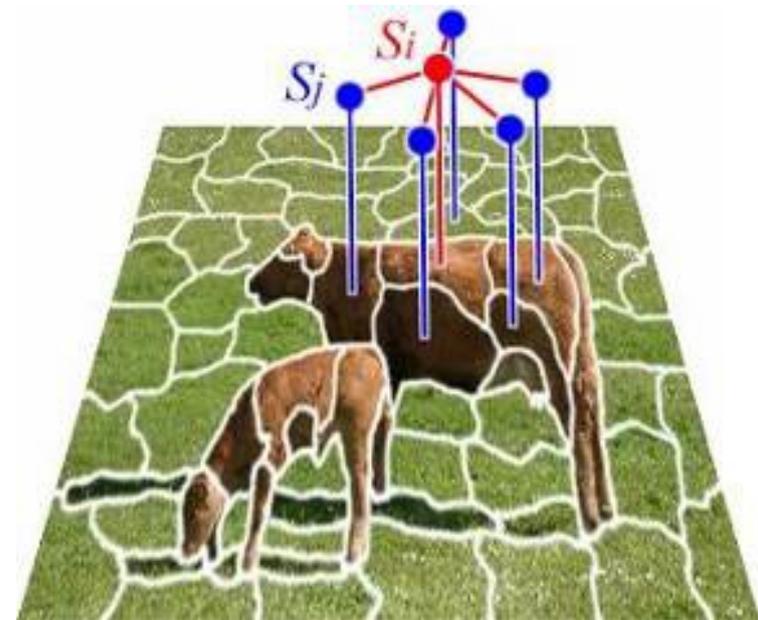
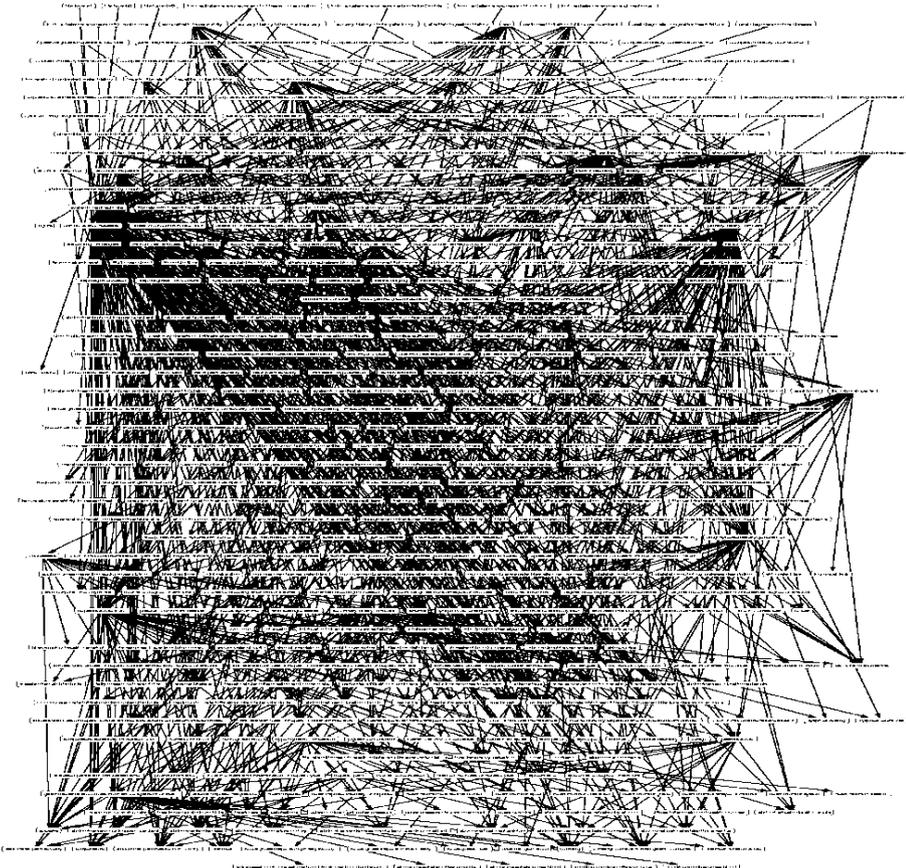
# Nicht-hierarchisches Clustering

- Parameter  $\epsilon \geq 0$
- Algorithmus:
  - Start: wähle beliebiges Objekt  $B_i$  und fasse alle Objekte in der  $\epsilon$ -Umgebung von  $B_i$  zur Klasse  $K_j$  hinzu
  - Wiederholung des Hinzufügens von Objekten in der  $\epsilon$ -Umgebung von allen bisherigen Objekten der Klasse  $K_j$
  - Wenn kein neues Element hinzufügbare, Wiederholen des gesamten Verfahrens für ein bisher noch nicht klassifiziertes Objekt  $B_k$  und neue Klasse  $K_{j+1}$
- Verfahren abhängig von geeigneter Wahl von  $\epsilon$

