

# MASCHINELLES LERNEN & DATAMINING

Vorlesung im Wintersemester 2017

Prof. E.G. Schukat-Talamazzini

Stand: 25. August 2017

## Teil VI

## Attributabhängigkeiten: graphische & kausale Modelle

### Analyse von Attributabhängigkeiten

Dependenzanalyse  $\stackrel{?}{=}$  Spaltengruppierung

	$\mathcal{X}_1$	$\mathcal{X}_2$	$\mathcal{X}_3$	$\mathcal{X}_4$
$o_1$	+	<i>low</i>	1.2	+4
$o_2$	−	<i>hi</i>	0.5	−3
$o_3$	+	<i>hi</i>	2.3	−3
$o_4$	+	<i>med</i>	2.1	−7

### Abhängigkeit $\neq$ Ähnlichkeit

- Lineare Abhängigkeiten  
 $E = m \cdot c^2$
- Skalenempfindlichkeit  
Temperatur in  $^{\circ}\text{C}$  oder  $^{\circ}\text{K}$
- Skalenübergreifend  
 $\mathcal{X}_i$  Geschlecht,  $\mathcal{X}_j$  Gehalt

### Struktur $\neq$ Partition

- keine Äquivalenzrelation  
Zeitreihen, Ortsgitter
- keine binäre Relation  
Alter, Geschlecht, Größe
- Kausalitätsrichtung ?  
 $\mathcal{X}_n$  Niederschlag,  $\mathcal{X}_m$  Ertrag

### Analyse von Attributabhängigkeiten

Mit welchem Ziel — zu welchem Zweck ?

#### Attributwerteprädiktion

Starke Abhängigkeit  $\Rightarrow \left\{ \begin{array}{l} \text{geringe a posteriori Streuung} \\ \text{kleiner Regressionsfehler} \end{array} \right\}$

- Voraussage · Imputation · Klassifikation

#### Strukturaufklärung

Lernen des am einfachsten strukturierten Datenmodells (*Occams Razor*)

- Interaktionen · Kausalitäten · Assoziationsregeln

#### Robuste Datenmodelle

Netzwerk ausgewählter Abhängigkeiten statt **saturiert**er W-Modelle

- geringe Kapazität · hohe Effizienz (Zeit/Speicher) · gute Induktivität

## Korrelation, Regression und Transinformation

Assoziationsregeln und Netzwerkanalyse

Bedingte statistische Unabhängigkeit

Graphische Modelle: ungerichtete Graphen

Kausale Modelle: gerichtete azyklische Graphen

Berechnen bedingter Wahrscheinlichkeiten

Parameterschätzung in Bayesnetzen und Loglinearmodellen

Aufdeckung der Abhängigkeitsstruktur

Kovarianzselektion

## Statistische Unabhängigkeit von Zufallsvariablen

### Statistische Unabhängigkeit

von Zufallsvariablen  $\mathbb{X}_1, \dots, \mathbb{X}_N$ , falls für alle  $x_1, \dots, x_N \in \mathbb{R}^N$  gilt:

$$P(\mathbb{X}_1 = x_1, \dots, \mathbb{X}_N = x_N) = \prod_{n=1}^N P(\mathbb{X}_n = x_n)$$

### Statistische Unkorreliertheit

von Zufallsvariablen  $\mathbb{X}_1, \dots, \mathbb{X}_N$ , falls für alle  $x_1, \dots, x_N \in \mathbb{R}^N$  gilt:

$$\mathcal{E}[\prod_{n=1}^N \mathbb{X}_n] = \prod_{n=1}^N \mathcal{E}[\mathbb{X}_n]$$

#### Bemerkungen

1. Aus der Unabhängigkeit folgt die Unkorreliertheit, aber nicht umgekehrt.
2. Für normalverteilte  $\mathbb{X} \sim \mathcal{N}(\mu, \mathbf{S})$  gilt:  $\mathbb{X}_i, \mathbb{X}_j$  korreliert gdw.  $\sigma_{ij} \neq 0$ .

## Statistische Unabhängigkeit von Ereignissen

### Paarweise statistische Unabhängigkeit

**Faktorisierbarkeit** oder (falls  $P(A) \neq 0$ ) **Neutralität**:

$$A \not\sim B \iff P(A, B) = P(A) \cdot P(B) \iff P(B|A) = P(B)$$

Beispiel: der Wurf zweier fairer Würfel

$$\begin{aligned} A &= \text{„gerade Augensumme“} & P(A, B) &= \frac{1}{12} = \frac{1}{2} \cdot \frac{1}{6} = P(A) \cdot P(B) \\ B &= \text{„erster Wurf ist sechs“} & P(B, C) &= \frac{1}{36} = \frac{1}{6} \cdot \frac{1}{6} = P(B) \cdot P(C) \\ C &= \text{„Augensumme ist sieben“} & P(A, C) &= 0 \neq \frac{1}{2} \cdot \frac{1}{6} = P(A) \cdot P(C) \end{aligned}$$

### Statistische Unabhängigkeit

der Ereignisse  $A_1, \dots, A_I$ , falls für alle Indexmengen  $\mathcal{I} \subseteq \{1, \dots, I\}$ :

$$P(\bigwedge_{i \in \mathcal{I}} A_i) = \prod_{i \in \mathcal{I}} P(A_i)$$

Stat. Unabhäng.  $\Rightarrow$  paarweise s.U.  
Stat. Unabhäng.  $\nRightarrow$  paarweise s.U.

$A$  = „erster Wurf hat gerade Augenzahl“  
 $B$  = „zweiter Wurf hat gerade Augenzahl“  
 $C$  = „Augensumme ist ungerade“

$$P(A, B, C) = 0 \neq \frac{1}{8} = P(A) \cdot P(B) \cdot P(C)$$

### Beweis.

1. Die uniforme Verteilungsdichte auf dem Träger

$$\{(x, y) \mid |x| + |y| \leq 1\} \subseteq \mathbb{R}^2$$

ist wegen

$$\mathcal{E}[\mathbb{X}\mathbb{Y}] = 0 = \mathcal{E}[\mathbb{X}] \cdot \mathcal{E}[\mathbb{Y}]$$

zwar unkorreliert, aus ihrer (hypothetischen) Unabhängigkeit folgt aber wegen

$$P(0, \cdot) \cdot P(\cdot, 0) = P(0, 0) = P(0, 1) = P(0, \cdot) \cdot P(\cdot, 1)$$

und  $P(0, \cdot) \neq 0$  sofort der Widerspruch  $P(\cdot, 0) = P(\cdot, 1)$ .

2. Es gilt nach Kovarianzdefinition

$$\sigma_{ij} = \text{Cov}[\mathbb{X}_i, \mathbb{X}_j] = \mathcal{E}[\mathbb{X}_i \mathbb{X}_j] - \mathcal{E}[\mathbb{X}_i] \cdot \mathcal{E}[\mathbb{X}_j] ;$$

daraus folgt die Behauptung — auch für nicht-normal verteilte Variablen.



## Korrelation und Kovarianz

### Definition

Es sei  $\omega \subset \mathbb{R}^N$  ein Datensatz mit der (empirischen) Kovarianzmatrix  $S = [\sigma_{ij}]$ . Die Zahlen

$$\rho_{ij} \stackrel{\text{def}}{=} \frac{\sigma_{ij}}{\sqrt{\sigma_{ii}} \cdot \sqrt{\sigma_{jj}}}, \quad 1 \leq i, j \leq N$$

heißen Pearsonsche **Korrelationskoeffizienten** der Attributpaare  $(x_i, x_j)$ .

### Bemerkungen

1. Betragsmäßig kleine (große) Werte  $\sigma_{ij}$  markieren einen schwachen (starken) Zusammenhang zwischen  $x_i$  und  $x_j$ .
2. Die **Kovarianzen** sind aber extrem skalierungsempfindlich ( $\sigma_{ii}, \sigma_{jj}$ ).
3. Die Korrelationskoeffizienten  $\rho_{ij}$  liegen stets im Intervall  $[-1, +1]$ .
4. Der Wert  $\rho_{ij} = 0$  markiert Unkorreliertheit, die Werte  $\rho_{ij} \in \{+1, -1\}$  hingegen **deterministische Abhängigkeit** (mit positiver/negativer Steigung).

## Korrelationsgruppierung

(Algorithmus)

GEGEBEN:

Daten  $\omega \subset \mathbb{R}^N$ , Schwelle  $\theta_\rho$ , „leere“ Gruppierung  $\gamma : i \mapsto \perp$ .

### 1 INITIALISIERUNG

Berechne alle Korrelationskoeffizienten  $\rho_{ij}$ .

### 2 ABSTEIGEND SORTIEREN

$$|\rho_{i_1 j_1}| \geq |\rho_{i_2 j_2}| \geq |\rho_{i_3 j_3}| \geq |\rho_{i_4 j_4}| \geq \dots \geq \dots \geq$$

### 3 FÜR ALLE $r = 1, 2, \dots, N(N-1)/2$ :

1. Wenn  $|\rho_{i_r j_r}| < \theta_\rho$  dann  $\rightsquigarrow$  ENDE.
2. Wenn  $\gamma(n) \neq \perp$  für alle  $n$  dann  $\rightsquigarrow$  ENDE.
3. Wenn  $\gamma(i_r) = \perp = \gamma(j_r)$  dann erzeuge neue Gruppe  $\{i_r, j_r\}$ .
4. Wenn  $\gamma(i_r) = \perp$  dann setze  $\gamma(i_r) \leftarrow \gamma(j_r)$ .
5. Wenn  $\gamma(j_r) = \perp$  dann setze  $\gamma(j_r) \leftarrow \gamma(i_r)$ .
6. Wenn  $\gamma(i_r) \neq \gamma(j_r)$  dann vereinige die Gruppen:  $\gamma(i_r) \cup \gamma(j_r)$ .

(zum3hlog(A))

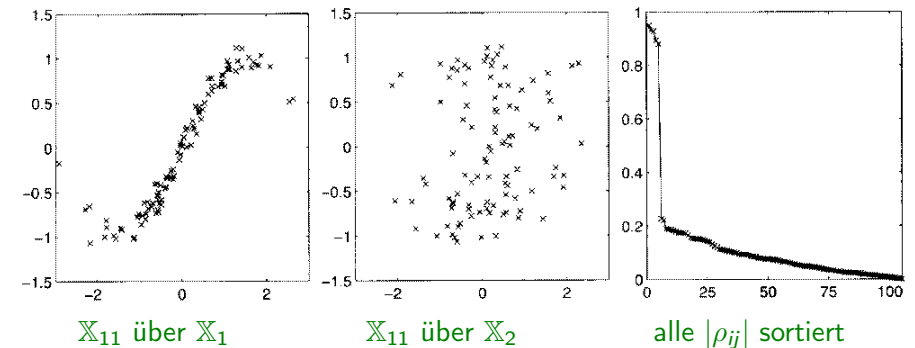
## Beispiel — Korrelationsanalyse synthetischer Daten

### Zufällig generierte Datenvektoren

$\omega = \{\mathbf{x}_1, \dots, \mathbf{x}_{100}\} \subset \mathbb{R}^{15}$  mit Wertetupeln der Zufallsvariablen

$$\mathbb{X}_n = \begin{cases} \mathcal{N}(0, 1) & n = 1, \dots, 10 \\ \sin(\mathbb{X}_{n-10}) + \mathcal{N}(0, 1/10) & n = 11, \dots, 15 \end{cases}$$

(10 Kanäle weißes Rauschen & 5 Kanäle verrauschte Sinuskopien)



## Korrelationsgruppierung

Was tut dieser Algorithmus ?

### Single-linkage Agglomeration — aber:

Terminiert bei Unterschreiten der Korrelationsschwelle.

Terminiert sobald alle Einermengen „verbraucht“ sind.

### Synthesedatenbeispiel

Für die Daten  $\omega \subset \mathbb{R}^{15}$  werden in den ersten fünf Schritten die Gruppen

$$\{1, 11\}, \{2, 12\}, \{3, 13\}, \{4, 14\}, \{5, 15\}$$

gebildet; anschließend gibt es jeweils drei gleichwahrscheinliche Möglichkeiten:

1. Eine der „alten“ Gruppen wird mit einem neuen Index aufgefüllt.
2. Zwei „alte“ Gruppen werden vereinigt.
3. Aus zwei „frischen“ Indizes wird eine neue Gruppe gebildet.

Mit der Ausnahme von 1. sind all diese Optionen höchst **unerwünscht**.

## Gestörte (lineare) Abhängigkeit

$\mathbb{Y} = f(\mathbb{X}) + \mathbb{E}$  mit Funktionsprototyp  $f : \mathbb{R} \rightarrow \mathbb{R}$  und Residuum  $\mathbb{E}$

### Lemma

Für zwei normalverteilte Zufallsvariablen  $\mathbb{X}, \mathbb{Y}$  mit

$$\mathbb{Y} = a\mathbb{X} + b + \mathbb{E}, \quad \mathbb{X} \sim \mathcal{N}(\mu_x, \sigma_x^2), \quad \mathbb{E} \sim \mathcal{N}(0, \sigma_e^2)$$

gehört  $\mathbb{Y}$  der Verteilungsaussage

$$\mathbb{Y} \sim \mathcal{N}(a\mu_x + b, a^2\sigma_x^2 + \sigma_e^2).$$

Die **Kovarianz** und die **Korrelation** zwischen  $\mathbb{X}$  und  $\mathbb{Y}$  betragen

$$\sigma_{xy} = a \cdot \sigma_x^2 \quad \text{bzw.} \quad \rho_{xy} = \text{sign}(a) \cdot \left(1 + \frac{\sigma_e^2}{a^2 \cdot \sigma_x^2}\right)^{-\frac{1}{2}}.$$

### Bemerkung

Die Korrelation  $\rho_{xy}$  erbt das Vorzeichen von  $a$ .

Der Betrag wächst und fällt mit  $\sigma_e^{-2}$  im Einheitsintervall.

## Beweis.

Berechnung der Kovarianz (o.B.d.A. ist  $\mu_x = 0$ ):

$$\begin{aligned} \sigma_{xy} &= \text{Cov}[\mathbb{X}, \mathbb{Y}] = \mathcal{E}[\mathbb{X}\mathbb{Y}] - \mu_x\mu_y \\ &= \mathcal{E}[a\mathbb{X}^2 + b\mathbb{X} + \mathbb{E}\mathbb{X}] - \mu_x\mu_y \\ &= a \cdot (\sigma_x^2 + \mu_x^2) + b\mu_x + \mu_e\mu_x - \mu_x\mu_y \\ &= a\sigma_x^2 + a\mu_x^2 + b\mu_x + \mu_e\mu_x - \mu_x\mu_y \\ &= a\sigma_x^2 \end{aligned}$$

Berechnung der Korrelation:

$$\begin{aligned} \rho_{xy} &= \frac{a\sigma_x^2}{\sqrt{\sigma_x^2 \cdot (a^2\sigma_x^2 + \sigma_e^2)}} \\ &= \text{sign}(a) \cdot \frac{a\sigma_x^2}{a\sigma_x^2 \cdot \sqrt{1 + \frac{\sigma_e^2}{a^2\sigma_x^2}}} \\ &= \text{sign}(a) \cdot \left(1 + \frac{\sigma_e^2}{a^2\sigma_x^2}\right)^{-\frac{1}{2}} \end{aligned}$$

□

## Kausalität und Scheinzusammenhang

Verursacht Diät-Cola wirklich Übergewicht ?

### Ursache und Wirkung

Korrelation und Abhängigkeit haben keine Vorzugsrichtung:

$$\left\{ \begin{array}{l} \mathbb{X}_i = \text{„Körpergewicht [kg]“} \\ \mathbb{X}_j = \text{„Konsum kalorienreduzierter Getränke [\ell]“} \end{array} \right\}$$

Hohe (positive) Korrelation(en)  $\rho_{ij} = \rho_{ji}$  ohne Hinweis auf Kausalrichtung.

### Versteckte gemeinsame Ursache oder Lederallergie ?

Das Korrelationsmaß hat keine Vorzugsrichtung:

$$\left\{ \begin{array}{l} \mathbb{X}_i = \text{„Kaffee gerührt?“} \\ \mathbb{X}_j = \text{„Kaffee schmeckt süß?“} \\ \dots \dots \dots \\ \mathbb{X}_k = \text{„Zuckerwürfel drin?“} \end{array} \right\} \quad \left\{ \begin{array}{l} \mathbb{X}_i = \text{„Lederschuhe an?“} \\ \mathbb{X}_j = \text{„Kopfschmerzen?“} \\ \dots \dots \dots \\ \mathbb{X}_k = \text{„Wagnergasse bis 3h?“} \end{array} \right\}$$

Hohe (positive) Korrelation  $\rho_{ij}$  ohne jeden (direkten) kausalen Zusammenhang.

## (Bi-) Partielle Korrelation

Vergleich nach Subtraktion der Ausgleichsgeraden

### Definition

Es seien Zufallsvariablen  $\mathbb{X}_1, \dots, \mathbb{X}_N$  gegeben; ferner bezeichne

$$\mathbb{X}_{i|k} = a_{i|k} \cdot \mathbb{X}_k + b_{i|k}, \quad (i, k \in \{1, \dots, N\}, i \neq k)$$

den linearen **Quadratmittelprädiktor** für  $\mathbb{X}_i$  aus  $\mathbb{X}_k$  („Ausgleichsgerade“).

Dann heißt

$$\rho_{ij|k} \stackrel{\text{def}}{=} \text{Corr}[\mathbb{X}_i - \mathbb{X}_{i|k}, \mathbb{X}_j - \mathbb{X}_{j|k}]$$

die **partielle Korrelation** zwischen  $\mathbb{X}_i$  und  $\mathbb{X}_j$  hinsichtlich  $\mathbb{X}_k$  und es heißt

$$\rho_{i|k,j|\ell} \stackrel{\text{def}}{=} \text{Corr}[\mathbb{X}_i - \mathbb{X}_{i|k}, \mathbb{X}_j - \mathbb{X}_{j|\ell}]$$

die **bipartielle Korrelation** zwischen  $\mathbb{X}_i$  und  $\mathbb{X}_j$  hinsichtlich  $\mathbb{X}_k$  und  $\mathbb{X}_\ell$ .

## (Bi-) Partielle Korrelation

Berechnung aus den gewöhnlichen Korrelationskoeffizienten

### Lemma

Es seien die Zufallsvariablen  $\mathbb{X}_1, \dots, \mathbb{X}_N$  und ihre Korrelationen  $\rho_{ij}$ ,  $i, j \in \{1, \dots, N\}$  gegeben.

1. Die partielle Korrelation zwischen  $\mathbb{X}_i$  und  $\mathbb{X}_j$  ohne den Einfluß von  $\mathbb{X}_k$  hat den Wert

$$\rho_{ij|k} = \frac{\rho_{ij} - \rho_{ik} \cdot \rho_{jk}}{\sqrt{(1 - \rho_{ik}^2) \cdot (1 - \rho_{jk}^2)}}.$$

2. Die bipartielle Korrelation zwischen  $\mathbb{X}_i$  und  $\mathbb{X}_j$  ohne den Einfluß von  $\mathbb{X}_k$  bzw.  $\mathbb{X}_\ell$  hat den Wert

$$\rho_{i|k,j|\ell} = \frac{\rho_{ij} - \rho_{ik}\rho_{jk} - \rho_{i\ell}\rho_{j\ell} + \rho_{i\ell}\rho_{k\ell}\rho_{j\ell}}{\sqrt{(1 - \rho_{ik}^2) \cdot (1 - \rho_{j\ell}^2)}}.$$

## Regressionsanalyse

### Definition

Eine Familie

$$\left[ f(\cdot | \mathbf{a}) : \mathbb{R}^N \rightarrow \mathbb{R} \right]_{\mathbf{a} \in \mathcal{M}}$$

von Abbildungen heißt **Funktionsprototyp** der Dimension  $N$ ; ein Element  $f(\cdot | \mathbf{a})$  der Familie heißt **Funktionsinstanz** zu  $\mathbf{a}$ .

Für einen Datensatz  $\omega \subset \mathbb{R}^N \times \mathbb{R}$  definieren wir den **Regressionsfehler**

$$\varepsilon(f, \mathbf{a}, \omega) \stackrel{\text{def}}{=} \sum_{(\mathbf{x}, y) \in \omega} (y - f(\mathbf{x}, \mathbf{a}))^2$$

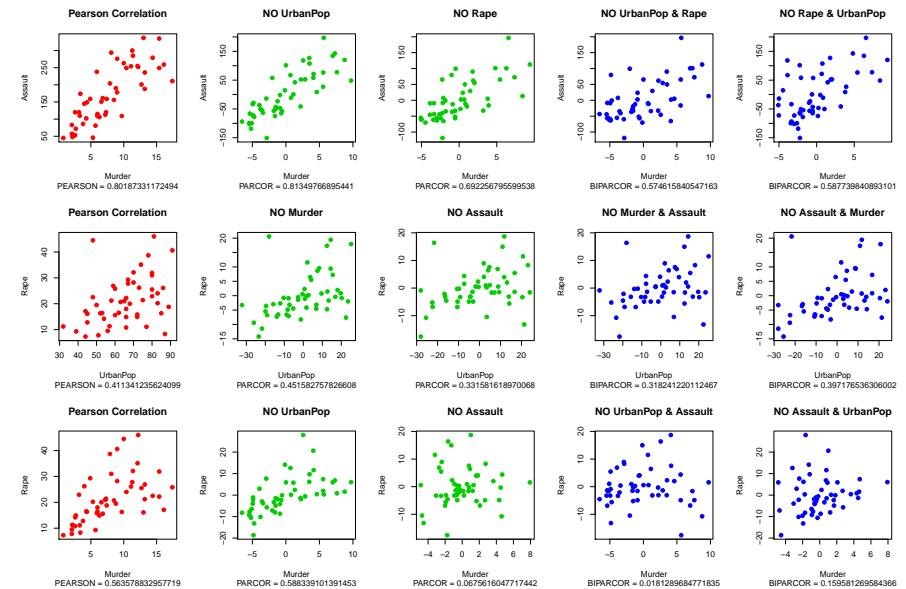
von  $f(\cdot | \mathbf{a})$  über  $\omega$ . Eine Funktionsinstanz  $f(\cdot | \mathbf{a}^*)$  mit minimalem Regressionsfehler heißt **Regressionsfunktion** von  $\omega$ , ihre Parameter  $\mathbf{a}^*$  heißen **Regressionsparameter**.

### Beispiel — lineare Regression

Die spezielle Familie der  $f(\cdot | \mathbf{a}) : (x_1, \dots, x_N) \mapsto a_0 + \sum_{n=1}^N a_n x_n$  mit  $\mathbf{a} \in \mathbb{R}^{N+1}$  heißt **affiner** oder — im Fall  $a_0 \equiv 0$  — **linearer** Funktionsprototyp.

## Beispiel — U.S. Arrests

Mord/Überfall · Metropol/Vergewaltigung · Mord/Vergewaltigung



## Beispiel — Ausgleichsgerade

Funktionsprototyp der Dimension  $N = 1$  ➡ Geradengleichungen  $y = a + bx$

### Regressionsparameter

für einen gegebenen Datensatz  $\omega \subset \mathbb{R} \times \mathbb{R}$

$$b = \frac{\sigma_{xy}}{\sigma_{xx}} \quad \text{und} \quad a = \mu_y - b\mu_x = \mu_y - \frac{\sigma_{xy}}{\sigma_{xx}} \cdot \mu_x$$

### Regressionsfehler

einer Geraden  $y = a + bx$  (Verschiebung  $\rightsquigarrow$  o.B.d.A.  $\mu_x = 0$ )

$$\begin{aligned} \frac{1}{T} \cdot \varepsilon(a, b, \omega) &= \frac{1}{T} \sum_t (y_t - a - bx_t)^2 = \dots \\ &= \sigma_{yy} + \mu_y^2 + a^2 + b^2 \sigma_{xx} - 2a\mu_y - 2b\sigma_{xy} \\ &= \sigma_{yy} \cdot (1 - \rho_{xy}^2) \quad (\text{Einsetzen } a = \mu_y \text{ und } b = \sigma_{xy}/\sigma_{xx}) \end{aligned}$$

### Aufgeklärte Varianz

Die quadrierte Korrelation  $\rho_{xy}^2 \in [0, 1]$  ist der proportionale Anteil der Varianz  $\sigma_{yy}$  von  $\mathbb{Y}$ , der durch die ZV  $\hat{\mathbb{Y}} = a + b \cdot \mathbb{X}$  aufgeklärt werden konnte.

## Lineare und nichtlineare Regression

### Kein Fall für Ausgleichsgeraden

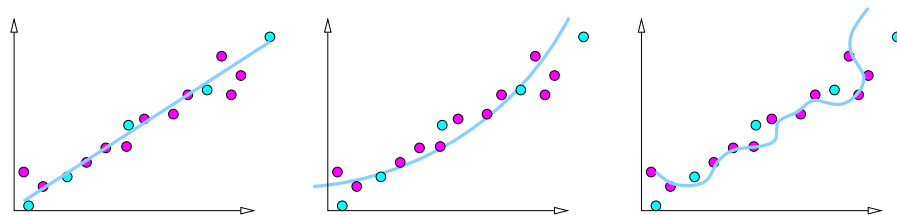
Betrachte die Taylorreihenentwicklung der **sinusoidalen** Abhängigkeit

$$y = \sin(x + \alpha) = \underbrace{\sin \alpha + x \cos \alpha}_{\text{linear}} - x^2 \frac{\sin \alpha}{2} - x^3 \frac{\cos \alpha}{6} \pm \dots$$

### Ausgleichspolynome

Affiner Regressionsansatz mit Termexpansion, z.B. polynomial für  $N = 3$ :

$$(x_1, x_2, x_3) \mapsto (1, x_1, x_2, x_3, x_1^2, x_2^2, x_3^2, x_1 x_2, x_1 x_3, x_2 x_3, \dots)$$



robuste Ausgleichsgerade

leichte Überanpassung

starke Überanpassung

## Lokale Regression

Eine Frage der guten Nachbarschaft

### Nächster-Nachbar-Regel

Belegmenge  $\omega^{(x)} = \{x_s\}$  ist einelementig.

$$\varepsilon(f, \mathbf{a}, \omega \mid \mathbf{x}) \stackrel{\text{def}}{=} (y_s - f(\mathbf{x}_s | \mathbf{a}))^2, \quad s = \underset{t=1..T}{\operatorname{argmin}} d(\mathbf{x}, \mathbf{x}_t)$$

### k-Nächste-Nachbarn-Regel

Scharfe Belegmenge  $\omega^{(x)}$  mit genau  $k$  Elementen.

$$\varepsilon(f, \mathbf{a}, \omega \mid \mathbf{x}) \stackrel{\text{def}}{=} \sum_{i=1}^k (y_{s_i} - f(\mathbf{x}_{s_i} | \mathbf{a}))^2$$

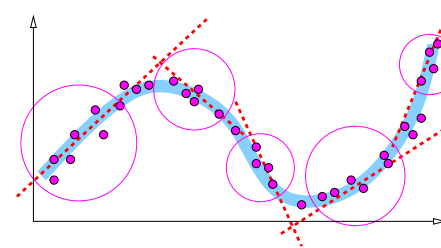
### Gewichtete Mittelung

Unschärfe Belegmenge  $\omega^{(x)}$  mit  $T$  Elementen.

$$\varepsilon(f, \mathbf{a}, \omega \mid \mathbf{x}) \stackrel{\text{def}}{=} \sum_{t=1}^T w_t \cdot (y_t - f(\mathbf{x}_t | \mathbf{a}))^2, \quad w_t \propto \exp \left\{ -\frac{1}{2\sigma^2} \cdot \|\mathbf{x} - \mathbf{x}_t\|^2 \right\}$$

## Lokale Regression

### Verzögertes Lernen



- lokales Modell „just in time“
- kein globales Modell

$$f(\mathbf{x}_t | \mathbf{a}^*) \approx y_t \quad (\forall t)$$

(Algorithmus)

GEGEBEN:

Lerndatenprobe  $\omega = [(\mathbf{x}_t, y_t)]_1^T \subset \mathbb{R}^N \times \mathbb{R}$  und Eingabevektor  $\mathbf{z} \in \mathbb{R}^N$

#### 1 NACHBARSCHAFT FIXIEREN

Berechne Nachbarschaftsmenge  $\omega^{(z)} \subset \omega$ , eventuell mit Gewichten  $\{w_t\}_1^T$ .

#### 2 LOKALE AUSGLEICHSCHEUNUNG

Schätze lokale Regressionsfunktion  $f(\cdot | \mathbf{a}^{(z)})$  für den  $\omega^{(z)}$ -Datensatz.

#### 3 VORHERSAGE TREFFEN

Setze  $\hat{y}(\mathbf{z}) := f(\mathbf{z} | \mathbf{a}^{(z)})$ .

(zumfithogIA)

## Lokale Regression

Konstante Funktionsprototypen · Disjunkte Nachbarschaften

### Konstanter Funktionsprototyp

$$f(\cdot | \mathbf{a}) : \begin{cases} \mathbb{R}^N & \rightarrow \mathbb{R} \\ \mathbf{x} & \mapsto \mathbf{a} \end{cases}, \quad \mathbf{a} \in \mathbb{R}$$

NN-Regel	k-NN-Regel	Distanzgewichte
$f_n(\mathbf{x}) = y_{t(\mathbf{x})}$	$f_k(\mathbf{x}) = \frac{1}{k} \sum_{i=1}^k y_{t_i(\mathbf{x})}$	$f_g(\mathbf{x}) = \mathbf{w}^\top \mathbf{y} / \ \mathbf{w}\ _1$
„Kopie“	„Ortsmittel“	„Schwerpunkt“

### Stückweise lineare Regression

(Algorithmus)

#### 1 GRUPPIERUNG

Lerne extensionale Partition  $\omega_1, \dots, \omega_K$  von  $\mathbf{x}_1, \dots, \mathbf{x}_T$  ( $K$ -means).

#### 2 STÜCKWEISE REGRESSION

Lerne lokale Regressionsfunktionen  $f(\cdot | \mathbf{a}_1), f(\cdot | \mathbf{a}_2), \dots, f(\cdot | \mathbf{a}_K)$ .

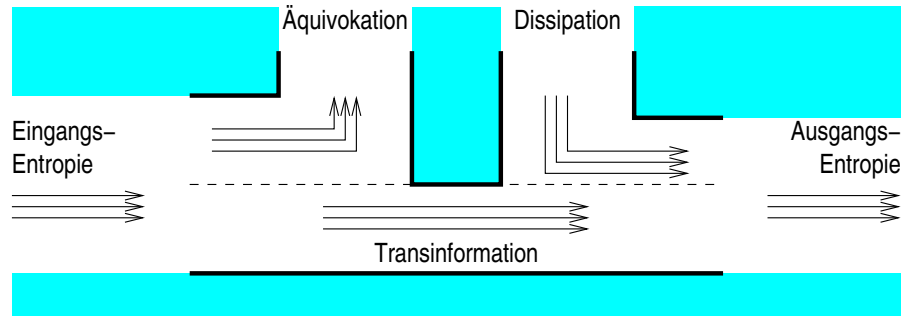
#### 3 VORHERSAGEPHASE

- Bestimme zu  $\mathbf{x} \in \mathbb{R}^N$  den Gruppenindex  $\lambda$ , also mit  $\mathbf{x} \in \Omega_\lambda \supset \omega_\lambda$ .
- Berechne den Vorhersagewert  $\hat{y}(\mathbf{x}) = f(\mathbf{x} | \mathbf{a}_\lambda)$ .

