

Grundlagen der Künstlichen Intelligenz

11 Bayessche Netze

Struktur, Semantik, Konstruktion, Inferenz

Volker Steinhage

Inhalt

- Motivation
- Bayessche Netze
- Bayessche Netze und Verbundwahrscheinlichkeitsverteilungen
- Bayessche Netze und bedingte Unabhängigkeiten
- Konstruktion von Bayesschen Netzen
- Inferenz in Bayesschen Netzen
- Exakte und approximative Inferenz

Motivation (1)

Die Verteilung der Verbundwahrscheinlichkeiten

	zahnschmerzen		\neg zahnschmerzen	
	verfangen	\neg verfangen	verfangen	\neg verfangen
loch	0.108	0.012	0.072	0.008
\neg loch	0.016	0.064	0.144	0.576

- erlaubt die Beantwortung aller Anfragen an die Domäne bzgl. unbedingter und bedingter W'keiten, indem die Anfragen als Disjunktion über den atomaren Ereignissen formuliert und die entspr. W'keiten aufaddiert werden
- wächst aber exponentiell in der Zahl der Zufallsvariablen und ist bzgl. der Erfassung der Verbundwahrscheinlichkeiten nicht naheliegend und einfach

Die Bayessche Regel mit der Annahme von absoluten und bedingten Unabhängigkeiten zwischen Zufallsvariablen erlaubt bereits eine effiziente Reduzierung der zur Anfragebeantwortung erforderlichen unbedingten und bedingten Wahrscheinlichkeiten

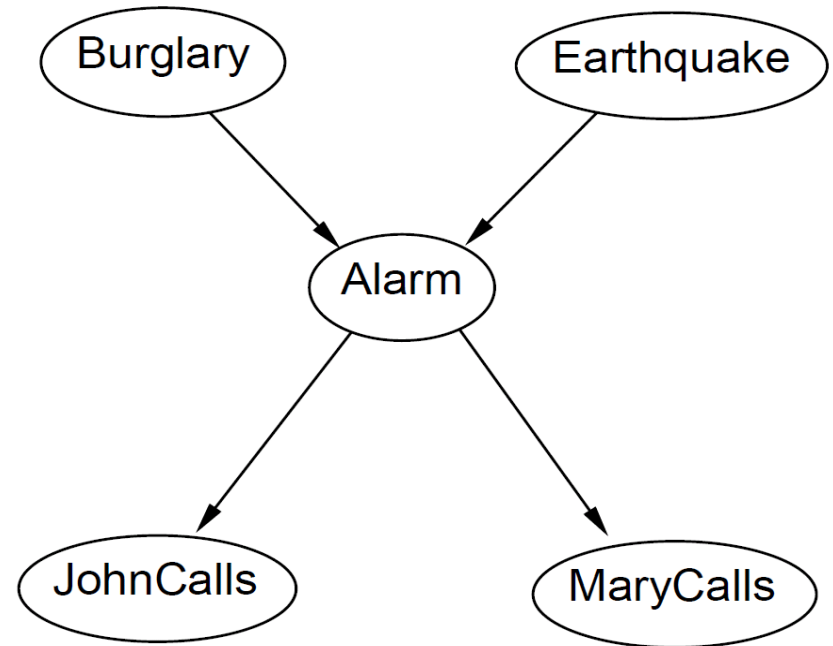
Motivation (2)

Bayessche Netze sind eine Repräsentationsform

- der Verteilung der Verbundwahrscheinlichkeiten
- zur Darstellung von Unabhängigkeiten zwischen Zufallsvariablen und erlauben so die effiziente Durchführung von Inferenz

Beispiel von Juda Pearl

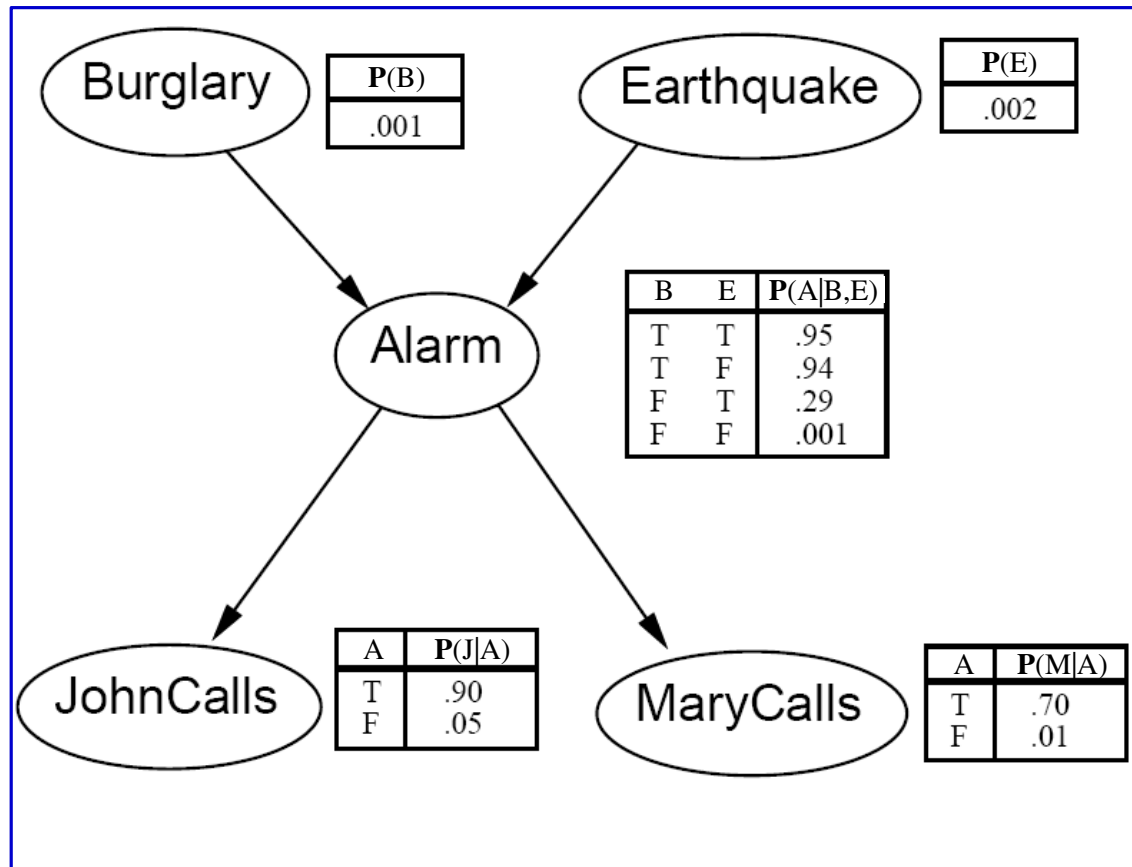
- Ihr Haus hat eine Einbruchsicherung (*Alarm*), die auf Einbrüche (*Burglary*), aber auch schon auf leichte Erdbeben* (*Earthquake*) reagiert
- Die Nachbarn John und Mary sagen zu, Sie bei Alarm im Büro anzurufen (*JohnCalls* bzw. *MaryCalls*)
- John verwechselt manchmal das Telefonläuten mit Alarm und ruft dann auch an
- Mary hört manchmal laute Musik und überhört dann den Alarm



* Juda Pearl lebt und lehrt in Los Angeles

Struktur von Bayesschen Netzen*

- 1) **Knoten**: Zufallsvariablen
- 2) **Gerichtete Kanten** zwischen Knoten: *direkte* Einflüsse bzw. Abhängigkeiten
- 3) Zu jedem Knoten gibt es eine **Tabelle von bedingten W'keiten** (*Conditional Probability Table, CPT*), die den Effekt der **Elternknoten** auf den Knoten quantifiziert



- 4) Der Graph ist **azyklisch**; also ein DAG

* auch *belief networks, probabilistic networks, causal networks*

Semantik von Bayesschen Netzen

Zwei Zugänge zum Verständnis von Bayesschen Netzen (BN):

1. BN repräsentiert die **Verbundwahrscheinlichkeitsverteilung** der Zufallsvariablen

~ Geeignete Sichtweise für die **Konstruktion des BNs**

2. BN kodiert eine Menge von **Unabhängigkeitsannahmen**

~ Geeignete Sichtweise zur **Konstruktion von Inferenzen**

Kodierung von Unabhängigkeitsannahmen in Bayesschen Netzen

- Allgemein: BNs repräsentieren *vollständig* die Abhängigkeiten von direkten Elternknoten