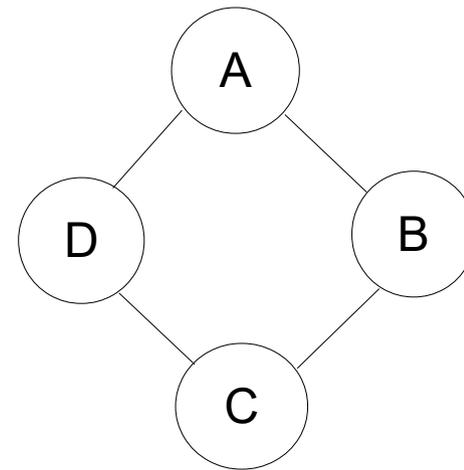
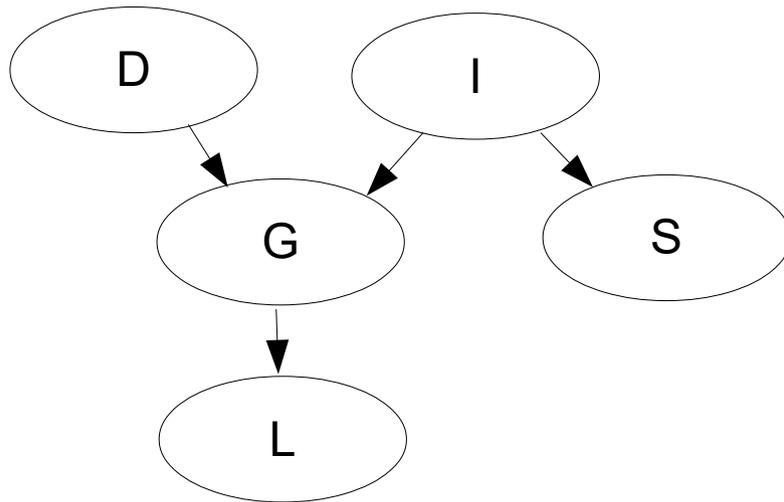
# Probabilistic Graphical Models (PGM)

- Kombination aus
  - Wahrscheinlichkeitstheorie
  - Graphentheorie
- Stellt allgemeine Frameworks zur Verfügung, subsumiert u.a.:
  - Kalman filters
  - Hidden Markov models
  - Ising models
  - ...

# Probabilistic Graphical Models (PGM)



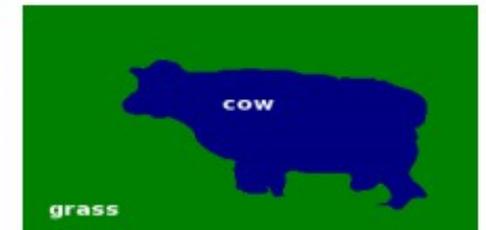
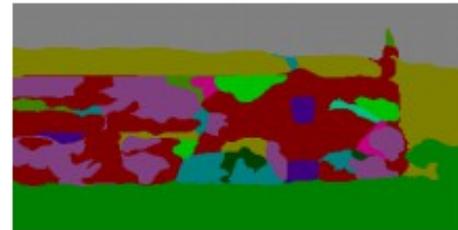
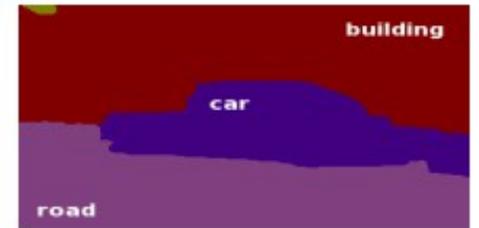
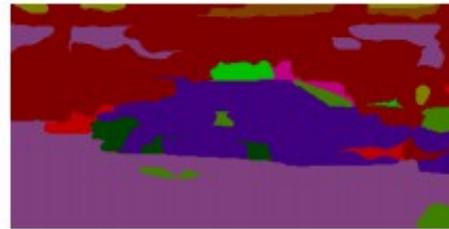
# Probabilistic Graphical Models (PGM)



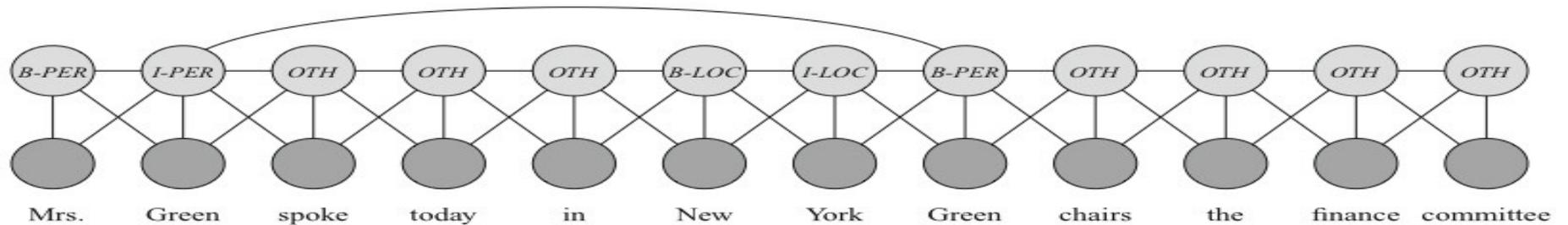
# Anwendungsbereiche

- Medizinische Diagnose
- Fehleranalyse
- Sprachverarbeitung
- Verkehrsanalyse
- Soziale Netzwerkanalyse
- Nachrichtendekodierung
- Computer Vision
- Bildsegmentation
- 3D Rekonstruktion
- Szenenanalyse
- Spracherkennung
- Roboter Lokalisierung
- ...