(zumfithogIA)

## Informationstheorie

Der gedächtnislose Informationskanal — Claude Shannon, 1949



### Der Informationskanal

ist durch die gemeinsame Verteilung  $f_{xy}(\cdot, \cdot)$  seiner **Eingangsvariablen**  $\mathbb{X}$  und seiner **Ausgangsvariablen**  $\mathbb{Y}$  charakterisiert.

### Kanalentropien

**Eingangsentropie**  $\mathcal{H}(\mathbb{X}) = \mathcal{E}[-\log f_x(\mathbb{X})]$  **Ausgangsentropie**  $\mathcal{H}(\mathbb{Y}) = \mathcal{E}[-\log f_y(\mathbb{Y})]$  **Gesamtentropie**  $\mathcal{H}(\mathbb{XY}) = \mathcal{E}[-\log f_{xy}(\mathbb{X}, \mathbb{Y})]$

## Transinformation normalverteilter Attribute

### Lemma

Für die (differentiellen) Entropien und die Transinformation normalverteilter Zufallsvariablen gelten die nachfolgenden Aussagen:

1. Wenn  $\mathbb{X} \sim \mathcal{N}(\mu, \sigma^2)$ , so gilt:

$$\mathcal{H}(\mathbb{X}) = \frac{1}{2} \cdot (\log \sigma^2 + 1 + \log(2\pi))$$

2. Wenn  $(\mathbb{X}_1, \dots, \mathbb{X}_N) \sim \mathcal{N}(\mu, \mathbf{S})$ , so gilt:

$$\begin{aligned} \mathcal{H}(\mathbb{X}_1 \dots \mathbb{X}_N) &= \frac{1}{2} \cdot (\log \det(\mathbf{S}) + N + N \log(2\pi)) \\ \mathcal{H}(\mathbb{X}_i \mathbb{X}_j) &= \frac{1}{2} \cdot \log(\sigma_{ii} \cdot \sigma_{jj} - \sigma_{ij}^2) + 1 + \log(2\pi) \end{aligned}$$

3. Für jedes bivariat normale Variablenpaar  $(\mathbb{X}_i, \mathbb{X}_j)$  gilt:

$$\mathfrak{I}(\mathbb{X}_i; \mathbb{X}_j) = -\frac{1}{2} \cdot \log(1 - \rho_{ij}^2)$$

## Bedingte Kanalentropien

Was Sie schon immer über Entropien wissen wollten, aber noch nie zu fragen wagten

### Definition

Der Informationskanal sei durch  $f_{xy}$  charakterisiert.

- $\mathcal{H}(\mathbb{X}|\mathbb{Y}) = \mathcal{E}[-\log f_{x|y}(\mathbb{X}|\mathbb{Y})]$  heißt **Äquivokation** des Kanals.
- $\mathcal{H}(\mathbb{Y}|\mathbb{X}) = \mathcal{E}[-\log f_{y|x}(\mathbb{Y}|\mathbb{X})]$  heißt **Dissipation** des Kanals.
- $\mathfrak{I}(\mathbb{X}; \mathbb{Y}) = \mathcal{E}[-\log \frac{f_x(\mathbb{X}) \cdot f_y(\mathbb{Y})}{f_{xy}(\mathbb{X}, \mathbb{Y})}]$  heißt **Transinformation** des Kanals.

### Lemma

In einem gedächtnislosen Informationskanal gelten die Aussagen:

### Divergenz

(Kullback-Leibler)

1.  $\mathcal{H}(\mathbb{X}|\mathbb{Y}) = \mathcal{H}(\mathbb{XY}) - \mathcal{H}(\mathbb{Y})$
2.  $\mathcal{H}(\mathbb{Y}|\mathbb{X}) = \mathcal{H}(\mathbb{XY}) - \mathcal{H}(\mathbb{X})$
3.  $\mathfrak{I}(\mathbb{X}; \mathbb{Y}) = \mathcal{H}(\mathbb{X}) + \mathcal{H}(\mathbb{Y}) - \mathcal{H}(\mathbb{XY})$
4.  $\mathfrak{I}(\mathbb{X}; \mathbb{Y}) = \mathcal{D}(f_{xy} \| f_x \cdot f_y)$

$$\mathcal{D}(f \| g) = \mathcal{E}_f[\log \frac{f}{g}]$$

### Beweis.

1. **Univariater Fall:**

$$\begin{aligned} \mathcal{H}(\mathbb{X}) &= \mathcal{E}[-\log \mathcal{N}(\mathbb{X} | \mu, \sigma^2)] = \mathcal{E}[\frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2} \left( \frac{\mathbb{X} - \mu}{\sigma} \right)^2] \\ &= \mathcal{E}[\frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2} \tilde{\mathbb{X}}^2] = \frac{1}{2} \cdot (\log \sigma^2 + 1 + \log(2\pi)) \end{aligned}$$

Beachte, daß  $\tilde{\mathbb{X}} = \frac{(\mathbb{X} - \mu)}{\sigma}$  standardnormalverteilt ist, d.h.  $\tilde{\mathbb{X}} \sim \mathcal{N}(0, 1)$ .

2. **Multivariater Fall:**

$$\begin{aligned} \mathcal{H}(\mathbb{X}) &= \mathcal{E}[-\log \mathcal{N}(\mathbb{X} | \mu, \mathbf{S})] = \mathcal{E}[\frac{1}{2} \log \det(2\pi\mathbf{S}) + \frac{1}{2} \cdot (\mathbb{X} - \mu)^\top \mathbf{S}^{-1} (\mathbb{X} - \mu)] \\ &= \mathcal{E}[\frac{1}{2} \log \det(2\pi\mathbf{S}) + \frac{1}{2} \cdot \tilde{\mathbb{X}}^\top \tilde{\mathbb{X}}] = \frac{1}{2} \cdot (\log \det(\mathbf{S}) + N + N \log(2\pi)) \end{aligned}$$

**Bivariater Fall:** gilt wegen  $\det \begin{pmatrix} \sigma_{ii} & \sigma_{ij} \\ \sigma_{ji} & \sigma_{jj} \end{pmatrix} = \sigma_{ii}\sigma_{jj} - \sigma_{ij}^2$ .

3. **Transformationen:**

$$\begin{aligned} \mathfrak{I}(\mathbb{X}_i; \mathbb{X}_j) &= \mathcal{H}(\mathbb{X}_i) + \mathcal{H}(\mathbb{X}_j) - \mathcal{H}(\mathbb{X}_i \mathbb{X}_j) \\ &= +\frac{1}{2} \cdot \log \left( \frac{\sigma_{ii}\sigma_{jj}}{\sigma_{ii}\sigma_{jj} - \sigma_{ij}^2} \right) = -\frac{1}{2} \cdot \log(1 - \rho_{ij}^2) \end{aligned}$$

$(1 - \rho_{ij}^2)$  ist der Anteil **unaufgeklärter Varianz**.

## Transinformation diskreter Attribute

### Wertebereiche und Verteilung

Es sei  $\mathbb{X} \in \{\xi_1, \dots, \xi_K\}$  und  $\mathbb{Y} \in \{\eta_1, \dots, \eta_L\}$  verteilt gemäß

$$p_{k\ell} = P(\mathbb{X} = \xi_k, \mathbb{Y} = \eta_\ell)$$

### Marginale und gemeinsame Entropien

$$\left. \begin{aligned} \mathcal{H}(\mathbb{X}\mathbb{Y}) &= - \sum_k \sum_\ell p_{k\ell} \cdot \log p_{k\ell} \\ \mathcal{H}(\mathbb{X}) &= - \sum_k \left( \sum_\ell p_{k\ell} \right) \cdot \log \left( \sum_\ell p_{k\ell} \right) \\ \mathcal{H}(\mathbb{Y}) &= - \sum_\ell \left( \sum_k p_{k\ell} \right) \cdot \log \left( \sum_k p_{k\ell} \right) \end{aligned} \right\} \Rightarrow \underbrace{- \sum_{k,\ell} p_{k\ell} \cdot \log \frac{p_{k\ell}}{p_{k\cdot} \cdot p_{\cdot\ell}}}_{\text{Transinformation } \mathfrak{I}(\mathbb{X}; \mathbb{Y})}$$

## Korrelation, Regression und Transinformation

### Assoziationsregeln und Netzwerkanalyse

### Bedingte statistische Unabhängigkeit

### Graphische Modelle: ungerichtete Graphen

### Kausale Modelle: gerichtete azyklische Graphen

### Berechnen bedingter Wahrscheinlichkeiten

### Parameterschätzung in Bayesnetzen und Loglinearmodellen

### Aufdeckung der Abhängigkeitsstruktur

### Kovarianzselektion

## Transinformation gemischtskaliger Attribute

$$\mathbb{X} \in \mathbb{R} \text{ und } \mathbb{Y} \in \{\eta_1, \dots, \eta_L\}$$

### Punktweise Transinformation

$$\mathfrak{I}(\mathbb{X}; \mathbb{Y}) = \int_{\mathbb{X}} \sum_y f(x, y) \cdot \mathfrak{I}(x; y)$$

Die „mutual information“ zwischen korrespondierenden Werten  $x$  und  $y$ :

$$\log \frac{f(x|y)}{f(x)} = \underbrace{\log \frac{f(x, y)}{f(x) \cdot f(y)}}_{\mathfrak{I}(x; y)} = \log \frac{f(y|x)}{f(y)}$$

### Faktor diskret

Gaußsche Mischverteilung

$$f(x, \eta_\ell) = \pi_\ell \cdot \mathcal{N}(x | \mu_\ell, \sigma_\ell^2)$$

### Schätzformel

$$\sum_{t=1}^T \frac{1}{T} \cdot \log \frac{\mathcal{N}(x_t | \mu_{\ell(t)}, \sigma_{\ell(t)}^2)}{\sum_\ell \pi_\ell \cdot \mathcal{N}(x_t | \mu_\ell, \sigma_\ell^2)}$$

### Faktor stetig

Diskriminantverteilung

$$f(x, \eta_\ell) = f_{\mathbb{X}}(x) \cdot p(\eta_\ell | x)$$

### Schätzformel

$$\sum_{t=1}^T \frac{1}{T} \cdot \log \frac{p(\eta_{\ell(t)} | x_t)}{\pi_{\ell(t)}}$$

## Assoziationsanalyse

Agrawal (SIGMOD Conference 1993) — mehr als 6.000× zitiert!

### Warenkorbdaten

Objekte = qualitative **Stücklisten**

$$\rightsquigarrow \Omega = \mathfrak{P}(\mathfrak{A})$$

$$\omega \subset \Omega = \{0, 1\}^N$$

über einem globalem **Artikelinventar**  $\mathfrak{A} = \{a_1, \dots, a_N\}$

### Assoziationsregeln

„Wer alle Produkte aus  $A$  kauft, der kauft auch alle Produkte aus  $B$ .“

$$\text{IF } A \text{ THEN } B, \quad A, B \in \Omega, \quad A \cap B = \emptyset$$

### Beispielregeln

IF {Windeln} THEN {Bier}

IF {Brot, Butter} THEN {Milch}

IF {Rosen, Wein, Goldbären} THEN {Kondome}

### Bemerkungen

1. Warenkorbdaten haben **binäre Attribute**.
2. Assoziationsregeln formulieren **multiple Abhängigkeiten**.



## Gute und schlechte Regeln

Abdeckungs- und Geltungsgrad einer Regel · Signifikanz ihrer Prämisse

### Definition

Es sei  $\omega \subset \Omega$  ein Datensatz,  $A, B \in \Omega$  zwei Stücklisten und

IF  $A$  THEN  $B$  (kürzer:  $A \rightarrow B$ ) eine Assoziationsregel. Die Größe

$$\text{supp}(A \rightarrow B) \stackrel{\text{def}}{=} \text{supp}(A \cup B), \quad \text{supp}(A) \stackrel{\text{def}}{=} \frac{|\{x \in \omega \mid x \supseteq A\}|}{|\omega|}$$

heißt **Support**,

$$\rightsquigarrow \hat{P}(A \cup B)$$

$$\text{conf}(A \rightarrow B) \stackrel{\text{def}}{=} \frac{\text{supp}(A \cup B)}{\text{supp}(A)}$$

heißt **Konfidenz** und

$$\rightsquigarrow \hat{P}(B|A)$$

$$\text{lift}(A \rightarrow B) \stackrel{\text{def}}{=} \frac{\text{supp}(A \cup B)}{\text{supp}(A) \cdot \text{supp}(B)}$$

heißt **Relevanz** der Assoziation  $A \rightarrow B$ .

$$\rightsquigarrow \frac{\hat{P}(B|A)}{\hat{P}(B)}$$

## Apriori-Basisalgorithmus

Schichtenweise Stücklisten- und Regelgenerierung

GEGEBEN

Warenkorbdaten  $\omega$ , Stückzahlgrenze  $N^*$ , Schwellen  $\theta_s, \theta_c, \theta_r$ .

### 1 INITIALISIERUNG

$$\mathcal{M}_1 \leftarrow \{\{i\} \mid \text{supp}(\{i\}) \geq \theta_s\}, \quad \mathcal{R} = \emptyset$$

### 2 SCHICHTEXPANSION ( $n = 2, \dots, N^*$ )

Erzeuge alle

$$A = B \cup \{i\} \quad \text{mit } B \in \mathcal{M}_{n-1}, \{i\} \in \mathcal{M}_1 \text{ und } i \notin B.$$

Bringe  $A$  nach  $\mathcal{M}_n$  falls  $\text{supp}(A) \geq \theta_s$ .

### 3 REGLERZEUGUNG

Für alle  $C \in \mathcal{M} = \bigcup_{n=1}^{N^*} \mathcal{M}_n$ :

Für alle Artikel  $j \in C$  teste

$$\text{conf}(C \setminus \{j\} \rightarrow \{j\}) \geq \theta_c \quad \wedge \quad \text{lift}(C \setminus \{j\} \rightarrow \{j\}) \geq \theta_r$$

und verbringe die Regel im Erfolgsfall nach  $\mathcal{R}$ .

## Extraktion nützlicher Assoziationsregeln

Eine Frage des Aufwandes

### Aufgabenstellung

Gesucht ist — bei gegebenen Warenkorbdaten — die Teilmenge solcher Regeln  $A \rightarrow B$

- mit **signifikantem Abdeckungsgrad**  $\text{supp}(A \rightarrow B) \geq \theta_s$
- und hohem **Geltungsgrad**  $\text{conf}(A \rightarrow B) \geq \theta_c$
- und erheblicher **Aussagekraft**.  $\text{lift}(A \rightarrow B) \geq \theta_r$

### Problem

Es gibt  $2^N$  kombinatorisch mögliche Stücklisten und es gibt  $3^N$  mögliche Assoziationsregeln ( $N = |\mathcal{A}|$ ).

### Lösungsansatz

Es gilt die **Antitonie**

$$A \supseteq B \Rightarrow \text{supp}(A) \leq \text{supp}(B)$$

und es gibt nur  $N$  **einelementige** Stücklisten.

## Regelformat und Artikelbeschreibung

### Assoziationsregeln mit multipler Conclusio

Statt einfacher Regeln  $C \setminus \{j\} \rightarrow \{j\}$  produziere

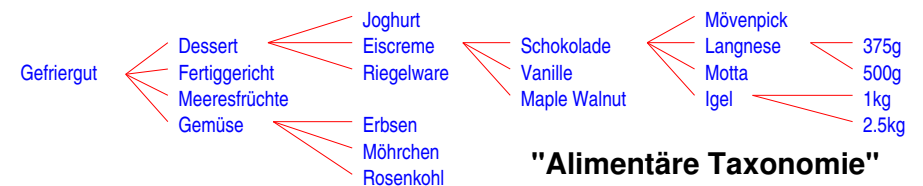
$$A \rightarrow B \quad \text{mit } A \cap B = \emptyset \text{ und } A \cup B = C.$$

Stufenweise Erzeugung („bottom-up“) unter Verwendung der Monotonie:

$$B_1 \subseteq B_2 \Rightarrow \begin{cases} \text{supp}(A_1 \rightarrow B_1) = \text{supp}(A_2 \rightarrow B_2) \\ \text{conf}(A_1 \rightarrow B_1) \geq \text{conf}(A_2 \rightarrow B_2) \\ \text{lift}(A_1 \rightarrow B_1) \stackrel{?}{\geq} \text{lift}(A_2 \rightarrow B_2) \end{cases}$$

### Aufwandsreduktion durch onthologische Gliederung

Artikeleinträge werden durch ihre Verallgemeinerungen aufgestockt.



## Assoziationsanalyse für (mehrwertige) Nominalskalen

### Verallgemeinerte Stücklisten

Listen von kontradiktionsfreien Attribut-Wert-Paaren:

$(windy = false, play = no, outlook = sunny, humidity = high)$

Es gibt  $\prod_n (L_n + 1)$  Stücklisten und  $\prod_n (2 \cdot L_n + 1)$  Assoziationsregeln.

### Beispielregeln

Klassisch:

- IF {Spaghetti} THEN {Rotwein, Tomaten, Basilikum}
- IF {Waits, Dylan, Bush} THEN {Spektor}

Mehrwertig:

- IF {humidity = high, windy = false} THEN {outlook = sunny}

Zweiwertig:

- IF {Pommes,  $\neg$  Ketchup} THEN {Mayonnaise}
- IF {E.Jelinek, Ch.Roche} THEN { $\neg$  U.Danella}

## Korrelation, Regression und Transinformation

## Assoziationsregeln und Netzwerkanalyse

## Bedingte statistische Unabhängigkeit

## Graphische Modelle: ungerichtete Graphen

## Kausale Modelle: gerichtete azyklische Graphen

## Berechnen bedingter Wahrscheinlichkeiten

## Parameterschätzung in Bayesnetzen und Loglinearmodellen

## Aufdeckung der Abhängigkeitsstruktur

## Kovarianzselektion

## Beispiel — Tennisdaten mit WEKA

5 Attribute · 14 Objekte · Apriori mit  $\theta_s = 15\%$ ,  $\theta_c = 90\%$

### Stücklistenaufstellung („itemsets“)

12 Einermengen · 47 Paare · 39 Tripel · 6 Quadrupel

### Beste Assoziationsregeln ( $\theta_c \equiv 100\%$ )

- 4 IF {humidity = normal, windy = false} THEN {play = yes}  
IF {temperature = cool} THEN {humidity = normal}  
IF {outlook = overcast} THEN {play = yes}
- 3 IF {temperature = cool, play = yes} THEN {humidity = normal}  
IF {outlook = rainy, windy = false} THEN {play = yes}  
IF {outlook = rainy, play = yes} THEN {windy = false}  
IF {outlook = sunny, humidity = high} THEN {play = no}  
IF {outlook = sunny, play = no} THEN {humidity = high}
- 2 IF {temp = cool, windy = false} THEN {humidity = normal, play = yes}  
IF {temp = cool, humidity = normal, windy = false} THEN {play = yes}

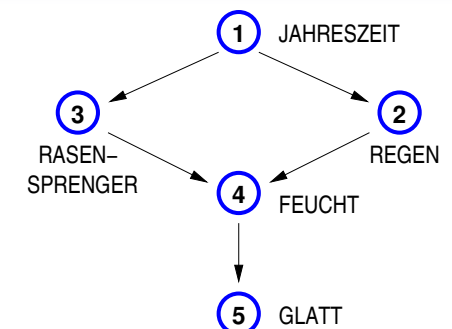
## Graphische Wahrscheinlichkeitsmodelle

### Regen $\leftarrow$ Jahreszeit

	$X_1$ w	$X_1$ f	$X_1$ s	$X_1$ h
$X_2 = 0$	0.2	0.3	0.1	0.7
$X_2 = 1$	0.8	0.7	0.9	0.3

### Feucht $\leftarrow$ Regen, Sprenger

	$X_2 X_3$	$\bar{X}_2 X_3$	$X_2 \bar{X}_3$	$\bar{X}_2 \bar{X}_3$
$X_4 = 0$	0.1	0.3	0.4	0.8
$X_4 = 1$	0.9	0.7	0.6	0.2



### Wozu Graphische Modelle ?

- Visualisierung quantitativer Zusammenhänge
- Inferenz von Abhängigkeitsbeziehungen
- Berechnung kausaler Effekte
- Effiziente Auswertung multivariater Modelle

### Jahreszeit $\leftarrow$

$X_1 = w$	0.25
$X_1 = f$	0.25
$X_1 = s$	0.25
$X_1 = h$	0.25

### Glatt $\leftarrow$ Feucht

	$X_4$	$\bar{X}_4$
$X_5 = 0$	0.3	0.9
$X_5 = 1$	0.7	0.1

## Wahrscheinlichkeit und Graphstruktur

### Datensatz

Wahrscheinlichkeits(dichte)werte

$$P : \Omega = \mathcal{X}_1 \times \dots \times \mathcal{X}_N \rightarrow \mathbb{R}_0^+$$

durch **Statistiken** des Datensatzes  
 $\omega \subset \Omega$  repräsentiert.

$$P(\mathbf{x}) = P(x_1) \cdot P(x_2|x_1) \cdot P(x_3|x_1) \cdot P(x_4|x_2, x_3) \cdot P(x_5|x_4)$$

### Modellformel

Saturiertes Modell

$$\prod_n L_n - 1 = 63$$

Naives Modell

$$\sum_n L_n - N = 7$$

Faktorisierung

$$3 + 4 + 4 + 4 + 2 = 17$$

### Dependenzmodell

Menge aller **bedingten**

Unabhängigkeiten zwischen

**Mengen von** Zufallsvariablen

$$\mathbb{S}(\mathbb{X}_2 | \mathbb{X}_1 | \mathbb{X}_3)$$

(„Regen“ unabhängig von  
 „Rasensprenger“, wenn „Jahreszeit“  
 gegeben)

### Graphisches Modell

- **Markovnetz**

ungerichteter Graph

„partielle Unabhängigkeit“

$$\mathbb{X}_i \not\leftrightarrow \mathbb{X}_j$$

- **Bayesnetz**

gerichteter azyklischer Graph

„kausale Abhängigkeit“

$$\mathbb{X}_i \rightarrow \mathbb{X}_j$$

## Simpsons Paradoxon #2

Eine farbenfrohe Mordstatistik für den Bundesstaat Florida

### Zweiwegetabelle: Hautfarbe & Strafmaß

FarbeMörder	#Todesurteil	#Haftstrafe	% T.U.
schwarz	17	149	11.4
weiß	19	141	12.5

kein Rassismus: ähnliche Todesurteilquote für Schwarz und Weiß

### Zusatzvariable: Hautfarbe des Opfers

FarbeOpfer	FarbeMörder	#Tod	#Haft	% T.U.
schwarz	schwarz	6	97	5.8
schwarz	weiß	0	9	0.0
weiß	schwarz	11	52	17.5
weiß	weiß	19	132	12.6

Der Mord an einem weißen Mitbürger kommt teurer zu stehen!

„marginal unabhängig“

„bedingt unabhängig“

## Simpsons Paradoxon #1

Geschlechtsspezifische Diskriminierung an der Universität

### Zweiwegetabelle: Geschlecht & Zulassungsquote

Geschlecht	#Bewerbung	#Zulassung	%
M	600	350	58.3
F	600	250	41.6

Frauen haben die geringeren Zulassungschancen!

### Zusatzvariable: Fakultätszugehörigkeit

Fakultät	Geschlecht	#Bewerbung	#Zulassung	%
TECH	M	100	25	25
TECH	F	300	75	25
PHIL	M	200	100	50
PHIL	F	200	100	50
THEO	M	300	225	75
THEO	F	100	75	75

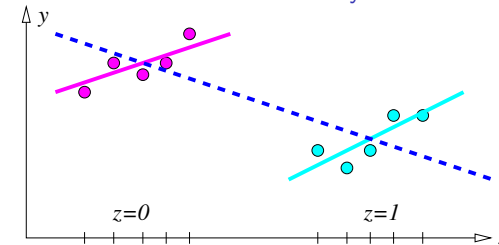
Männer tendieren zu Fakultäten mit hoher Zulassungsquote!

„marginal abhängig“

„bedingt abhängig“

## Simpsons Paradoxon #3

Das Maradonna-Syndrom: Fußballspielen ist ungesund!



### Drei Attribute

$\mathbb{X}$  = „Fußballaktivität“

$\mathbb{Y}$  = „Lebenserwartung“

$\mathbb{Z}$  = „Geschlecht“

### Bedingte Abhängigkeit

Weibliche wie männliche

Regressionsgeraden

$$Y = f(X|0) \quad \text{bzw.} \quad Y = f(X|1)$$

besitzen **positive** Steigung.

### Grund:

Frauen sind tendenziell **langlebig** und stehen eher auf **Volleyball+Ayurveda**.

### Marginale Abhängigkeit

Geschlechtsneutrale

Regressionsgerade

$$Y = f(X)$$

besitzt **negative** Steigung.

## Bedingte statistische Unabhängigkeit

von Mengen von Zufallsvariablen

### Definition

Es seien  $A, B, Z$  drei paarweise disjunkte Teilmengen der Zufallsvariablen  $\{\mathbb{X}_1, \dots, \mathbb{X}_N\}$ . Dann heißt  $A$  **bedingt statistisch unabhängig** von  $B$  bezogen auf  $Z$  genau dann, wenn gilt

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

und wir schreiben

$$\mathfrak{S}(A | Z | B) .$$

Ferner heißen  $A$  und  $B$  **bedingt faktorisierbar** bezogen auf  $Z$ , falls es zwei geeignete Funktionen (sic!)  $f$  und  $g$  gibt mit

$$P(A, B, Z) = f(A, Z) \cdot g(B, Z) .$$

### Marginale statistische Unabhängigkeit

Der Spezialfall „gewöhnlicher“ statistischer Unabhängigkeit ergibt sich für UA-Postulate der Form  $\mathfrak{S}(A | Z | B)$  mit  $Z = \emptyset$ .

### Beweis.

- (1)⇒(2)⇒(3)⇒(1)

$$\begin{aligned} P(b | a, z) &= \frac{P(a, b | z)}{P(a | z)} = \frac{P(a | b, z) \cdot P(b | z)}{P(a | z)} \\ &= \frac{P(a | z) \cdot P(b | z)}{P(a | z)} = P(b | z) \end{aligned}$$

- (1)⇒(2)⇒(3)⇒(1)

$$P(a, b, z) = P(a, z) \cdot P(b | a, z) = P(a, z) \cdot P(b | z) =: f(a, z) \cdot g(b, z)$$

- (1)⇒(2)⇒(3)⇒(1)

$$\begin{aligned} P(a | b, z) &= \frac{P(a, b, z)}{P(b, z)} = \frac{P(a, b, z)}{\sum_a P(a, b, z)} \\ &= \frac{f(a, z) \cdot g(b, z)}{\sum_a f(a, z) \cdot g(b, z)} = \frac{f(a, z)}{\sum_a f(a, z)} \end{aligned}$$

Der letzte Ausdruck ist offenbar unabhängig von  $b$ .

□

## Rechenregeln für bedingte Unabhängigkeiten

### Lemma

Die folgenden Allaussagen über die Werte  $a, b$  und  $z$  dreier Zufallsvariablen  $\mathbb{X}_a, \mathbb{X}_b, \mathbb{X}_z$  sind äquivalent:

1.  $P(a | b, z) = P(a | z)$  ( $a \not\sim b$  wenn  $z$  bekannt)
2.  $P(b | a, z) = P(b | z)$  ( $b \not\sim a$  wenn  $z$  bekannt)
3.  $P(a, b, z) = f(a, z) \cdot g(b, z)$  (Faktorisierbarkeit)

Diese Äquivalenz gilt entsprechend für **Mengen** von Zufallsvariablen.

### Weitere äquivalente Formulierungen

für die bedingte statistische Unabhängigkeit zwischen drei Zufallsvariablen:

1.  $P(a, b, z) = \frac{P(a, z) \cdot P(b, z)}{P(z)}$
2.  $P(a, b | z) = P(a | z) \cdot P(b | z)$
3.  $P(a, b, z) = P(a | z) \cdot P(b, z)$

### Diskrete ZV

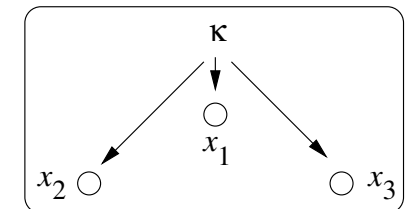
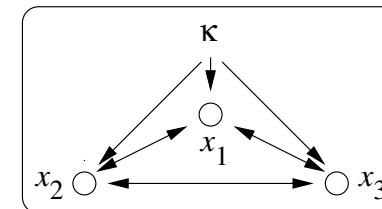
$P(\mathbb{Y} = y | \mathbb{I} = i)$  ist konstant bzgl.  $i$ .

### Stetige ZV

Lineare Regression  
 $\mathbb{Y} | \mathbf{x} \sim \mathcal{N}(a + b\mathbf{x}, \sigma^2)$   
mit  $b = 0$ .

## Beispiel — Numerische Klassifikation

Normale und naive Bayesregel



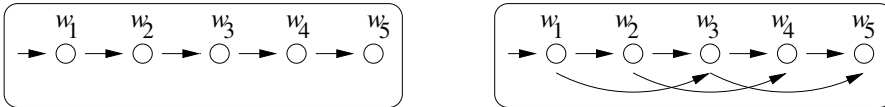
### Datenerzeugungsmodell

für Merkmale  $x_1, \dots, x_N \in \mathbb{R}$  und Klassenvariable  $y \in \{\Omega_1, \dots, \Omega_K\}$ :

$$f(\mathbf{x}, \Omega_\kappa) = P(\Omega_\kappa) \cdot f(\mathbf{x} | \Omega_\kappa)$$

- **Multivariate Normalverteilungsdichte** (saturiertes Modell):  
 $x_i \leftarrow \{x_j | j \neq i\}$  für alle  $i$
- **Klassenbedingte Unabhängigkeit** (ausgedünntes Modell):  
 $f(\mathbf{x} | \Omega_\kappa) = \prod_i \mathcal{N}(x_i | \mu_i, \sigma_i^2)$  ergibt  $x_i \leftarrow \emptyset$  für alle  $i$

## Beispiel — $N$ -Gramm-Grammatiken



### Datenerzeugungsmodell

für eine Symbolfolge (Wortfolge)  $\mathbf{w} = w_1 \dots w_M$  ist die **Kettenregel**

$$P(\mathbf{w}) = \prod_{m=1}^M P(w_m | w_1, \dots, w_{m-1}) \simeq \prod_{m=1}^M P(w_m | w_{m-2}, w_{m-1}) \\ \simeq \prod_{m=1}^M P(w_m | w_{m-1})$$

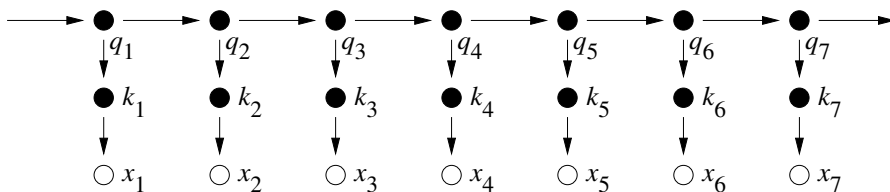
mit den statistischen Abhängigkeiten  $\left\{ \begin{array}{l} w_m \leftarrow w_{m-1} \text{ (Bigramme)} \\ w_m \leftarrow \{w_{m-2}, w_{m-1}\} \text{ (Trigramme)} \end{array} \right\}$ .