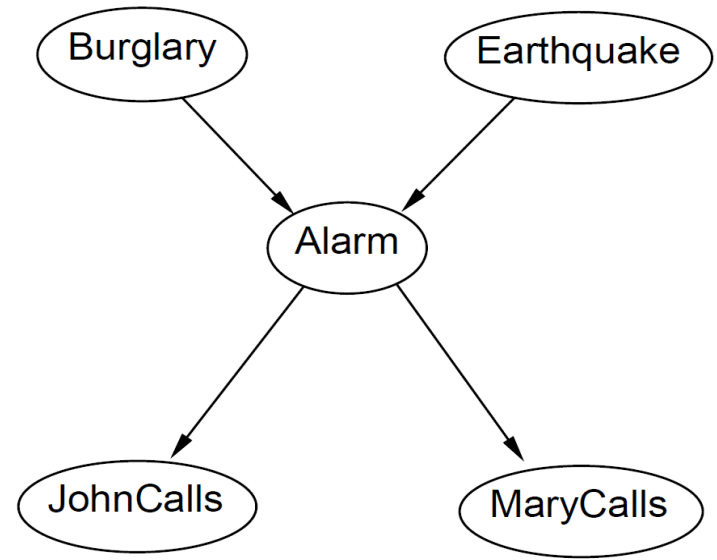
- am Beispiel:

- *Alarm* hängt *nur* von *Burglary* und *Earthquake* ab
- *MaryCalls* hängt *nur* von *Alarm* ab

~ *MaryCalls* ist also unabhängig von *JohnCalls*, *Earthquake* und *Burglary*:

$$\mathbf{P}(\text{MarryCalls} | \text{JohnCalls}, \text{Alarm}, \text{Earthquake}, \text{Burglary}) = \mathbf{P}(\text{MaryCalls} | \text{Alarm})$$

- ~ Bayessche Netze kodieren damit also auch *Unabhängigkeitsannahmen*



Bayessche Netze und Verbundwahrscheinlichkeit

Ein BN ist eine **kompakte Repräsentation einer Verbundw'keitsverteilung**:

Jedes atomare Ereignis der Verteilung ist eine Konjunktion einer bestimmten Wertebelegung $P(X_1 = x_1 \wedge \dots \wedge X_n = x_n)$, abgekürzt: $P(x_1, \dots, x_n)$.

Für jedes atomare Ereignis gilt nach der **Produktregel**:

$$\begin{aligned} P(x_1, \dots, x_n) &= P(x_n \mid x_{n-1}, \dots, x_1) \cdot \dots \cdot P(x_2 \mid x_1) P(x_1) \\ &= \prod_{i=1}^n P(x_i \mid x_{i-1}, \dots, x_1). \end{aligned}$$

Wegen der **Unabhängigkeitsannahmen** ist dies äquivalent zu:

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i \mid \text{parents}(X_i)).$$

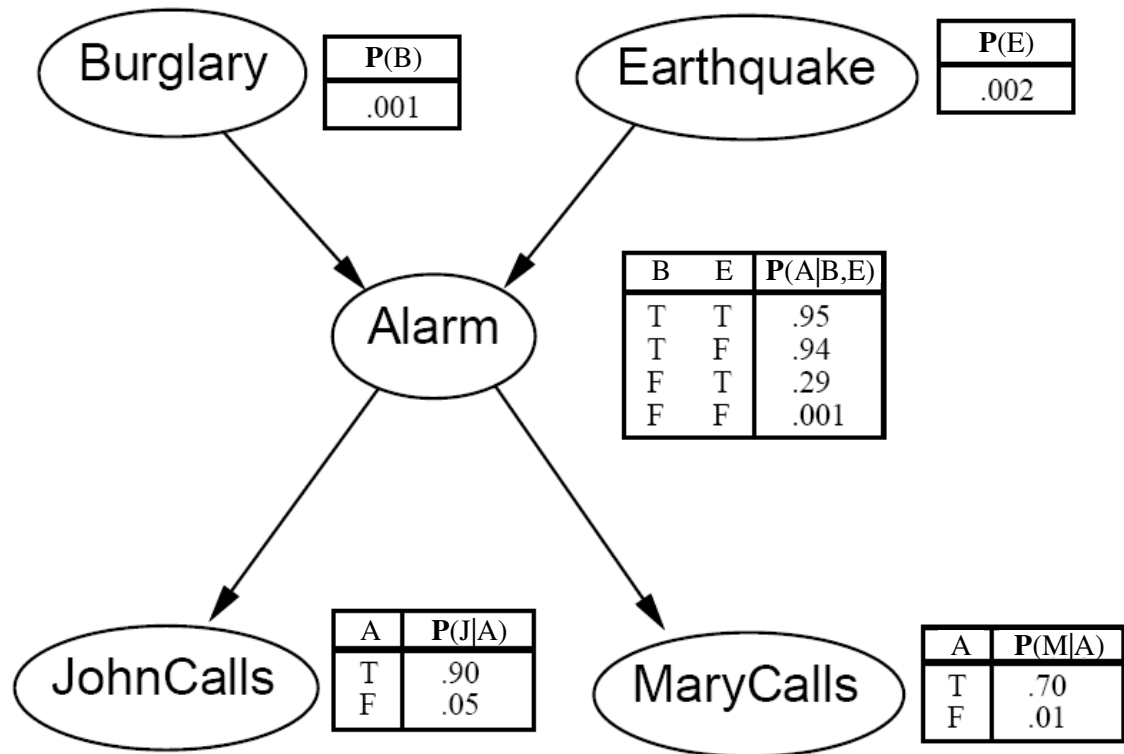
→ aus **Netztopologie** und **CPTs** sind alle Verbundw'keiten ableitbar

Beispiel für Berechnung eines atomaren Ereignisses

- W'keiten für negative Ereignisse ergeben sich als

$$P(\neg x) = 1 - P(x).$$

- W'keiten für atomare Ereignisse ergeben sich durch Faktorisierung über die Produktregel:



$$\begin{aligned}
 P(j, m, a, \neg b, \neg e) &= P(j | a) \cdot P(m | a) \cdot P(a | \neg b, \neg e) \cdot P(\neg b) \cdot P(\neg e) \\
 &= 0.9 \cdot 0.7 \cdot 0.001 \cdot 0.999 \cdot 0.998 \\
 &= 0.00062
 \end{aligned}$$

Kompaktheit Bayesscher Netze

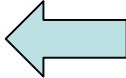
- Zur expliziten Repräsentation der **Verbundwahrscheinlichkeitsverteilung** brauchen wir eine Tabelle der Größe 2^n bei n Booleschen Variablen
- Falls in einem BN mit n Knoten jeder Knoten maximal k Eltern hat, brauchen wir nur n **Tabellen der Größe 2^k** bei booleschen Variablen

Beispiel: $n = 20$ und $k = 5$

$$\rightarrow 2^n = 2^{20} = 1.048.576 \quad \text{vs.} \quad n \cdot 2^k = 20 \cdot 2^5 = 640$$

- Im ungünstigsten Fall kann natürlich auch ein BN exponentiell groß werden (wenn jede Variable von jeder anderen direkt beeinflusst wird)
- Abhängigkeit von der *Strukturiertheit* der Anwendungsdomäne (lokale vs. globale Interaktion) und dem Geschick des Designs

Entwurf eines Bayesschen Netzes (1)

1. Wähle Menge von relevanten Variablen, welche die Domäne beschreiben
2. Ordne alle Variablen 
3. Nimm erste Variable in der Liste
4. Gib alle direkten Einflüsse von Knoten, die schon im Netz sind, auf den Knoten für diese Variable an: Kanten + CPT
5. Streiche die Variable aus der Liste
6. Solange Liste nicht leer: gehe zu Schritt 3

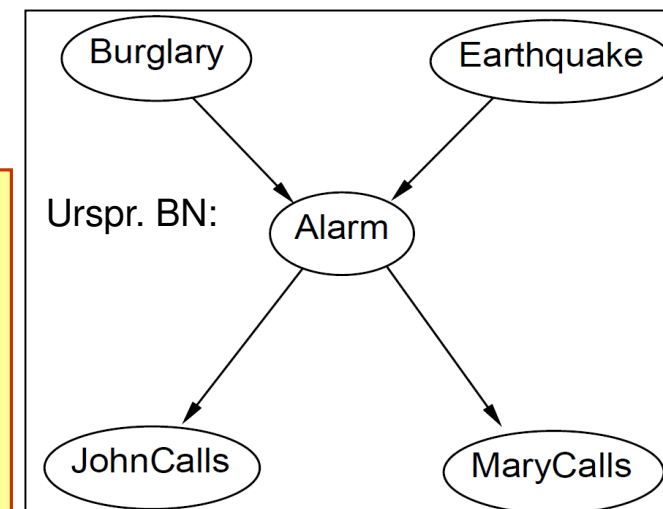
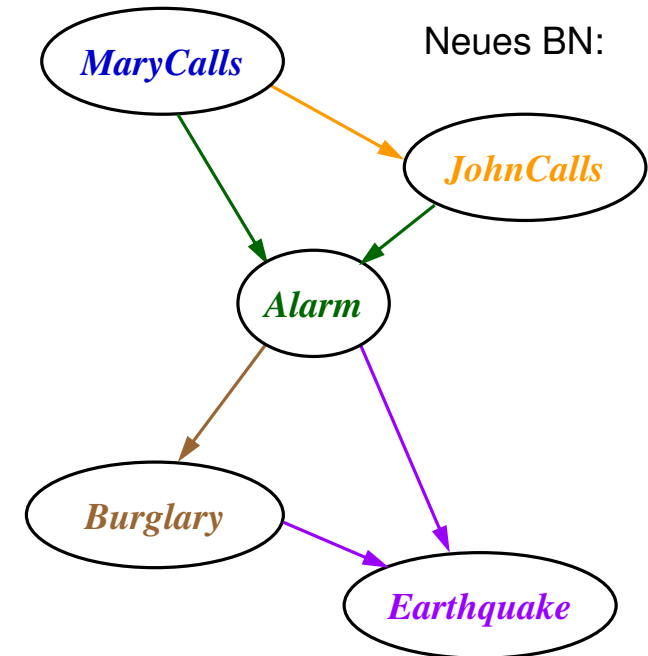
Frage: Welche Ordnung der Liste ist geeignet?