# Beispiel: Bildsegmentierung



# Beispiel: Sprachverarbeitung

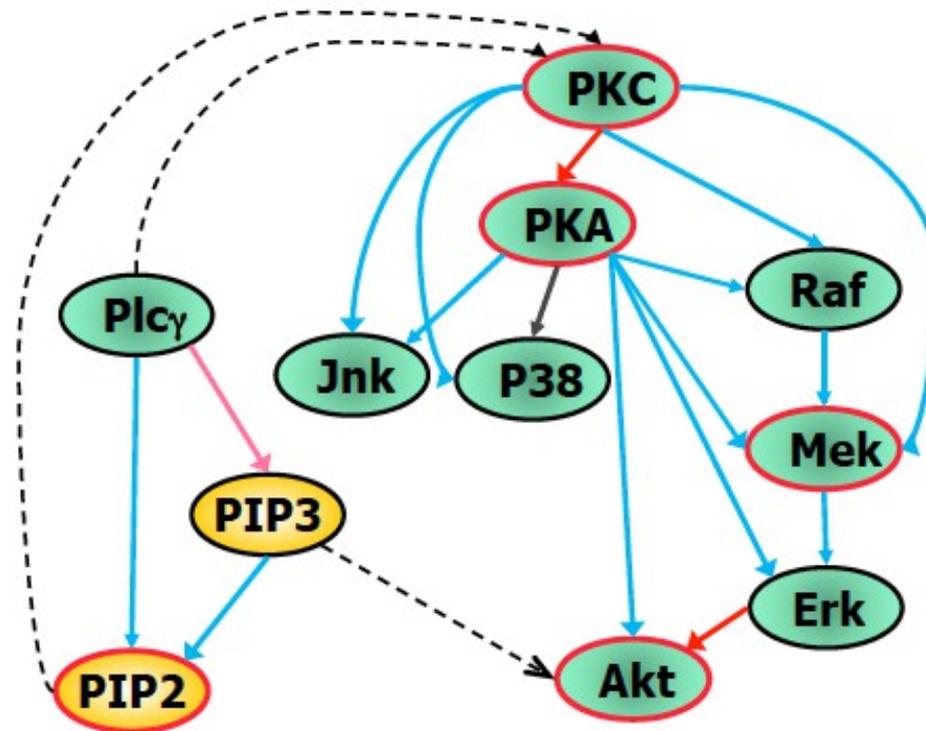


## KEY

<i>B-PER</i>	Begin person name	<i>I-LOC</i>	Within location name
<i>I-PER</i>	Within person name	<i>OTH</i>	Not an entity
<i>B-LOC</i>	Begin location name		

(a)

# Beispiel: Protein-Signaling Network



# Themen bei PGMs

- Repräsentation
  - Gerichtete, ungerichtete Modelle
  - Temporale Modelle
  - „Plate“ Modelle
- Inferenz, Schlussfolgern
  - Exakt vs. Approximativ
  - Entscheidungsfindung
- Lernen
  - Parameter vs. Struktur
  - Vollständige vs. Unvollständige Daten

# 2 Arten von PGMs

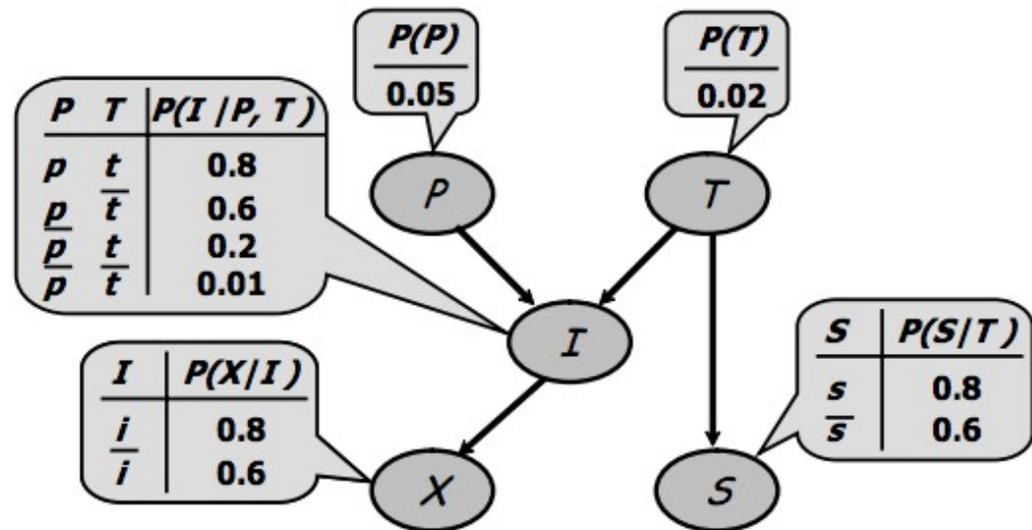
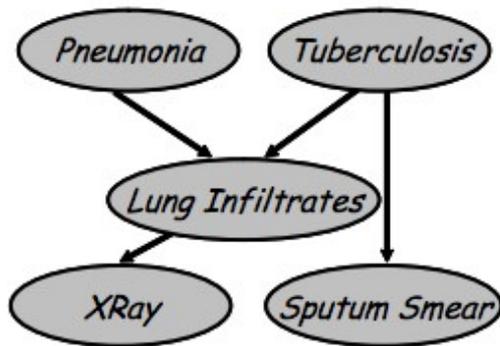
- Bayesian networks  
(auch Belief networks oder kausale Netzwerke genannt)
  - Gerichtete graphische Modelle
- Markov networks  
(auch Markov random fields (MRFs) genannt).
  - Ungerichtete graphische Modelle

# Bayes Network

Gerichteter azyklischer Graph (directed acyclic graph – DAG)

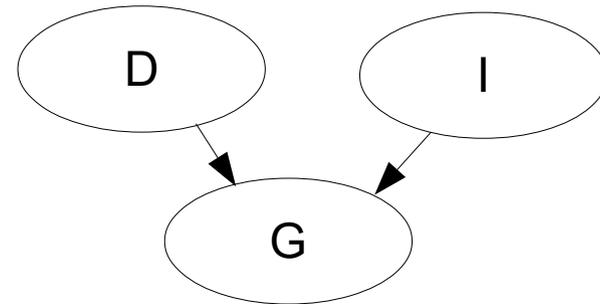
- Knoten von  $G$  sind Zufallsvariablen  $X_i$
- Kanten: direkter Einfluss von einem Knoten auf den anderen
- Jeder Zufallsvariable  $X_i$  ist eine bedingte Wahrscheinlichkeitsverteilung (conditional probability distribution – CPD).
- Die CPD für  $X_i$  gegeben die Elternknoten ( $Pa_{X_i}$ ) ist  $P(X_i | Pa_{X_i})$ .