### Unabhängigkeitspostulate der Bigramm-Grammatik

$\mathfrak{I}(\{w_m\} | \{w_{m-1}\} | \{w_1, \dots, w_{m-2}\})$  für alle  $m = 2, \dots, M$ .

## Beispiel — (Semi-)kontinuierliches HMM

mit eindimensionalen Ausgabewerten



### Datenerzeugungsmodell

**Beobachtbare Wertefolge**  $\mathbf{x} = x_1, \dots, x_T$  mit  $x_t \in \mathbb{R}$

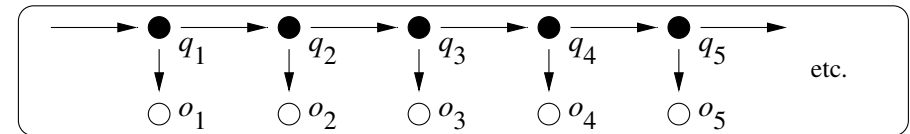
**Verborgene Komponentenfolge**  $\mathbf{k} = k_1 \dots k_T$  mit  $k_t \in \mathcal{K}$

**Verborgene Zustandsfolge**  $\mathbf{q} = q_1 \dots q_T$  mit  $q_t \in \mathcal{Q}$

$$P(\mathbf{X}) = P(\mathbf{X} | \lambda) = \sum_{\mathbf{q} \in \mathcal{Q}^T} \sum_{\mathbf{k} \in \mathcal{K}^T} P(\mathbf{X}, \mathbf{k}, \mathbf{q} | \lambda)$$

mit statistischen Abhängigkeiten  $q_t \leftarrow q_{t-1}$ ,  $k_t \leftarrow q_t$  und  $x_t \leftarrow k_t$ .

## Beispiel — Hidden Markov Modelle



### Datenerzeugungsmodell

**Beobachtbare Zeichenfolge**  $\mathbf{o} = o_1 \dots o_T$  mit  $o_t \in \mathcal{O}$

**Verborgene Zustandsfolge**  $\mathbf{q} = q_1 \dots q_T$  mit  $q_t \in \mathcal{Q}$

$$P(\mathbf{o}) = P(\mathbf{o} | \lambda) = \sum_{\mathbf{q} \in \mathcal{Q}^T} P(\mathbf{o}, \mathbf{q} | \lambda) = \sum_{\mathbf{q} \in \mathcal{Q}^T} \prod_{t=1}^T P(q_t | q_{t-1}) \cdot P(o_t | q_t)$$

mit statistischen Abhängigkeiten  $q_t \leftarrow q_{t-1}$  und  $o_t \leftarrow q_t$ .

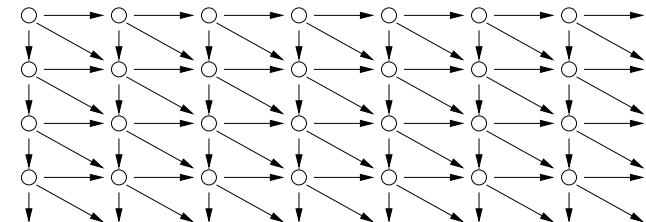
### Unabhängigkeitspostulate des HMM

$\mathfrak{I}(\{q_{t+1}\} | \{q_t\} | \{q_1, \dots, q_{t-1}; o_1, \dots, o_t\})$  und

$\mathfrak{I}(\{o_{t+1}\} | \{q_{t+1}\} | \{q_1, \dots, q_t; o_1, \dots, o_t\})$

## Beispiel — 2D Markov Random Field

Texturmodelle in der Grauwertbildanalyse



### Datenerzeugungsmodell

**Beobachtbare** Zufallsvariablen  $x_{i,j}$

auf dem **Ortsgitter**  $i = 1, \dots, I$  und  $j = 1, \dots, J$

mit statistischen Abhängigkeiten  $x_{i,j} \leftarrow \{x_{i-1,j-1}, x_{i,j-1}, x_{i-1,j}\}$ .

### Unabhängigkeitspostulate des MRF

Für alle Gitterpunkte  $(n, m) \in \mathbb{Z} \times \mathbb{Z}$  ist gefordert:

$\mathfrak{I}(\{\mathbb{X}_{n,m}\} | \{\mathbb{X}_{n,m-1}, \mathbb{X}_{n-1,m}, \mathbb{X}_{n-1,m-1}\} | \{\mathbb{X}_{i,j} | i < n, j < m\})$

## Dependenzmodelle

### Algebraische Charakterisierung von Abhängigkeitsstrukturen

#### Definition

Sei  $V = \{X_1, \dots, X_N\}$  eine Menge von Zufallsvariablen und  $P(\cdot)$  eine Verteilung über  $V$ . Die Relation  $\mathfrak{S} = \mathfrak{S}_P$  mit

$$\mathfrak{S} : \mathfrak{P}X \times \mathfrak{P}X \times \mathfrak{P}X \rightarrow \{0, 1\}$$

heißt **Dependenzmodell von  $P(\cdot)$** , wenn für alle (disjunkten) Variablenmengen  $A, B, Z \subset V$  gilt:

$$\mathfrak{S}(A | Z | B) \Leftrightarrow P(A | B, Z) = P(A | Z)$$

#### Bemerkungen

- Es gibt  $4^N$  viele Variablenkombinationen  $A, B, Z$ .  
Es gibt  $2^{4^N}$  viele dreistellige Mengenrelationen  $\mathfrak{S}$  über  $V$ .  
Wieviele  $\mathfrak{S}$  davon sind ein **valides Dependenzmodell  $\mathfrak{S}_P$** ?
- Simpsons Paradoxa:  $\mathfrak{S}(A|Z|B) \not\Rightarrow \mathfrak{S}(A|\emptyset|B)$  und  $\mathfrak{S}(A|Z|B) \not\Rightarrow \mathfrak{S}(A|\emptyset|B)$

#### Beweis.

**SYM** Symmetrie

$$\mathfrak{S}(A|Z|B) \Rightarrow P(A, Z, B) = \underbrace{f(A, Z) \cdot g(B, Z)}_{P(B, Z, A)} \Rightarrow \mathfrak{S}(B|Z|A)$$

**DEC** Dekomposition

$$P(A, Z, B) = \sum_C P(A, Z, B, C) = \sum_C f(A, Z) \cdot \overbrace{g(B, C, Z)}^{\tilde{g}(B, Z)} = f(A, Z) \cdot \sum_C \tilde{g}(B, C, Z)$$

beweist  $\mathfrak{S}(A|Z|B)$ ; analoge Herleitung von  $\mathfrak{S}(A|Z|C)$ .

**WUN** Schwache Vereinigung

$$\begin{aligned} \mathfrak{S}(A | Z | B, C) \Rightarrow P(A, Z, B, C) &= f(A, Z) \cdot g(B, C, Z) \\ &= \tilde{f}(A, Z, C) \cdot \tilde{g}(B, Z, C) \Rightarrow \mathfrak{S}(A | Z, C | B) \end{aligned}$$

**CON** Kontraktion

$$P(A | Z, B, C) = \underbrace{P(A | Z, B)}_{\mathfrak{S}(A|Z,B|C)} = \underbrace{P(A | Z)}_{\mathfrak{S}(A|Z|B)} \Rightarrow \mathfrak{S}(A | Z | B, C)$$

**INT** Durchschnitt (Beweis zu äquivalenter Formulierung INT\* folgt)

□

## Pearlsche Dependenzaxiome

### Axiomatische Charakterisierung aller „erlaubten“ $\mathfrak{S}$ -Relationen

#### Satz (Judea Pearl)

Es sei  $P(\cdot)$  eine Wahrscheinlichkeitsverteilung über  $X_1, \dots, X_N$  und  $\mathfrak{S}(\cdot | \cdot | \cdot)$  das zugehörige Dependenzmodell.

Dann gelten für alle (paarweise disjunkten) Variablenmengen  $A, B, C, Z$  die folgenden vier Aussagen:

**SYM Symmetrie**  $\mathfrak{S}(A | Z | B) \Leftrightarrow \mathfrak{S}(B | Z | A)$

**DEC Dekomposition**  $\mathfrak{S}(A | Z | BUC) \Rightarrow \mathfrak{S}(A | Z | B) \wedge \mathfrak{S}(A | Z | C)$

**WUN Schwache Vereinigung**  $\mathfrak{S}(A | Z | BUC) \Rightarrow \mathfrak{S}(A | ZUC | B)$

**CON Kontraktion**  $\mathfrak{S}(A | Z | B) \wedge \mathfrak{S}(A | ZUB | C) \Rightarrow \mathfrak{S}(A | Z | BUC)$

Falls  $P(\cdot)$  zudem streng positiv ( $\forall \mathbf{x} \in \Omega : P(\mathbf{x}) > 0$ ) ist, gilt sogar:

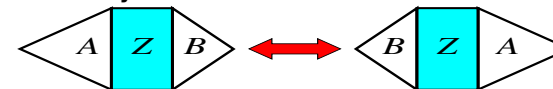
**INT Durchschnitt**

$$\mathfrak{S}(A | ZUC | B) \wedge \mathfrak{S}(A | ZUB | C) \Rightarrow \mathfrak{S}(A | Z | BUC)$$

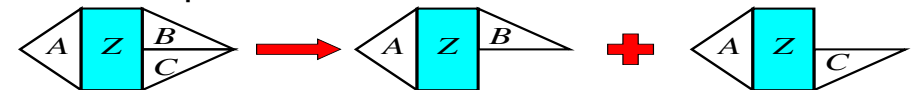
## Pearlsche Dependenzaxiome

### Beweis durch angestrenktes Hingucken

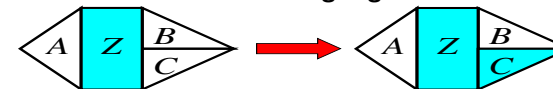
**SYM – Symmetrie**



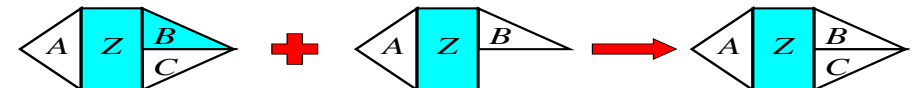
**DEC – Dekomposition**



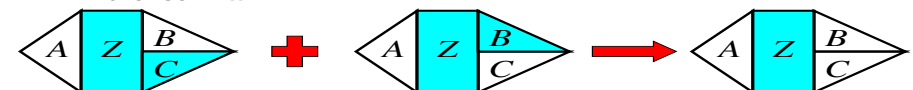
**WUN – Schwache Vereinigung**



**CON – Kontraktion**



**INT – Durchschnitt**



## Bemerkungen zu Pearls Axiomen

1. Die logische Umkehrung der Implikation CON folgt aus den Axiomen DEC und WUN.
2. Die logische Umkehrung von INT folgt mit zweimaliger Anwendung von WUN.
3. Die Axiomatisierung kann auf nichtdisjunkte ZV-Mengen ausgedehnt werden. Aus den obengenannten Axiomen sowie der zusätzlichen Forderung  $\Im(A \mid Z \mid Z)$  beweist man die Aussage

$$\Im(A \mid Z \mid B) \Leftrightarrow \Im(A, Z \mid Z \mid B, Z)$$

4. Die fünf Axiome sind voneinander logisch unabhängig. Beweis durch Gegenbeispiele.
5. Das Axiom INT findet sich auch in der folgenden, äquivalenten Fassung **INT\*** (Lauritzen,  $Z = \emptyset$ ) wieder:

$$\Im(A \mid C \mid B) \wedge \Im(A \mid B \mid C) \Rightarrow \Im(A \mid \emptyset \mid B, C)$$

## Vollständigkeitsvermutung (Pearl & Paz, 1985)

### Trügerische Hoffnung

Wenn  $\Im$  die Axiome SYM, DEC, WUN & CON erfüllt, so heißt  $(V, \Im)$  **Semigraphoid** und es gibt eine Wahrscheinlichkeitsverteilung  $P(\cdot)$  mit

$$P(A \mid B, Z) = P(A \mid Z) \iff \Im(A \mid Z \mid B).$$

Wenn zusätzlich das Durchschnittsaxiom (INT) erfüllt ist, so kann für das **Graphoid**  $(V, \Im)$  sogar ein streng positives  $P(\cdot)$  gefunden werden.

### Satz (Studeny, 1992)

Weder für die Relationenmenge

$$\{\Im_P \mid P \text{ Wahrscheinlichkeitsverteilung}\}$$

noch für deren Teilmenge

$$\{\Im_P \mid P \text{ streng positive Wahrscheinlichkeitsverteilung}\}$$

gibt es ein korrektes und vollständiges **endliches Axiomensystem**.

## Durchschnittsaxiom INT

Garantiert ausschließlich für streng positive Verteilungen!

### Herleitung für streng positive $P(\cdot)$

Auf Grund der Prämissen von INT\* gelten die Faktorisierungen

$$P(a, b, c) = k(a, c) \cdot \ell(b, c) = g(a, b) \cdot h(b, c)$$

und für beliebige Werte  $c$  — also zum Beispiel für  $c_0$  beliebig aber fest — gilt

$$g(a, b) = k(a, c_0) \cdot \frac{\ell(b, c_0)}{h(b, c_0)} =: \pi(a) \cdot \rho(b).$$

Dann gilt die marginale Unabhängigkeit  $\{a\} \not\sim \{b, c\}$  wegen der Faktorisierung

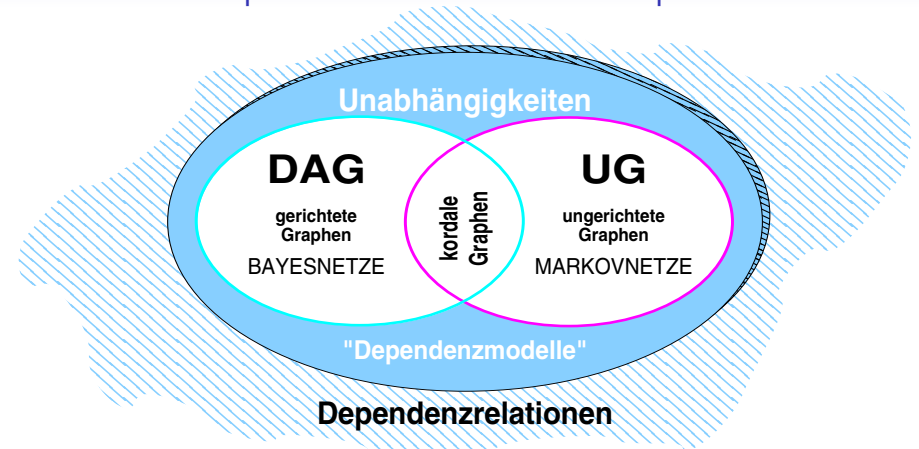
$$P(a, b, c) = \pi(a) \cdot [\rho(b) \cdot h(b, c)].$$

### Gegenbeispiel

Die drei binärwertige Zufallsvariablen mit  $\mathbb{X}_1 = \mathbb{X}_2 = \mathbb{X}_3$  und  $P(\mathbb{X}_i = 1) = \frac{1}{2}$  für alle  $i \in \{1, 2, 3\}$  sind nicht streng positiv (z.B.  $P(1, 1, 0) = 0$ ) und widerlegen INT:

$$\Im(\mathbb{X}_1 \mid \mathbb{X}_2 \mid \mathbb{X}_3), \quad \Im(\mathbb{X}_1 \mid \mathbb{X}_3 \mid \mathbb{X}_2), \quad \neg \Im(\mathbb{X}_1 \mid \emptyset \mid \mathbb{X}_2, \mathbb{X}_3)$$

## Dependenzmodelle und Graphen



- ? Welche Dependenzmodelle sind durch UG charakterisierbar
- ? Welche Dependenzmodelle sind durch DAG charakterisierbar
- ? Welche Dependenzmodelle liegen gleichzeitig in beiden Klassen
- ? Welche Dependenzmodelle sind komplexer als jede Graphstruktur



Korrelation, Regression und Transinformation

Assoziationsregeln und Netzwerkanalyse

Bedingte statistische Unabhängigkeit

Graphische Modelle: ungerichtete Graphen

Kausale Modelle: gerichtete azyklische Graphen

Berechnen bedingter Wahrscheinlichkeiten

Parameterschätzung in Bayesnetzen und Loglinearmodellen

Aufdeckung der Abhängigkeitsstruktur

Kovarianzselektion

## Graphische Verteilungen und Dependenzmodelle

Überrepräsentation & Unterrepräsentation von  $\mathfrak{I}(\cdot | \cdot | \cdot)$  durch  $\text{sep}(\cdot | \cdot | \cdot)$

### Definition

Es sei  $P(\cdot)$  eine Wahrscheinlichkeitsverteilung auf  $V$  und  $\mathfrak{I}$  ihr Dependenzmodell. Der ungerichtete Graph  $\mathcal{G} = (V, \mathcal{E})$  heißt

- **Abhängigkeitsbild** von  $P$  gdw.

$$\mathfrak{I}(A | Z | B) \Rightarrow \text{sep}(A | Z | B)$$

- **Unabhängigkeitsbild** von  $P$  gdw.

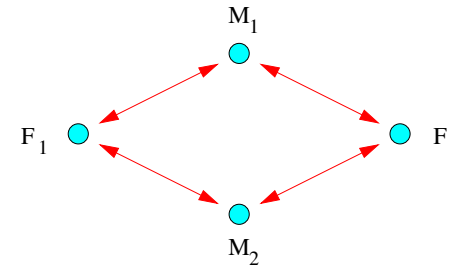
$$\mathfrak{I}(A | Z | B) \Leftarrow \text{sep}(A | Z | B)$$

- **perfektes Bild** von  $P$  gdw.

$$\mathfrak{I}(A | Z | B) \Leftrightarrow \text{sep}(A | Z | B)$$

Die Verteilung  $P(\cdot)$  (und das Modell  $\mathfrak{I}$ ) heißen **graphisch**, wenn ein ungerichteter Graph existiert, der  $\mathfrak{I}$  perfekt abbildet.

## Trennungsrelation im ungerichteten Graphen



Attributwerte:  $\pm$ infiziert

### Partnertauschmodell

Wegen

$$P(f_2 | f_1, m_1, m_2) = P(f_2 | m_1, m_2)$$

gilt  $\mathfrak{I}(F_1 | M_1, M_2 | F_2)$ .

### Partnertauschgraph

$\text{sep}(F_1 | M_1, M_2 | F_2)$  und  $\text{sep}(M_1 | F_1, F_2 | M_2)$

### Definition

Es sei  $\mathcal{G} = (V, \mathcal{E})$  ein ungerichteter Graph und  $A, B, Z \subset V$  disjunkte Knotenmengen. Die Menge  $Z$  **trennt**  $A$  **von**  $B$  genau dann, wenn alle Pfade zwischen Elementen  $a \in A$  und  $b \in B$  mindestens einen Knoten  $z \in Z$  enthalten. Wir schreiben dafür:

$$\text{sep}(A | Z | B)$$

⇒ „ $Z$  blockiert alle Verbindungen zwischen Knoten aus  $A$  und  $B$ “

## Über A-Bilder, U-Bilder und P-Bilder

### Bemerkungen

1. Die Trennungsrelation im UG ist **monoton** in der Barriere  $Z$ :

$$\text{sep}(A | Z | B) \text{ und } \tilde{Z} \supset Z \Rightarrow \text{sep}(A | \tilde{Z} | B)$$

2. Es gilt die „marginale Trennung“  $\text{sep}(\{a\} | \emptyset | \{b\})$  genau dann, wenn  $a, b \in V$  zu verschiedenen Zusammenhangskomponenten gehören.
3. A-Bild  $\rightsquigarrow$  für adjazente Knoten gilt keinerlei Unabhängigkeit (der diskrete Graph ist A-Bild jedes  $P$ )
4. U-Bild  $\rightsquigarrow$  für nichtadjazente Knoten gilt  $\geq 1$  Unabhängigkeit (der vollständige Graph ist U-Bild jedes  $P$ )
5. Nicht alle Dependenzmodelle  $\mathfrak{I}$  besitzen ein perfektes Bild. Für das nichtmonotone Modell mit zwei Würfeln  $\mathbb{W}_1, \mathbb{W}_2$  und die Signalglocke  $\mathbb{G}$  für Pasch gilt nämlich

$$\mathfrak{I}(\mathbb{W}_1 | \emptyset | \mathbb{W}_2) \text{ und nicht } \mathfrak{I}(\mathbb{W}_1 | \mathbb{G} | \mathbb{W}_2).$$



## Die drei Markoveigenschaften

### Definition

Es sei  $P(\cdot)$  eine Wahrscheinlichkeitsverteilung auf  $V$  und  $\mathfrak{S}$  ihr Dependenzmodell. Der ungerichtete Graph  $\mathcal{G} = (V, \mathcal{E})$  erfüllt die

- paarweise Markoveigenschaft**

gdw. für alle nichtadjazenten  $a, b \in V$  gilt:

$$\mathfrak{S}(a \mid V \setminus \{a, b\} \mid b)$$

- lokale Markoveigenschaft**

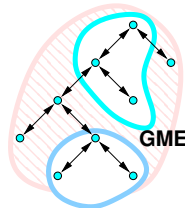
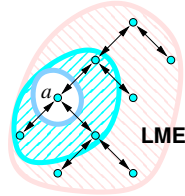
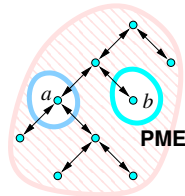
gdw. für alle jede Variable  $a \in V$  gilt:

$$\mathfrak{S}(a \mid \text{bd}(a) \mid V \setminus \text{cl}(a))$$

- globale Markoveigenschaft**

gdw. für alle  $A, B, Z \subset V$  mit  $\text{sep}(A|Z|B)$  gilt:

$$\mathfrak{S}(A \mid Z \mid B)$$



### Beweis.

- GME  $\Rightarrow$  LME**

Es sei  $a \in V$ .

Offensichtlich werden die beiden Mengen  $\{a\}$  und  $V \setminus \text{cl}(a)$  durch den Rand  $\text{bd}(a)$  von  $a$  separiert.

Damit folgt die Behauptung aus der Anwendung von GME.

- LME  $\Rightarrow$  PME**

Zunächst gilt wegen der Voraussetzung LME die Aussage

$$\mathfrak{S}(a \mid \text{bd}(a) \mid V \setminus \text{cl}(a))$$

Wegen der Teilmengenbeziehung

$$V \setminus \{a, b\} = \text{bd}(a) \cup ((V \setminus \text{cl}(a)) \setminus \{b\})$$

kann mittels Axiom WUN

$$\mathfrak{S}(a \mid V \setminus \{a, b\} \mid V \setminus \text{cl}(a))$$

gefolgert werden und mittels Axiom DEC wird verkürzt zu

$$\mathfrak{S}(a \mid V \setminus \{a, b\} \mid b).$$

□

## Die Markoveigenschaften für „Semigraphoide“

Markovnetze  $\hat{=}$  minimale Unabhängigkeitsbilder

### Definition

Der Graph  $\mathcal{G}$  heißt **Markovnetz** von  $P(\cdot)$ , wenn er minimal mit der globalen Markoveigenschaft für  $P(\cdot)$  ist.

Das Markovnetz  $\mathcal{G}$  ignoriert keine Abhängigkeiten, höchstens Unabhängigkeiten, aber davon so wenige wie möglich.

### Satz

Sei  $\mathcal{G} = (V, \mathcal{E})$  und  $P(\cdot)$  auf  $V$  gegeben. Dann gilt

globale ME  $\Rightarrow$  lokale ME  $\Rightarrow$  paarweise ME,

aber es gilt im allgemeinen weder die Umkehrrichtung

paarweise ME  $\Rightarrow$  lokale ME

noch die Umkehrrichtung

lokale ME  $\Rightarrow$  globale ME

### Beweis.

- LME  $\nRightarrow$  GME** (Gegenbeispiel)

$$U - W - X - Y - Z \quad (U = W, Y = Z, X = W \cdot Y)$$

mit binärwertigen, gleichverteilten Variablen.

Es gilt zwar die lokale ME, aber  $\mathfrak{S}(U, W \mid X \mid Y, Z)$  scheitert wegen

$$\begin{aligned} P(U = W = Y = Z = 1 \mid X = 0) &= 0 \\ P(U = W = 1 \mid X = 0) \cdot P(Y = Z = 1 \mid X = 0) &\neq 0 \end{aligned}$$

- PME  $\nRightarrow$  LME** (Gegenbeispiel)

$$X - Y - Z \quad (X = Y = Z)$$

mit binärwertigen, gleichverteilten Variablen.

Dann sind  $\mathfrak{S}(X \mid Z \mid Y)$  und  $\mathfrak{S}(X \mid Y \mid Z)$  trivialerweise erfüllt, überflüssigerweise sogar auch  $\mathfrak{S}(Y \mid X \mid Z)$ . Aber es gilt keineswegs

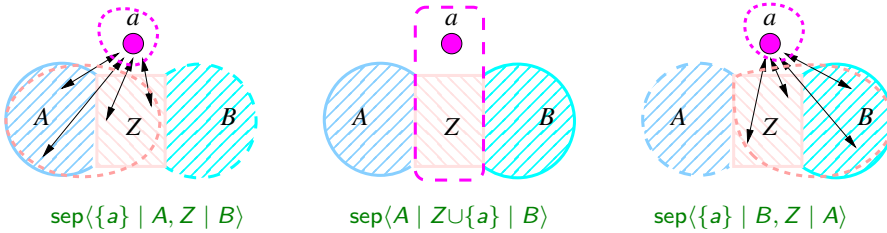
$$\mathfrak{S}(X \mid \text{bd}(X) \mid V \setminus \text{cl}(X))$$

denn  $\text{bd}(X) = \emptyset$  und  $V \setminus \text{cl}(X) = \{Y, Z\}$ , und es ist  $X$  natürlich nicht marginal unabhängig von  $\{Y, Z\}$ .

□

## Die Markoveigenschaften für „Graphoide“

Äquivalenz für strikt positive Wahrscheinlichkeitsverteilungen



### Satz

Sei  $\mathcal{G}$  ein UG. Erfüllt die Dependenzrelation  $\mathfrak{S}$  von  $P(\cdot)$  für alle disjunkten Mengen  $A, B, C, Z \subset V$  die Eigenschaft

- INT Durchschnitt**

$$\mathfrak{S}(A \mid Z \cup C \mid B) \wedge \mathfrak{S}(A \mid Z \cup B \mid C) \Rightarrow \mathfrak{S}(A \mid Z \mid B \cup C),$$

so gilt

globale ME  $\Leftrightarrow$  lokale ME  $\Leftrightarrow$  paarweise ME.

### Beweis.

Es ist nur die Implikation  $\text{PME} \Rightarrow \text{GME}$  zu zeigen, die wir durch absteigende Induktion über die Größe  $n = |Z|$  beweisen.

- Induktionsanfang:**

Für  $n = N - 2$  liefert PME die Behauptung (o.B.d.A. sei  $|A| = |B| = 1$ ).

- Induktionsschluß:**

Wir unterscheiden die beiden Fälle  $A \cup B \cup Z = V$  und  $A \cup B \cup Z \neq V$ .

- Fall 1:** Sei o.B.d.A.  $|A| > 1$  und  $a \in A$ . Dann gelten nach WUN die beiden Trennungsaussagen

$$\text{sep}(A \setminus \{a\} \mid Z \cup \{a\} \mid B), \quad \text{sep}(\{a\} \mid Z \cup A \setminus \{a\} \mid B).$$

Nach I.V. übersetzen diese in die korrespondierenden Unabhängigkeiten und mit Axiom INT folgt  $\mathfrak{S}(A \mid Z \mid B)$ .

- Fall 2:** Für jedes  $a \in V \setminus (A \cup B \cup Z)$  gilt  $\text{sep}(A \mid Z \cup \{a\} \mid B)$  und mindestens eine der beiden Trennungsaussagen

$$\text{sep}(\{a\} \mid A, Z \mid B), \quad \text{sep}(\{a\} \mid B, Z \mid A).$$

Im ersten Fall folgt das Resultat  $\mathfrak{S}(A \mid Z \mid B)$  nach den Axiomen INT, DEC und im zweiten Fall nach den Axiomen SYM, INT, DEC aus den übersetzten Trennungsaussagen (I.V.). □

## Markovnetzkonstruktion

(1:1)-Abbildung aller partiellen (Un-)Abhängigkeiten

### Lemma

Erfüllt die Dependenzrelation  $\mathfrak{S}$  von  $P(\cdot)$  die Axiome SYM, DEC und INT, so gibt es ein **eindeutiges Markovnetz**  $\mathcal{G} = (V, \mathcal{E})$  zu  $\mathfrak{S}$ .

Für alle Variablenpaare  $a, b \in V$  gilt:

$$\{a, b\} \notin \mathcal{E} \Leftrightarrow \mathfrak{S}(a \mid V \setminus \{a, b\} \mid b)$$

### Satz (Pearl & Paz, 1985)

Die Dependenzrelation  $\mathfrak{S}$  ist **graphisch genau** dann, wenn sie die Axiome SYM, DEC, INT, SUN und TRA erfüllt.

- SUN Starke Vereinigung**

$$\mathfrak{S}(A \mid Z \mid B) \Rightarrow \mathfrak{S}(A \mid Z \cup C \mid B)$$

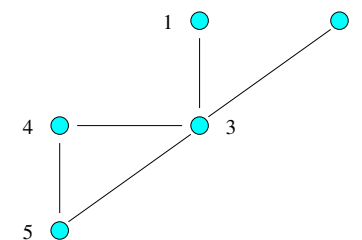
- TRA Transitivität**

Für alle Variablen  $x \in V$  gilt:

$$\mathfrak{S}(A \mid Z \mid B) \Rightarrow \mathfrak{S}(A \mid Z \mid \{x\}) \vee \mathfrak{S}(\{x\} \mid Z \mid B)$$

## Beispiel — qualitative graphische Inferenz

„Vorhersage einer Reiseankunftszeit“



### Uhrzeitwertige Zufallsvariable

Zwei Passanten — zwei Armbanduhren

- $X_1$  Zeit auf Armanduhr I
- $X_2$  Zeit auf Armanduhr II
- $X_3$  die wahre Uhrzeit
- $X_4$  die Fahrtzeit „Jena–Weimar“
- $X_5$  die Ankunftszeit in Weimar

### Markovnetzerzeugung

Kantenlöschverfahren mit den Vorbehalten:

$$\left\{ \begin{array}{l} \neg \mathfrak{S}(X_1 \mid X_2, X_4, X_5 \mid X_3) \\ \neg \mathfrak{S}(X_2 \mid X_1, X_4, X_5 \mid X_3) \end{array} \right\}, \quad \left\{ \begin{array}{l} \neg \mathfrak{S}(X_3 \mid X_1, X_2, X_5 \mid X_4) \\ \neg \mathfrak{S}(X_3 \mid X_1, X_2, X_4 \mid X_5) \\ \neg \mathfrak{S}(X_4 \mid X_1, X_2, X_3 \mid X_5) \end{array} \right.$$

### Inferenz durch Ablesen von Trennungseigenschaften

Bedingte, aber nicht partielle Unabhängigkeiten:  $\mathfrak{S}(X_1, X_2 \mid X_3 \mid X_5)$

## Beispiel — Würfelpaar und Glocke

### Nichtgraphische Verteilungen

Viele interessante Verteilungen liegen außerhalb der Klasse ungerichteter graphischer Modelle.