Beispiel (1)

Ordnung: MaryCalls, JohnCalls, Alarm, Burglary, Earthquake

- 1) Wähle *MaryCalls*: \sim keine Elternknoten
- 2) Wähle *JohnCalls*: Bei Evidenz *MaryCalls* sei *JohnCalls* wahrscheinlicher
 - $\sim P(\text{JohnCalls} \mid \text{MaryCalls}) \neq P(\text{JohnCalls})$
 - \sim *MaryCalls* wird Elternknoten von *JohnCalls*
- Wähle *Alarm*: Bei Evidenzen *MaryCalls* und *JohnCalls* sei *Alarm* wahrscheinlicher
 - \sim *MaryCalls* und *JohnCalls* werden Elternknoten von *Alarm*
- Wähle *Burglary*: Hierfür sei die Evidenz *Alarm* alleine hinreichend
 - \sim *Alarm* wird Elternknoten von *Burglary*
- Wähle *Earthquake*: Bei alleiniger Evidenz *Alarm* sei *Earthquake* wahrscheinlich; bei gemeinsamer Evidenz *Alarm* und *Burglary* sei *Earthquake* weniger wahrscheinlich
 - \sim *Burglary* und *Alarm* werden Elternknoten von *Earthquake*

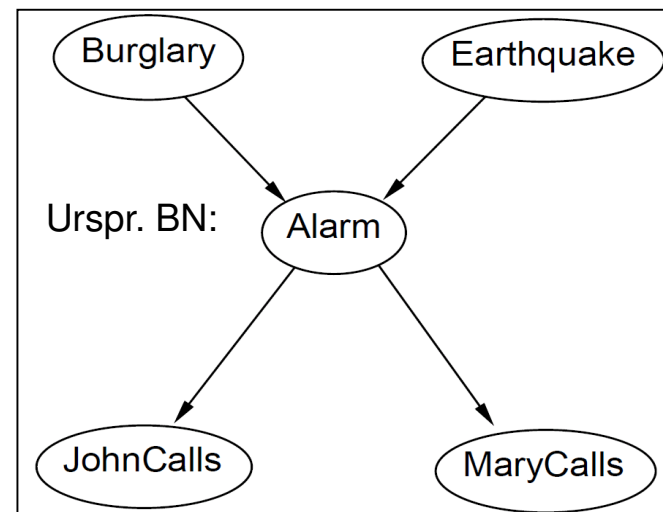
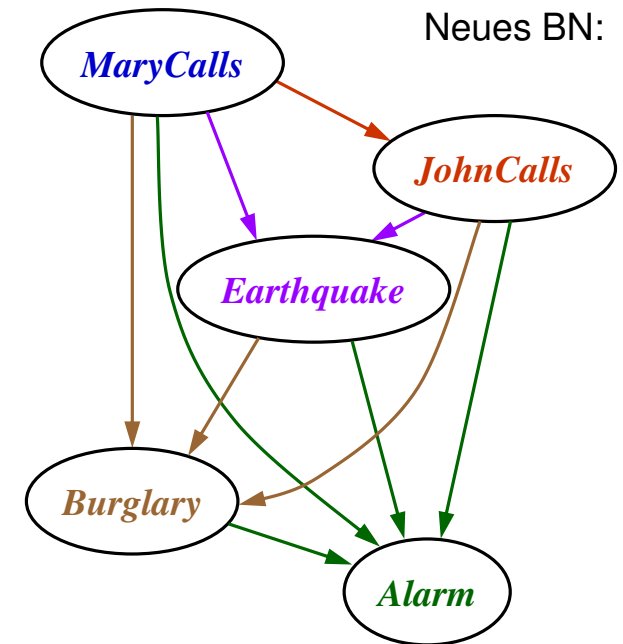
- Die Zahl der Abhängigkeiten (Kanten) ist gegenüber dem urspr. Netzentwurf von vier auf sechs um zwei gestiegen
- Die Zahl der zu ermittelnden Wahrscheinlichkeiten ist um drei gestiegen
- Gravierend ist die Qualität der neuen Abhängigkeiten: Wie soll z.B. $P(\text{Earthquake} \mid \text{Alarm}, \text{Burglary})$ erfasst werden?



Beispiel (2)

- Ordnung: *MaryCalls*, *JohnCalls*, *Eathquake*,
Burglary, *Alarm*

- ~ Die Zahl der Abhängigkeiten (Kanten) hat sich gegenüber dem ursprünglichen Netzentwurf mehr als verdoppelt
- ~ Die Zahl der zu ermittelnden Wahrscheinlichkeiten ist auf 31 gestiegen und entspricht somit der vollständigen Verbundwahrscheinlichkeitsverteilung
- ~ Wiederum nur schwer zu erfassende neue Abhängigkeiten



Entwurf eines Bayesschen Netzes (2)

- Der Aufbau eines „*diagnostischen Netzes*“ führt zu Bedingungen von Symptomen zu Ursachen und damit zu neuen Abhängigkeiten zwischen ansonsten unabhängigen Ursachen – und oft auch zwischen unabhängigen Symptomen
- Besser ist der Aufbau eines „*kausalen Netzes*“:
 - starte mit den grundlegenden Ursachen (*root causes*)
 - erweitere schrittweise jeweils um die direkten Auswirkungen
 - bis zu den Blattknoten, die ohne Auswirkungen auf andere Variable sind
- Bemerkung: alle drei Netze des Beispiels repräsentieren dieselbe Verteilung der Verbundw'keiten, berücksichtigen jedoch in unterschiedlichem Maße Unabhängigkeiten!

Bayessche Netze und Graphische Modelle

Bayessche Netze zählen zu den sogenannten *Graphischen Modellen*:

- Graphische Modelle werden als Kombination von Wahrheitstheorie und Graphentheorie betrachtet
- Andere Formen von *Graphischen Modellen* sind Markov Random Fields, Conditional Random Fields, Faktorgraphen, u.a.

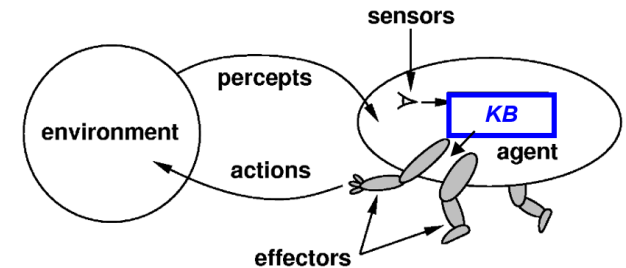
-
- Bisher: Betrachtung von BNs zur **effizienten Kodierung der Wahrheiten**
 - Jetzt: Betrachtung des operationellen Teils, nämlich die **Inferenz** in BNs

Inferenz in Bayesschen Netzen

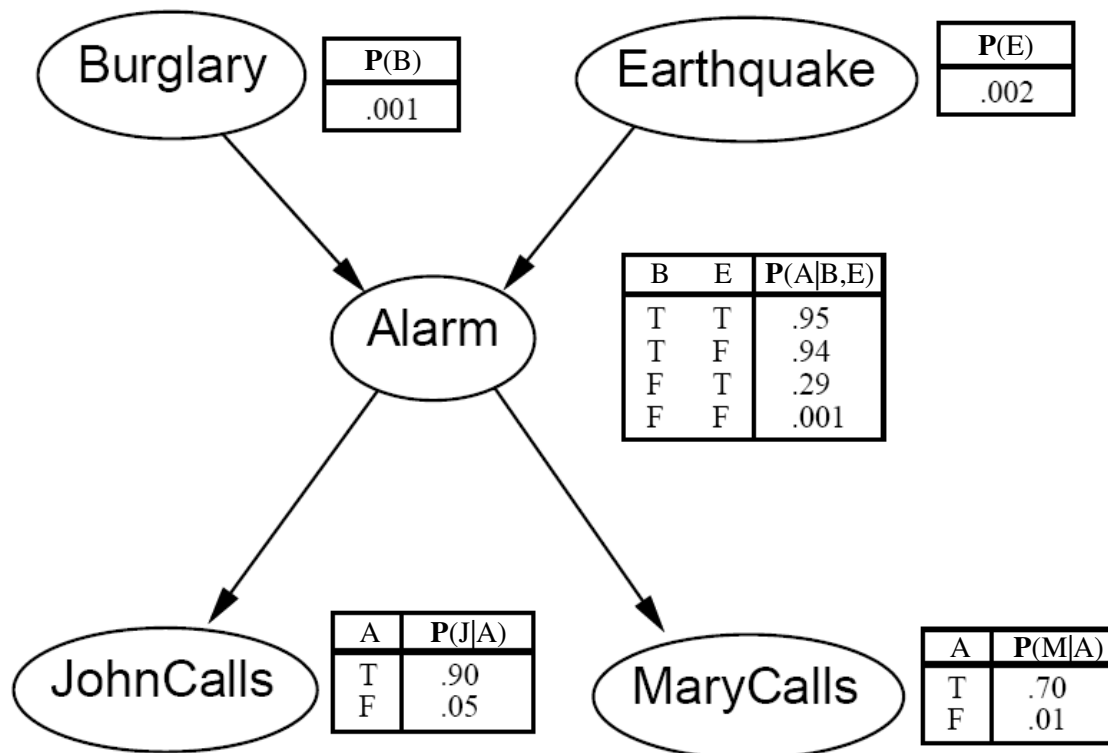
Probabilistische Inferenz:

Gegeben: **Instanzierte Evidenzvariablen**

Gesucht: Wahrscheinlichkeitsverteilung von **Anfragevariablen**:



$$P(\text{Query} \mid \text{Evidence})$$



Exakte Inferenz durch Aufzählen (1)

Aufgabe: Anfrage $P(X|e)$ bei Belegung e der Evidenzvariablenmenge E

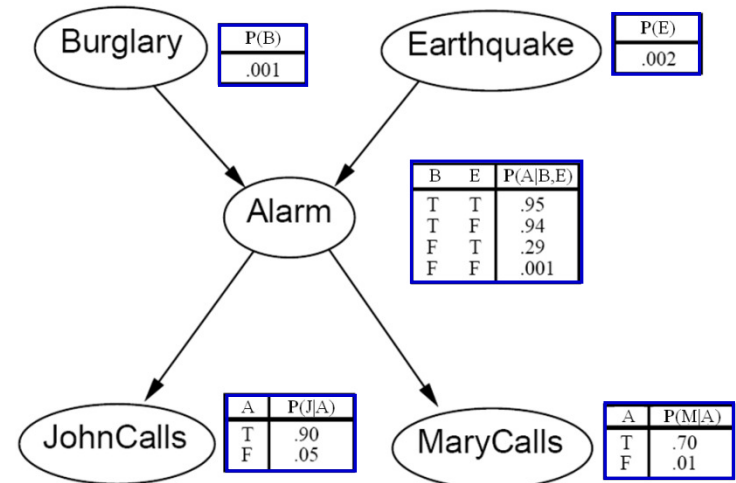
- ~ Aus letzter Vorlesung: Summation der Verbundw'keiten über der Menge Y aller unbeobachteten Variablen:

$$P(X | e) = \alpha \cdot P(X, e) = \alpha \cdot \sum_y P(X, e, y)$$

- Wegen Unabhängigkeitsannahmen der BNs ist die Anfrage durch die Summierung über Produkten von bedingten Wahr'keiten im Netz zu beantworten:

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i)).$$

- Die bedingten W'keiten sind in den CPTs des BNs notiert



Exakte Inferenz durch Aufzählen (2)

- Eine systematische Methode, unbeobachtete Variablen („hidden variables“) aus der Verbundverteilung heraus zu marginalisieren

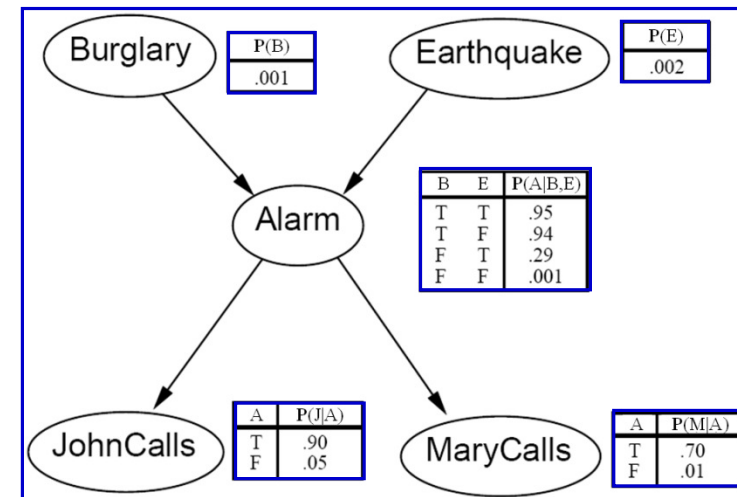
- Einfache Anfrage im „Burglary“-Netzwerk sei:

$P(\text{Burglary} \mid \text{JohnCalls} = \text{true}, \text{MaryCalls} = \text{true})?$

- Abgekürzt: $P(B|j,m) = P(B,j,m) / P(j,m)$ // per Def.

$$= \alpha \cdot P(B,j,m)$$

$$= \alpha \cdot \sum_e \sum_a P(B,e,a,j,m)$$



- Faktorisiere die Verbundverteilung als Produkt von CPT-Einträgen:

$$P(B|j,m) = \alpha \cdot \sum_e \sum_a P(B) \cdot P(e) \cdot P(a|B,e) \cdot P(j|a) \cdot P(m|a)$$

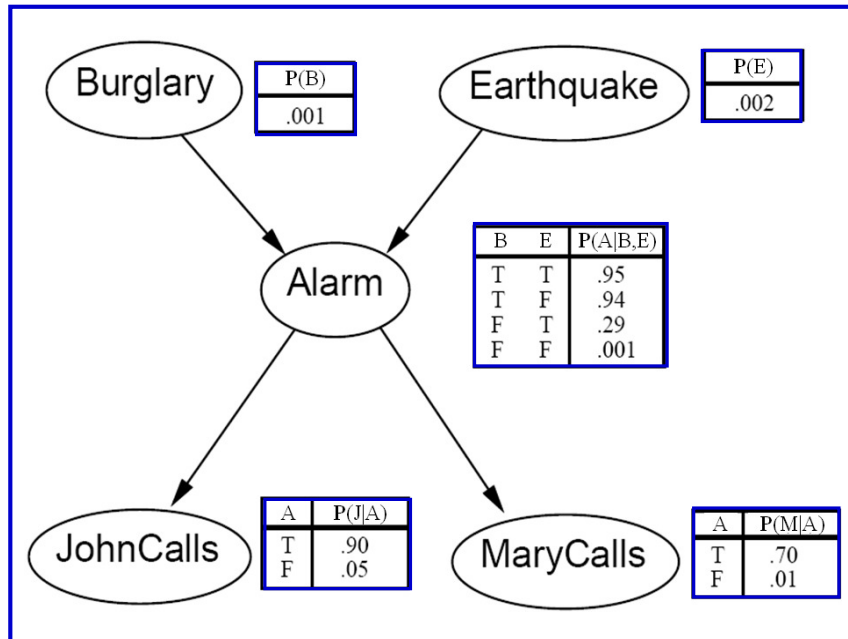
$$= \alpha \cdot P(B) \cdot \sum_e P(e) \cdot \sum_a P(a|B,e) \cdot P(j|a) \cdot P(m|a)$$

- Rekursive *depth-first*-Aufzählung der Faktoren für n Boolesche Variable in $O(n)$ Platz- und $O(2^n)$ Zeitkomplexität

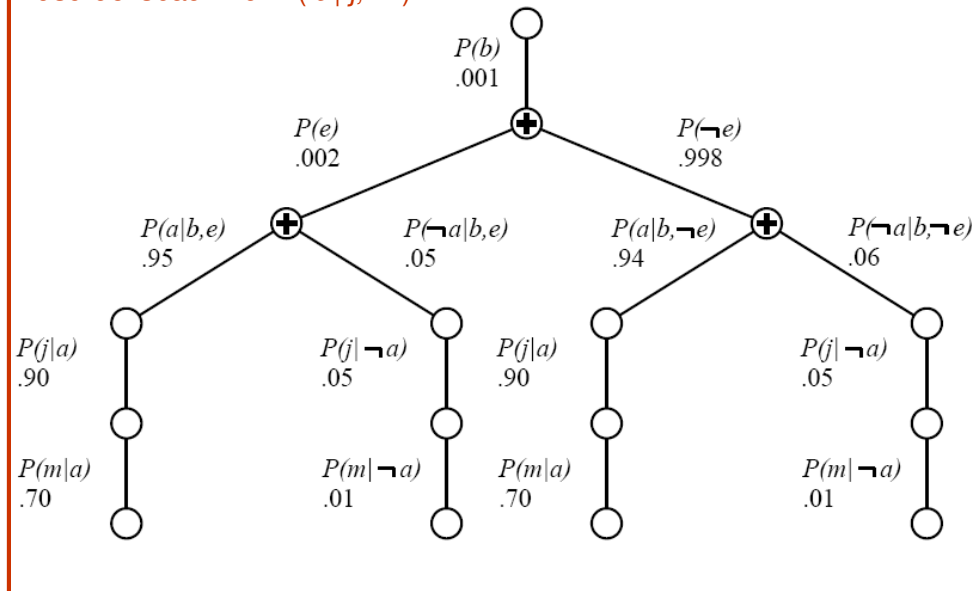
Exakte Inferenz durch Aufzählen (3)