# Bayes Network mit CPDs



# Wahrscheinlichkeitsverteilung

I	D	G	P(I,D,G)
i0	d0	g1	0.126
i0	d0	g2	0.168
i0	d0	g3	0.126
i0	d1	g1	0.009
i0	d1	g2	0.045
i0	d1	g3	0.126
i1	d0	g1	0.252
i1	d0	g2	0.0224
i1	d0	g3	0.0056
i1	d1	g1	0.06
i1	d1	g2	0.036
i1	d1	g3	0.024



# Konditionierung auf g1

$P(I,D,g1)$

I	D	G	$P(I,D,G)$
i0	d0	g1	0.126
i0	d1	g1	0.009
i1	d0	g1	0.252
i1	d1	g1	0.06

# Reduktion und Normalisierung

$P(I,D,g1)$

→

$P(I,D|g1)$

I	D	G	$P(I,D,g1)$
i0	d0	g1	0.126
i0	d1	g1	0.009
i1	d0	g1	0.252
i1	d1	g1	0.06
			<b>0.447</b>

I	D	$P(I,D g1)$
i0	d0	0.282
i0	d1	0.02
i1	d0	0.564
i1	d1	0.134
		<b>1.000</b>

# Marginalisierung

$P(I,D)$

→

$P(I)$

I	D	P(I,D)
i0	d0	0.282
i0	d1	0.02
i1	d0	0.564
i1	d1	0.134
		<b>1.000</b>

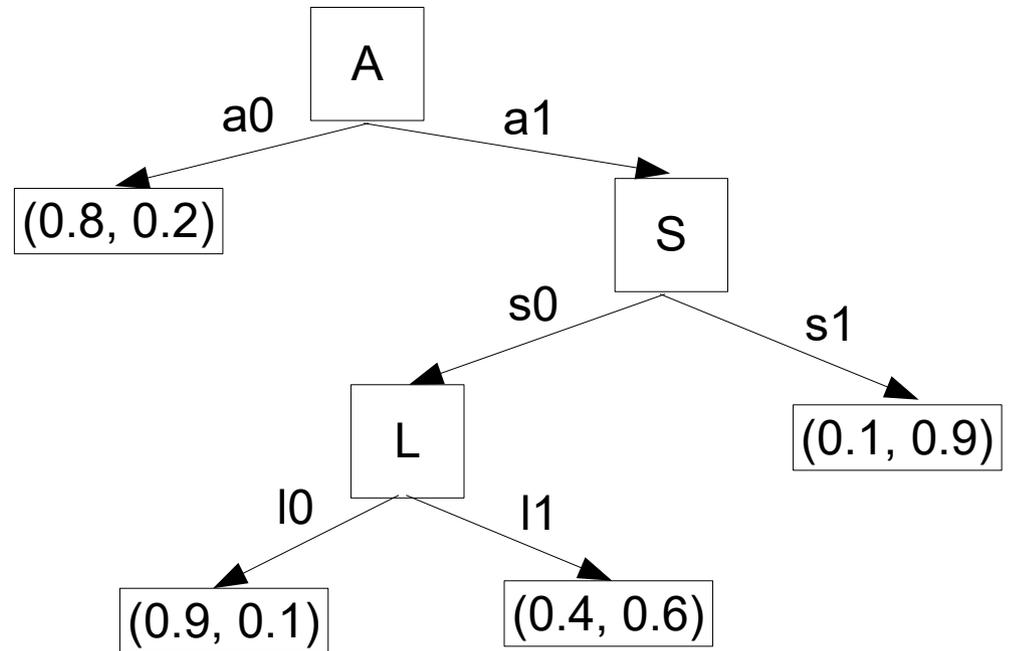
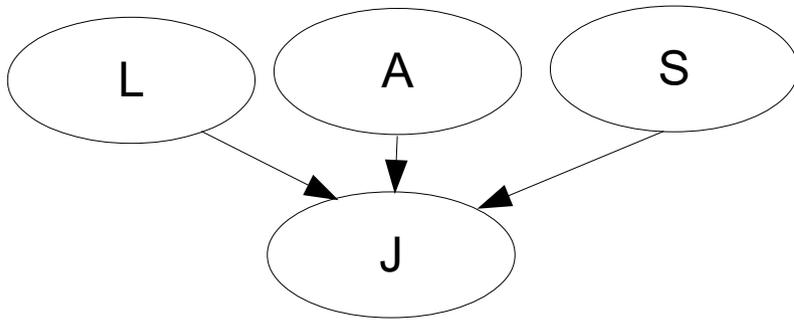
I	P(I)
i0	0.282+0.02 = 0.302
i1	0.564+0.134 = 0.698
	<b>1.000</b>

# Darstellung von CPDs

- Tabellarische CPDs
- Baum-strukturierte CPDs
- Lokale Wahrscheinlichkeitsmodelle
- Logistische CPDs & Generalisierungen
- Noisy OR / AND
- Lineare Gauss Verteilungen & Generalisierungen
- ...