- Selbst ein streng positives  $P(\cdot)$  garantiert lediglich die Axiome SYM, DEC und INT, nicht aber SUN oder TRA.

### Würfel-Glocken-Experiment

Es schlägt die starke Vereinigung (SUN) fehl:

$$\mathfrak{S}(W_1 \mid \emptyset \mid W_2) \quad \text{aber} \quad \neg \mathfrak{S}(W_1 \mid G \mid W_2)$$

Bei **unfairen Würfeln** gilt auch keine Transitivität (TRA) mehr:

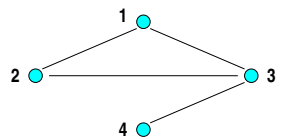
$$\mathfrak{S}(W_1 \mid \emptyset \mid W_2) \quad \text{aber weder} \quad \mathfrak{S}(W_1 \mid \emptyset \mid G) \quad \text{noch} \quad \mathfrak{S}(G \mid \emptyset \mid W_2)$$

### Bemerkung

Die drei Axiome DEC, INT, SUN liefern eine beachtliche Äquivalenz:

$$\mathfrak{S}(A \mid Z \mid B) \quad \Leftrightarrow \quad \forall a \in A, b \in B : \mathfrak{S}(\{a\} \mid Z \mid \{b\})$$

## Beispiel — pathologische Verteilung ohne Markovnetz



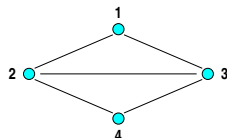
Kantenlöschverfahren

### Dependenzstruktur

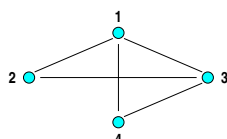
Gegeben sind bedingte Unabhängigkeiten

$$\mathfrak{S}(1 \mid 2, 3 \mid 4) \quad \mathfrak{S}(2 \mid 1, 3 \mid 4)$$

zuzüglich aller Symmetrien.



minimales U-Bild #1

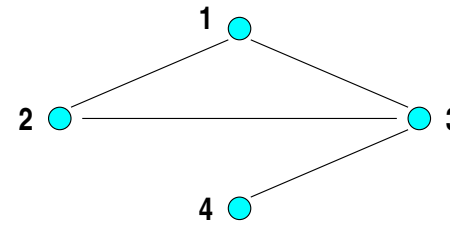


minimales U-Bild #2

### Eigenschaften

- $\mathfrak{S}$  erfüllt die Axiome SYM, DEC, WUN, CON.
- $\mathfrak{S}$  widerspricht dem Axiom INT, weil  $\mathfrak{S}(1, 2 \mid 3 \mid 4)$  fehlt.
- $\mathfrak{S}$  gehorcht einer Verteilung  $P$ , aber  $P$  ist wegen  $\neg \text{INT}$  nicht streng positiv!
- Das Kantenlöschverfahren ergibt kein Unabhängigkeitsbild, weil  $\text{sep}(1|3|4)$  gilt, aber nicht  $\mathfrak{S}(1|3|4)$ .
- Es gibt kein eindeutiges Markovnetz!

## Beispiel — eine Unverteilung mit Markovnetz



Eindeutiges Markovnetz

### Dependenzstruktur

Gegeben sind die bedingten „Unabhängigkeiten“

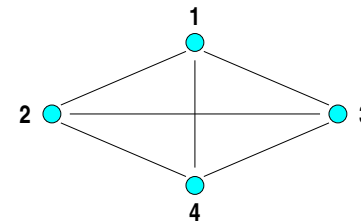
$$\begin{array}{ll} \mathfrak{S}(1 \mid 2 \mid 3) & \mathfrak{S}(1, 2 \mid 3 \mid 4) \\ \mathfrak{S}(1 \mid 3 \mid 4) & \mathfrak{S}(1 \mid 2, 3 \mid 4) \\ \mathfrak{S}(2 \mid 3 \mid 4) & \mathfrak{S}(2 \mid 1, 3 \mid 4) \end{array}$$

zuzüglich aller Symmetrien.

### Eigenschaften

- $\mathfrak{S}$  erfüllt die Axiome SYM, DEC, WUN und INT.
- $\mathfrak{S}$  widerspricht dem Axiom CON, denn es gelten zwar  $\mathfrak{S}(1 \mid 2 \mid 3)$  und  $\mathfrak{S}(1 \mid 2, 3 \mid 4)$ , aber keineswegs  $\mathfrak{S}(1 \mid 2 \mid 3, 4)$ .
- $\mathfrak{S}$  besitzt wegen  $\neg \text{CON}$  kein Wahrscheinlichkeitsmodell mit  $\mathfrak{S} = \mathfrak{S}_P$ .
- $\mathfrak{S}$  besitzt aber wegen SYM, DEC, INT ein eindeutiges Markovnetz.

## Beispiel — Unverteilung mit Monsternetz



vollständiger Graph

### Dependenzstruktur

Gegeben sind die Postulate

$$\begin{array}{ll} \mathfrak{S}(1 \mid 3 \mid 4) & \mathfrak{S}(2 \mid 3 \mid 4) \\ \mathfrak{S}(1, 2 \mid 3 \mid 4) & \mathfrak{S}(1 \mid 2 \mid 4) \end{array}$$

zuzüglich aller Symmetrien.

### Eigenschaften

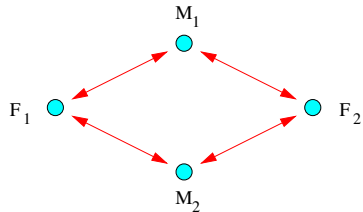
- $\mathfrak{S}$  erfüllt die Axiome SYM, DEC, INT.
- $\mathfrak{S}$  erfüllt nicht das Axiom WUN.
- Wegen SYM, DEC, INT gibt es ein eindeutiges Markovnetz  $\mathcal{G}$ .
- Der Graph  $\mathcal{G}$  ist offenbar (keine Löschung) **vollständig**.
- Der Graph  $\mathcal{G}$  „hilft uns nicht sparen“ ...

## Faktorisierung von $P(\cdot)$

über den Cliques eines ungerichteten Graphen

### Definition

Die Menge  $C \subseteq V$  heißt **Clique** von  $\mathcal{G} = (V, \mathcal{E})$ , wenn  $(C, \mathcal{E}|_C)$  einen maximalen zusammenhängenden Teilgraphen von  $\mathcal{G}$  bildet.



### Der diamantene Graph

besitzt genau vier Cliques:

$$C_1 = \{M_1, F_1\}, \quad C_2 = \{M_1, F_2\},$$

$$C_3 = \{M_2, F_1\}, \quad C_4 = \{M_2, F_2\}$$

### Gibbs-Verteilung des PT-Modells über dem Diamanten

$$P(m_1, m_2, f_1, f_2) = \frac{1}{Z} \cdot \phi_1(m_1, f_1) \cdot \phi_2(m_1, f_2) \cdot \phi_3(m_2, f_1) \cdot \phi_4(m_2, f_2)$$

mit den **Kompatibilitäts-** oder **Kernfunktionen** (keine Wahrscheinlichkeiten!)

$$\phi_i(\xi_{i_1}, \xi_{i_2}) = \begin{cases} \alpha_i & \xi_{i_1} = \xi_{i_2} \text{ gleicher Gesundheitszustand} \\ \beta_i & \xi_{i_1} \neq \xi_{i_2} \text{ genau ein Partner infiziert} \end{cases}$$

## Faktorisierungs- und Markoveigenschaften

### Lemma

Für alle ungerichteten Graphen  $\mathcal{G} = (V, \mathcal{E})$  und für alle Wahrscheinlichkeitsmodelle  $P : \mathcal{X}_V \rightarrow \mathbb{R}$  gilt:

$$\boxed{\text{FAK} \Rightarrow \text{GME}} \Rightarrow \text{LME} \Rightarrow \text{PME}$$

### Satz (Hammersley & Clifford, 1971)

Für jede streng positive Wahrscheinlichkeitsverteilung  $P(\cdot)$  und jeden ungerichteten Graphen  $\mathcal{G}$  gilt:

$$\boxed{\text{FAK} \Leftrightarrow \text{GME}} \Leftrightarrow \text{LME} \Leftrightarrow \text{PME}$$

### Bemerkung

Im Falle numerischer Zufallsvariablen ist als Voraussetzung des HC-Satzes auch die **Existenz und Stetigkeit** der Dichtefunktion  $f : \mathcal{X}_V \rightarrow \mathbb{R}$  zu fordern.

## FAK — die Faktorisierungseigenschaft

### Definition

Die Wahrscheinlichkeitsverteilung  $P(\cdot)$  **zerfällt über dem Graphen**  $\mathcal{G} = (V, \mathcal{E})$ , wenn es für jede vollständige Menge  $A \subset V$  eine nichtnegative **Kernfunktion**

$$\phi_A : \bigotimes_{a \in A} \mathcal{X}_a \rightarrow \mathbb{R}_0^+$$

über dem kartesischen Produkt aller  $A$ -Wertebereiche gibt mit

$$P(\mathbf{x}) = \prod_{A \text{ vollständig}} \phi_A(\mathbf{x}_A)$$

O.B.d.A. können wir diese Faktorisierungseigenschaft (FAK) aber auch unter Beschränkung auf die Menge  $\mathcal{C}(\mathcal{G})$  der **Cliques** von  $\mathcal{G}$  definieren:

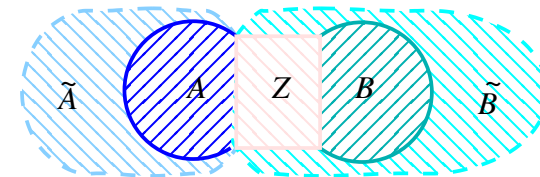
$$P(\mathbf{x}) = \prod_{A \in \mathcal{C}(\mathcal{G})} \phi_A(\mathbf{x}_A)$$

### Beweis.

FAK  $\Rightarrow$  GME

Es seien  $A, B, Z \subset V$  disjunkt mit  $\text{sep}(A \mid Z \mid B)$ .

Sei  $\tilde{A}$  die Zusammenhangshülle von  $A$  in  $\mathcal{G}_{V \setminus Z}$  und sei  $\tilde{B} = V \setminus (\tilde{A} \cup Z)$ .



$A, B$  gehören sicherlich zu verschiedenen Zusammenhangskomponenten im Restgraphen  $\mathcal{G}_{V \setminus Z}$ , also gilt für jede Clique  $C \subset V$  genau eine der beiden Bedingungen

$$C \subseteq \tilde{A} \cup Z \quad \text{oder} \quad C \subseteq \tilde{B} \cup Z.$$

Die (garantierte: FAK) Faktorisierung gewinnt damit das folgendes Aussehen:

$$P(\mathbf{x}) = \prod_{C \in \mathcal{C}} \phi_C(\mathbf{x}_C) = \prod_{C \in \mathcal{C}_A} \phi_C(\mathbf{x}_C) \cdot \prod_{C \in \mathcal{C}_B} \phi_C(\mathbf{x}_C) = g(\mathbf{x}_{\tilde{A} \cup Z}) \cdot h(\mathbf{x}_{\tilde{B} \cup Z})$$

Nach Definition der bedingten Unabhängigkeit folgt daraus  $\mathfrak{I}(\tilde{A} \mid Z \mid \tilde{B})$  und nach zweimaliger Anwendung des Axioms DEC auch die GME-Behauptung  $\mathfrak{I}(A \mid Z \mid B)$ . □

FAK  $\nleftrightarrow$  GME für pathologische  $P(\cdot)$ 

Moussouris (1974)

## Beweis.

GME  $\Rightarrow$  FAK

Aus völlig trivialen Gründen (auch  $V \subseteq V$ ) gibt es eine Mammut-Faktorisierung à la

$$P(\mathbf{x}) = \prod_{A \subseteq V} \phi_A(\mathbf{x}_A).$$

Wegen der Eigenschaft  $P(\mathbf{x}) > 0$  strenger Positivität läßt sich diese Darstellung schmerzfrei logarithmieren:

$$\log P(\mathbf{x}) = \sum_{A \subseteq V} \log \phi_A(\mathbf{x}_A)$$

Nach einer sogenannten „*Möbius-Inversion*“ (sehr schwierig!) lassen sich in obigem Ausdruck durch Faktorisierung nach partiellen Unabhängigkeiten Zug um Zug alle Nicht-Cliquen-Summanden eliminieren.  $\square$

## Gegenbeispiel

Betrachte  $V = \{\mathbb{X}_1, \mathbb{X}_2, \mathbb{X}_3, \mathbb{X}_4\}$  und die Verteilung

$$P(\mathbf{x}) = \begin{cases} 1/8 & \mathbf{x} \in \begin{pmatrix} 0 & 1 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \end{pmatrix} \\ 0 & \text{sonst} \end{cases}$$

Für alle  $(x_2, x_4)$  ist entweder  $P(\mathbb{X}_1 | x_2, x_4)$  oder  $P(\mathbb{X}_3 | x_2, x_4)$  eine degenerierte Abbildung, also besteht trivialerweise keinerlei Abhängigkeit von  $\mathbb{X}_3$  bzw.  $\mathbb{X}_1$ . Gleiches gilt auch für alle  $(x_1, x_3)$ , also gilt insgesamt

$$\mathfrak{I}(\mathbb{X}_1 | \mathbb{X}_2, \mathbb{X}_4 | \mathbb{X}_3) \quad \wedge \quad \mathfrak{I}(\mathbb{X}_2 | \mathbb{X}_1, \mathbb{X}_3 | \mathbb{X}_4)$$

Der **Diamant** besitzt die GME, aber  $P(\cdot)$  ist nicht  $\diamond$ -faktorisierbar:

$$\begin{aligned} 0 \neq 1/8 &= P(0, 0, 0, 0) = \phi_{1,2}(0, 0) \cdot \phi_{2,3}(0, 0) \cdot \phi_{3,4}(0, 0) \cdot \phi_{4,1}(0, 0) \\ 0 &= P(0, 0, 1, 0) = \phi_{1,2}(0, 0) \cdot \phi_{2,3}(0, 1) \cdot \phi_{3,4}(1, 0) \cdot \phi_{4,1}(0, 0) \\ 0 \neq 1/8 &= P(0, 0, 1, 1) = \phi_{1,2}(0, 0) \cdot \phi_{2,3}(0, 1) \cdot \phi_{3,4}(1, 1) \cdot \phi_{4,1}(1, 0) \\ 0 \neq 1/8 &= P(1, 1, 1, 0) = \phi_{1,2}(1, 1) \cdot \phi_{2,3}(1, 1) \cdot \phi_{3,4}(1, 0) \cdot \phi_{4,1}(0, 1) \end{aligned}$$

## Zerlegbare graphische Modelle

## Gibbs-Verteilungen

Verteilungen in Cliquenproduktform:

$$P(\mathbf{x}) = P(\mathbf{x}_V) = \prod_{C \in \mathcal{C}} \phi_C(\mathbf{x}_C) / \sum_{\mathbf{x} \in \Omega} \prod_{C \in \mathcal{C}} \phi_C(\mathbf{x}_C)$$

Die Potentialfunktionen  $\phi_C(\cdot)$  sind i.a. **keine** (⚡) Wahrscheinlichkeiten.

## Zerlegbarkeit

Wann zerfällt  $P(\mathbf{x})$  in ein Produkt **bedingter Randverteilungen** ?

- Wenn die Cliquen des Modellgraphen als Baum angeordnet sind!
- Die Baumstruktur regelt die Abhängigkeitsrichtungen.

Es besteht Freiheit in der Wahl, welche **Außencliquen** ein Blatt und welche eine Wurzel werden.

## Beispiel — Markovketten I

Faktorisierung mit unterschiedlicher Variablenordnung

$$\mathbb{X}_1 \longleftrightarrow \mathbb{X}_2 \longleftrightarrow \mathbb{X}_3 \longleftrightarrow \mathbb{X}_4$$

## Faktorisierung = Kettenregel + Unabhängigkeiten

$$P(x_1, x_2, x_3, x_4) = P(x_1) \cdot P(x_2 | x_1) \cdot \underbrace{P(x_3 | x_1, x_2)}_{P(x_3 | x_2)} \cdot \underbrace{P(x_4 | x_1, x_2, x_3)}_{P(x_4 | x_3)}$$

Jede Variable kann als **Baumwurzel** nominiert werden — so auch  $\mathbb{X}_3$ :

$$P(x_3, x_2, x_4, x_1) = P(x_3) \cdot P(x_2 | x_3) \cdot \underbrace{P(x_4 | x_3, x_2)}_{P(x_4 | x_3)} \cdot \underbrace{P(x_1 | x_3, x_2, x_4)}_{P(x_1 | x_2)}$$

Aber **nicht jede Variablenfolge** ist mit der Baumstruktur verträglich:

$$P(x_1, x_4, x_2, x_3) = P(x_1) \cdot \underbrace{P(x_4 | x_1)}_{\text{⚡}} \cdot \underbrace{P(x_2 | x_1, x_4)}_{\text{⚡}} \cdot \underbrace{P(x_3 | x_1, x_4, x_2)}_{\text{⚡}}$$

## Beispiel — Markovketten II

Faktorisierung mit unterschiedlichen Cliquenbäumen

$$(\mathbb{X}_1, \mathbb{X}_2) \longleftrightarrow (\mathbb{X}_2, \mathbb{X}_3) \longleftrightarrow (\mathbb{X}_3, \mathbb{X}_4)$$

Faktorisierung = Cliquen + Baum + Wurzelauswahl

$$P(x_1, x_2, x_3, x_4) = f(x_1, x_2) \cdot g(x_2, x_3) \cdot h(x_3, x_4)$$

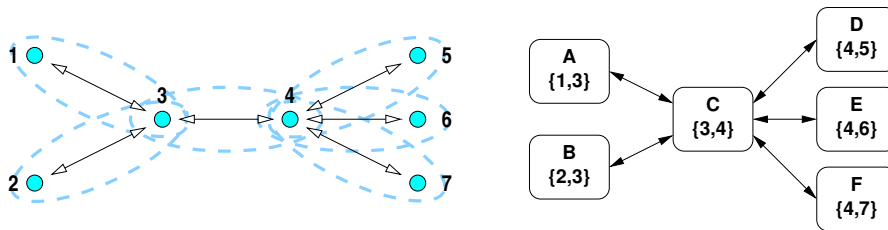
$$P(x_1, x_2) \cdot P(x_3|x_2) \cdot P(x_4|x_3)$$

$$P(x_1|x_2) \cdot P(x_2, x_3) \cdot P(x_4|x_3)$$

$$P(x_1|x_2) \cdot P(x_2|x_3) \cdot P(x_3, x_4)$$

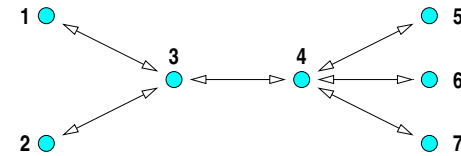
➡ Jede Wurzelnominierung definiert eine **valide** Modellformel.

## Beispiel — Markovbäume II



$$\begin{aligned} P(\mathbf{x}) &= \frac{\text{Cliquenwahrscheinlichkeit}}{\text{Cliquenschnittwahrscheinlichkeit}} \\ &= \frac{P(A) \cdot P(B) \cdot P(C) \cdot P(D) \cdot P(E) \cdot P(F)}{P(A \cap C) \cdot P(B \cap C) \cdot P(C \cap D) \cdot P(C \cap E) \cdot P(C \cap F)} \\ &= \frac{P(1, 3) \cdot P(2, 3) \cdot P(3, 4) \cdot P(4, 5) \cdot P(4, 6) \cdot P(4, 7)}{P(3) \cdot P(3) \cdot P(4) \cdot P(4) \cdot P(4)} \\ &= P(3) \cdot P(1|3) \cdot P(2|3) \cdot P(4|3) \cdot P(5|4) \cdot P(6|4) \cdot P(7|4) \end{aligned}$$

## Beispiel — Markovbäume I



**Fakt**

Ist  $\mathcal{G}$  ein Baum, so sind alle Cliquen zweielementig.

Die  $N - 1$  Cliquen bilden selbst wieder einen Baum.

Faktorisierung im Beispiel

Kettenregel & Variablenbaumtraversierung

$$P(\mathbf{x}) = P(3) \cdot P(1|3) \cdot P(2|3) \cdot P(4|3) \cdot P(5|4) \cdot P(6|4) \cdot P(7|4)$$

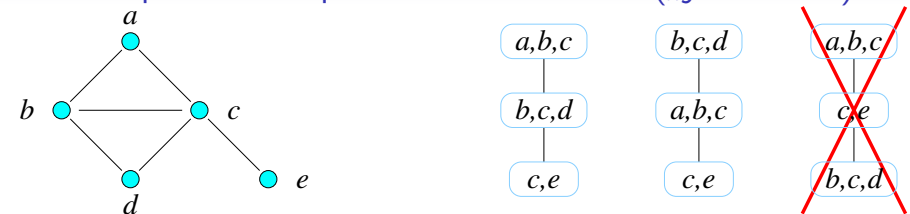
Faktorisierung allgemein

Traversieren  $\rightsquigarrow$  konsistente Variablenordnung  $\rightsquigarrow$  Einfachbedingungen

$$P(\mathbf{x}) = \prod_n P(x_n | \cdot) = \prod_n P(x_n | x_{\pi(n)})$$

Denn für alle  $\mathbb{X}_n \in V$  gilt:  $\text{sep}(\mathbb{X}_n | \mathbb{X}_{\pi(n)} | V \setminus \text{off}(\mathbb{X}_n))$

## Beispiel — Cliquenverbundbaum („join tree“)



Drei Cliquen — aber welche Baumstruktur?

Die Cliquen  $C_1 = \{a, b, c\}$ ,  $C_2 = \{b, c, d\}$ ,  $C_3 = \{c, e\}$  bilden paarweise einen nichtleeren Schnitt.

➡  $\mathfrak{I}(C_1|C_2|C_3) \rightsquigarrow C_1 - C_2 - C_3$  ist U-Bild von  $P(\cdot)$

$\mathfrak{I}(C_2|C_1|C_3) \rightsquigarrow C_2 - C_1 - C_3$  ist U-Bild von  $P(\cdot)$

➡  $C_1 - C_3 - C_2$  ist **nicht** U-Bild von  $P(\cdot)$ , da  $\neg \mathfrak{I}(C_1|C_3|C_2)$

$\mathfrak{I}(C_1|C_2|C_3)$  und  $\mathfrak{I}(C_2|C_1|C_3) \rightsquigarrow \mathfrak{I}(C_1, C_2 | \emptyset | C_3)$  (INT) verletzt, also  $\geq 2$  minimale U-Bilder.

Zerlegung

Beide konsistenten Verbundbäume ergeben nach Traversierung:

$$P(a, b, c, d, e) = P(a) \cdot P(b|a) \cdot P(c|a, b) \cdot P(d|b, c) \cdot P(e|c)$$

## Kordalität und Zerlegbarkeit

Äquivalente Eigenschaften ungerichteter Graphen

### Definition

Das Mengentripel  $(A, Z, B)$  heißt **Zerlegung des ungerichteten Graphen**  $\mathcal{G} = (V, \mathcal{E})$ , falls gilt:

Partition

$$A \uplus Z \uplus B = V$$

Trennung

$$\text{sep}\langle A \mid Z \mid B \rangle$$

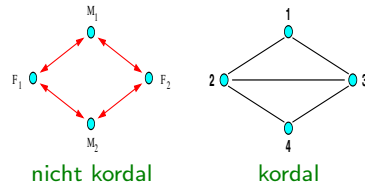
Vollständigkeit

$\mathcal{G}_Z$  ist vollständig

Der Graph  $\mathcal{G}$  selbst heißt **zerlegbar**, wenn er vollständig ist oder aber eine Zerlegung mit zerlegbaren Teilgraphen  $\mathcal{G}_{A \cup Z}$  und  $\mathcal{G}_{B \cup Z}$  besitzt.

### Definition

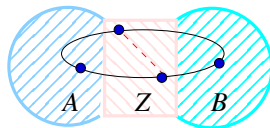
Ein ungerichteter Graph  $\mathcal{G} = (V, \mathcal{E})$  heißt **kordal** oder **trianguliert** genau dann, wenn jeder Zyklus der Länge  $\geq 4$  mindestens eine **Sehne** besitzt.



### Beweis.

Wir zeigen die Implikation „zerlegbar  $\Rightarrow$  kordal“

- **Induktionsanfang:**  
 $|V| \leq 3$  impliziert trivialerweise die Kordalität.
- **Induktionsschritt:** Sei also  $(A, Z, B)$  eine Zerlegung von  $\mathcal{G}$ . Nach Induktionsvoraussetzung sind dann  $\mathcal{G}_{A \cup Z}$  und  $\mathcal{G}_{B \cup Z}$  kordal.



Angenommen,  $\mathcal{G}$  besitzt einen Zyklus  $\geq 4$  ohne Sehne. Dieser muß wegen der I.V. Knoten in  $A$  und auch in  $B$  haben, passiert also mindestens  $2 \times$  die Menge  $Z$  und teilt deshalb  $\geq 2$  Knoten mit  $Z$ . Diese sind aber wegen der Vollständigkeit von  $Z$  verbunden — ⚡

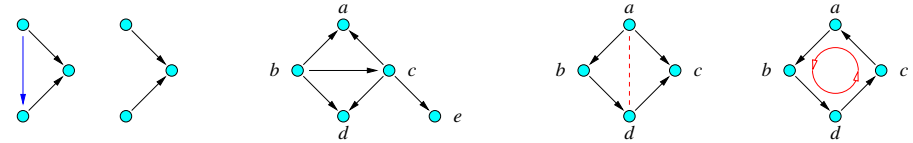
□

## Moralische Graphen

„Alle Elternpaare sind (miteinander!) verheiratet.“

### Definition

Ein gerichteter Graph heißt **moralisch**, wenn jedes konvergierende Kantenpaar aus zwei adjazenten Knoten entspringt.



### Satz

Für einen ungerichteten Graphen  $\mathcal{G}$  sind die Eigenschaften äquivalent:

1.  $\mathcal{G}$  ist zerlegbar.
2.  $\mathcal{G}$  ist kordal.
3.  $\mathcal{G}$  läßt sich azyklisch und moralisch richten.
4.  $\mathcal{G}$  besitzt die Cliqueneliminationseigenschaft.
5. Es gibt einen verträglichen Verbundbaum für  $\mathcal{G}$ .

## Cliqueneliminationseigenschaft

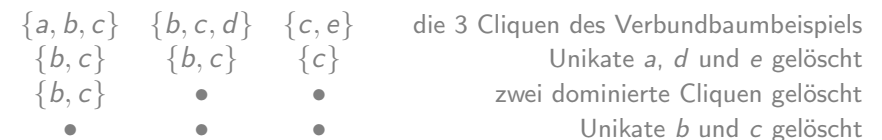
### Definition

Der ungerichtete Graph  $\mathcal{G}$  besitzt die **Cliqueneliminationseigenschaft**, wenn alle Knoten aller seiner Cliques durch wiederholte Anwendung folgender Operationen eliminiert werden können:

- **Unikatknoten**  
Lösche einen Knoten, der nur in einer einzigen Clique auftaucht.
- **Dominierte Mengen**  
Lösche eine Clique, die Teilmenge einer anderen Clique ist.

### Der Schlüsselgraph besitzt die CEP

Schachmatt in sieben Zügen:





## Verträgliche Verbundbäume

### Definition

Sei  $\mathcal{G} = (V, \mathcal{E})$  ein ungerichteter Graph. Der Graph  $\mathcal{G}^*$  ist ein **mit  $\mathcal{G}$  verträglicher Verbundbaum**, falls gilt:

1. Die Knoten  $\mathcal{G}^*$  sind genau die Cliques  $\mathcal{C}(\mathcal{G})$ .
2.  $\mathcal{G}^*$  ist zusammenhängend und zyklensfrei.
3. Für jeden Knoten  $a \in V$  gilt:  
Je zwei  $a$  enthaltende Cliques besitzen einen Verbindungspfad, der ausschließlich Cliques  $C$  mit  $a \in C$  enthält.

### Für den Schlüsselgraphen ist VB #3 nicht verträglich

Im dritten Verbundbaum



gilt  $b \in C_1$  und  $b \in C_2$ , aber es ist  $b \notin C_3$ , obwohl  $C_3$  auf dem einzigen verfügbaren Pfad von  $C_1$  nach  $C_2$  liegt.

### Beweis.

**GME  $\Rightarrow$  FAK**

(die Umkehrung gilt ja sowieso)

Induktion über die Zerlegungshierarchie von  $\mathcal{G}$ :

Sei  $\text{sep}\langle A|Z|B \rangle$  eine Zerlegung von  $\mathcal{G}$ . Dann gilt

$$\begin{aligned} P(\mathbf{x}_V) &= P(\mathbf{x}_{A \cup Z}) \cdot P(\mathbf{x}_B | \mathbf{x}_{A \cup Z}) \\ &= P(\mathbf{x}_{A \cup Z}) \cdot P(\mathbf{x}_B | \mathbf{x}_Z) \\ &= \frac{P(\mathbf{x}_{A \cup Z}) \cdot P(\mathbf{x}_{B \cup Z})}{P(\mathbf{x}_Z)} \end{aligned}$$

wegen  $\mathfrak{S}(A|Z|B)$  nach GME.

Der Nenner  $P(\mathbf{x}_Z)$  ist bereits ein Cliquenfaktor, weil  $Z$  vollständig ist.

Die beiden Zählerterme sind nach Induktionsvoraussetzung über  $\mathcal{G}_{A \cup Z}$  bzw.  $\mathcal{G}_{B \cup Z}$  faktorisiert, bestehen also ausschließlich aus Cliquentermen.

Die behauptete Faktorisierung ergibt sich durch Zusammenfassen und Umgruppieren nach  $\mathcal{G}$ -Cliques. □

## Zerlegbarkeit und Faktorisierung

### Lemma (Cliquenschnittformel)

Sei  $\mathcal{G}$  ein zerlegbarer ungerichteter Graph. Dann gilt für alle  $P(\cdot)$

**FAK**  $\Leftrightarrow$  **GME**

und diese Faktorisierung besteht aus cliquenbezogenen Randverteilungen:

$$P(\mathbf{x}) = \prod_{C \in \mathcal{C}} \frac{P(\mathbf{x}_C)}{P(\mathbf{x}_{C \cap \pi(C)})}$$

Dabei bezeichnet  $\pi(C)$  die eindeutig bestimmte Vorgängerclique von  $C$  im (festen, aber beliebigen) verträglichen Verbundbaum.