Beispiel: $P(\text{Burglary} \mid \text{JohnCalls} = \text{true}, \text{MaryCalls} = \text{true})$?

$$\sim P(B|j,m) = \alpha \cdot P(B) \cdot \sum_e P(e) \cdot \sum_a P(a|B,e) \cdot P(j|a) \cdot P(m|a)$$



Ausdrucksbaum für $P(b \mid j, m)$:



- Führt zu $P(b \mid j, m) = \alpha \cdot 0.00059224$ und analog $P(\neg b \mid j, m) = \alpha \cdot 0.0014919$
- Normierung: $P(B \mid j, m) = \alpha \cdot \langle 0.00059224, 0.0014919 \rangle = \langle 0.284, 0.716 \rangle$
- Bei beiden Anrufen besteht die Wahr'keit von 28,4% für einen Einbruch

Aufzählungsalgorithmus

function **ENUMERATION-ASK**(X, \mathbf{e}, bn) **returns** a distribution over X

inputs: X , the query variable
 \mathbf{e} , observed values for variables \mathbf{E}
 bn , a Bayesian network with variables $\{X\} \cup \mathbf{E} \cup \mathbf{Y}$

$Q(X) \leftarrow$ a distribution over X , initially empty

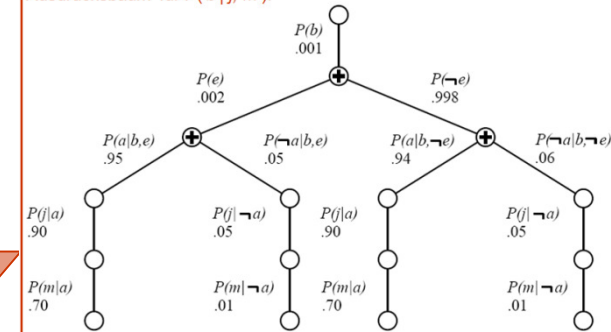
for each value x_i of X **do**

 extend \mathbf{e} with value x_i for X

$Q(x_i) \leftarrow \text{ENUMERATE-ALL}(\text{VARS}[bn], \mathbf{e})$

return $\text{NORMALIZE}(Q(X))$

Ausdrucksbaum für $P(b|j, m)$:



function **ENUMERATE-ALL**($vars, \mathbf{e}$) **returns** a real number

if $\text{EMPTY?}(vars)$ **then return** 1.0

$Y \leftarrow \text{FIRST}(vars)$

if Y has value y in \mathbf{e}

then return $P(y | Pa(Y)) \times \text{ENUMERATE-ALL}(\text{REST}(vars), \mathbf{e})$

else return $\sum_y P(y | Pa(Y)) \times \text{ENUMERATE-ALL}(\text{REST}(vars), \mathbf{e}_y)$

 where \mathbf{e}_y is \mathbf{e} extended with $Y = y$

$Pa(Y)$ für Eltern(Y)

then für Evidenzvariable & Anfragevariable

else für unbeob. Variable

Redundante Berechnungen und Variablenelimination

Erhebliche Beschleunigung durch Vermeidung **wiederholter oder unnötiger Berechnungen** möglich:

- 1) **Redundante** Teilberechnungen sind nur einmal auszuführen und die entspr. Zwischenergebnisse zu speichern. Dazu sind die Ausdrücke von rechts nach links (bzw. im Baum von unten nach oben) auszuwerten.
- 2) **Irrelevant** heißt eine Variable V bzgl. einer Anfragevariablen X , wenn gilt: V ist weder Vorfahre der Anfragevariablen X noch Vorfahre einer Evidenzvariablen.

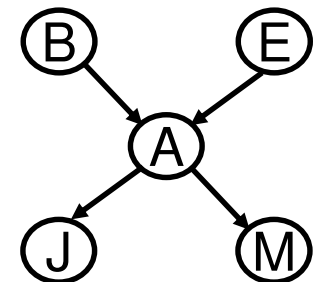
Bspl.: $P(\text{JohnCalls} \mid \text{Burglary} = \text{true})$

$$P(J \mid b) = \alpha \cdot P(b) \cdot \sum_e P(e) \cdot \sum_a P(a \mid b, e) \cdot P(J \mid a) \cdot \sum_m P(m \mid a)$$

$X = \text{JohnCalls}$, $\mathbf{E} = \{\text{Burglary}\}$

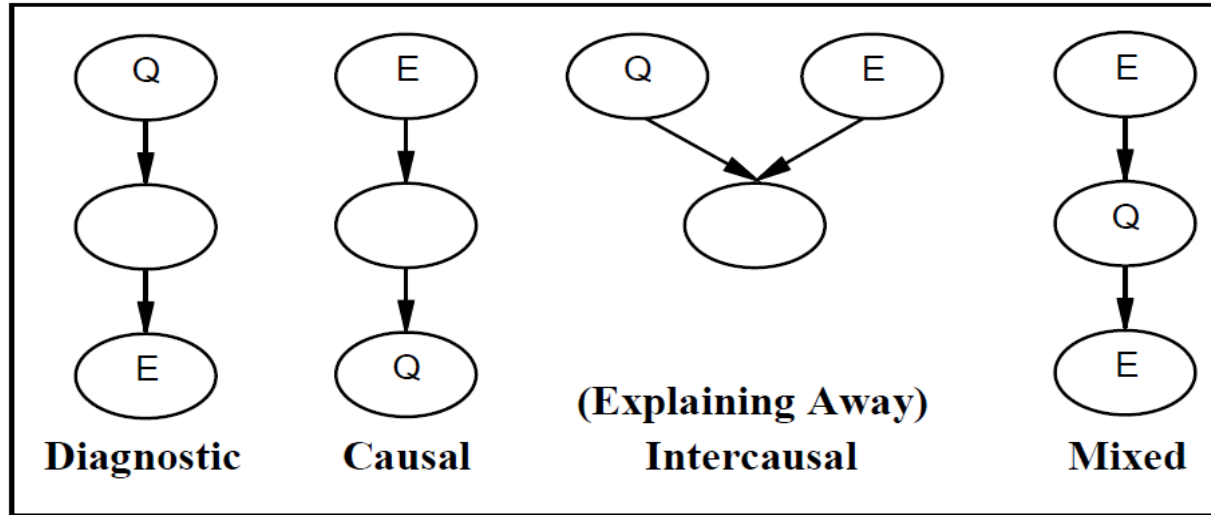
$\text{Ancestors}(\{X\} \cup \mathbf{E}) = \text{Ancestors}(\{J, B\}) = \{\text{Alarm}, \text{Earthquake}, \text{Burglary}\}$

$\rightarrow \text{MaryCalls}$ ist irrelevant



Alle Blattknoten, die für irrelevante Variablen stehen, sind eliminierbar. Nach deren Entfernung gibt es neue Blattknoten, die möglicherweise auch irrelevant sind, ...

Typen von Inferenzen



(1) **Diagnostisch:** Von Effekten zu Ursachen

Bspl.: $P(\text{burglary} \mid \text{johnCalls}) = 0.016$

(2) **Kausal:** Von Ursachen zu Effekten

$P(\text{johnCalls} \mid \text{burglary}) = 0.86$

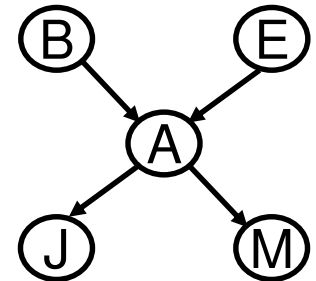
(3) **Interkausal:** Zwischen Ursachen eines gemeinsamen Effektes

$P(\text{burglary} \mid \text{alarm}) = 0.376$, aber

$P(\text{burglary} \mid \text{alarm}, \text{earthquake}) = 0.003$ („Explaining away“)

(4) **Gemischt:** Kombination von (1) – (3)

$P(\text{alarm} \mid \text{johnCalls}, \neg \text{earthquake}) = 0.03$



Komplexität der exakten Inferenz

Einfach verbundene Netzwerke (sog. „*Polytrees*“):

- Def.[Polytree]: Unter Vernachlässigung der Kantenrichtung sind jeweils zwei beliebige Knoten des Netzes durch *maximal* einen Pfad verbunden.
- Zeit- und Platzkomplexität *linear* in der Größe des Netzes:
 - mit Größe = Zahl der CPT-Einträge: $O(d^k \cdot n)$ mit n Knoten mit max. d Variablenwerten und max. k Elternknoten
 - Zeit- und Platzkomplexität linear in der Knotenzahl des Netzes
 - bei max. k Elternknoten

Mehrfach verbundene Netze:

- mindestens NP-hart (man kann z.B. 3-SAT leicht nachbilden)
- daher Betrachtung approximativer Inferenz

Approximative Inferenz

Idee: Betrachte *BN als Beschreibung eines Zufallsprozesses* und *führe eine stochast. Simulation des Prozesses durch*, um gewünschte *W'keiten* zu schätzen

Stochastische Simulation:

1. Ziehe N Beispiele von einer Stichprobenverteilung S
2. Berechne eine *approximative* Wahr'keit \hat{P} auf der Grundlage der Verteilung der relativen Häufigkeiten
3. Zeige, dass diese gegen die wahre gesuchte Wahr'keit P konvergiert

Verfahren, die Wahr'keiten durch eine stochastische Simulation, nämlich das Ziehen von zufällig angeordneten *Stichproben* (*Beispielen*, *Samples*), abzuschätzen, werden auch als *Monte-Carlo-Methoden* bezeichnet

Stichprobe von einem leeren Netzwerk (1)

Einfachste Sampling-Methode: Erzeugung von Ereignissen,
denen keine Evidenz zugeordnet ist:

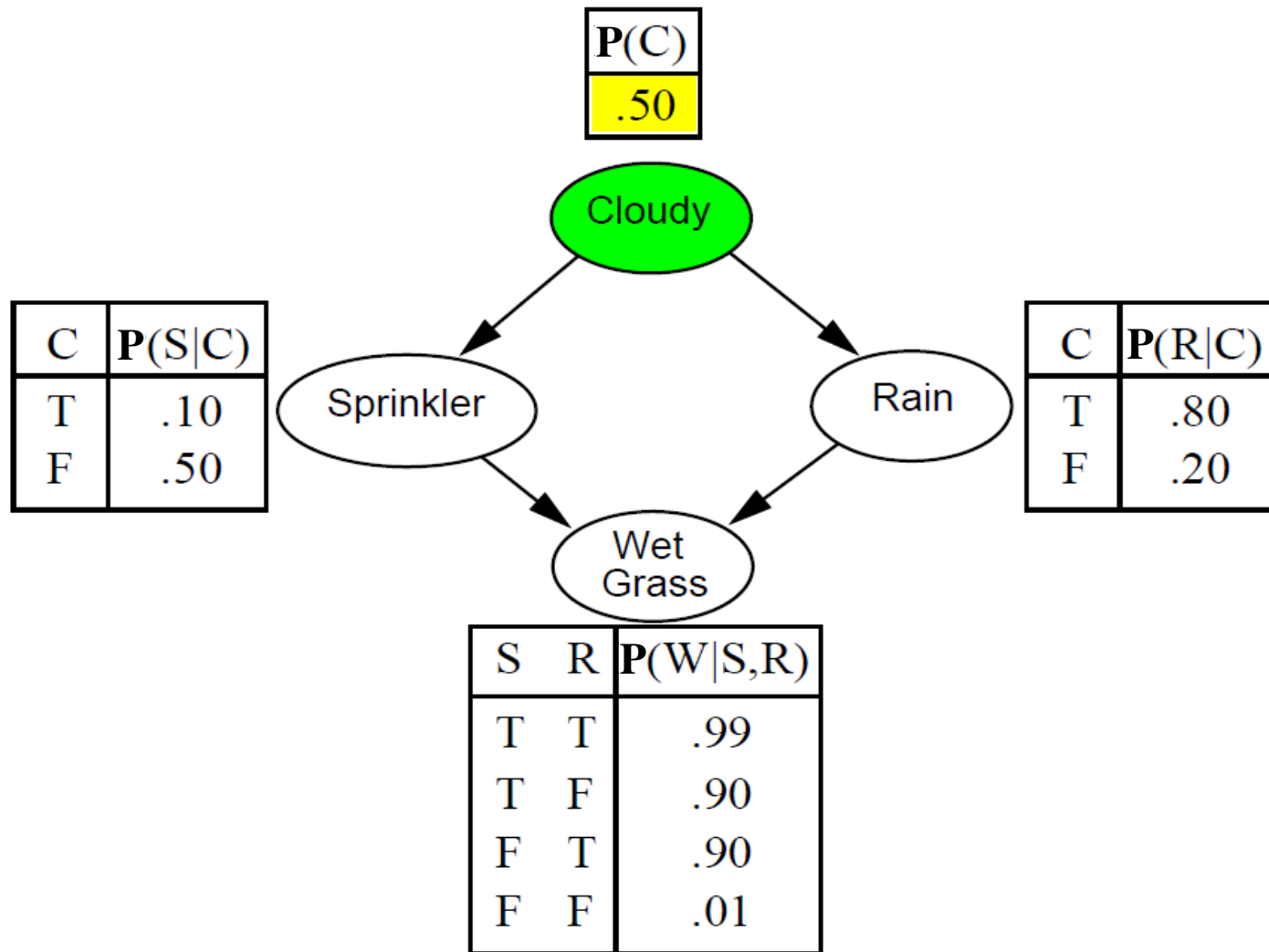
~ Sampling aller Variablen in ihrer topologischen Reihenfolge

Keine
Evidenz-
variablen

```
function PRIOR-SAMPLE( $bn$ ) returns an event sampled from  $bn$   
inputs:  $bn$ , a belief network specifying joint distribution  $\mathbf{P}(X_1, \dots, X_n)$   
 $\mathbf{x} \leftarrow$  an event with  $n$  elements  
for  $i = 1$  to  $n$  do  
     $x_i \leftarrow$  a random sample from  $\mathbf{P}(X_i \mid \text{parents}(X_i))$   
    given the values of  $\text{parents}(X_i)$  in  $\mathbf{x}$   
return  $\mathbf{x}$ 
```

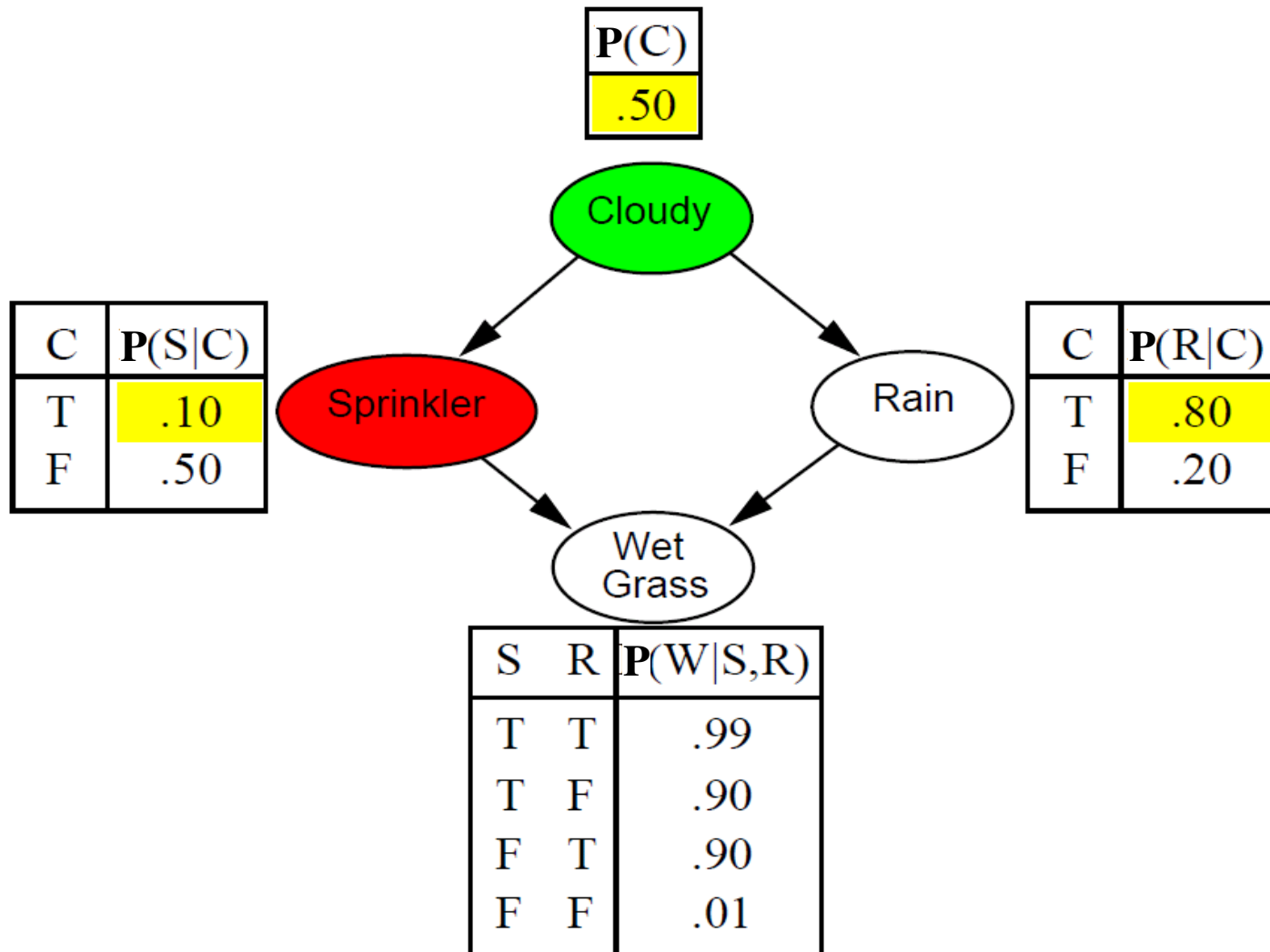
Beispiel zu Prior-Sample (1)