# Baum CPD

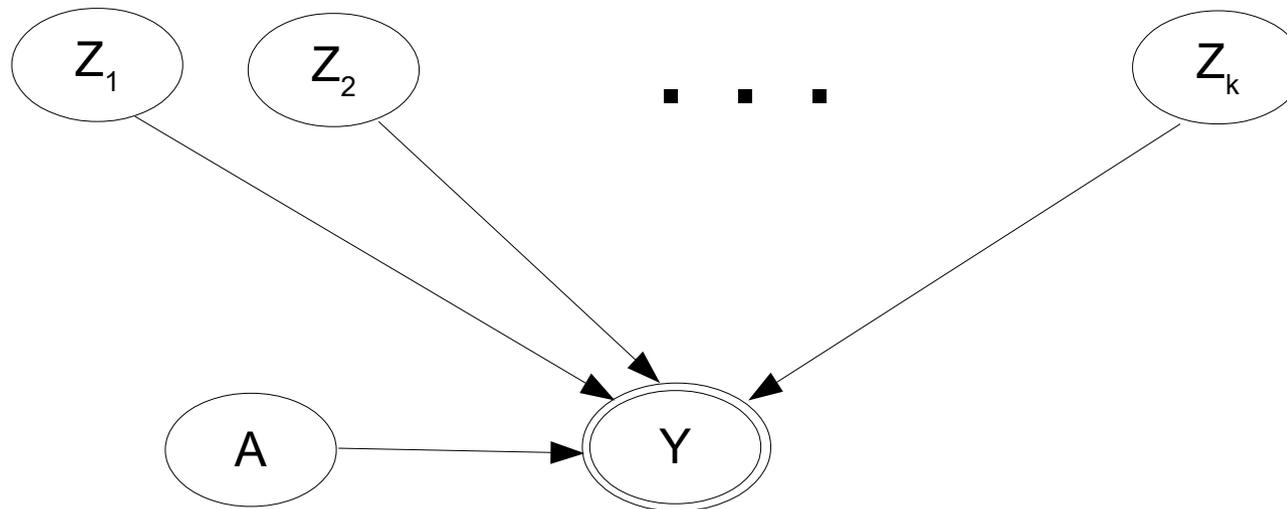
$P(J | L, A, S)$



# Multiplexer CPD

$$A = \{1, k\}$$

$$A = a \Rightarrow Y = Z_a$$



$$P(Y | A, Z_1, \dots, Z_k) = \begin{cases} 1 & \dots Y = Z \\ 0 & \dots \text{ansonsten} \end{cases}$$

0 ..... ansonsten

# Faktoren

Faktor  $\Phi(X_1, \dots, X_n)$

Funktion  $\Phi : \text{Val}(X_1, \dots, X_n) \rightarrow \mathbb{R}$

Scope =  $\{X_1, \dots, X_n\}$

- Jede JPD und CPD ist ein Faktor
- Unnormalisierte Funktion  $P(I, D, g1)$  ist Faktor über  $\{I, D\}$

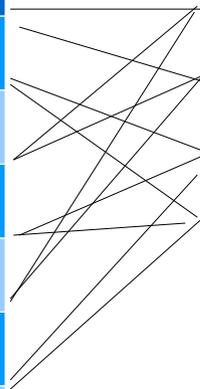
I	D		$P(I, D, g1)$
i0	d0	g1	0.126
i0	d1	g1	0.009
i1	d0	g1	0.252
i1	d1	g1	0.06
			<b>0.447</b>

# Allgemeine Faktoren

A	B	$\phi$
a0	b0	30
a0	b1	5
a1	b0	1
a1	b1	30

# Faktor Produkt

A	B	$\Phi(A,B)$
a0	b0	0.5
a0	b1	0.8
a1	b0	0.1
a1	b1	0
a2	b0	0.3
a2	b1	0.9



B	C	$\Phi(B,C)$
b0	c0	0.5
b0	c1	0.7
b1	c0	0.1
b1	c1	0.2



A	B	C	$\Phi(A,B,C)$
a0	b0	c0	$0.5 \cdot 0.5 = 0.25$
a0	b0	c1	$0.5 \cdot 0.7 = 0.35$
a0	b1	c0	$0.8 \cdot 0.1 = 0.08$
a0	b1	c1	$0.8 \cdot 0.2 = 0.16$
a1	b0	c0	$0.1 \cdot 0.5 = 0.05$
a1	b0	c1	$0.1 \cdot 0.7 = 0.07$
a1	b1	c0	$0 \cdot 0.1 = 0$
a1	b1	c1	$0 \cdot 0.2 = 0$
a2	b0	c0	$0.3 \cdot 0.5 = 0.15$
a2	b0	c1	$0.3 \cdot 0.7 = 0.21$
a2	b1	c0	$0.9 \cdot 0.1 = 0.09$
a2	b1	c1	$0.9 \cdot 0.2 = 0.18$

# Faktor Marginalisierung

A	B	C	$\Phi(A,B,C)$
a0	b0	c0	0.25
a0	b0	c1	0.35
a0	b1	c0	0.08
a0	b1	c1	0.16
a1	b0	c0	0.05
a1	b0	c1	0.07
a1	b1	c0	0
a1	b1	c1	0
a2	b0	c0	0.15
a2	b0	c1	0.21
a2	b1	c0	0.09
a2	b1	c1	0.18

A	C	$\Phi(A,C)$
a0	c0	0.33
a0	c1	0.51
a1	c0	0.05
a1	c1	0.07
a2	c0	0.3
a2	c1	0.9

# Faktor Reduktion

A	B	C	$\Phi(A,B,C)$
a0	b0	c0	0.25
a0	b0	c1	0.35
a0	b1	c0	0.08
a0	b1	c1	0.16
a1	b0	c0	0.05
a1	b0	c1	0.07
a1	b1	c0	0
a1	b1	c1	0
a2	b0	c0	0.15
a2	b0	c1	0.21
a2	b1	c0	0.09
a2	b1	c1	0.18



A	B	C	$\Phi(A,B)$
a0	b0	c0	0.25
a0	b1	c0	0.08
a1	b0	c0	0.05
a1	b1	c0	0
a2	b0	c0	0.15
a2	b1	c0	0.09

# Warum Faktoren?

- Grundlegende Bestandteile zur Definition von Wahrscheinlichkeitsverteilungen in hoch-dimensionalen Räumen
- Menge von grundlegenden Operationen zur Manipulation der Wahrscheinlichkeitsverteilungen