### Die beiden CSF für den Schlüsselgraphen

Die Verbundbäume  $\{a, b, c\} - \{b, c, d\} - \{c, e\}$  und  $\{b, c, d\} - \{a, b, c\} - \{c, e\}$  liefern die äquivalenten Faktorisierungen

$$P(\cdot) = \frac{P(a, b, c) \cdot P(b, c, d) \cdot P(c, e)}{P(b, c) \cdot P(c)} \quad \text{und} \quad P(\cdot) = \frac{P(b, c, d) \cdot P(a, b, c) \cdot P(c, e)}{P(b, c) \cdot P(c)}.$$

### Beweis.

#### Cliquenschnittformel

Es sei  $C_1, \dots, C_M$  eine mit der Nachfolgerrelation eines Cliquenverbundbaums von  $\mathcal{G}$  verträgliche Cliquenordnung.

Für jede Clique  $C_i$  bezeichne  $C_{\pi(i)}$  die eindeutig bestimmte Elterclique ( $\pi(i) < i$ ). Dann gilt (für alle  $i$ ) die Trennungsrelation

$$\text{sep}\langle C_i | C_{\pi(i)} | C_1, \dots, C_{i-1} \rangle$$

und wegen GME auch die entsprechende bedingte Unabhängigkeit.

$$\begin{aligned} P(\mathbf{x}) &= P(x_1, \dots, x_N) = \prod_{i=1}^M P(\mathbf{x}_{C_i} | \mathbf{x}_{C_1}, \dots, \mathbf{x}_{C_{i-1}}) \\ &= \prod_{i=1}^M P(\mathbf{x}_{C_i} | \mathbf{x}_{C_{\pi(i)}}) \\ &= \prod_{i=1}^M P(\mathbf{x}_{C_i} | \mathbf{x}_{C_i \cap C_{\pi(i)}}) \\ &= \prod_{i=1}^M \frac{P(\mathbf{x}_{C_i})}{P(\mathbf{x}_{C_i \cap C_{\pi(i)}})} \end{aligned}$$



# Graphtriangulierung & Verbundbaumkonstruktion

(Algorithmus)

## 1 KNOTENORDNUNG

Ordne Knoten nach maximalem Rang; setze sukzessiv:

$$v_{i+1} \stackrel{\text{def}}{=} \operatorname{argmax}_{v \notin V(i)} |\{v' \in V \mid (v, v') \in \mathcal{E}, v' \in V(i)\}|$$

## 2 KANTENERZEUGUNG

Für  $i = N, \dots, 1$

$$\mathcal{E} \leftarrow \mathcal{E} \cup \{v', v''\}$$

falls  $v', v'' \in V(i-1)$  und falls  $\{v_i, v'\}, \{v_i, v''\} \in \mathcal{E}$ .

## 3 CLIQUENORDNUNG

Fixiere Reihenfolge  $C_1, \dots, C_M$  nach dem maximalen Knotenrang.

## 4 KANTENERZEUGUNG

Für  $i = 2, \dots, M$  erzeuge neue Kante  $C_{\pi(i)} \rightarrow C_i$  mit

$$\pi(i) < i \quad \text{und} \quad |C_{\pi(i)} \cap C_i| \text{ ist maximal.}$$

(sumitriogIA)

## Beispiel

Knotenfolge:  
 $a^0 b^1 c^2 d^2 e^1$

Neue  
Kanten:  
(keine)

Cliquenfolge:  
 $C_1 : abc^{012}$   
 $C_2 : bcd^{122}$   
 $C_3 : ce^{21}$   
(beliebig)

VB-Kanten:  
 $1 \rightarrow 2,$   
 $1 \rightarrow 3$

# Zwischenbilanz

für ungerichtete graphische Modelle

1. Nicht jede Verteilung ist graphisch.
2. Streng positive Verteilungen erlauben aber, mit dem Kantenlöschverfahren ein Markovnetz (minimales U-Bild) zu erzeugen.
3. Markovnetze faktorisieren gemäß ihrer Cliquenstruktur, aber nicht zwingend in Wahrscheinlichkeiten.
4. Durch Triangulieren des Markovnetzes werden einige Unabhängigkeiten außer Gefecht gesetzt, aber dafür gewinnen wir eine Kettenregel (CSF).

## Korrelation, Regression und Transinformation

## Assoziationsregeln und Netzwerkanalyse

## Bedingte statistische Unabhängigkeit

## Graphische Modelle: ungerichtete Graphen

## Kausale Modelle: gerichtete azyklische Graphen

## Berechnen bedingter Wahrscheinlichkeiten

## Parameterschätzung in Bayesnetzen und Loglinearmodellen

## Aufdeckung der Abhängigkeitsstruktur

## Kovarianzselektion

# Ursache und Wirkung

Gerichtete azyklische Graphen

## Kausalrichtung

Drei Attribute · zwei Interaktionen · drei Wirkkonfigurationen:

kaskadierend „Wetter“  $\rightarrow$  „Ernte“  $\rightarrow$  „Preis“  $\Im(\mathbb{X}_1|\mathbb{X}_2|\mathbb{X}_3)$

divergierend „Größe“  $\leftarrow$  „Alter“  $\rightarrow$  „Lesefähigkeit“  $\Im(\mathbb{X}_1|\mathbb{X}_2|\mathbb{X}_3)$

konvergierend „Würfel<sub>1</sub>“  $\rightarrow$  „Glocke“  $\leftarrow$  „Würfel<sub>2</sub>“  $\neg \Im(\mathbb{X}_1|\mathbb{X}_2|\mathbb{X}_3)$

Modelle kausaler Beziehungen:  $\left\{ \begin{array}{l} \text{erklärende} \\ \text{vermittelnde} \\ \text{diagnostische} \end{array} \right\}$  Variablen.

## Lemma (Erinnerung)

Ein gerichteter Graph  $\mathcal{G} = (V, \mathcal{E})$ ,  $\mathcal{E} \subseteq V \times V$ , ist **azyklisch** genau dann, wenn es eine **kantenverträgliche Knotenordnung**  $V = \{v_1, \dots, v_N\}$  gibt:

$$(v_i, v_j) \in \mathcal{E} \Rightarrow i < j \quad \text{für alle } i, j \in \{1, \dots, N\}$$

M.a.W.: Ein DAG („directed acyclic graph“) besitzt keine gerichteten Zyklen (**Pfade**); ungerichtete Zyklen (**Ketten**) sind hingegen erlaubt.

## $\delta$ -Trennungsrelation für gerichtete azyklische Graphen

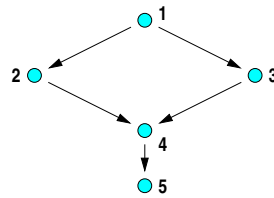
### Definition

Es sei  $\mathcal{G} = (V, \mathcal{E})$  ein gerichteter azyklischer Graph und  $A, B, Z \subset V$  disjunkte Knotenmengen. Eine **Kette** zwischen den Knoten  $a$  und  $b$  heißt **blockiert von  $Z$**

- wenn sie einen nichtkonvergierenden Knoten  $c \in Z$  enthält
- oder wenn sie einen konvergierenden Knoten  $c \notin Z$  enthält, der auch keinen Nachfolger in  $Z$  besitzt.

Die Menge  $Z$  **trennt  $A$  von  $B$**  genau dann, wenn alle Ketten zwischen Elementen  $a \in A$  und  $b \in B$  von Elementen aus  $Z$  blockiert werden. Wir schreiben dafür:

$$\text{sep}_\delta \langle A \mid Z \mid B \rangle$$



### Beispiel

Es gilt  $\text{sep}_\delta \langle 2 \mid 1 \mid 3 \rangle$ , denn:

Die Kette  $2 \leftarrow 1 \rightarrow 3$  ist von  $\mathbb{X}_1$  blockiert wegen  $1 \in Z = \{1\}$

Die Kette  $2 \rightarrow 4 \leftarrow 3$  ist von  $\mathbb{X}_4$  blockiert wg.  $4, 5 \notin Z = \{1\}$

## Kausale Verteilungen und Dependenzmodelle

Überrepräsentation & Unterrepräsentation von  $\mathfrak{S}(\cdot \mid \cdot \mid \cdot)$  durch  $\text{sep}_\delta \langle \cdot \mid \cdot \mid \cdot \rangle$

### Definition

Es sei  $P(\cdot)$  eine Wahrscheinlichkeitsverteilung auf  $V$  und  $\mathfrak{S}$  ihr Dependenzmodell. Der gerichtete azyklische Graph  $\mathcal{G} = (V, \mathcal{E})$  heißt

- Abhängigkeitsbild** von  $P$  gdw.

$$\mathfrak{S}(A \mid Z \mid B) \Rightarrow \text{sep}_\delta \langle A \mid Z \mid B \rangle$$

- Unabhängigkeitsbild** von  $P$  gdw.

$$\mathfrak{S}(A \mid Z \mid B) \Leftarrow \text{sep}_\delta \langle A \mid Z \mid B \rangle$$

- perfektes Bild** von  $P$  gdw.

$$\mathfrak{S}(A \mid Z \mid B) \Leftrightarrow \text{sep}_\delta \langle A \mid Z \mid B \rangle$$

Die Verteilung  $P(\cdot)$  (und das Modell  $\mathfrak{S}$ ) heißen **kausal**, wenn ein gerichteter azyklischer Graph existiert, der  $\mathfrak{S}$  perfekt abbildet.

## Ein Trennungskriterium

### Satz

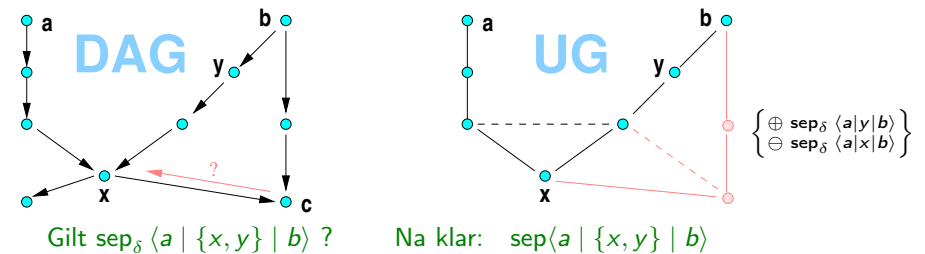
In einem DAG  $\mathcal{G} = (V, \mathcal{E})$  gilt für alle disjunkten Mengen  $A, B, Z \subset V$ :

$$\text{sep}_\delta \langle A \mid Z \mid B \rangle_{\mathcal{G}} \iff \text{sep} \langle A \mid Z \mid B \rangle_{\mathcal{G}^*}$$

Dabei bezeichne  $W$  die Vorgängerhülle von  $A \cup Z \cup B$  und  $\mathcal{G}^*$  sei der **moralische Graph**  $(\mathcal{G}_W)^m$  von  $\mathcal{G}_W$ .

### Beispiel

- obere Kette:**  $y$  blockiert, aber  $x$  blockiert nicht!
- untere Kette:**  $x$  blockiert und  $c$  blockiert.



## Rekursive Faktorisierung

### Definition

Die Wahrscheinlichkeitsverteilung  $P(\cdot)$  **zerfällt rekursiv** über dem gerichteten azyklischen Graphen  $\mathcal{G} = (V, \mathcal{E})$ , wenn es für jede Variable  $a \in V$  eine nichtnegative **Kernfunktion**

$$\phi_a : \mathcal{X}_a \times \bigotimes_{v \in \text{pa}(a)} \mathcal{X}_v \rightarrow \mathbb{R}_0^+$$

gibt mit

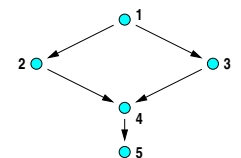
$$P(\mathbf{x}) = \prod_{a \in V} \phi_a(x_a, \mathbf{x}_{\text{pa}(a)})$$

Es bezeichnet  $\text{pa}(a)$  die **Eltermenge**  $\{v \mid (v, a) \in \mathcal{E}\}$  von  $a$ .

### Beispiel

Im Rasensprengergraphen zerfällt die Verteilung, falls es Potentialfunktionen gibt mit:

$$P(\mathbf{x}) = \phi_1(x_1) \cdot \phi_2(x_1, x_2) \cdot \phi_3(x_1, x_3) \cdot \phi_4(x_2, x_3, x_4) \cdot \phi_5(x_4, x_5)$$



## Die reduzierte Kettenregel

### Lemma

Wenn  $P(\cdot)$  über  $\mathcal{G}$  rekursiv zerfällt, können die Kernfunktionen  $\phi_a$  o.B.d.A. gemäß

$$\phi_a(x_a, \mathbf{x}_{pa(a)}) = P_{a|pa(a)}(x_a | \mathbf{x}_{pa(a)})$$

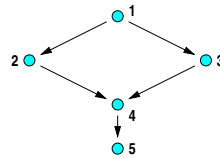
als bedingte Einzelwertwahrscheinlichkeiten gestaltet werden und es gilt — bei kantenverträglicher Variablenordnung — die **reduzierte Kettenregel**:

$$P(\mathbf{x}) = \prod_{i=1}^N P(x_i | \mathbf{x}_{pa(\mathbb{X}_i)})$$

### Beispiel

Im Rasensprengergraphen kann die Faktorisierung wie folgt gewählt werden:

$$P(\mathbf{x}) = P(x_1) \cdot P(x_2 | x_1) \cdot P(x_3 | x_1) \cdot P(x_4 | x_2, x_3) \cdot P(x_5 | x_4)$$



## Faktorisierung und Markoveigenschaft

### Satz

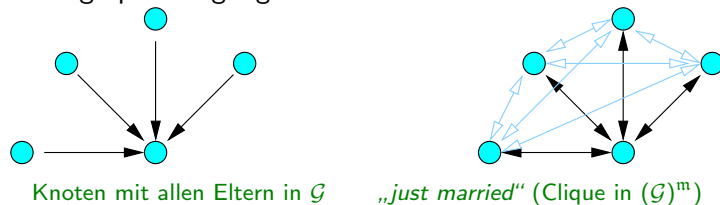
Wenn die Verteilung  $P(\cdot)$  über dem DAG  $\mathcal{G}$  rekursiv zerfällt, dann zerfällt  $P(\cdot)$  auch über dem moralischen Graphen  $(\mathcal{G})^m$  von  $\mathcal{G}$ .

$\mathcal{G}$  ist dann sicherlich ein Unabhängigkeitsbild von  $P(\cdot)$ , d.h. es gilt:



### Beweisidee

Moralgrapherzeugung



### Beweis.

Wir vereinbaren eine verträgliche Ordnung  $V = \{\mathbb{X}_1, \dots, \mathbb{X}_N\}$  und wir wissen, daß nun die Kausalitätsbeziehung gilt:

$$\mathbb{X}_j \in pa(\mathbb{X}_i) \Rightarrow j < i$$

Wir berechnen nun die Randverteilung der ersten  $n$  Variablen:

$$\begin{aligned} P(\mathbf{x}_{1..n}) &= \sum_{\mathbf{x}_{n+1}} \dots \sum_{\mathbf{x}_N} P(\mathbf{x}) = \sum_{\mathbf{x}_{n+1}} \dots \sum_{\mathbf{x}_N} \prod_{i=1}^N \phi_i(x_i, \mathbf{x}_{pa(\mathbb{X}_i)}) \\ &= \prod_{i=1}^n \phi_i(x_i, \mathbf{x}_{pa(\mathbb{X}_i)}) \cdot \prod_{i=n+1}^N \underbrace{\left( \sum_{\mathbf{x}_i \in \mathcal{X}_i} \phi_i(x_i, \mathbf{x}_{pa(\mathbb{X}_i)}) \right)}_{\sigma_i} \end{aligned}$$

Daraus folgt für die bedingte Wahrscheinlichkeit  $P(x_n | \mathbf{x}_{1..n-1})$ :

$$\dots = \frac{P(\mathbf{x}_{1..n})}{P(\mathbf{x}_{1..n-1})} = \frac{\prod_{i=1}^n \phi_i(x_i, \mathbf{x}_{pa(\mathbb{X}_i)}) \cdot \prod_{i=n+1}^N \sigma_i}{\prod_{i=1}^{n-1} \phi_i(x_i, \mathbf{x}_{pa(\mathbb{X}_i)}) \cdot \prod_{i=n}^N \sigma_i} = \frac{\phi_n(x_n, \mathbf{x}_{pa(\mathbb{X}_n)})}{\sigma_n}$$

Wenn wir also normierte Faktoren verwenden ( $\sigma_n \equiv 1$ ), entsprechen die  $\phi_n$  gerade den klassischen Kettenregellgliedern  $P(x_n | \cdot)$ . Daß diese tatsächlich nur von  $\mathbb{X}_n$  und deren Eltervariablen abhängen, ergibt sich aus der Argumentstruktur von  $\phi_n$ . □

### Beweis.

#### • FME\* $\Rightarrow$ FME

(moralische Faktorisierung)

Für jede Variable  $a \in V$  ist die Menge  $\{a\} \cup pa(a)$  im moralischen Graphen  $(\mathcal{G})^m$  von  $\mathcal{G}$  vollständig, denn  $a$  ist mit jedem Elter adjazent und alle Eltern wurden miteinander verheiratet.

Damit bilden die Potentialfunktionen  $\phi_a$  auch eine Cliquenfaktorisierung auf  $(\mathcal{G})^m$ .

#### • FME $\Rightarrow$ GME

(für  $(\mathcal{G})^m$ ; gilt immer)

#### • GME $\Rightarrow$ GME\*

Gilt nun  $\text{sep}_\delta \langle A | Z | B \rangle$  in  $\mathcal{G}$ , so ist auch  $\text{sep} \langle A | Z | B \rangle$  in  $(\mathcal{G})^m$ .

Es besitzt  $(\mathcal{G})^m$  die globale ME für  $P(\cdot)$ , also ist auch  $\mathfrak{I} \langle A | Z | B \rangle$ . □

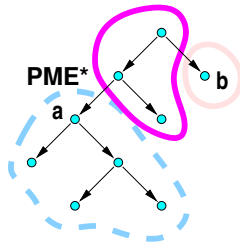
## Die drei Markoveigenschaften

### Definition

Es sei  $P(\cdot)$  eine Wahrscheinlichkeitsverteilung auf  $V$  und  $\mathfrak{S}$  ihr Dependenzmodell. Der gerichtete azyklische Graph  $\mathcal{G} = (V, \mathcal{E})$  erfüllt die

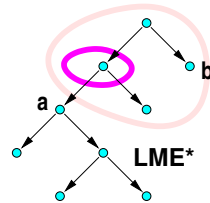
- **paarweise Markoveigenschaft**  
gdw. für alle nichtadjazenten  $a, b \in V$  gilt:

$$\mathfrak{S}(a \mid V \setminus \text{off}(a) \setminus \{b\} \mid b)$$



- **lokale Markoveigenschaft**  
gdw. für jede Variable  $a \in V$  gilt:

$$\mathfrak{S}(a \mid \text{pa}(a) \mid V \setminus \text{off}(a))$$



- **globale Markoveigenschaft**  
gdw. für alle  $A, B, Z \subset V$  mit  $\text{sep}_\delta \langle A|Z|B \rangle$  gilt:

$$\mathfrak{S}(A \mid Z \mid B)$$

### Beweis.

$\text{PME}^* \not\Rightarrow \text{LME}^*$

(alle anderen Richtungen nur im alten Vorlesungsskriptum)

Als Gegenbeispiel betrachte die vier binärwertigen, uniform verteilten Zufallsvariablen  $X = Y = Z$  und  $W$  und den DAG mit den Kanten

$$Z \rightarrow W \rightarrow X \quad \text{und} \quad Z \rightarrow Y \rightarrow W.$$

Der Graph besitzt die paarweise ME, denn von den insgesamt vier nichtadjazenten Variablenpaaren erfüllen nur  $(X, Y)$  und  $(X, Z)$  die Nachkommenbedingung. Damit sind

$$\mathfrak{S}(X \mid W, Z \mid Y) \quad \text{und} \quad \mathfrak{S}(X \mid W, Y \mid Z)$$

zu überprüfen — die Faktorzerlegung ergibt sich aber wie folgt:

$$P(x, y \mid z, w) = \begin{cases} 1 & x = y = z \\ 0 & \text{sonst} \end{cases} = \delta_{xz} \cdot \delta_{yz}$$

Ganz analog ergibt sich auch  $P(x, z \mid y, w) = \delta_{xy} \cdot \delta_{zy}$ . Der Graph besitzt aber nicht die lokale ME, denn die Unabhängigkeit

$$\mathfrak{S}(X \mid \underbrace{\text{pa}(X)}_W \mid \underbrace{V \setminus \text{off}(X)}_{\{W, Y, Z\}})$$

bedingt nach Axiom DEC auch  $\mathfrak{S}(X \mid W \mid Y, Z)$ , was die Verteilung  $P(\cdot)$  offensichtlich nicht hergibt.  $\square$

## Die Markoveigenschaften für Semi-/Graphoide

Bayesnetze  $\triangleq$  minimale Unabhängigkeitsbilder

### Definition

Der Graph  $\mathcal{G}$  heißt **Bayesnetz** von  $P(\cdot)$ , wenn er minimal mit der globalen Markoveigenschaft für  $P(\cdot)$  ist.

Das Bayesnetz  $\mathcal{G}$  ignoriert keine Abhängigkeiten, höchstens Unabhängigkeiten, aber davon so wenige wie möglich.

### Satz

Sei  $\mathcal{G} = (V, \mathcal{E})$  und  $P(\cdot)$  auf  $V$  gegeben. Dann gilt

$$\text{FAK}^* \Leftrightarrow \text{GME}^* \Leftrightarrow \text{LME}^* \Rightarrow \text{paarweise ME},$$

aber es gilt im allgemeinen nicht die Umkehrrichtung

$$\text{PME}^* \Rightarrow \text{LME}^*.$$

Für streng positive Verteilungen  $P(\cdot)$  gilt sogar die Äquivalenz

$$\text{LME}^* \Leftrightarrow \text{PME}^*.$$

## Axiomatisierung kausaler Dependenzmodelle ?

### Satz

Ist das Dependenzmodell  $\mathfrak{S}$  kausal, so gelten die folgenden sieben unabhängigen Axiome:

$$\text{SYM} \quad \text{Symmetrie} \quad \mathfrak{S}(A \mid Z \mid B) \Leftrightarrow \mathfrak{S}(B \mid Z \mid A)$$

$$\text{C/D} \quad \text{Komposition/Dekomposition} \quad \mathfrak{S}(A \mid Z \mid BUC) \Leftrightarrow \mathfrak{S}(A \mid Z \mid B) \wedge \mathfrak{S}(A \mid Z \mid C)$$

$$\text{INT} \quad \text{Durchschnitt} \quad \mathfrak{S}(A \mid ZUC \mid B) \wedge \mathfrak{S}(A \mid ZUB \mid C) \Rightarrow \mathfrak{S}(A \mid Z \mid BUC)$$

$$\text{WUN} \quad \text{Schwache Vereinigung} \quad \mathfrak{S}(A \mid Z \mid BUC) \Rightarrow \mathfrak{S}(A \mid ZUC \mid B)$$

$$\text{CON} \quad \text{Kontraktion} \quad \mathfrak{S}(A \mid Z \mid B) \wedge \mathfrak{S}(A \mid ZUB \mid C) \Rightarrow \mathfrak{S}(A \mid Z \mid BUC)$$

$$\text{WTR} \quad \text{Schwache Transitivität} \quad \mathfrak{S}(A \mid Z \mid B) \wedge \mathfrak{S}(A \mid ZU\{x\} \mid B) \Rightarrow \mathfrak{S}(A \mid Z \mid \{x\}) \vee \mathfrak{S}(\{x\} \mid Z \mid B)$$

$$\text{CHO} \quad \text{Kordialität} \quad \mathfrak{S}(a \mid c, d \mid b) \wedge \mathfrak{S}(c \mid a, b \mid d) \Rightarrow \mathfrak{S}(a \mid c \mid b) \vee \mathfrak{S}(a \mid d \mid b)$$

## Markovdecken und Markovgrenzen

### Definition

Sei  $\mathfrak{S}$  ein Dependenzmodell auf  $V$  und  $V = \{\mathbb{X}_1, \dots, \mathbb{X}_N\}$  eine **Variablenordnung**.

- Eine Menge  $B \subset V$  heißt **Markovdecke** von  $c \in V$  bezüglich  $A \subset V$  genau dann, wenn gilt:

$$B \subseteq A \quad \wedge \quad \mathfrak{S}(\{c\} \mid B \mid A \setminus B)$$

- Ist  $B$  minimal mit dieser Eigenschaft, so heißt  $B$  eine **Markovgrenze**.
- Die Folge  $B_1, \dots, B_N$  heißt **Grenzsystem** von  $\mathfrak{S}$  bezüglich Variablenordnung  $\mathbb{X}_1, \dots, \mathbb{X}_N$  genau dann, wenn jede Menge  $B_n$  eine Markovgrenze von  $\mathbb{X}_n$  bezüglich  $V_n = \{\mathbb{X}_1, \dots, \mathbb{X}_{n-1}\}$  ist.
- Ein gerichteter azyklischer Graph  $\mathcal{G}$ , dessen Eltermengen  $pa(\mathbb{X}_n)$  ein Grenzsystem von  $\mathfrak{S}$  bilden, heißt **Grenzengraph** von  $\mathfrak{S}$ .

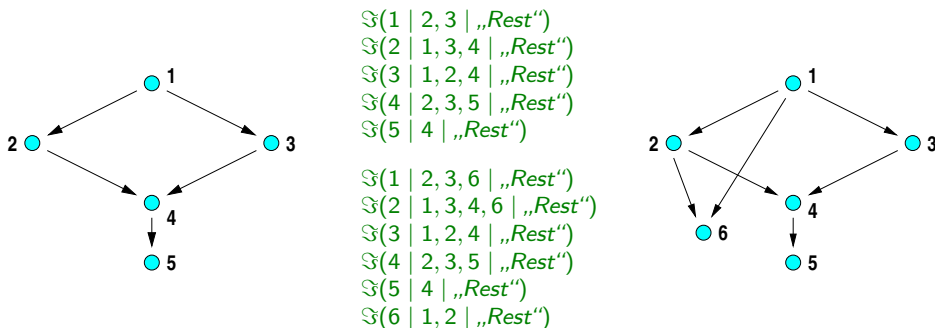
Markovdecken einer Markovkette:  $\mathfrak{S}(\mathbb{X}_n \mid \{\mathbb{X}_{n-1}, \mathbb{X}_{n+1}\} \mid V \setminus \{\mathbb{X}_{n-1}, \mathbb{X}_n, \mathbb{X}_{n+1}\})$

## Markovdecken gegen den Rest der Welt

### Fragestellung

In ungerichteten Graphen fallen die beiden folgenden Fragestellungen zusammen:

- LME** Lokale Markoveigenschaft: Gegen welche Variablen wird  $a \in V$  durch seine unmittelbaren Nachbarn  $bd(a)$  abgeschirmt?
- AME** Allgemeine Markoveigenschaft: Durch welche Menge wird  $a \in V$  gegen den „Rest der Welt“ abgeschirmt?



## Bayesnetzkonstruktion

### Lemma (Verma 1986)

Ist  $\mathfrak{S}$  ein Semigraphoid, so ist jeder Grenzengraph von  $\mathfrak{S}$  ein Bayesnetz von  $\mathfrak{S}$ .

Ist  $\mathfrak{S}$  ein Graphoid, so ist der Grenzengraph von  $\mathfrak{S}$  **bei gegebener Variablenordnung** eindeutig.

$\mathcal{G}$  ist ein Bayesnetz für die Verteilung  $P(\cdot)$  genau dann, wenn er die LME\* für  $\mathfrak{S} = \mathfrak{S}_P$  besitzt und die Eltermengen  $pa(\mathbb{X}_n)$  minimal mit dieser Eigenschaft sind (Markovgrenzen von  $\mathbb{X}_n$  bzgl.  $V \setminus off(\mathbb{X}_n)$ ).

(Algorithmus)

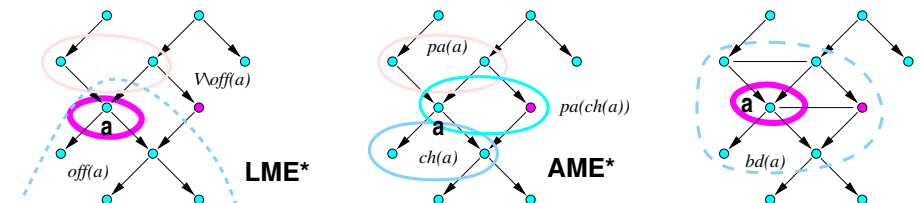
- Wähle eine Variablenordnung  $V = \{\mathbb{X}_1, \dots, \mathbb{X}_N\}$  aus.
- Wähle  $\mathbb{X}_1$  als Wurzel und ordne die Randverteilung  $P_1(x_1)$  zu.
- Für alle  $i \geq 2$  berechne ein minimales  $B_i$  mit

$$B_i \subseteq \{\mathbb{X}_1, \dots, \mathbb{X}_{i-1}\} \quad \text{und} \quad P(x_i \mid x_1, \dots, x_{i-1}) = P(x_i \mid \mathbf{x}_{B_i})$$

und kreierte Knoten  $\mathbb{X}_i$  mit der Vorgängermenge  $pa(\mathbb{X}_i) = B_i$  und der lokalen Verteilung  $P_i(x_i \mid \mathbf{x}_{B_i})$ .

(zumthog1A)

## Allgemeine Markoveigenschaft



### Lemma (AME\*)

Es sei  $\mathcal{G}$  ein Bayesnetz für  $\mathfrak{S}$ . Für jedes  $a \in V$  bildet die Vereinigung der folgenden Variablenmengen eine Markovdecke bzgl.  $V$ :

- die Menge  $pa(a)$  der direkten Vorfahren von  $a$ ,
- die Menge  $ch(a)$  der direkten Nachkommen von  $a$ ,
- die Menge der direkten Vorfahren der direkten Nachkommen von  $a$ .

Mit anderen Worten:

$$\mathfrak{S}(a \mid pa(a) \cup ch(a) \cup pa(ch(a)) \setminus \{a\} \mid \text{„Rest“})$$

## Beweis.

- $\mathcal{G}$  ist ein Bayesnetz von  $\mathfrak{S}$ , also insbesondere ein Unabhängigkeitsbild; folglich gilt die GME\*.
- Wir haben also nur die Trennungseigenschaft

$$\text{sep}_\delta \langle a \mid B_a \mid \text{„Rest“} \rangle_{\mathcal{G}}, \quad B_a \stackrel{\text{def}}{=} \text{pa}(a) \cup \text{ch}(a) \cup \text{pa}(\text{ch}(a)) \setminus \{a\}$$

zu zeigen.