1) Stichprobe aus $\mathbf{P}(\text{Cloudy}) = \langle 0.5, 0.5 \rangle$ liefere *true* :



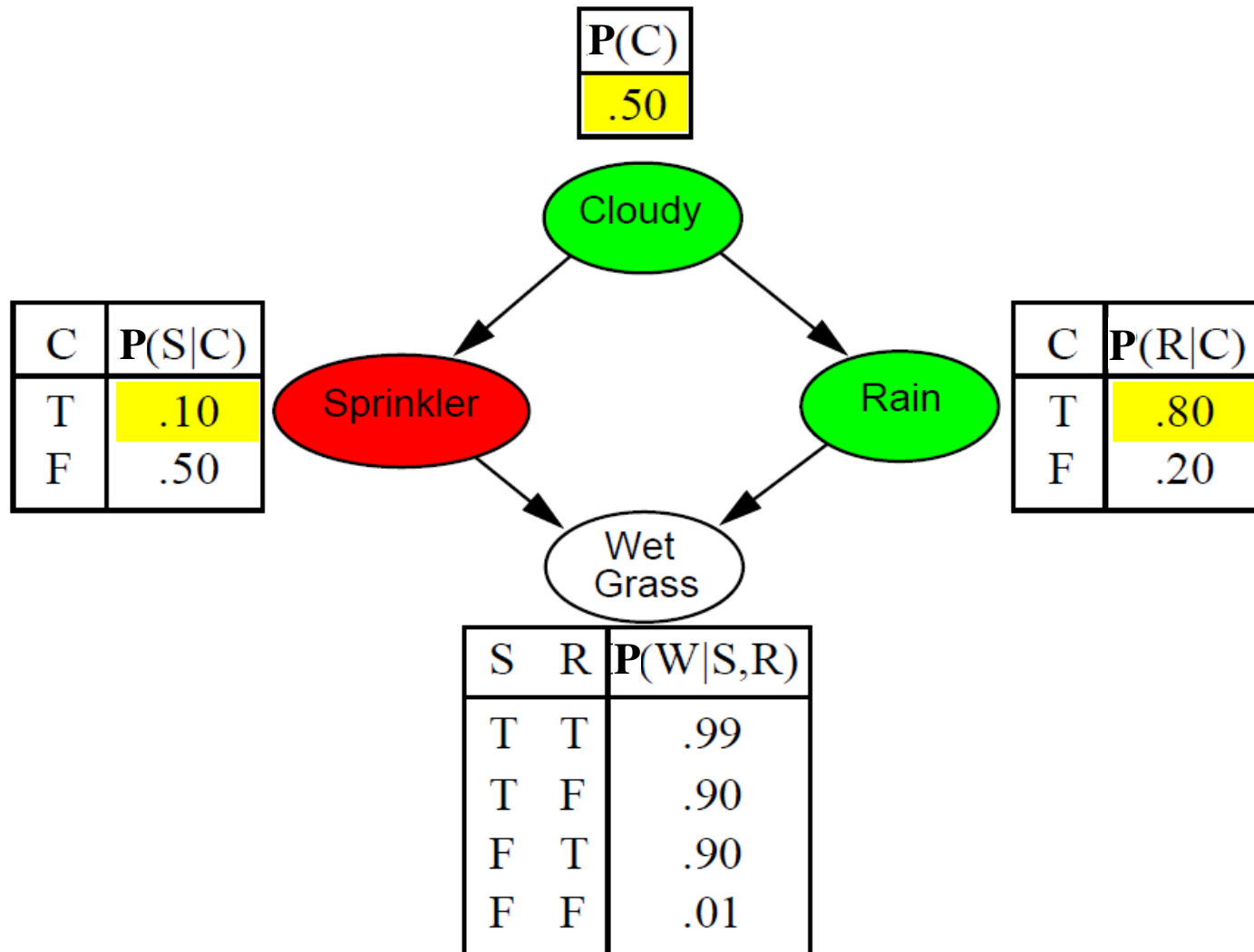
Beispiel zu Prior-Sample (2)

2) Stichprobe aus $\mathbf{P}(\text{Sprinkler}|\text{cloudy}) = \langle 0.1, 0.9 \rangle$ liefere *false* :



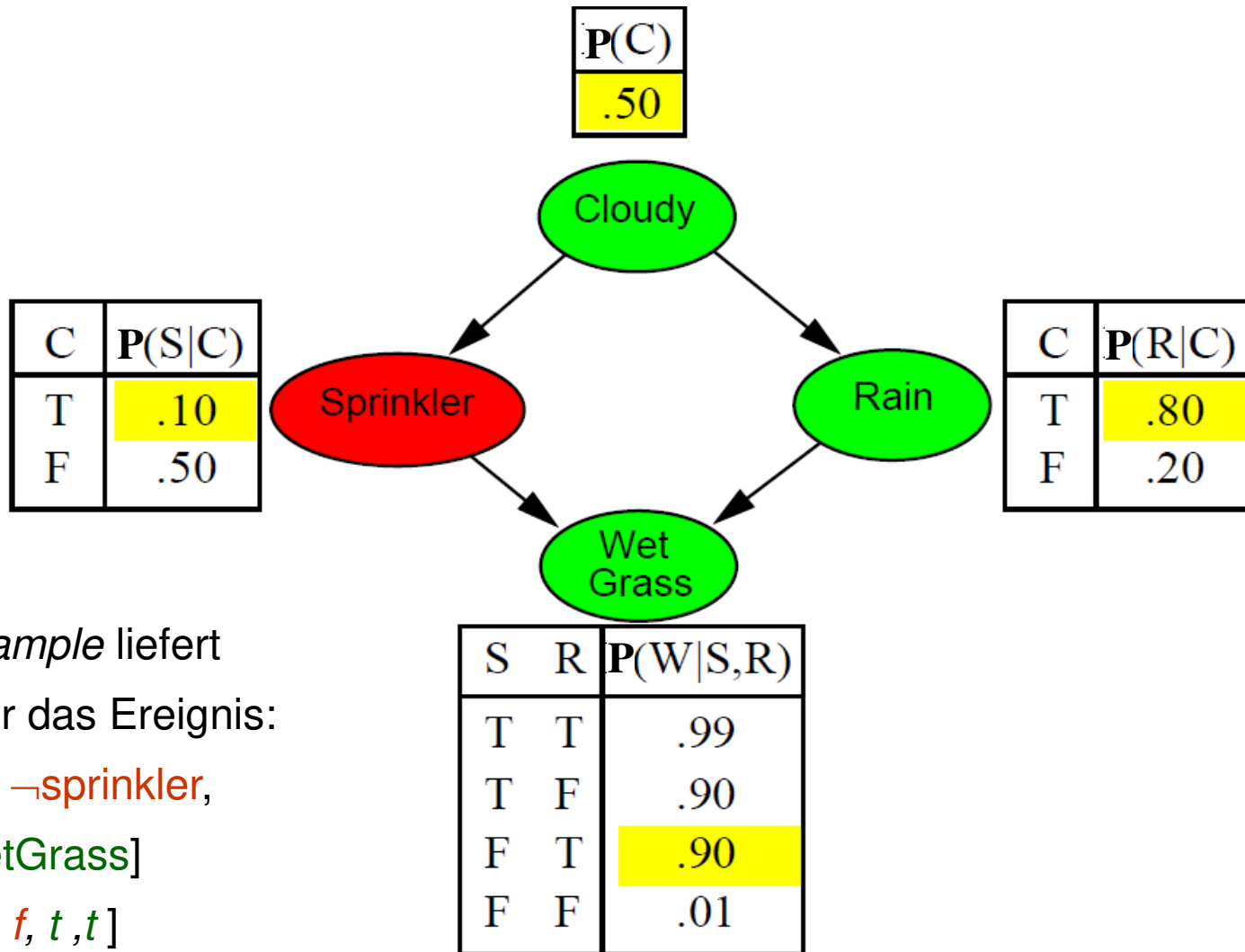
Beispiel zu Prior-Sample (3)

3) Stichprobe aus $\mathbf{P}(\text{Rain}|\text{cloudy}) = \langle 0.8, 0.2 \rangle$ liefere *true* :



Beispiel zu Prior-Sample (4)

4) Stichprobe aus $P(\text{WetGrass} | \neg\text{sprinkler}, \text{rain}) = \langle 0.9, 0.1 \rangle$ liefere *true* :



Prior-Sample liefert
also hier das Ereignis:
[cloudy, \neg sprinkler,
rain, wetGrass]
bzw. [*t*, *f*, *t*, *t*]

Stichprobe von einem leeren Netzwerk (2)

Die W'keit, dass *Prior-Sample* ein bestimmtes atomares Ereignis erzeugt ist

Sampling nach
Bayes-Netz

$$S_{PS}(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i \mid \text{parents}(X_i)) = P(x_1, \dots, x_n),$$

Def. Verbundwk.
nach Bayes-Netz

also die **korrekte Verbundwahrscheinlichkeit** des atomaren Ereignisses.