# Unabhängigkeit

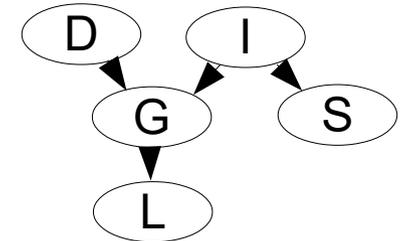
- Unabhängigkeit

$$P(X = x, Y = y) = P(X = x)P(Y = y)$$

- Bedingte Unabhängigkeit

$$P(X = x, Y = y | Z = z) = P(X = x | Z = z)P(Y = y | Z = z)$$

# Faktorisierung



- Kettenregel

$$P(L, G, S, D, I) = P(L|G, S, D, I) * P(G|S, D, I) * P(S|D, I) * P(D|I) * P(I)$$

- Kettenregel für Bayes'sche Netze

$$P(L, G, S, D, I) = P(L|G) * P(G|D, I) * P(S | I) * P(D) * P(I)$$

$$P(X_1, \dots, X_n) = \prod_i P(X_i | \text{Par}_G(X_i))$$

→ Distribution = Product of Factors

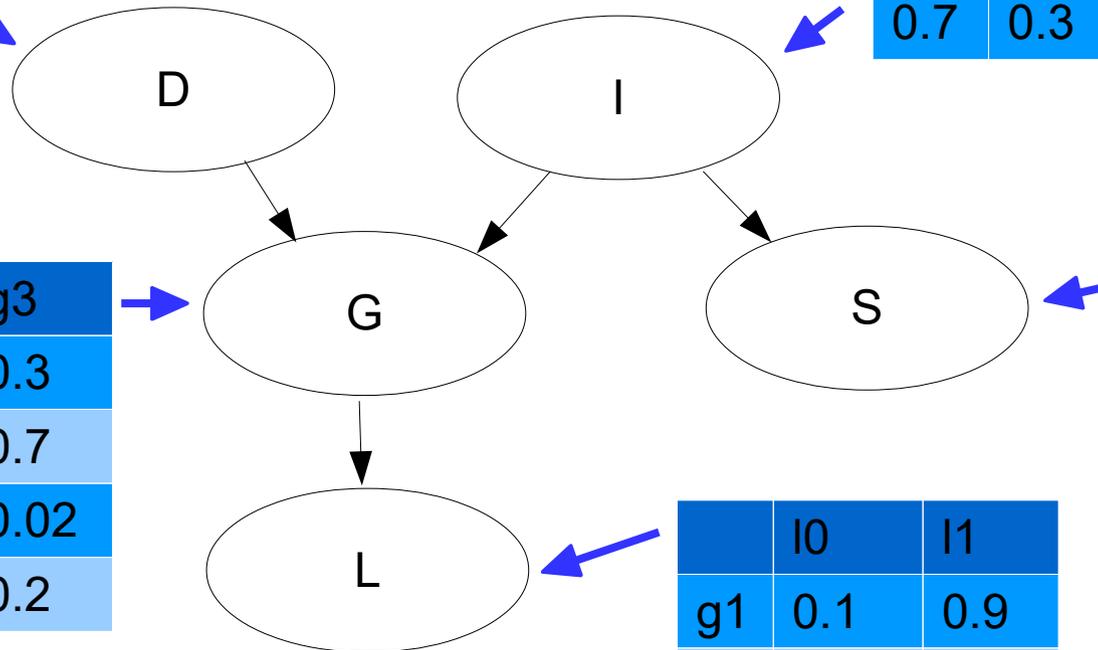
- Bedingte Unabhängigkeitsannahme,  
lokale Markov Annahme:  $(X_i \perp \text{NonDescendants}X_i | \text{Par}_G(X_i))$

- P faktorisierung über G, wenn  
 $P(X_1, \dots, X_n) = \prod_i P(X_i | \text{Par}_G(X_i))$

# Schlussfolgern in BN

d0	d1
0.6	0.4

i0	i1
0.7	0.3



	s0	s1
i0	0.95	0.05
i1	0.2	0.8

	g1	g2	g3
i0,d0	0.3	0.4	0.3
i0,d1	0.05	0.25	0.7
i1,d0	0.9	0.08	0.02
i1,d1	0.5	0.3	0.2

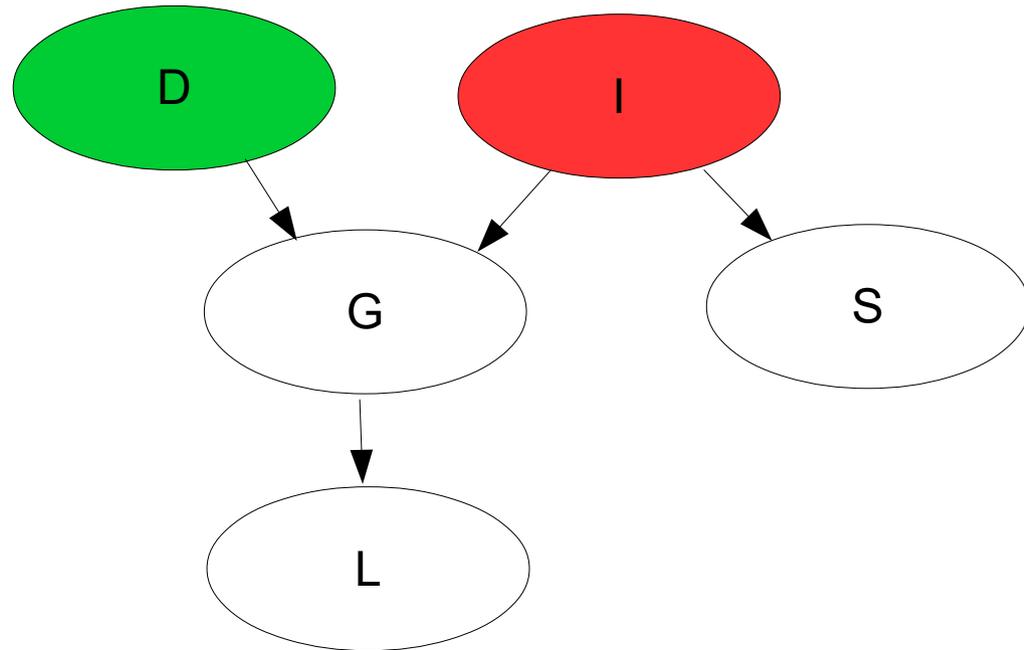
	l0	l1
g1	0.1	0.9
g2	0.4	0.6
g3	0.99	0.01