- Die Trennungseigenschaft beweisen wir im Moralgraphen  $(\mathcal{G})^m$ .  
Dort hat Knoten  $a \in V$  als Nachbarn genau alle ehemaligen Eltern und Kinder des DAG sowie zusätzlich all jene Knoten, zu denen gemeinsame Kinder in  $\mathcal{G}$  existieren, mit anderen Worten gilt:

$$\text{bd}_{(\mathcal{G})^m}(a) = B_a$$

- Selbstverständlich wird  $a$  im Moralgraphen  $(\mathcal{G})^m$  — wie in jedem UG wegen der LME — durch seinen Rand  $\text{bd}_{(\mathcal{G})^m}(a)$  von allen Restknoten getrennt:

$$\text{sep}_\delta \langle a \mid B_a \mid \text{„Rest“} \rangle_{(\mathcal{G})^m}$$

Damit ist die Behauptung gezeigt.

□

## Beweis.

- Jedes P mit dem Diamant-UG  $f_1 \overset{m_1}{\diamond} f_2 \overset{m_2}{\diamond}$  als perfektem Bild.
- Jedes P mit dem Konvergenz-DAG  $w_1 \rightarrow g \leftarrow w_2$  als perfektem Bild.
- Jede nichtkausale loglineare Verteilung mit der Modellformel

$$P(a, b, c) = \phi_1(b, c) \cdot \phi_2(a, c) \cdot \phi_3(a, b)$$

denn der **vollständige UG** ist das eindeutige Markovnetz zu P, enthält aber die  $\{b, c\}$ ,  $\{a, c\}$ ,  $\{a, b\}$  nicht als Cliquen.

- Wegen des Spezialfalls partieller Unabhängigkeiten besitzen  $\mathcal{G}_{UG}$  und  $(\mathcal{G}_{DAG})^m$  identische Kanten, das Markovnetz ist also der Moralgraph des Bayesnetzes.

$\mathcal{G}_{UG}$  muß dann aber auch kordal sein, denn jeder Kreis  $\geq 4$  muß im (azyklischen) Bayesnetz einen konvergierenden Knoten besitzen, folglich (aus Gründen der Moral) auch eine Sehne.

- Im Falle der Unmoral gäbe es  $a \rightarrow z \leftarrow b$ , aber weder  $a \rightarrow b$  noch  $a \leftarrow b$ . Für die „historischen Abschlüsse“  $Z$  von  $\{z\}$  und  $W$  von  $\{a, b\}$  gilt dann aber

$$\text{sep}(\{a\} \mid W \setminus \{a, b\} \mid \{b\})_{(\mathcal{G}_W)^m},$$

aber nicht

$$\text{sep}(\{a\} \mid Z \setminus \{a, b\} \mid \{b\})_{(\mathcal{G}_Z)^m},$$

ein eklatanter ⚡ zum SUN-Axiom (P graphisch!), da  $W \setminus \{a, b\} \subset Z \setminus \{a, b\}$  gilt.

□

## Graphische versus kausale Verteilungen

### Lemma

- Es gibt graphische Verteilungen, die nicht kausal sind.
- Es gibt kausale Verteilungen, die nicht graphisch sind.
- Es gibt Verteilungen, die weder graphisch noch kausal sind.
- Ist  $P : \mathcal{X} \rightarrow \mathbb{R}$  sowohl graphisch als auch kausal, so ist jedes Markovnetz von  $\mathfrak{S}_P$  kordal/trianguliert.
- Ist  $P : \mathcal{X} \rightarrow \mathbb{R}$  sowohl graphisch als auch kausal, so ist jedes Bayesnetz von  $\mathfrak{S}_P$  moralisch.

### Beweisidee

$$\begin{array}{ccccc} & & \mathfrak{S}(\{a\} \mid \text{„Rest“} \mid \{b\})_{P(\cdot)} & & \\ & & \uparrow & & \\ \text{sep} \langle A|Z|B \rangle_{\mathcal{G}_{UG}} & \Leftrightarrow & \mathfrak{S} \langle A|Z|B \rangle_{P(\cdot)} & \Leftrightarrow & \text{sep}_\delta \langle A|Z|B \rangle_{\mathcal{G}_{DAG}} \\ & & & & \updownarrow \\ & & & & \text{sep} \langle A|Z|B \rangle_{(\mathcal{G}_{DAG})^m} \end{array}$$

## Beispiele

### Markovnetze mit 3, 4, 5 oder 6 Variablen

• • •	$P(x) \cdot P(y) \cdot P(z)$	3-diskret	⊕
• — • •	$P(x, y) \cdot P(z)$	2+1-diskret	⊕
• — • — •	$P(x, y) \cdot P(y, z) / P(y)$	kaskadiert	⊕
△	$P(x, y, z)$	saturiert	⊕
◇	$\phi(x, y) \cdot \phi(y, z) \cdot \phi(z, w) \cdot \phi(w, x)$	Diamant	⊖
◁ ▷	$P(x, y, z) \cdot P(y, z, w) / P(y, z)$	3/3-Cliquen	⊕
▷ ▷	$P(x, y, z) \cdot P(v, w, z) / P(z)$	3/3-Cliquen	⊕
△ ≡ △	$\phi(x_1, x_2, x_3) \cdot \phi(y_1, y_2, y_3) \cdot \phi(x_1, y_1) \cdot \phi(x_2, y_2) \cdot \phi(x_3, y_3)$	Toblerone	⊖

## Beispiele

### Bayesnetze mit 3 oder 4 Variablen

• • •	$P(x) \cdot P(y) \cdot P(z)$	3-diskret	⊕
•→• •	$P(x) \cdot P(y x) \cdot P(z)$	2+1-diskret	⊕
•→•→•	$P(x) \cdot P(y x) \cdot P(z y)$	kaskadiert	⊕
•←•→•	$P(x y) \cdot P(y) \cdot P(z y)$	divergent	⊕
•→•←•	$P(x) \cdot P(y x, z) \cdot P(z)$	konvergent	⊖
△	$P(x) \cdot P(y x) \cdot P(z x, y)$	saturiert	⊕
◁▷	$P(x) \cdot P(y x) \cdot P(z x, y) \cdot P(w y, z)$	3-3-Cliquen	⊕
◁▷	$P(x) \cdot P(y x) \cdot P(w y) \cdot P(z x, y, w)$	unmoralisch!	⊖

## Berechnung bedingter Wahrscheinlichkeiten

### Verbundverteilung

Gemeinsame Verteilung  $P(x_1, \dots, x_n)$  **aller** Variablen in **Produktform**.

### Randverteilungen

Gemeinsame Verteilung für eine Teilmenge  $A \subset V$ :

$$\begin{aligned}
 P(V \setminus \{x_i\}) &= \sum_{x_i} P(x_1, \dots, x_n) \\
 P(V \setminus \{x_i, x_j\}) &= \sum_{x_i} \sum_{x_j} P(x_1, \dots, x_n) \\
 P(V \setminus \{x_{i_1}, \dots, x_{i_m}\}) &= \sum_{x_{i_1}} \dots \sum_{x_{i_m}} P(x_1, \dots, x_n)
 \end{aligned}$$

### Bedingte Verteilungen

Einfluß einer Zufallsvariablen  $x_j$  auf eine andere  $x_i$ :

$$P(x_i|x_j) = \frac{P(x_i, x_j)}{P(x_j)} = \frac{\sum \dots \sum P(x_1, \dots, x_n)}{\sum \sum \dots \sum P(x_1, \dots, x_n)}$$

## Korrelation, Regression und Transinformation

## Assoziationsregeln und Netzwerkanalyse

## Bedingte statistische Unabhängigkeit

## Graphische Modelle: ungerichtete Graphen

## Kausale Modelle: gerichtete azyklische Graphen

## Berechnen bedingter Wahrscheinlichkeiten

## Parameterschätzung in Bayesnetzen und Loglinearmodellen

## Aufdeckung der Abhängigkeitsstruktur

## Kovarianzselektion

## Warum Bayesnetze ?

Weil sie in Wahrscheinlichkeiten faktorisieren !

## Was ist Inferenz ?

**Logik** Axiome, Schlußregeln ➡ neue Sätze

**Arithmetik** Parameterwerte, Operationen ➡ Funktionswerte

**Stochastik** Observablen, W-Modell ➡ bedingte W'keiten

## A posteriori Verteilungen

**Evidenz**  $E = \{e_1, \dots, e_m\}$  („*instanziierte*“ Variablen)

$$P(x_i|E) = P(x_i = \xi \mid e_1 = \eta_1, \dots, e_m = \eta_m)$$

Rand- und Rückschlußverteilungen sind aufwendig zu berechnen!

- Eliminiere Variablen in ökonomischer Reihenfolge — gemäß Dependenzstruktur bzw. Modellformel.
- Propagationsalgorithmen, Marker-Passing, Sampling ...



## Notation der Rechengrößen

für baumförmige Bayesnetze

### Wahrscheinlichkeitsparametermatrix

Jeder Knoten  $y$  im DAG hat **genau einen** Elterknoten  $x$ .

$$\begin{aligned} \mathbf{M}_{y|x} &= P(y|x) = [P(y = \eta_j | x = \xi_i)]_{ij} \\ &= \begin{Bmatrix} P(y = \eta_1 | x = \xi_1) & \cdots & P(y = \eta_k | x = \xi_1) \\ \vdots & & \vdots \\ P(y = \eta_1 | x = \xi_m) & \cdots & P(y = \eta_k | x = \xi_m) \end{Bmatrix} \end{aligned}$$

### Evidenz

Instanziierte Variablen  $e \in V$  bzw.  $E \subseteq V$ .

### Belief-Funktion

Subjektive Einschätzung von  $x$  auf Grundlage von  $E$  (Wahrheitsfeld):

$$\begin{aligned} \text{bel}(x) &\stackrel{\text{def}}{=} P(x|E) \\ \text{bel}(x) &= P(x | z = \zeta) = (P(x = \xi_1 | z = \zeta), \dots, P(x = \xi_\ell | z = \zeta))^T \end{aligned}$$

## Unidirektionale Fortpflanzung in Ketten

Drei Variablen

**Beispiel:**  $x \rightarrow y \rightarrow z$ , Evidenz  $\{z = \zeta\}$

Nach der Bayesformel gilt wiederum:

$$\text{bel}(x) = P(x | z = \zeta) = \frac{P(x) \cdot P(z = \zeta | x)}{P(z = \zeta)} \propto P(x) \cdot \lambda(x)$$

Der **diagnostische Vektor** lautet nunmehr

$$\begin{aligned} \lambda(x) &= P(z = \zeta | x) = \sum_y P(z = \zeta, y | x) \\ &= \sum_y P(z = \zeta | y) \cdot P(y|x) \\ &= \mathbf{M}_{y|x} \bullet \lambda(y) \end{aligned}$$

$\mathbf{M}_{y|x} \bullet \lambda(y)$  bezeichnet das Vektor-Matrix-Produkt über die Variable  $y$ .

## Unidirektionale Fortpflanzung in Ketten

Zwei Variablen

**Beispiel:**  $x \rightarrow y$ , Evidenz  $\{y = \eta\}$

Nach der Bayesformel gilt:

$$\text{bel}(x) = P(x | y = \eta) = \frac{P(x) \cdot P(y = \eta | x)}{P(y = \eta)} \propto P(x) \cdot \lambda(x)$$

mit der a priori Wahrsch'keit  $P(x)$  und dem **diagnostischen Vektor**

$$\lambda(x) = P(y = \eta | x) \quad (\eta\text{-te Spalte der Matrix } \mathbf{M}_{y|x}).$$

$P(x) \cdot \lambda(x)$  bezeichnet das komponentenweise Produkt.

## Unidirektionale Fortpflanzung in Ketten

Mehr als drei Variablen

**Beispiel:**  $x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_n$ , Evidenz  $\{x_n = \xi\}$

Nach der Bayesformel gilt wiederum:

$$\text{bel}(x_1) = P(x_1 | x_n = \xi) = \frac{P(x_1) \cdot P(x_n = \xi | x_1)}{P(x_n = \xi)} \propto P(x_1) \cdot \lambda(x_1)$$

Der **diagnostische Vektor** gehorcht der Rekursion:

$$\begin{aligned} \lambda(x_1) &= \mathbf{M}_{x_2|x_1} \bullet \lambda(x_2) \\ &= \mathbf{M}_{x_2|x_1} \bullet \mathbf{M}_{x_3|x_2} \bullet \lambda(x_3) \\ &= \mathbf{M}_{x_2|x_1} \bullet \mathbf{M}_{x_3|x_2} \bullet \mathbf{M}_{x_4|x_3} \bullet \lambda(x_4) \\ &= \mathbf{M}_{x_2|x_1} \bullet \mathbf{M}_{x_3|x_2} \bullet \dots \bullet \mathbf{M}_{x_{n-1}|x_{n-2}} \bullet \underbrace{P(x_n = \xi | x_{n-1})}_{\mathbf{M}_{\xi|x_{n-1}}} \end{aligned}$$



## Bidirektionale Fortpflanzung in Ketten

Beispiel:  $e^+ \rightarrow v \rightarrow w \rightarrow x \rightarrow y \rightarrow z \rightarrow e^-$

A posteriori Wahrscheinlichkeiten nach Bayesformel:

$$\begin{aligned} \text{bel}(x) &= P(x | e^+, e^-) \propto P(e^- | x, e^+) \cdot P(x | e^+) \\ &= P(e^- | x) \cdot P(x | e^+) = \lambda(x) \cdot \pi(x) \end{aligned}$$

Diagnostische Evidenz

Kausale Evidenz

$$\lambda(x) = P(e^- | x)$$

$$\pi(x) = P(x | e^+)$$

Fortpflanzung rückwärts

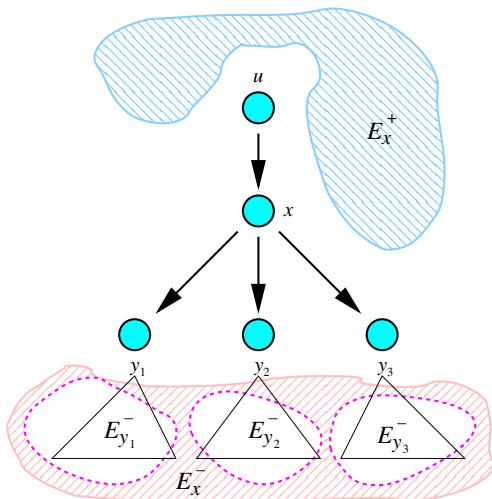
$$\begin{aligned} \pi(x) &= P(x | e^+) \\ &= \sum_w P(x | w, e^+) \cdot P(w | e^+) \\ &= \sum_w P(x | w) \cdot P(w | e^+) \\ &= \pi(w) \bullet M_{x|w} \end{aligned}$$

Fortpflanzung vorwärts

$$\begin{aligned} \lambda(x) &= P(e^- | x) \\ &= \sum_y P(e^-, y | x) \\ &= \sum_y P(e^- | y) \cdot P(y | x) \\ &= M_{y|x} \bullet \lambda(y) \end{aligned}$$

## Bidirektionale Fortpflanzung in Bäumen

Zerlegung der Evidenz



Vertikale Zerlegung

kausal/diagnostisch:

$$E_x = E_x^+ \uplus E_x^-$$

Horizontale Zerlegung

des diagnostischen Teils

$$E_x^- = \biguplus_{y_\ell \in \text{ch}(x)} E_{y_\ell}^-$$

Horizontale Zerlegung

des kausalen Teils

$$E_{y_\ell}^+ = E_x^+ \uplus \biguplus_{k \neq \ell} E_{y_k}^-$$

## Bidirektionale Fortpflanzung in Bäumen

Zerlegung der Belief-Funktion

Zerlegung der Evidenz

Für  $x \in V$  unterscheiden wir zwei Quellgebiete:

$$E = E_x^+ \uplus E_x^- \quad \text{mit} \quad \begin{cases} E_x^- \subset \text{ch}(x) & \text{„flußabwärts“} \\ E_x^+ \subset V \setminus \text{ch}(x) & \text{„flußaufwärts“} \end{cases}$$

Belief-Funktion

Nach Kettenregel und  $\text{sep}_\delta \langle E_x^- | \{x\} | E_x^+ \rangle$  folgt:

$$\begin{aligned} \text{bel}(x) &= P(x | E_x^+, E_x^-) \\ &\propto P(E_x^-, x | E_x^+) \\ &= P(E_x^- | x, E_x^+) \cdot P(x | E_x^+) = \lambda(x) \cdot \pi(x) \end{aligned}$$

$\pi(x)$  = **kausale** Unterstützung von  $x$  durch die Vorgänger

$\lambda(x)$  = **diagnostische** Unterstützung von  $x$  durch die Nachfolger

## Bidirektionale Fortpflanzung in Bäumen

Diagnostische und prädiktive Wahrscheinlichkeiten

Diagnostische Komponente

Seien  $u_1, \dots, u_r$  die Nachfolger von  $x$ :

$$\lambda(x) = P(E_x^- | x) = P(E_{u_1}^-, \dots, E_{u_r}^- | x) = \prod_{s=1}^r \underbrace{P(E_{u_s}^- | x)}_{\lambda_{u_s}(x)}$$

Falls  $\{x = \xi\}$  selbst instanziiert, so erzeuge Dummyknoten  $d$  mit  $\lambda_d(x) = \mathbf{1}_{x=\xi}$ .

Prädiktive Komponente

Sei  $u \in V$  der Vater (die Mutter) von  $x$ :

$$\begin{aligned} \pi(x) &= P(x | E_x^+) = \sum_u P(x, u | E_x^+) \\ &= \sum_u P(x | u) \cdot P(u | E_x^+) =: M_{x|u} \bullet \pi_x(u) \end{aligned}$$

## Bidirektionale Fortpflanzung in Bäumen

Variablenversetzte diagnostische und prädiktive Komponenten

### Berechnung von $\lambda_x(u)$

für  $u \rightarrow x$

$$\begin{aligned}\lambda_x(u) &= \sum_x P(E_x^- | u, x) \cdot P(x|u) \\ &= \sum_x P(E_x^- | x) \cdot P(x|u) \\ &= \sum_x \lambda(x) \cdot P(x|u) \\ &= M_{x|u} \bullet \lambda(x)\end{aligned}$$

### Berechnung von $\pi_y(x)$

für  $u \rightarrow x$  und  $y \leftarrow x \rightarrow z$

$$\begin{aligned}\pi_y(x) &= P(x | E_y^+) \\ &= P(x | E_x^+, E_z^-) \\ &\propto P(E_z^- | x, E_x^+) \cdot P(x | E_x^+) \\ &= \lambda_z(x) \cdot \pi(x) \\ &= \lambda_z(x) \cdot M_{x|u} \bullet \pi_x(u)\end{aligned}$$

Spezialfall:  $x = \xi$  evident

$$\lambda_x(u) = P(x = \xi | u)$$

( $\xi$ -te Spalte von Matrix  $M_{x|u}$ )

Allgemeinfall:  $\geq 3$  Kinder

$$\pi_y(x) = \pi(x) \cdot \sum_{z \neq y} \lambda_z(x)$$

## Inferenz in moralischen Bayesnetzen

Vorwärts-Rückwärts-Algorithmus über Variablenkomplexen

1. **ENTFERNE ALLE KANTENRICHTUNGEN**  
↔ äquivalentes kordales Markovnetz
2. **BILDE VETRÄGLICHEN VERBUNDBAUM**  
mit Cliquensequenz  $\mathcal{C} = \{C_1, \dots, C_K\}$
3. **KONSTRUIERE VARIABLENKOMPLEXE**  
 $\mathbb{Y}_k := \bigotimes_{\mathbb{X}_j \in B_k} \mathbb{X}_j$  mit  $B_k := C_k \setminus C_{k-1}$
4. **EXPANDIERE VERTEILUNGSPARAMETER**  
 $M_{k|\ell} = (P(\mathbb{Y}_k = \eta | \mathbb{Y}_\ell = \zeta) | \eta \in \mathcal{Y}_k, \zeta \in \mathcal{Y}_\ell)$
5. **EHEKUTIERE VR-ALGORITHMUS AUF VB**

(sumfting/A)

### Moralische Bayesnetze

**Verbundbaum:**  
spezielle x/y-Kombinationen

**Markovnetz:**  
nur Imputation!

### Unmoralische Bayesnetze

**Schummeln:**  
Verheiraten aller Elternpaare

**Monte Carlo:**  
Iteratives Auswürfeln und Neuschätzen

## Vorwärts-Rückwärts-Algorithmus

in baumförmigen Bayesnetzen

### Start

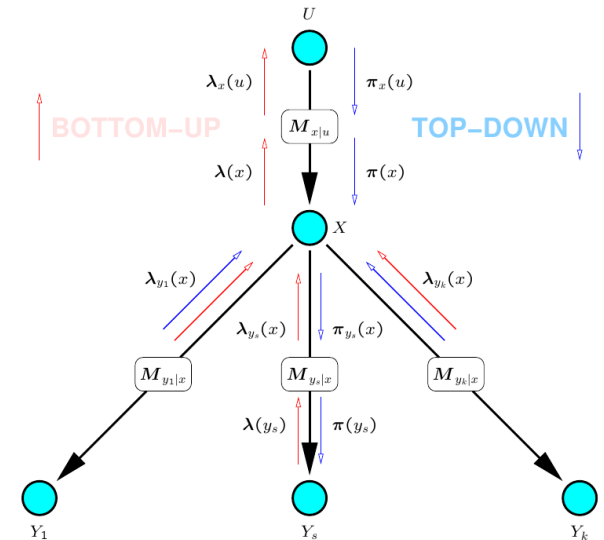
$$\begin{aligned}\pi(x_0) &= M_{x_0| \cdot} \quad (\text{Wurzel}) \\ \lambda(x_\ell) &= \mathbf{1} \quad (\text{Blatt}) \\ \lambda(x_e) &= e^{(\zeta)} \quad (\text{Evidenz})\end{aligned}$$

### Bottom-up

$$\begin{aligned}\lambda(y_k) & \quad (\text{I.V.}) \\ \lambda_{y_k}(x) &= M_{y_k|x} \bullet \lambda(y_k) \\ \lambda(x) &= \prod_k \lambda_{y_k}(x)\end{aligned}$$

### Top-down

$$\begin{aligned}\pi(x) & \quad (\text{I.V.}) \\ \pi_{y_k}(x) &\propto \pi(x) \cdot \prod_{\ell \neq k} \lambda_{y_\ell}(x) \\ \pi(y_k) &= M_{y_k|x} \bullet \pi_{y_k}(x)\end{aligned}$$



## Spezialfall Imputation

Vorhersage eines Attributwertes aus allen anderen

### Belief-Funktion

mit Zielvariable  $x_k$  und Evidenzvariablen  $E = V \setminus \{x_k\}$ :

$$\text{bel}(x_k)_\xi = P(x_k = \xi | \mathbf{x}_E) = \frac{P(x_k = \xi, \mathbf{x}_E)}{P(\mathbf{x}_E)} = \frac{P(\mathbf{x}_{|x_k=\xi})}{P(\mathbf{x}_E)}$$

$\mathbb{X}_k$  ist diskretes Attribut

1. Für alle  $\xi_\ell \in \mathcal{X}_k$  berechne  $q_\ell = P(\mathbf{x}_{|x_k=\xi_\ell})$  mit

$$\mathbf{x}_{|x_k=\xi_\ell} = (x_1, \dots, x_{k-1}, x_k = \xi_\ell, x_{k+1}, \dots, x_n)^\top \in \mathbb{R}^n.$$

2. Setze  $\text{bel}(x_k)_\ell = q_\ell / \sum_i q_i$ .

$\mathbb{X}_k$  ist stetiges Attribut

Effiziente Lösungsmöglichkeit trotz  $|\mathcal{X}_k| = \infty$  ?

## Spezialfall Imputation

Vorhersage eines normalverteilten Attributwertes

### Die Geheimfunktion

Es ist  $P(x_k = \xi \mid \mathbf{x}_E) = \mathcal{N}(\xi \mid \mu, \sigma^2)$ , also erhalten wir Resultate der Form

$$P(\mathbf{x}_{|x_k=\xi}) = P(\mathbf{x}_E) \cdot P(\xi \mid \mathbf{x}_E) = \underbrace{c \cdot \mathcal{N}(\xi \mid \mu, \sigma^2)}_{g_{c,\mu,\sigma}(\xi)}$$

durch Auswertung des Bayesnetzes an der Stelle  $\xi \in \mathbb{R}$ .

### Lemma

Die unbekannten Parameter  $c > 0$ ,  $\mu \in \mathbb{R}$  und  $\sigma > 0$  der skalierten univariaten Gaußdichte

$$g_{c,\mu,\sigma}(\xi) \stackrel{\text{def}}{=} c \cdot \mathcal{N}(\xi \mid \mu, \sigma^2)$$

können aus den Funktionswerten von  $g(\cdot)$  an vier reellen Stützstellen bestimmt werden.

Diese Entschlüsselungstechnik läßt sich auf **multivariate** Gaußdichten verallgemeinern.

### Beweis.

Wir definieren die *lograt*-Funktion  $\ell(x, y) = -2 \cdot \log(g(x)/g(y))$  und folgern die Identität

$$\ell(x, y) = \frac{1}{\sigma^2} \cdot (x^2 - y^2 - 2\mu \cdot (x - y)) .$$

Wir definieren nun die Differentiale

$$\ell_h^-(x) \stackrel{\text{def}}{=} \ell(x, x-h) = \frac{1}{\sigma^2} \cdot (+2hx - h^2 - 2h\mu)$$

$$\ell_h^+(x) \stackrel{\text{def}}{=} \ell(x, x+h) = \frac{1}{\sigma^2} \cdot (-2hx - h^2 + 2h\mu)$$

für  $h > 0$  und finden nach deren Addition einen Lösungsausdruck

$$\hat{\sigma}^2 = -2 \cdot \frac{h^2}{\ell_h^+(x) + \ell_h^-(x)}$$

für die gesuchte Varianz. Anschließend können wir aus jeder der Differentialformeln den Erwartungswert berechnen, z.B.:

$$\hat{\mu} = \frac{\hat{\sigma}^2 \cdot \ell_h^+(x) + 2hx + h^2}{2h} = \frac{\hat{\sigma}^2}{2h} \cdot \ell_h^+(x) + x + \frac{h}{2}$$

Schließlich bestimmen wir noch den Skalierungsfaktor  $c$ ; die numerisch stabilste Methode besteht in einer weiteren Auswertung der Geheimfunktion  $g(\cdot)$ , und zwar am Dichtegipfel:

$$\hat{c} = \frac{g(x)}{\mathcal{N}(x \mid \hat{\mu}, \hat{\sigma}^2)} = \frac{g(\hat{\mu})}{\mathcal{N}(\hat{\mu} \mid \hat{\mu}, \hat{\sigma}^2)} = \frac{g(\hat{\mu})}{\mathcal{N}(0 \mid 0, \hat{\sigma}^2)} = \sqrt{2\pi} \cdot \hat{\sigma} \cdot g(\hat{\mu})$$

□

## Imputation in (nichtkordalen) Markovnetzen

Effiziente Berechnung als  $\text{bd}(x_k)$ -ausgedünntes Cliquesprodukt

### Bedingte Wahrscheinlichkeit nach Faktorisierung

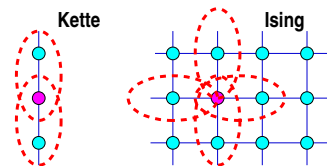
$$P(\mathbb{X}_k = x_k \mid \mathbb{X}_{V \setminus k} = \mathbf{x}') = \frac{P(x_k, \mathbf{x}')}{P(\mathbf{x}')} = \frac{\prod_{C \in \mathcal{C}} \phi_C(x_k, \mathbf{x}')}{\sum_{\xi \in \mathcal{X}_k} \prod_{C \in \mathcal{C}} \phi_C(x_k, \mathbf{x}')} .$$

### Ausklammern & Kürzen aller Gibbspotenziale $\phi_C$ mit $x_k \notin C$

Reduzierte Faktorisierung über  $\mathcal{C}_{(k)} := \{C \mid x_k \in C\} = \{C \mid C \subseteq \text{cl}(x_k)\}$  (wegen  $\mathfrak{I}(x_k \mid \text{bd}(x_k) \mid \text{„Rest“})$  nicht ganz unerwartet!)

### Binäres Zielattribut $|\mathcal{X}_k| = 2$

$$\log \text{odds}(\mathbf{x}') = \log \frac{P(1 \mid \mathbf{x}')}{P(0 \mid \mathbf{x}')} = \sum_{C \in \mathcal{C}_{(k)}} \log \frac{\phi_C(\mathbf{x}_{C \setminus k}, 1)}{\phi_C(\mathbf{x}_{C \setminus k}, 0)}$$



## Korrelation, Regression und Transinformation

### Assoziationsregeln und Netzwerkanalyse

### Bedingte statistische Unabhängigkeit

### Graphische Modelle: ungerichtete Graphen

### Kausale Modelle: gerichtete azyklische Graphen

### Berechnen bedingter Wahrscheinlichkeiten

### Parameterschätzung in Bayesnetzen und Loglinearmodellen

### Aufdeckung der Abhängigkeitsstruktur

### Kovarianzselektion

## Diskrete loglineare Modelle

Spezialfall: drei Variablen ( $N = 3$ )

### Dreiwegetabellen

Drei diskrete Zufallsvariablen  $V = \{\mathbb{X}_1, \mathbb{X}_2, \mathbb{X}_3\} = \{a, b, c\}$

- Endliche Wertebereiche  $\mathcal{X}_a, \mathcal{X}_b, \mathcal{X}_c$
- Endlich viele Zellen  $(j, k, l) \in \mathcal{X}_a \times \mathcal{X}_b \times \mathcal{X}_c$
- Würfel  $\{p_{jkl}\}$  von Wahrscheinlichkeiten  $\sum_j \sum_k \sum_l p_{jkl} = 1$
- Würfel  $\{n_{jkl}\}$  von (absoluten) Häufigkeiten  $\sum_j \sum_k \sum_l n_{jkl} = T$

### Loglineares Verteilungsmodell

Produktform

$$p_{jkl} = \prod_{A \in \Delta} \underbrace{\phi_A(\mathbf{x}_A)}_{z_{jkl}^A} \quad \Delta \subset \mathfrak{P}V$$

Summenform

$$\log p_{jkl} = \sum_{A \in \Delta} \underbrace{\log \phi_A(\mathbf{x}_A)}_{u_{jkl}^A}$$

## Schätzung der kanonischen Modellparameter

### Normierungseigenschaft

$$1 \stackrel{!}{=} \sum_{jkl} p_{jkl} = \sum_{jkl} \exp \left\{ \sum_{A \in \Delta} u_{jkl}^A \right\} = e^u \cdot \sum_{jkl} \exp \left\{ \sum_{A \neq \emptyset} u_{jkl}^A \right\}$$

### Multinomial gezogene Stichprobe

$$P(\mathbf{n}|\mathbf{p}) = P(\{n_{jkl}\} | \{p_{jkl}\}) = \frac{T!}{\prod_{j,k,l} n_{jkl}!} \cdot \prod_{j,k,l} p_{jkl}^{n_{jkl}}$$

### Logarithmierte Likelihood-Funktion

$$\ell_{\text{ML}}(\mathbf{n}|\mathbf{p}) = \log \frac{T!}{\prod_{j,k,l} n_{jkl}!} + \sum_{j,k,l} n_{jkl} \log p_{jkl}$$

### Maximum-Likelihood-Schätzwerte

Kanonische Verteilungsparameter für das saturierte Modell

$$\hat{\mathbf{p}} = \underset{\mathbf{p}}{\operatorname{argmax}} \ell_{\text{ML}}(\mathbf{n}|\mathbf{p}) \quad \Rightarrow \quad \hat{p}_{jkl} = \frac{n_{jkl}}{T}$$

## Beispiele — Dreiwegemodelle

Menge der (maximalen) Interaktionsterme · „Generatoren“

### Unabhängiges Modell

$a, b, c$

$$\begin{aligned} \log p_{jkl} &= u + u_j^a + u_k^b + u_l^c \\ p_{jkl} &= P(a = \alpha_j) \cdot P(b = \beta_k) \cdot P(c = \chi_l) = p_{j..} \cdot p_{..k} \cdot p_{..l} \end{aligned}$$

### Kettenförmiges Modell

$ab, ac$

$$\begin{aligned} \log p_{jkl} &= u + u_j^a + u_k^b + u_l^c + u_{jk}^{ab} + u_{jl}^{ac} \\ p_{jkl} &= \frac{p_{jk..} \cdot p_{j..l}}{p_{j..}} \quad \text{bzw.} \quad \frac{p_{jkl}}{p_{j..}} = \frac{p_{jk..}}{p_{j..}} \cdot \frac{p_{j..l}}{p_{j..}} \end{aligned}$$

### Saturiertes Modell

$abc$

$$\log p_{jkl} = u + u_j^a + u_k^b + u_l^c + u_{jk}^{ab} + u_{jl}^{ac} + u_{kl}^{bc} + u_{jkl}^{abc}$$

## Diskrete Loglinearmodelle

### Definition

Die Familie diskreter Wahrscheinlichkeitsfunktionen  $\{p_{\mathbf{x}}\}_{\mathbf{x} \in \Omega}$  der Gestalt

$$\log p_{\mathbf{x}} = \sum_{A \in \Delta} u_{\mathbf{x}}^A, \quad \mathbf{x} \in \Omega, \quad \Delta \subseteq \mathfrak{P}V$$

heißt **Loglinearmodell** mit der Menge  $\Delta$  von **Interaktionstermen**.

1. Das Loglinearmodell heißt **hierarchisch**, falls gilt:

$$A \subseteq B \quad \text{und} \quad B \in \Delta \quad \Rightarrow \quad A \in \Delta$$

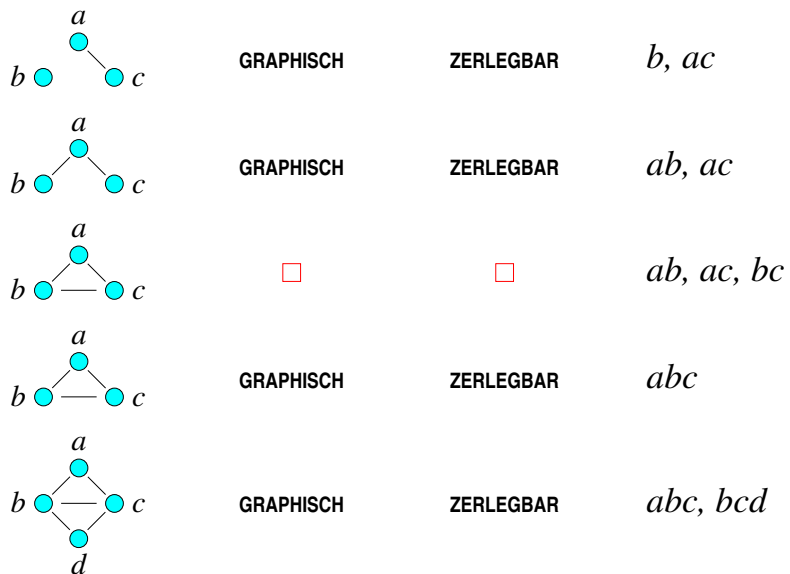
2. Das Loglinearmodell heißt **graphisch**, wenn gilt:

$$C \in \Delta \quad \Leftrightarrow \quad \forall a, b \in C : \{a, b\} \in \Delta$$

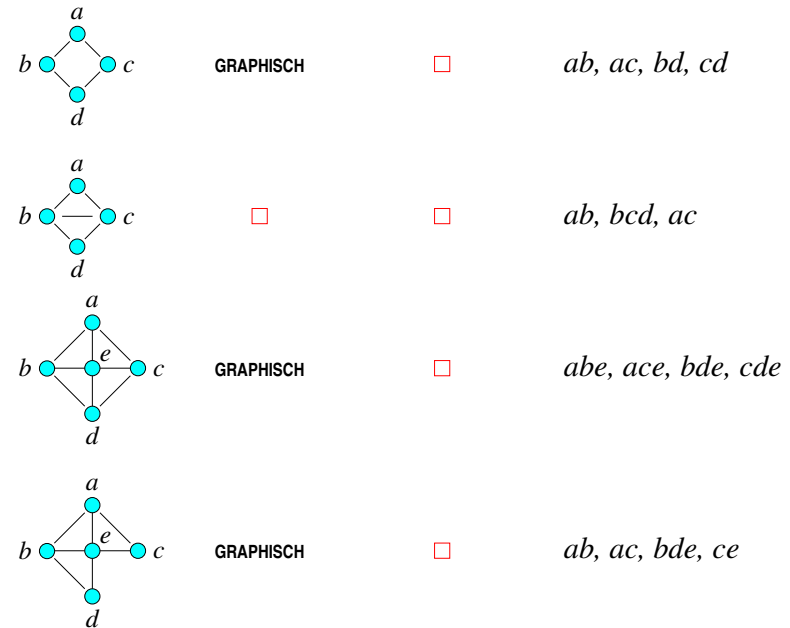
3. Ein graphisches LLM heißt **zerlegbar**, wenn sein Graph kordal ist.

Die maximalen Interaktionsterme eines hierarchischen Loglinearmodells heißen **Generatoren**. Die Generatorenmenge wird auch als **Modellformel** bezeichnet.

## Beispiele — Loglinearmodelle I



## Beispiele — Loglinearmodelle II



## Elementare und marginale Ereignisse

Häufigkeit und charakteristische Funktion

### Definition

Es sei  $\{n_x\}_{x \in \Omega}$  die Tafel elementarer Ereignishäufigkeiten über  $V$ .  
Das Zahlenfeld  $\{n_{x_A}\}_{x_A \in \Omega_A}$  für eine Variablenmenge  $A \subseteq V$  mit Einträgen

$$n_{x_A} \stackrel{\text{def}}{=} \sum_{x_{A'} \in \Omega_{A'}} n_x, \quad A' = V \setminus A$$

heißt **marginale Tafel** oder Tabelle für  $A$ .

### Definition

Sei  $V = \{\mathbb{X}_1, \dots, \mathbb{X}_N\}$ ,  $A \subseteq V$  und  $x_A \in \Omega_A$  ein marginales Ereignis.  
Die zweiwertige Abbildung

$$\varphi_{x_A} : \begin{cases} \Omega & \rightarrow \{1, 0\} \\ \xi & \mapsto \begin{cases} 1 & \xi_j = x_j \text{ für alle } j \in A \\ 0 & \text{sonst} \end{cases} \end{cases}$$

heißt **charakteristische Funktion** von  $x_A$ .

## Maximum-Entropie-Prinzip

Edwin Thompson Jaynes, 1957

### Satz (ML $\triangleq$ ME)

Es sei ein hierarchisches loglineares Modell

$$\log p_x = \sum_{A \in \Delta} u_x^A$$

gegeben sowie die Häufigkeitstafel  $\{n_x\}_{x \in \Omega}$  der Daten  $\omega \subset \Omega$ .

1. Die **Maximum-Likelihood**-Parameter  $\{u_x^A\}_{x_A}$  des Modells erfüllen die Bedingungsgleichungen (\*)

$$\mathcal{E}[\varphi_{x_A}(\mathbb{X}) \mid \mathbf{u}] = \frac{n_{x_A}}{T}, \quad A \in \Delta, \mathbf{x} \in \Omega.$$

2. Unter allem Wahrscheinlichkeitsverteilungen, die das Gleichungssystem (\*) erfüllen, hat obiges loglineare Modell mit Parametern  $\{u_x^A\}_{x_A}$  die **maximale Entropie**.

„Unter allen Zuständen eines physikalischen Systems, die kompatibel mit dem vorhandenen Wissen sind, ist der zu wählen, welcher die Entropie maximiert.“

## Lernen der Loglinearparameter

GIS-Algorithmus — Generalized Iterative Scaling

### Satz (Deming & Stephan, 1940)

Mit der abkürzenden Schreibweise  $z_x^A = \exp(u_x^A)$  gilt:

Das Iterationsverfahren

$$z_x^A \leftarrow z_x^A \cdot \left( \frac{n_{x_A}/T}{\mathcal{E}[\varphi_{x_A}(\mathbb{X})]} \right)^{1/|\Delta|} = z_x^A \cdot \left( \frac{\sum_{y \in \Omega} \varphi_{x_A}(y) \cdot \frac{n_y}{T}}{\sum_{y \in \Omega} \varphi_{x_A}(y) \cdot \prod_{B \in \Delta} z_y^B} \right)^{1/|\Delta|}$$

mit den Startwerten  $z_x^A \equiv 1$  konvergiert gegen die Maximum-Likelihood-Schätzwerte des loglinearen Modells.

#### Bemerkung

Die Gleichung für  $\emptyset \in \Delta$  garantiert  $\sum p_x = 1$ .

Das Bedingungssystem ist konsistent: alle  $C_y = \sum \varphi_{x_A}(y)$  sind gleich  $|\Delta|$ .

Beweis  $\rightsquigarrow$  Skriptum „Stochastische Grammatikmodelle“

## Beispiel — Tafelobst im Tetrapack

Daten = 100 Obstkörbe

	0:4	1:3	2:2	3:1	4:0	
$n_{00}$	60	2	1	0	0	63
$n_{01}$	0	1	2	1	0	4
$n_{10}$	0	1	2	1	0	4
$n_{11}$	0	0	1	2	26	29
	60	4	6	4	26	100

Generalized Iterative Scaling

$i$	Loglinearparameter				Wahrscheinlichkeiten in Promille				
	$z_{00}$	$z_{01}$	$z_{11}$	$1/z$	$\begin{pmatrix} 00 \\ 00 \end{pmatrix}$	$\begin{pmatrix} 10 \\ 00 \end{pmatrix}$	$\begin{pmatrix} 11 \\ 00 \end{pmatrix}$	$\begin{pmatrix} 11 \\ 10 \end{pmatrix}$	$\begin{pmatrix} 11 \\ 11 \end{pmatrix}$
0	1	1	1	16	62.5	62.5	62.5	62.5	62.5
1	1.2	0.693	1.03	10.9	192	63.9	54.7	46.8	103
2	1.33	0.509	1.06	9.05	351	51.1	40.5	32.2	139
3	1.4	0.402	1.1	8.39	460	37.8	29.5	23.1	172
4	1.43	0.339	1.13	8.15	517	28.9	22.8	18	201
6	1.46	0.279	1.17	8.03	561	20.6	16.6	13.4	235
9	1.47	0.252	1.19	8.01	579	17.1	13.9	11.3	251
12	1.47	0.245	1.19	8	584	16.3	13.2	10.7	255
16	1.47	0.244	1.2	8	585	16	13	10.6	256
20	1.47	0.243	1.2	8	585	16	13	10.6	256
saturiertes Modell:					600	10	10	10	260

Iterationsanfang

$$z_{00} = z_{01} = z_{11} = 1$$

Iterationsschritt

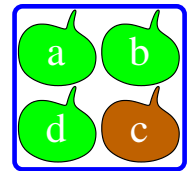
$$z_{\xi\eta} \leftarrow z_{\xi\eta} \cdot \left( \frac{n_{\xi\eta}/100}{\mathcal{E}[\varphi_{\xi\eta}(\mathbb{X})]} \right)^{1/5}$$

## Beispiel — Tafelobst im Tetrapack

### Diamantenes Verteilungsmodell

für die vier frischen/faulen Äpfel:

$$P : \begin{cases} \{0, 1\}^4 & \rightarrow [0, 1] \\ (\alpha, \beta, \zeta, \delta) & \mapsto z \cdot z_{\alpha\beta}^{ab} \cdot z_{\beta\zeta}^{bc} \cdot z_{\zeta\delta}^{cd} \cdot z_{\delta\alpha}^{da} \end{cases}$$



### Datensammlung und Statistiken

Absolute Häufigkeiten

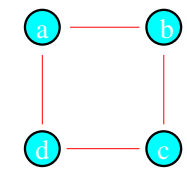
$$\{n_{\alpha\beta\zeta\delta} \mid \alpha, \beta, \zeta, \delta \in \{0, 1\}\}$$

Minimale suffiziente Statistiken, z.B. für  $ab \in \Delta$ :

$$n_{\alpha\beta}^{ab} = \sum_{a,b,c,d} \delta_{a=\alpha} \cdot \delta_{b=\beta} \cdot n_{abcd}$$

### Symmetrie I

$$n_{\xi\eta}^{ab} = n_{\xi\eta}^{bc} = n_{\xi\eta}^{cd} = n_{\xi\eta}^{da} = n_{\xi\eta}$$



### Symmetrie II

$$n_{01} = n_{10}$$

## MLS $\hat{=}$ relative Ereignishäufigkeiten

Happy End — für alle kausalen Verteilungen

### Zerlegbare Loglinearmodelle

Cliquen  $\mathcal{C} = \{C_1, \dots, C_M\}$

$$P(\mathbf{x}) = \prod_{C \in \mathcal{C}} z_{\mathbf{x}}^C = \prod_{C \in \mathcal{C}} \frac{P(\mathbf{x}_C)}{P(\mathbf{x}_{C \cap \pi(C)})}$$

### Bayesnetze

Ordnung  $V = \{V_1, \dots, V_N\}$

$$P(\mathbf{x}) = \prod_{n=1}^N P(x_n \mid \mathbf{x}_{\text{pa}(n)}) = \prod_{n=1}^N M_{x_n \mid \text{pa}(n)}(\mathbf{x})$$

### Maximum-Likelihood

$$\hat{z}_{\mathbf{x}}^C = \frac{n_{\mathbf{x}_C}}{n_{\mathbf{x}_{C \cap \pi(C)}}}$$

### Maximum-Likelihood

$$\hat{M}_{x_n \mid \text{pa}(n)}(\mathbf{x}) = \frac{n_{\mathbf{x}_{\{x_n\} \cup \text{pa}(n)}}}{n_{\mathbf{x}_{\text{pa}(n)}}}$$

Korrelation, Regression und Transinformation

Assoziationsregeln und Netzwerkanalyse

Bedingte statistische Unabhängigkeit

Graphische Modelle: ungerichtete Graphen

Kausale Modelle: gerichtete azyklische Graphen

Berechnen bedingter Wahrscheinlichkeiten

Parameterschätzung in Bayesnetzen und Loglinearmodellen

Aufdeckung der Abhängigkeitsstruktur

Kovarianzselektion

## Gütemaße für Modellstrukturen

### Definition

Mit den **Maximum-Likelihood-Parametern**

$$\hat{\theta}_{\Delta}(\omega) = \operatorname{argmax}_{\theta} \ell_{\omega}(\Delta, \theta) = \operatorname{argmax}_{\theta} \sum_{\mathbf{x} \in \omega} \log P(\mathbf{x} \mid \Delta, \theta)$$

und der ML-bezogenen Bewertung  $\hat{\ell}_{\omega}(\Delta) = \ell_{\omega}(\Delta, \hat{\theta}_{\Delta}(\omega))$  heißt

$$\operatorname{dev}(\Delta) \stackrel{\text{def}}{=} 2 \cdot \left( \hat{\ell}(\mathfrak{P}V) - \hat{\ell}(\Delta) \right), \quad \mathfrak{P}V = \{V\} = \text{saturiertes Modell}$$

die **Devianz** der Modellstruktur  $\Delta$  für die Daten  $\omega$ .

### Lemma

Das Devianzmaß besitzt die folgenden Eigenschaften:

1. Gilt  $\omega \sim P(\cdot \mid \Delta)$ , so ist die Devianz asymptotisch  $\chi_d^2$ -verteilt, wobei  $d$  die Differenz der Freiheitsgrade von  $\Delta$  und saturiertem Modell bezeichne.
2. Es gilt  $\operatorname{dev}(\mathfrak{P}V) = 0$  und  $\mathcal{E}[\operatorname{dev}(\Delta)] = d$ .

## Welches ist die beste Modellstruktur ?

Wahrscheinlichste Kombination aus Struktur und Parametern

### Gegeben

Datenprobe  $\omega$  aus der Objektmenge  $\Omega$  über den Variablen  $V$

$$\Omega = \bigotimes_{a \in V} \mathcal{X}_a$$

### Gesucht

Das bestpassende graphische/kausale/kordale/loglineare Modell

$$\begin{aligned} \hat{\Delta} &= \operatorname{argmax}_{\Delta \subset \mathfrak{P}V} J_{\omega}(\Delta) = \operatorname{argmax}_{\Delta \subset \mathfrak{P}V} \frac{f_{\text{prior}}(\Delta) \cdot P(\omega \mid \Delta)}{P(\omega)} \\ P(\omega \mid \Delta) &= \sum_{\theta \in \mathcal{M}(\Delta)} f_{\text{prior}}(\theta \mid \Delta) \cdot P(\omega \mid \Delta, \theta) \end{aligned}$$

### Markovnetze

$\binom{N}{2}$  Kanten & insgesamt  
 $2^{\binom{N}{2}}$  ungerichtete Graphen

### Bayesnetze

$N!$  Ordnungen & jeweils  
 $2^{\binom{N}{2}}$  zyklensfreie Graphen

### Loglinear

$2^N$  Terme  
 $2^{2^N}$  Modelle

## Einige regularisierte Gütemaße

### Kreuzvalidierung

Datenpartition  $\omega = \omega_a \uplus \omega_b$

$$J(\Delta) = \ell_{\omega_b}(\Delta, \hat{\theta}_{\Delta}(\omega_a))$$

### Rotationsvalidierung ( $L^1O$ )

„leave-one-out“  $\omega^{(x)} = \omega \setminus \{\mathbf{x}\}$

$$J(\Delta) = \sum_{\mathbf{x} \in \omega} \ell_{\{\mathbf{x}\}}(\Delta, \hat{\theta}_{\Delta}(\omega^{(x)}))$$

### ML-Bewertung + Strafterm

$$J(\Delta) = \hat{\ell}_{\omega}(\Delta) - \psi(N) \cdot |\theta_{\Delta}|$$

### AIC $\Rightarrow \psi(N) \equiv 1$

„Akaike Information Criterion“

### BIC $\Rightarrow \psi(N) = \frac{1}{2} \log N$

„Bayesian Information Criterion“

### Entropie

Bedingte Entropien  $\mathcal{H}(\mathbb{X}_n \mid \mathbf{x}) = - \sum_{\xi \in \mathcal{X}_n} P(\xi \mid \mathbf{x}) \cdot \log P(\xi \mid \mathbf{x})$

$$J(\Delta) = \mathcal{H}(\Delta) = \sum_{n=1}^N \sum_{\mathbf{x} \in \mathcal{X}_{\mathbf{pa}(n)}} P(\mathbf{x}) \cdot \mathcal{H}(\mathbb{X}_n \mid \mathbb{X}_{\mathbf{pa}(n)} = \mathbf{x})$$

## Die K2-Metrik für Bayesnetze

Cooper & Herskovitz, 1991

### Fakt

Eine perfekte Gütefunktion wäre die a posteriori Wahrscheinlichkeit  $P(\Delta|\omega)$  der Modellstruktur auf Basis der Datenprobe.

### Gleich- und Dirichletverteilungsannahme

für Bayesnetzstruktur  $\Delta$  und -parameter  $\mathbf{M}_{n|pa(n)}$ :

$$P(\Delta|\omega) \propto P(\omega|\Delta) = \int \underbrace{\mathcal{D}(\theta|\Psi) \cdot P(\omega|\theta, \Delta)}_{P(\omega^{(\Psi)}|\theta, \Delta)} d\theta$$

### K2-Metrik

Geschlossene Darstellung der a posteriori Wahrscheinlichkeit:

$$J(\Delta) = \prod_{n=1}^N \prod_{x \in \mathcal{X}_{pa(n)}} \frac{(L_n - 1)!}{(n_x^{pa(n)} + L_n - 1)!} \cdot \prod_{\xi \in \mathcal{X}_n} n_{x, \xi}^{pa(n), \{n\}}$$

## SFS — Sequential Forward Selection

Gierige bottom-up Suche (Whitney 1971 · Buntine 1991)

### 1 INITIALISIERUNG

$$\mathcal{G} = (V, \emptyset)$$

### 2 AUSWAHL

einer nützlichsten neuen Kante

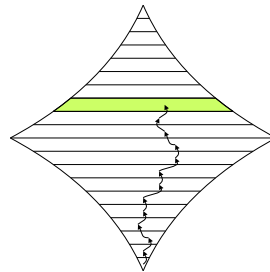
$$e^* = \operatorname{argmax} \{J(E, e) \mid e \in \mathcal{E}_V \setminus E\}$$

### 3 TERMINIERUNG

Wenn  $J(E, e^*) \leq J(E)$

dann  $\rightsquigarrow$  **ENDE**

sonst  $E \leftarrow E \cup \{e^*\}$  und  $\rightsquigarrow$  2.



#### Bemerkung

SFS trifft voreilige Entscheidungen (Horizont=1) und verfehlt i.a. die Optimallösung.

$$E^{(1)} \subset E^{(2)} \subset E^{(3)} \subset \dots$$

## Suchverfahren

Wer findet die Stecknadel im Heuhaufen vor Anbruch des jüngsten Tages ?

### Ungerichtete Graphen — Markovnetze

Gesucht ist eine J-optimale Teilmenge von

$$\mathcal{E}_V = \{\{a, b\} \mid a, b \in V, a \neq b\}$$

➔ Jedes  $\mathcal{E} \subseteq \mathcal{E}_V$  ist „erlaubt“!

### Kombinatorische Merkmalauswahl

Alle „wrapper“-Verfahren sind sinngemäß anwendbar:

- $\left\{ \begin{array}{l} \text{backward} \\ \text{forward} \end{array} \right\}$  selection: sukzessiv Kanten  $\left\{ \begin{array}{l} \text{entfernen} \\ \text{einfügen} \end{array} \right\}$
- pulsierende Suche · geordnete Suche · evolutionäre Suche

## SBE — Sequential Backward Elimination

Gierige top-down Suche (Marill & Green 1963 · Edwards/MIM 1995)

### 1 INITIALISIERUNG

$$\mathcal{G} = (V, \mathcal{E}_V)$$

### 2 AUSWAHL

einer nutzlosesten alten Kante

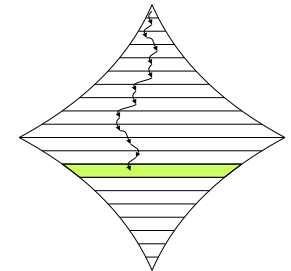
$$e^* = \operatorname{argmax} \{J(E \setminus e) \mid e \in E\}$$

### 3 TERMINIERUNG

Wenn  $J(E \setminus e^*) \leq J(E)$

dann  $\rightsquigarrow$  **ENDE**

sonst  $E \leftarrow E \setminus \{e^*\}$  und  $\rightsquigarrow$  2.



#### Bemerkung

SBE aufwändiger als SFS:

Start mit umfangreicheren E  
Längerer Weg zum Ziel

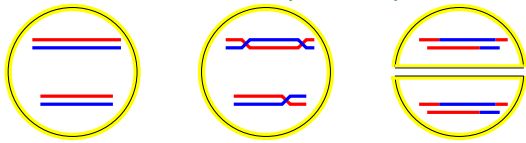




## Beispiel — Chromosomensequenzierung

Wir sortieren das Erbgut von „*barley powdery mildew fungus*“

### Haploide Vererbung (Meiose)



nach der Verschmelzung Crossover-Operation nach der Zellteilung

### Erhebung

70 Geschwisterindividuen

6 binäre phänotypische Attribute

$\mathbb{X}_i \stackrel{1:1}{=} \text{unbekannter Genlocus}$

### Hypothese über Genloci $\mathbb{X}_1, \dots, \mathbb{X}_6$

- unterschiedliche Chromosomen  $\rightsquigarrow$  unabhängig
- gleiches Chromosom  $\rightsquigarrow$  distanzabhängig korreliert
- Sequenz von Genen  $g_1, g_2, g_3 \Rightarrow \mathfrak{S}(g_1 \mid g_2 \mid g_3)$

### Resultat

$$d \longleftrightarrow a \longleftrightarrow b \longleftrightarrow e \longleftrightarrow c \longleftrightarrow f$$

## K2-Algorithmus

Elternsuchverfahren (Cooper & Herskovits, 1992)

(Algorithmus)

### 1 INITIALISIERUNG

Eine Variablenordnung ist a priori vorzugeben:

$$V = \{x_1, \dots, x_N\}, \quad n = 2$$

### 2 ELTERNAUSWAHL

Triff eine Vorwärtsauswahl (SFS) bezüglich K2-Bewertung:

$$\text{pa}(x_n) = \arg\max \{J(A) \mid A \subseteq \{x_1, \dots, x_{n-1}\}\}$$

### 3 TERMINIERUNG

Wenn  $n < N$  dann  $n \leftarrow n + 1$  und  $\rightsquigarrow$  2 sonst  $\rightsquigarrow$  ENDE.

(sumf3nog(A))

## Suchverfahren

Die Stecknadel piekt jetzt nur noch auf einer Seite !

### Gerichtete azyklische Graphen — Bayesnetze

Optimale Teilmenge von  $\mathfrak{C}_{\{x_i\}} = \{(a, b) \mid a \neq b\}$

- UG-Kantenselektion — Test auf Kordalität
- DAG-Kantenselektion — Test auf Zyklen
- DAG-Kantenselektion — Test auf Zyklen und Moralität

Optimale Teilmengen von  $\mathfrak{C}_{(x_i)} = \{(x_i, x_j) \mid 1 \leq i < j \leq N\}$

- Lineare Variablenordnung  $V = \{\mathbb{X}_1, \dots, \mathbb{X}_n\}$  vorlegen
- Optimale Eltermenge  $B_n \subseteq V_n$  für jedes  $\mathbb{X}_n$  berechnen
- (zulässiges Verfahren sofern  $J(\cdot)$  in „Familienterme“ zerfällt)

Exakte Suche für eingeschränkte Netzstrukturen

- Bäume & Fallschirme
- Minimaler Spannbaum (Cormen, Leiserson, Rivest 1990)

## Tetrad III Algorithmus

UG-Kantenselektion (Scheines 1996)

(Algorithmus)

### 1 KORRELATIONSTEST

Lösche Verbindungskanten für marginal unabhängige ZV

$$\mathfrak{S}(\mathbb{X}_i \mid \emptyset \mid \mathbb{X}_j)$$

### 2 PARTIELLE UNABHÄNGIGKEIT

Lösche Verbindungskanten für bedingt unabhängige ZV

$$\mathfrak{S}(\mathbb{X}_i \mid \text{bd}(\mathbb{X}_i) \cup \text{bd}(\mathbb{X}_j) \mid \mathbb{X}_j)$$

### 3 Teste $m$ -elementige Teilmengen von $C_{ij} = \text{bd}(\mathbb{X}_i) \cup \text{bd}(\mathbb{X}_j)$ .

### 4 ORIENTIERUNGSPHASE

1. Wähle eine Variablenordnung  $V = \{\mathbb{X}_1, \dots, \mathbb{X}_N\}$ .
2. Orientiere alle Kanten gemäß Ordnungsindex.
3. Ergänze Kanten für „unshielded collider“ und orientiere sie.

(sumf3nog(A))

## Bayesian Network SFFS

Pulsierende DAG-Kantenselektion für Bayesnetze (Blanco & Inza, 2002)

(Algorithmus)

### 1 INITIALISIERUNG

$$\mathcal{G} = (V, \emptyset), \quad n = 0, \quad \iota_0 = J(\emptyset)$$

### 2 VORWÄRTSAUSWAHL

$$\mathbf{e}^* = \operatorname{argmax} \{J(E, \mathbf{e}) \mid \mathbf{e} \in \mathfrak{E}_{\{x_i\}} \setminus E \wedge \text{DAG}(E, \mathbf{e})\}$$

Setze  $E \leftarrow E \cup \{\mathbf{e}^*\}$ ,  $n \leftarrow |n| + 1$  und  $\iota_n = J(E)$  und dann  $\rightsquigarrow$  2.

### 3 RÜCKWÄRTSAUSWAHL

$$\mathbf{e}^* = \operatorname{argmax} \{J(E \setminus \mathbf{e}) \mid \mathbf{e} \in \mathfrak{E}_{\{x_i\}} \setminus E\}$$

Wenn  $J(E \setminus \mathbf{e}^*) \leq \iota_{n-1}$  dann  $\rightsquigarrow$  2.

Sonst setze  $E \leftarrow E \setminus \{\mathbf{e}^*\}$ ,  $n \leftarrow n - 1$  und  $\iota_n = J(E)$  und  $\rightsquigarrow$  3.

### 4 TERMINIERUNG

Wenn Zielkardinalität  $n = n_0$  erreicht dann  $\rightsquigarrow$  ENDE.

(zumft3h0g1A)

## TBN — Baumförmige Bayesnetze

Erinnerung: moralische Bayesnetze  $\hat{=}$  zerlegbare Markovnetze

### Modellformel für ein TBN mit Wurzel $\mathbb{X}_{i_0}$

$V = \{x_1, \dots, x_N\}$  und  $\pi : V \setminus \{i_0\} \rightarrow V$  mit  $\text{pa}(x_j) = \{x_{\pi_j}\}$  für  $j \neq i_0$ :

$$P(\mathbf{x}) = P(x_{i_0}) \cdot \prod_{j \neq i_0} P(x_j | x_{\pi_j}) = \prod_{i=1}^n P(x_i) \cdot \prod_{j \neq i_0} \underbrace{\frac{P(x_j, x_{\pi_j})}{P(x_j) \cdot P(x_{\pi_j})}}_{\exp(\mathfrak{I}(x_j; x_{\pi_j}))}$$

Nur die **punktweisen Transinformationen** sind abhängig von der Baumstruktur!