# Kausalität

$$P(I1) \approx 0.5$$

$$P(I1 \mid i0) \approx 0.39$$

$$P(I1 \mid i0, d0) \approx 0.51$$



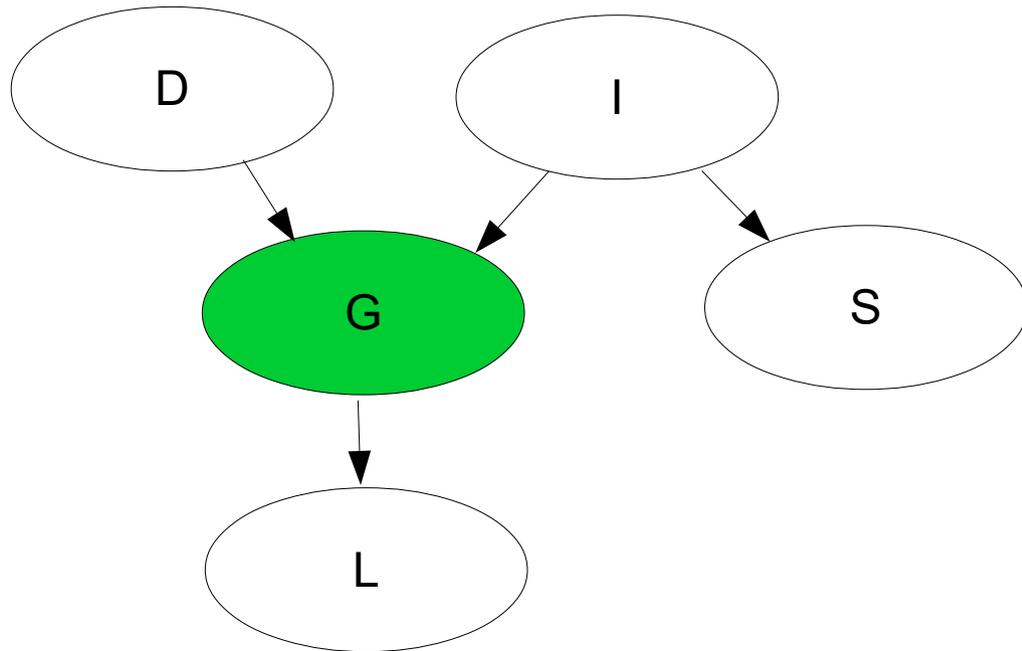
# Evidenz

$$P(d1) = 0.4$$

$$P(d1 \mid g3) \approx 0.63$$

$$P(i1) = 0.3$$

$$P(i1 \mid g3) \approx 0.08$$



# Interkausalität

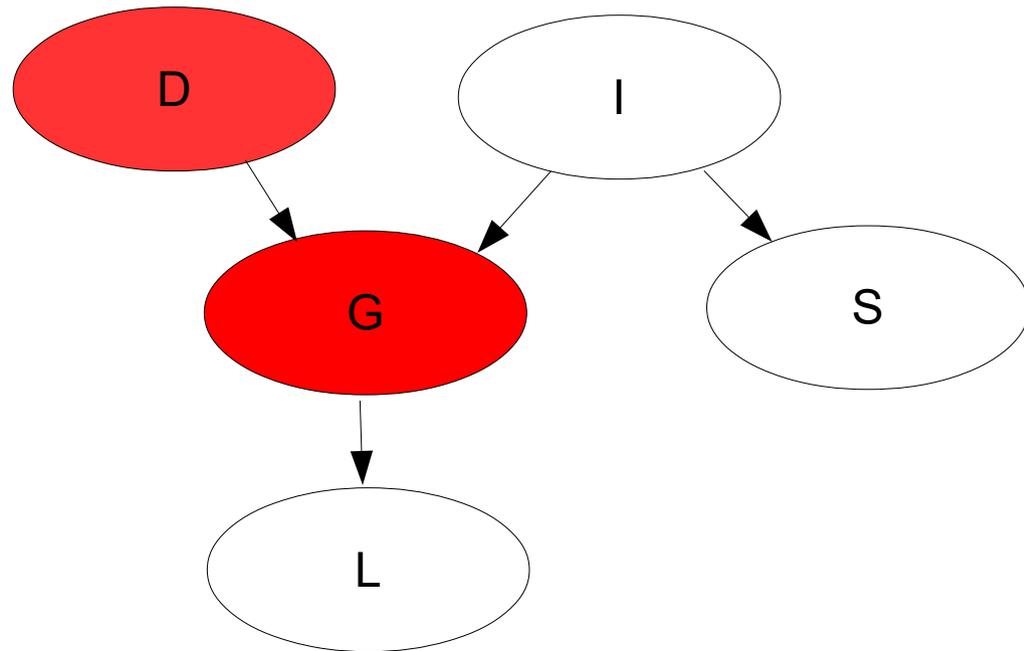
$$P(d1)=0.4$$

$$P(d1|g3)\approx 0.63$$

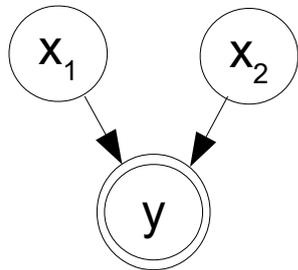
$$P(i1)=0.3$$

$$P(i1|g3)\approx 0.08$$

$$P(i1|g3,d1)\approx 0.11$$



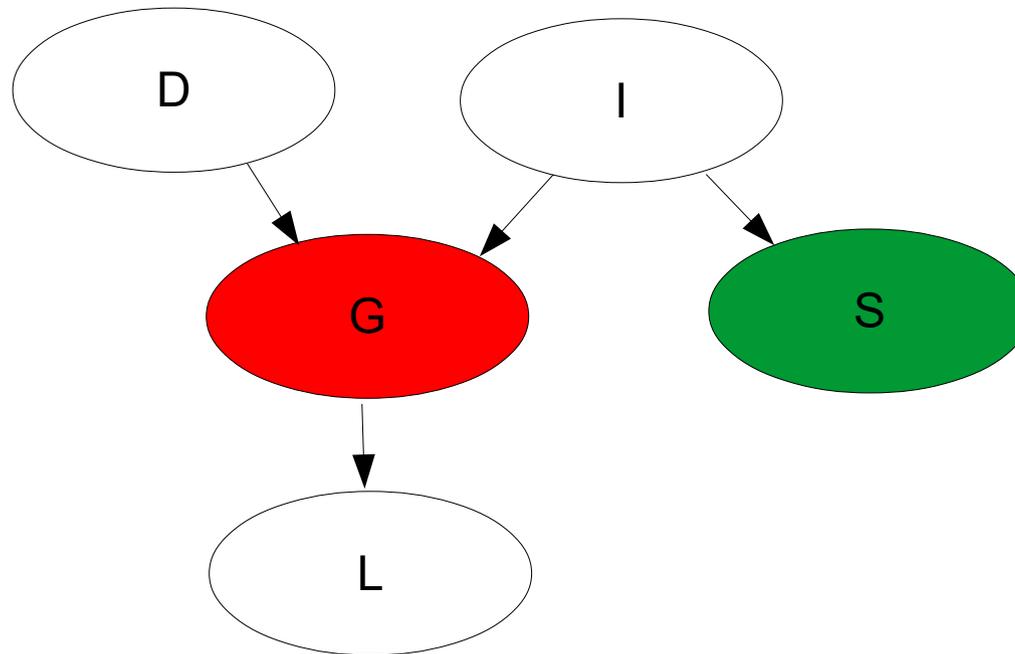
# Explaining away



$P(x_2 \mid y=1, x_1=1)$

$x_1$	$x_2$	$y$	$P$
0	0	0	0.25
0	1	1	0.25
1	0	1	0.25
1	1	1	0.25

# Interkausalität (2)



# Flow of Influence

Verändert die Konditionierung auf X  
die Wahrscheinlichkeit von Y?

- $X \rightarrow Y$
- $X \leftarrow Y$
- $X \rightarrow W \rightarrow Y$  ja
- $X \leftarrow W \leftarrow Y$  ja
- $X \leftarrow W \rightarrow Y$  ja
- $X \rightarrow W \leftarrow Y$  nein  
(v-structure)

# Flow of Influence

Verändert die Konditionierung auf X  
die Wahrscheinlichkeit von Y, gegeben Evidenz Z?

Evidence Z	W nicht in Z	W in Z
• $X \rightarrow W \rightarrow Y$	ja	nein
• $X \leftarrow W \leftarrow Y$	ja	nein
• $X \leftarrow W \rightarrow Y$	ja	nein
• $X \rightarrow W \leftarrow Y$ (v-structure)	?	?

- Ja, wenn W oder ein Nachkomme von W in Z  
- Sonst nein

# d-Separierbarkeit

d (=directed)-Separierbarkeit. Definition von aktiven Pfaden:

- Kausaler Pfad  $X \rightarrow Z \rightarrow Y$ : aktiv  $\Leftrightarrow Z$  nicht beobachtet
- Evidenz Pfad  $X \leftarrow Z \leftarrow Y$ : aktiv  $\Leftrightarrow Z$  nicht beobachtet
- Gemeinsamer Ursache  $X \leftarrow Z \rightarrow Y$ : aktiv  $\Leftrightarrow Z$  nicht beobachtet
- Gemeinsamer Effekt  $X \rightarrow Z \leftarrow Y$ : aktiv  $\Leftrightarrow Z$  oder einer seiner Nachkommen beobachtet

# d-Separierbarkeit

Ein Pfad  $X_1 \text{---} \dots \text{---} X_n$  ist aktiv gegeben Evidenz  $E$

- Wenn für jede v-Struktur  $X_{i-1} \rightarrow X_i \leftarrow X_{i+1}$  entweder  $X_i$  oder ein Nachkomme von  $X_i$  in  $E$  ist;
- Kein anderer Knoten entlang des Pfades in  $E$  ist.

Für drei Mengen  $X, Y, Z$  von Knoten aus  $G$  gilt.  $X$  und  $Y$  sind d-separiert gegeben  $Z$   $d\text{-sep}_G(X; Y \mid Z)$ , wenn es keinen aktiven Pfaden zwischen zwei Knoten  $x \in X$  und  $y \in Y$  gegeben  $Z$  gibt.

# Zusammenfassung: Bayesian Network

Zwei äquivalente Betrachtungsweisen der Graphenstruktur

- Faktorisierung:  $G$  ermöglicht die Representation von  $P$
- I-map: In  $G$  kodierte Unabhängigkeiten halten in  $P$

Wenn  $P$  über  $G$  faktorisiert, dann können aus dem Graphen  $G$  die Unabhängigkeiten extrahiert werden, die in  $P$  halten müssen (Independence Map).

# Ausblick

- Nächster Termin:

**Donnerstag, 15.12.2016 13.15-13.45 (s.t.)**

Strukturelle, syntaktische Mustererkennung