### Relevanter Anteil der logarithmierten Likelihood-Zielgröße

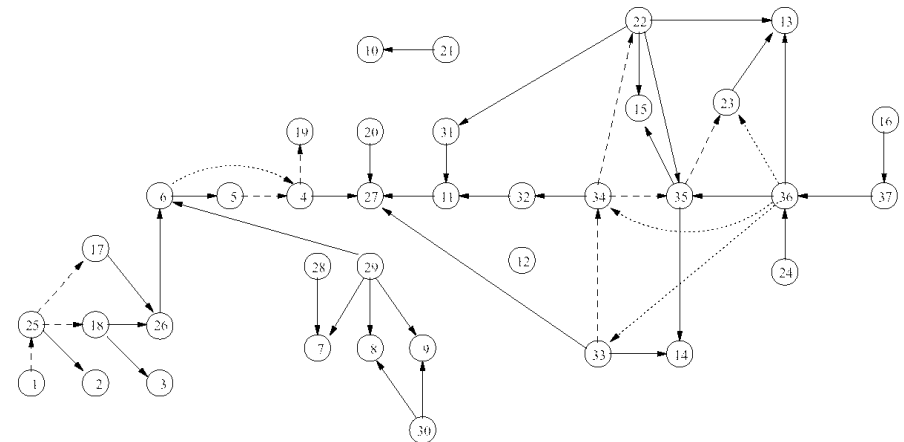
für Lerndatenprobe  $\omega$  und Baumkantenmenge  $E = \{\{j, \pi_j\} \mid j \neq i_0\}$ :

$$\ell_{\text{ML}}(\omega | E) = \sum_{(i,j) \in E} \underbrace{\{\mathcal{H}(\omega, P_{\mathbb{X}_i}) + \mathcal{H}(\omega, P_{\mathbb{X}_j}) - \mathcal{H}(\omega, P_{\mathbb{X}_i \mathbb{X}_j})\}}_{\mathfrak{I}_{\omega}(\mathbb{X}_i; \mathbb{X}_j)}$$

➔ Berechnung aller **empirischen Transinformationen**

## Beispiel — Alarmkette

37 Attribute · 46 → 45 Kanten · 370 Lernbeispiele



Gelernt: 45 Kanten

## Kausalfpfadanalyse mit baumförmigen Bayesnetzen

Suche nach dem minimalen Spannbaum (Chow & Liu 1968)

(Algorithmus)

### 1 INITIALISIERUNG

Berechne alle Transinformationswerte  $(i, j = 1, \dots, N)$ :

$$\text{TI}(\mathbb{X}_i, \mathbb{X}_j) \stackrel{\text{def}}{=} \sum_{x_i \in \mathcal{X}_i} \sum_{x_j \in \mathcal{X}_j} P(x_i, x_j) \cdot \log \frac{P(x_i, x_j)}{P(x_i) \cdot P(x_j)}$$

### 2 BEWERTETER GRAPH

Erzeuge  $\tilde{\mathcal{G}} = (V, V^2, \beta)$  mit der Kantengewichtung

$$\beta : \begin{cases} V^2 & \rightarrow \mathbb{R} \\ \{x_i, x_j\} & \mapsto -\text{TI}(\mathbb{X}_i, \mathbb{X}_j) \end{cases}$$

### 3 SPANNBAUM ( $O(N^2 \log N)$ ) SLAC — „single-linkage agglomerative clustering“

Konstruiere den minimalen spannenden Baum  $\mathcal{G} \subset \tilde{\mathcal{G}}$ .

### 4 ORIENTIERUNG VON $\mathcal{G}$

Wähle eine beliebige Wurzelvariable  $v_0 \in V$ .

Alle Kanten von  $\mathcal{G}$  werden „wurzelwärts“ gerichtet.

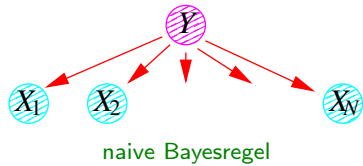
(zumft3h0g1A)

## Klassifizieren mit Bayesnetzen

$\mathbb{X}_1, \dots, \mathbb{X}_N \rightarrow \mathbb{Y} \in \{1, 2, \dots, K\}$

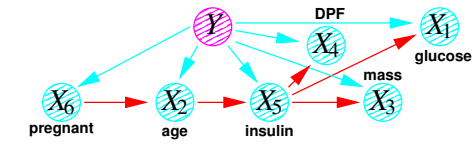
### Naives Bayesnetz

$$P(y) \cdot \prod_n P(x_n | y)$$



### Tree Augmented Bayesnet

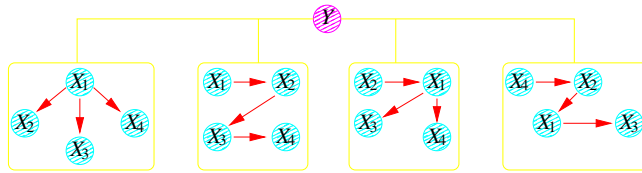
$$P(x_1 | y) \cdot \prod_{n=2}^N P(x_n | y, x_{pa(n)})$$



TABN (Friedman, Geiger, Goldszmidt 1998)

### Bayes-Multinetz

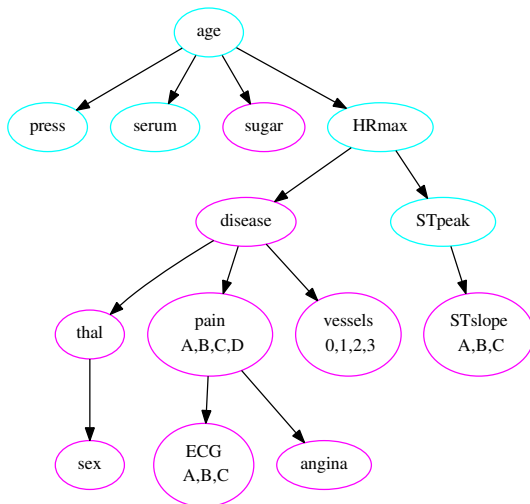
$$P(y) \cdot P(\mathbf{x} | E_y, \theta_y)$$



„Bayes wood“ (Heckerman 1991)

## Beispiel — Statlog Herzdatensammlung

13 Attribute · 270 Objekte · Klassifikation: „disease“



$\mathcal{G}(\mathbb{X}_{14} | \mathbb{X}_3, \mathbb{X}_8, \mathbb{X}_{12}, \mathbb{X}_{13} | \text{„Rest“})$

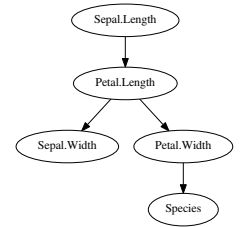
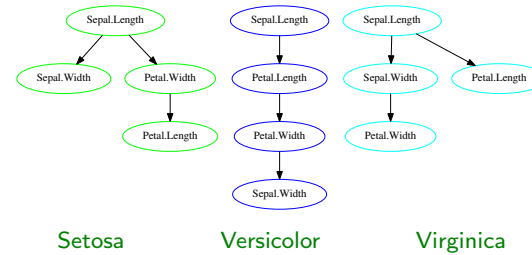
1. age (IR)
2. sex {male, female}
3. chest pain {A, B, C, D}
4. blood pressure (IR)
5. serum cholestoral (IR)
6. fasting blood sugar {T, F}
7. resting ECG results [0 : 2]
8. maximum heart rate achieved (IR)
9. exercise induced angina {T, F}
10. ST depression (exercise:rest) (IR)
11. slope of peak exercise ST {A, B, C}
12. vessels colored by flourosopy [0 : 3]
13. thal {normal, fixed, defect}
14. heart disease {T, F}

## Beispiel — Fishers Irisdatensatz

5 Attribute ( $\mathbb{R}^4 \times \{1, 2, 3\}$ ) · 150 Objekte (50 je Spezies)

### Transinformatiionsmatrix

		$X_1$	$X_2$	$X_3$	$X_4$	$X_5$
$l_{sepal}$	$X_1$	0	1.33	2.04	1.88	0.69
$w_{sepal}$	$X_2$	1.33	0	1.43	1.40	0.38
$l_{petal}$	$X_3$	2.04	1.43	0	2.64	1.37
$w_{petal}$	$X_4$	1.88	1.40	2.64	0	1.43
species	$X_5$	0.69	0.38	1.37	1.43	0



Bayesbaum für alle fünf Attribute

Bayeswald ein Baum je Spezies

Korrelation, Regression und Transinformation

Assoziationsregeln und Netzwerkanalyse

Bedingte statistische Unabhängigkeit

Graphische Modelle: ungerichtete Graphen

Kausale Modelle: gerichtete azyklische Graphen

Berechnen bedingter Wahrscheinlichkeiten

Parameterschätzung in Bayesnetzen und Loglinearmodellen

Aufdeckung der Abhängigkeitsstruktur

Kovarianzselektion

## Stetige Loglinearmodelle

Motivation: multivariate Normalverteilungsdichte

### Definition

Es sei die  $N$ -dimensionale multivariate Normalverteilungsdichte

$$\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \mathbf{S}) = |\mathbf{S}|^{-1/2} \cdot \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{S}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

mit dem Mittelwertvektor  $\boldsymbol{\mu}$  und der Kovarianzmatrix  $\mathbf{S}$  gegeben.  
Die Werte  $\alpha$ ,  $\beta_i$ ,  $\kappa_{ij}$  der Darstellung

$$\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \mathbf{S}) = \exp\left(\alpha + \sum_i \beta_i \cdot x_i + \sum_{i,j} \kappa_{ij} \cdot x_i x_j\right)$$

heißen **kanonische Parameter** der exponentiellen Familie; die Matrix  $\mathbf{K} = [\kappa_{ij}]$  heißt **Konzentrationsmatrix** oder **Präzisionsmatrix**.

### Kanonische Parameter & Standardparameter

$$\alpha = -\frac{1}{2} \cdot (\log |2\pi \mathbf{S}| + \boldsymbol{\mu}^\top \mathbf{S}^{-1} \boldsymbol{\mu}), \quad \beta = \mathbf{S}^{-1} \boldsymbol{\mu}, \quad \mathbf{K} = \mathbf{S}^{-1}.$$

### Loglinearmodelle

Die Kovarianzterme in  $\Delta$  sind die nichtnegativen Koeffizienten der Summationsterme

$$u^{(i_1, \dots, i_N)} \cdot x_1^{i_1} x_2^{i_2} x_3^{i_3} \dots x_n^{i_n} \dots x_{N-1}^{i_{N-1}} x_N^{i_N}$$

des Dichtefunktionsexponenten. Insbesondere fällt dem Term

$$u^{(0, \dots, 0)} \cdot x_1^0 x_2^0 \dots x_N^0 = u^{(0, \dots, 0)} \cdot 1 = u^{(0, \dots, 0)}$$

wieder die Rolle des Normierungsfaktors zu.

Die Vektoren  $\mathbf{i}$  können wir auch als *Multimengen* von Zufallsvariablen auffassen.

### Gaußsche graphische Modelle

Hier werden ausschließlich Kovarianzterme  $\mathbf{i} \in \Delta$  zugelassen mit

$$\sum_{n=1}^N i_n = i_1 + i_2 + i_3 + \dots + i_N \leq 2.$$

### Beweis.

Das Modell ist auch graphisch, denn es gibt grundsätzlich keinerlei Interaktion zwischen mehr als zwei Variablen. Es gilt  $\mathfrak{S}(A \mid Z \mid B)$  genau dann, wenn  $\Delta$  ausschließlich  $AU\mathcal{Z}$ -Terme und  $BU\mathcal{Z}$ -Terme enthält, also wenn es keine  $\{a, b\}$ -Terme mit  $a \in A$  und  $b \in B$  gibt.  $\square$

## Stetige Loglinearmodelle

### Definition

Sei  $\Delta \subset \mathbb{N}^N$  eine (endliche) Menge von Exponenten- $N$ -Tupeln.  
Die Familie stetiger Wahrscheinlichkeitsdichtefunktionen der Gestalt

$$\log f_\Delta(\mathbf{x}) = \sum_{\mathbf{i} \in \Delta} u^{\mathbf{i}} \cdot \prod_{n=1}^N x_n^{i_n}, \quad \mathbf{x} \in \mathbb{R}^N$$

heißt **stetiges Loglinearmodell** über  $V$  mit der Menge  $\Delta$  von **Kovarianztermen**.

### Lemma

Für Loglinearmodelle  $\Delta$ , die Normalverteilungen sind, gilt:

$$\Delta \text{ hierarchisch} \quad \Rightarrow \quad \Delta \text{ graphisch}$$

Wir nennen diese Familien **Gaußsche Graphische Modelle** oder **Kovarianzselektionsmodelle**.

## Wissenswertes über multivariate Normalverteilungsdichten

### Lemma

Für normalverteilte Zufallsvariablen  $V = \{\mathbb{X}_1, \dots, \mathbb{X}_N\}$  gelten die folgenden Aussagen:

1. **Summenbildung:**  $\mathbb{X} = \mathbb{X}' + \mathbb{X}'' \sim \mathcal{N}(\boldsymbol{\mu}' + \boldsymbol{\mu}'', \mathbf{S}' + \mathbf{S}'')$
2. **Affine Abbildung:**  $\mathbf{A}\mathbb{X} + \mathbf{b} \sim \mathcal{N}(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\mathbf{S}\mathbf{A}^\top)$
3. **Marginalisierung:**  $\mathbb{X}_C \sim \mathcal{N}(\boldsymbol{\mu}_C, \mathbf{S}_{CC})$
4. **Konditionierung:**  $\mathbb{X}_{A|\mathbf{x}_B} \sim \mathcal{N}(\boldsymbol{\mu}_{A|\mathbf{x}_B}, \mathbf{S}_{A|\mathbf{x}_B})$

Dabei gelte  $A \uplus B = V$  und es sind definiert:

$$\begin{aligned} \boldsymbol{\mu}_{A|\mathbf{x}_B} &= \boldsymbol{\mu}_A + \mathbf{S}_{AB} \cdot \mathbf{S}_{BB}^{-1} \cdot (\mathbf{x}_B - \boldsymbol{\mu}_B) & \boldsymbol{\mu} &= (\boldsymbol{\mu}_A^\top, \boldsymbol{\mu}_B^\top)^\top \\ \mathbf{S}_{A|\mathbf{x}_B} &= \mathbf{S}_{AA} - \mathbf{S}_{AB} \cdot \mathbf{S}_{BB}^{-1} \cdot \mathbf{S}_{BA} & \mathbf{S} &= \begin{pmatrix} \mathbf{S}_{AA} & \mathbf{S}_{AB} \\ \mathbf{S}_{BA} & \mathbf{S}_{BB} \end{pmatrix} \end{aligned}$$

5. Für die bedingte Kreuzkovarianzmatrix gilt der Zusammenhang:

$$(\mathbf{S}_{A|\mathbf{x}_B})^{-1} = \mathbf{K}_{AA} = (\mathbf{S}^{-1})_{AA}$$

## Marginalisierung

Ähnlich wie schon zuvor für Vektoren definieren wir Matrixausschnitte durch

$$M_{A,B} \stackrel{\text{def}}{=} (M_{ab} \mid a \in A, b \in B) .$$

Die Matrix  $S_{CC}$  insbesondere enthält also alle Varianzen von und Kovarianzen zwischen Variablen aus  $C$ .

Die Matrix  $S_{AB}$  heißt übrigens auch „Kreuzkovarianzmatrix“ der Variablenmengen  $A$  und  $B$ .

## Konditionierung

Bei geeigneter Variablennummerierung gilt in der Situation  $A \uplus B = V$ :

$$S = \begin{pmatrix} S_{AA} & S_{AB} \\ S_{BA} & S_{BB} \end{pmatrix} = \begin{pmatrix} S_{AA} & S_{AB} \\ S_{AB}^\top & S_{BB} \end{pmatrix}, \quad K = \begin{pmatrix} K_{AA} & K_{AB} \\ K_{BA} & K_{BB} \end{pmatrix}$$

Die Matrixgleichung ergibt sich aus der (unschönen!) Formel zur Blockmatrixinvertierung.

## Beweis.

Die partielle Unabhängigkeit, d.h. die Frage nach einer Kante oder keiner Kante zwischen zwei Variablen im Markovnetz, läßt sich ganz einfach aus der inversen Kovarianzmatrix  $K$  ablesen.

### Beweisidee 1:

Betrachte die bedingte Verteilung mit  $A = \{a, b\}$  und  $B = V \setminus \{a, b\}$ .

Das Variablenpaar  $(X_a, X_b)$  ist, bei gegebenem  $x_B$ , mit der Kovarianzmatrix  $S_{ab|}$  normalverteilt.

$X_a, X_b$  sind unabhängig genau dann, wenn  $S_{ab|}$  eine Diagonalmatrix ist; dies aber ist genau dann der Fall, wenn ihre Inverse, also  $K_{\{a,b\}}$  diagonal ist, also falls  $\kappa_{ab} = \kappa_{ba} = 0$  ist.

### Beweisidee 2:

Die Normalverteilungsdichte ist faktorisiert in Gibbs-Komponenten mit maximal zwei Variablen.

Sie läßt sich also in zwei Faktoren  $g_{V \setminus \{a\}}$  und  $h_{V \setminus \{b\}}$  genau dann zerlegen, wenn mindestens die Gibbs-Komponente für  $\{a, b\}$  fehlt.

□

## Kovarianz und Konzentration

Nulleinträge  $\hat{=}$  marginale & partielle Unabhängigkeiten

### Lemma (Wermuth 1976)

Für normalverteilte Variablen  $a, b \in V \sim \mathcal{N}(\mu, S)$  mit  $a \neq b$  gilt:

- **Marginale Unabhängigkeit:**  $\mathfrak{I}(a \mid \emptyset \mid b) \iff s_{ab} = 0$
- **Partielle Unabhängigkeit:**  $\mathfrak{I}(a \mid \text{Rest} \mid b) \iff \kappa_{ab} = 0$

## Gaußscher Diamant

$$\begin{pmatrix} \kappa_{11} & \kappa_{12} & 0 & \kappa_{14} \\ \kappa_{21} & \kappa_{22} & \kappa_{23} & 0 \\ 0 & \kappa_{32} & \kappa_{33} & \kappa_{34} \\ \kappa_{41} & 0 & \kappa_{43} & \kappa_{44} \end{pmatrix}$$



$$\begin{array}{cc} X_1 & - & X_2 \\ | & & | \\ X_4 & - & X_3 \end{array}$$

## Gaußscher Schlüssel

$$\begin{pmatrix} \kappa_{11} & \kappa_{12} & 0 & 0 & 0 \\ \kappa_{21} & \kappa_{22} & \kappa_{23} & \kappa_{24} & 0 \\ 0 & \kappa_{32} & \kappa_{33} & \kappa_{34} & \kappa_{35} \\ 0 & \kappa_{42} & \kappa_{43} & \kappa_{44} & \kappa_{45} \\ 0 & 0 & \kappa_{53} & \kappa_{54} & \kappa_{55} \end{pmatrix}$$



$$X_1 - X_2 \triangleleft \begin{array}{c} X_4 \\ X_3 \end{array} \triangleright X_5$$

## Charakterisierung bedingter Unabhängigkeiten

### Lemma

Für Variablenmengen  $A, B, Z$  mit  $V = A \uplus B \uplus Z$  und die bedingte Kreuzkovarianzmatrix

$$S_{AB|Z} = [s_{ab}]_{a \in A, b \in B}^{b \in B}, \quad s_{ab} \stackrel{\text{def}}{=} \text{Cov}[X_a, X_b \mid X_Z = x_Z]$$

gilt die Beziehung:

$$S_{AB|Z} = S_{AB} - S_{AZ} \cdot (S_{ZZ})^{-1} \cdot S_{ZB}$$

### Satz (Speed & Kiiveri 1986)

Für normalverteilte Variablen  $V = A \uplus B \uplus Z$  mit Kovarianzmatrix  $S$  sind äquivalent:

1.  $S_{AB} = S_{AZ} \cdot (S_{ZZ})^{-1} \cdot S_{ZB}$
2.  $(S^{-1})_{AB} = 0$  beziehungsweise  $K_{AB} = 0$
3.  $\mathfrak{I}(A \mid Z \mid B)$

# Maximum-Likelihood-Schätzung

für Gaußsche Graphische Modelle

## Satz (Dempster 1972)

Es sei  $\mathcal{G} = (V, \mathcal{E})$  ein Gaußsches Graphisches Modell mit der Generatorenmenge  $\mathcal{C} \subset \mathfrak{P}V$  und sei  $\omega \subset \mathbb{R}^N$  ein Datensatz mit den Statistiken  $\mathbf{m}$  und  $\Sigma$ . Dann bilden  $\mathbf{m}$  und  $\{\Sigma_{CC} \mid C \in \mathcal{C}\}$  eine minimale suffiziente Statistik des Modells für  $\omega$ .

Die Maximum-Likelihood-Parameter  $\mu$  und  $\mathbf{S}$  bzw.  $\mathbf{K}$  gehorchen den Bedingungen

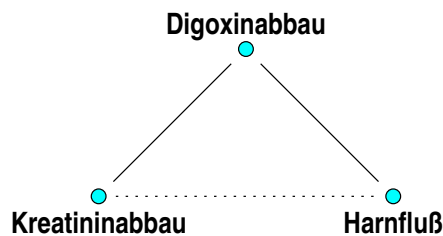
$$\begin{aligned} \mu &= \mathbf{m} \\ s_{ab} &= \sigma_{ab} & \{a, b\} \in \mathcal{E} \vee a = b \\ \kappa_{ab} &= 0 & \{a, b\} \notin \mathcal{E} \wedge a \neq b \end{aligned}$$

## Bemerkung

Da  $\mathcal{N}(\mu, \Sigma)$  das saturierte Modell ist, beträgt die Devianz:

$$\text{dev}(\mathcal{C}) = 2 \cdot (\ell(\mu, \Sigma) - \ell(\mu, \mathbf{S})) = T \cdot \log(\det \mathbf{S} / \det \Sigma)$$

## Beispiel — Digoxin-Abbau



## Datensammlung

$$\omega \subset \mathbb{R}^3$$

$$|\omega| = 35 \text{ Patienten}$$

$\mathbb{X}$  = Abbau von Kreatinin

$\mathbb{Y}$  = Abbau von Digoxin

$\mathbb{Z}$  = Harnflußrate

$$\chi^2\text{-Test} \Rightarrow \{\mathbb{X}, \mathbb{Z}\} \notin \mathcal{E}$$

# Maximum-Likelihood-Schätzung

Existenz & Eindeutigkeit

## Datenkovarianzmatrix

$$\Sigma = \begin{pmatrix} 3.023 & 1.258 & 1.004 \\ 1.258 & 1.709 & 0.842 \\ 1.004 & 0.842 & 1.116 \end{pmatrix}$$

## ML-Kovarianzmatrix

$$\mathbf{S} = \begin{pmatrix} 3.023 & 1.258 & 0.620 \\ 1.258 & 1.709 & 0.842 \\ 0.620 & 0.842 & 1.116 \end{pmatrix}$$

## ML-Konzentrationsmatrix

$$\mathbf{K} = \begin{pmatrix} 0.477 & -0.351 & 0.000 \\ -0.351 & 1.190 & -0.703 \\ 0.000 & -0.703 & 1.426 \end{pmatrix}$$

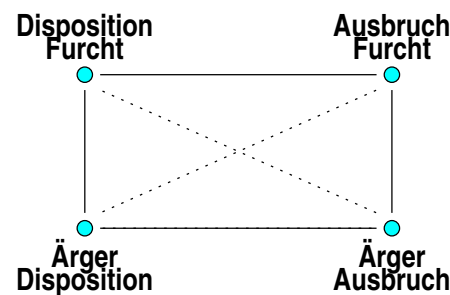
## Satz (Dempster 1972)

Es seien  $\mathbf{A}, \mathbf{B}$  zwei symmetrische, positiv-definite  $(N \times N)$ -Matrizen. Ferner sei  $\mathcal{E} \subseteq \{1, \dots, N\} \times \{1, \dots, N\}$  symmetrisch mit  $(i, i) \in \mathcal{E}$  für alle  $i$ . Dann gibt es eine symmetrische, positiv-definite Matrix  $\mathbf{S}$  mit

$$\begin{aligned} s_{ij} &= a_{ij} & \forall (i, j) \in \mathcal{E} \\ (\mathbf{S}^{-1})_{ij} &= b_{ij} & \forall (i, j) \notin \mathcal{E} \end{aligned}$$

und  $\mathbf{S}$  ist eindeutig mit diesen Eigenschaften.

## Beispiel — Furcht versus Ärger



## Datensammlung

$$\omega \subset \mathbb{R}^4$$

$$|\omega| = 684 \text{ Versuchspersonen}$$

### • Augenblickszustand:

$\mathbb{X}$  = Furcht

$\mathbb{W}$  = Ärger

### • mentale Prägung:

$\mathbb{Z}$  = Furcht

$\mathbb{Y}$  = Ärger

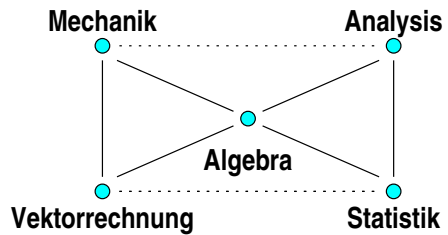
$\chi^2$ -Test ergibt:

$$\mathfrak{S}(\mathbb{X} \mid \mathbb{W}, \mathbb{Z} \mid \mathbb{Y})$$

$$\mathfrak{S}(\mathbb{W} \mid \mathbb{X}, \mathbb{Y} \mid \mathbb{Z})$$



## Beispiel — Punktezah in Übungsserien



### Datensammlung

$\omega \subset \{0, 1, 2, \dots, 100\}^5$   
 $|\omega| = 88$  Studierende

- 5 Übungsgruppen in 5 Fächern
- je 100 Punkte erzielbar

### Inferenz

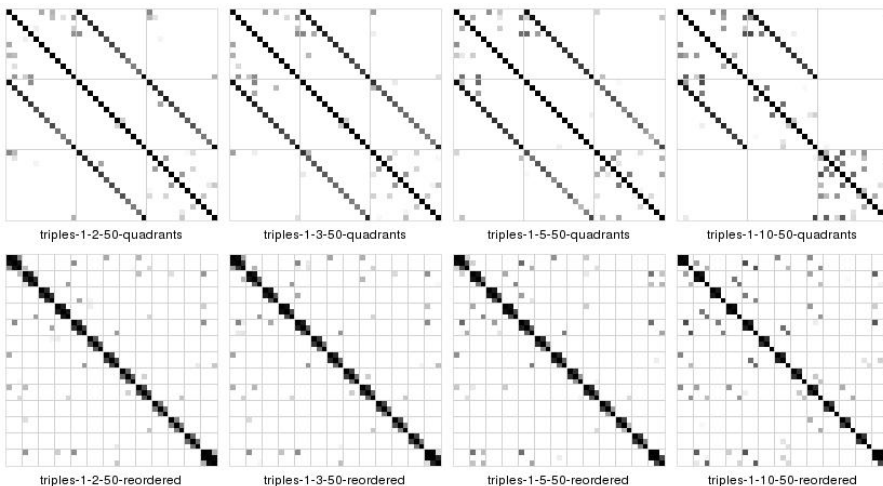
$\mathbb{S}(\text{Stat} \mid \text{Alg}, \text{Analysis} \mid \cdot)$

$\mathbb{S}(\text{Mech} \mid \text{Alg}, \text{Vektor} \mid \cdot)$

Zentrale Befähigung: **Algebra**

## Beispiel — Sprachsignalparameter

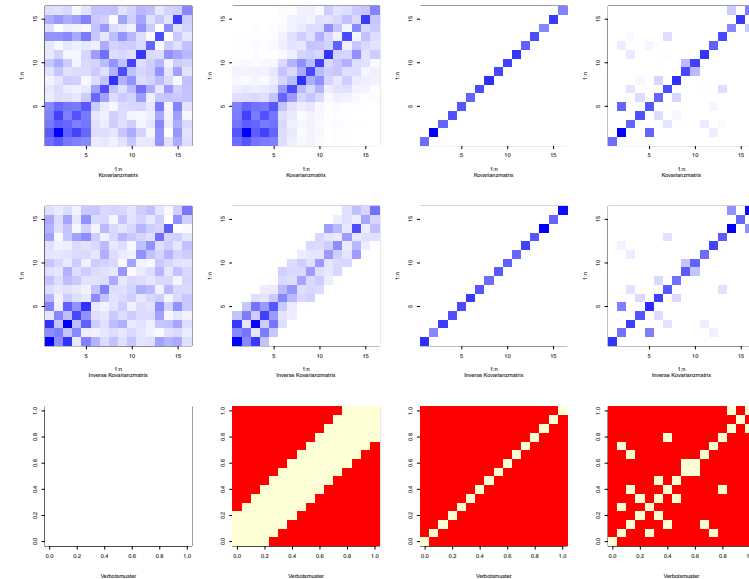
12 MFCC-Parameter · drei Zeitpunkte



### Dünne Abhängigkeitsstruktur

Die Vergangenheit wird durch unmittelbare Vorgänger „maskiert“.

## Beispiel — Schriftzeichenklassifikation



### Beispiel

Datensatz  
 letter  
 16 Merkmale  
 alle Klassen

oben:  
 Kovarianz  
 $\hat{S} = C^{-1}$

Mitte:  
 Konzentration  
**C** erfüllt **A**

unten:  
 Adjazenz **A**  
 Abhängigkeits-  
 muster  
 (gegeben)

Korrelation, Regression und Transinformation

Assoziationsregeln und Netzwerkanalyse

Bedingte statistische Unabhängigkeit

Graphische Modelle: ungerichtete Graphen

Kausale Modelle: gerichtete azyklische Graphen

Berechnen bedingter Wahrscheinlichkeiten

Parameterschätzung in Bayesnetzen und Loglinearmodellen

Aufdeckung der Abhängigkeitsstruktur

Kovarianzselektion

## Zusammenfassung (6)

1. **Kovarianz** und **Korrelation** sind **quantitative** Charakterisierungen der **linearen** Aspekte („*Regression*“) statistischer Abhängigkeit.
2. Die **Transinformation** quantifiziert statistische Abhängigkeit in allgemeiner Form, setzt aber die Kenntnis der **wahren Verteilung** voraus.
3. Die **Warenkorbanalyse** sucht **Assoziationsregeln** mit gleichermaßen hohen Werten für **Support**, **Konfidenz** und **Relevanz** (z.B. **Apriori-Algorithmus**).
4. Das **Dependenzmodell** charakterisiert die **bedingten Unabhängigkeiten**  $\mathfrak{I}(A \mid Z \mid B)$  je dreier Attributmengen einer Verteilung.
5. Verteilungen heißen **graphisch (kausal)**, wenn ihr DM durch die ( $\delta$ -)**Separation** eines **UG (DAG)** gegeben ist.
6. Kausale Modelle faktorisieren **attributweise** in **bedingte Wahrsch'keiten**, graphische Modelle faktorisieren **cliquenweise** in **Gibbspotenziale**.
7. **Kordale Modelle** besitzen einen **triangulierten** UG, einen **moralischen** DAG und eine **Cliquenschnittfaktorisierung**.
8. Statistische **Inferenz** ist nur für **Ketten** und (Verbund-)**Bäume** effizient.
9. Die **ML-Schätzung** der **Modellparameter** aus Daten beruht auf **relativen Häufigkeiten** (DAG) oder dem **Maximum-Entropie-Prinzip** (UG).
10. Die Aufdeckung der **Modellstruktur** basiert auf **Unabhängigkeitstests** (Kantenelimination/Grenzengraph) oder **gieriger Suche** mit **Strafterm**.