Sei N die Anzahl aller gezogenen Stichproben durch *Prior-Sample* und sei $N_{PS}(x_1, \dots, x_n)$ die Anzahl der durch *Prior-Sample* gezogenen Samples des atomaren Ereignisses x_1, \dots, x_n .

Es lässt sich zeigen, dass die **relative Häufigkeit** $N_{PS}(x_1, \dots, x_n)/N$ nach folg. Stichprobenw'keit konvergiert:

$$\lim_{N \rightarrow \infty} \hat{P}(x_1, \dots, x_n) = \lim_{N \rightarrow \infty} N_{PS}(x_1, \dots, x_n)/N = S_{PS}(x_1, \dots, x_n) = P(x_1, \dots, x_n).$$

Eine Schätzung mit dieser Grenzwerteigenschaft heißt **konsistent**.

Schreibweise: $\hat{P}(x_1, \dots, x_n) \approx P(x_1, \dots, x_n)$

Stichprobe von einem leeren Netzwerk (3)

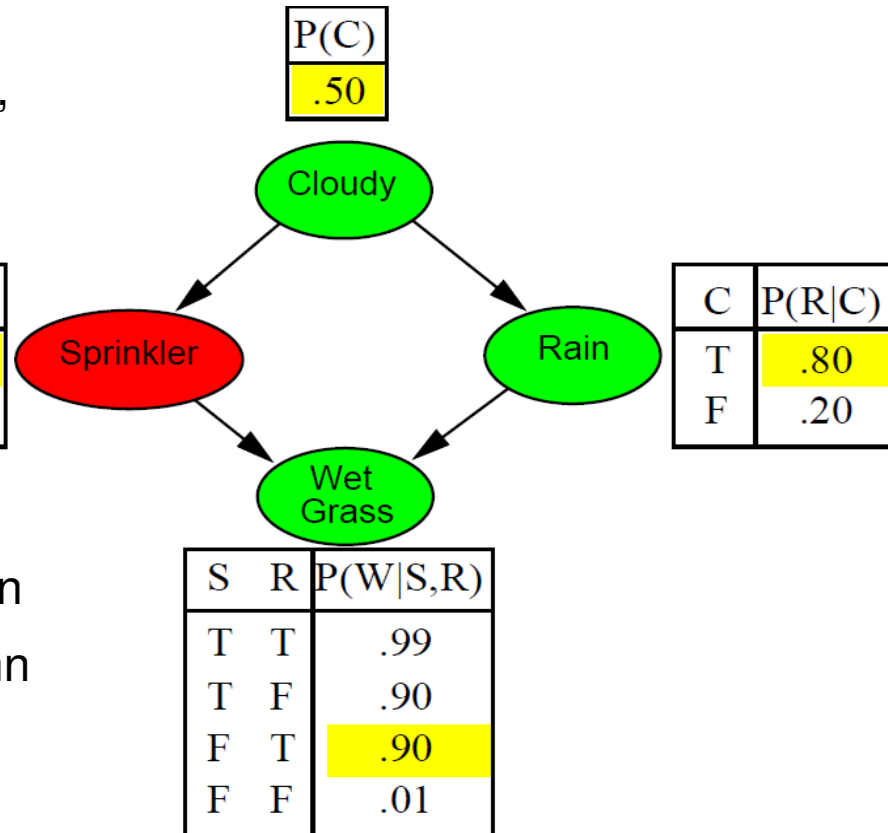
Am Beispiel: die W'keit für (cloudy,¬sprinkler,rain,wetGrass)

$$P(\text{cloudy}, \neg \text{sprinkler}, \text{rain}, \text{wetGrass}) = 0.5 \cdot 0.9 \cdot 0.8 \cdot 0.9 = 0.324$$

Im Grenzwert großer N erwarten wir, dass ca. 32.4% aller Stichproben diesem Ereignis entsprechen: $\lim_{N \rightarrow \infty} N_{PS}(\text{cloudy}, \neg \text{sprinkler}, \text{rain}, \text{wetGrass})/N \approx 0.324$

Die Wahrscheinlichkeit $P(\text{cloudy}, \neg \text{sprinkler}, \text{rain}, \text{wetGrass})$ wird also als Bruchteil aller von dem Sampling-Prozess gezogenen atom. Ereignisse geschätzt, die mit diesem Ereignis übereinstimmen

C	P(S C)
T	.10
F	.50



Wenn wir also 1000 Stichproben generieren und 511 davon ergeben $Rain = true$, dann schätzen wir $\hat{P}(Rain = true) = 0.511$.

Rejection Sampling

Berechne $\hat{P}(X | \mathbf{e})$ durch Samples, die zu \mathbf{e} passen.

$N = \# \text{ Samples}$

```
function REJECTION-SAMPLING( $X, \mathbf{e}, bn, N$ ) returns an estimate of  $P(X | \mathbf{e})$ 
  local variables:  $\mathbf{N}$ , a vector of counts over  $X$ , initially zero
  for  $j = 1$  to  $N$  do
     $\mathbf{x} \leftarrow \text{PRIOR-SAMPLE}(bn)$ 
    if  $\mathbf{x}$  is consistent with  $\mathbf{e}$  then
       $\mathbf{N}[x] \leftarrow \mathbf{N}[x] + 1$  where  $x$  is the value of  $X$  in  $\mathbf{x}$ 
  return NORMALIZE( $\mathbf{N}[X]$ )
```

$\mathbf{N} = \text{Vektor der Zählungen der } X\text{-Werte}$

Beispiel: Schätze $\mathbf{P}(\text{Rain} | \text{Sprinkler} = \text{true})$ unter Verwendung von 100 Samples

Ergebnis sei: 27 Stichproben mit $\text{Sprinkler} = \text{true}$, von diesen

- 8 mit $\text{Rain} = \text{true}$
- 19 mit $\text{Rain} = \text{false}$

$$\leadsto \hat{P}(\text{Rain} | \text{Sprinkler} = \text{true}) = \text{Normalize}(\langle 8, 19 \rangle) = \langle 0.296, 0.704 \rangle$$

Bewertung von Rejection Sampling

Rejection Sampling erzeugt also konsistente A-Posteriori-Schätzungen

Problem: Hoffnungslos teuer bei hoher Ablehnungsrate von Stichproben

Im Beispiel: 73 abgelehnte Stichproben von insgesamt 100 Stichproben

Der Bruchteil der mit der Evidenz \mathbf{e} inkonsistenten Stichproben

- ist hoch bei kleinem $P(\mathbf{e})$ und
- wächst exponentiell mit steigender Zahl der Evidenzvariablen

Besser ist die folg. *W'keitsgewichtung* bzw. Likelihood-Gewichtung, die nur Ereignisse *erzeugt*, die konsistent zur Evidenz \mathbf{e} sind

Der folgende Algorithmus *Likelihood Weighting* wird illustriert mit der Abfrage $P(\text{Rain} \mid \text{Sprinkler} = \text{true}, \text{WetGrass} = \text{true})$

Likelihood: eine nicht normalisierte bedingte Wahrscheinlichkeit

Likelihood-Gewichtung (1)

Idee:

- 1) Fixiere Werte der Evidenzvariablen **E** und sample nur die Nichtevidenzvariablen $\mathbf{Z} = \{X\} \cup \mathbf{Y}$ mit Anfragevariable **X** und unbeobachteten Variablen **Y**
 - 2) Gewichte jedes Sample mit der Likelihood, die ihm nach der Evidenz zukommt
- Damit ist jedes erzeugte Ereignis mit der Beobachtung **E** konsistent
 - Nicht alle Ereignisse sind gleich wichtig
 - Ereignisse, deren Evidenz wahrscheinlich erscheint (gemessen über dem Produkt der bedingt W'keiten jeder Evidenzvariablen bei bekannten Eltern), gehen mit höherem Gewicht ein
 - Ereignisse, deren tatsächliche Evidenz unwahrscheinlich erscheint, erhalten weniger Gewicht

Likelihood-Gewichtung (2)

$N = \# \text{ Samples}$

```
function LIKELIHOOD-WEIGHTING( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
  local variables:  $W$ , a vector of weighted counts over  $X$ , initially zero
  for  $j = 1$  to  $N$  do
    Ereignis  $x$  mit Gewicht  $w$   $\rightarrow x, w \leftarrow \text{WEIGHTED-SAMPLE}(bn, e)$ 
     $W[x] \leftarrow W[x] + w$  where  $x$  is the value of  $X$  in  $x$ 
  return NORMALIZE( $W[X]$ )
```

W = Vektor der gewichteten Zählungen der X -Werte

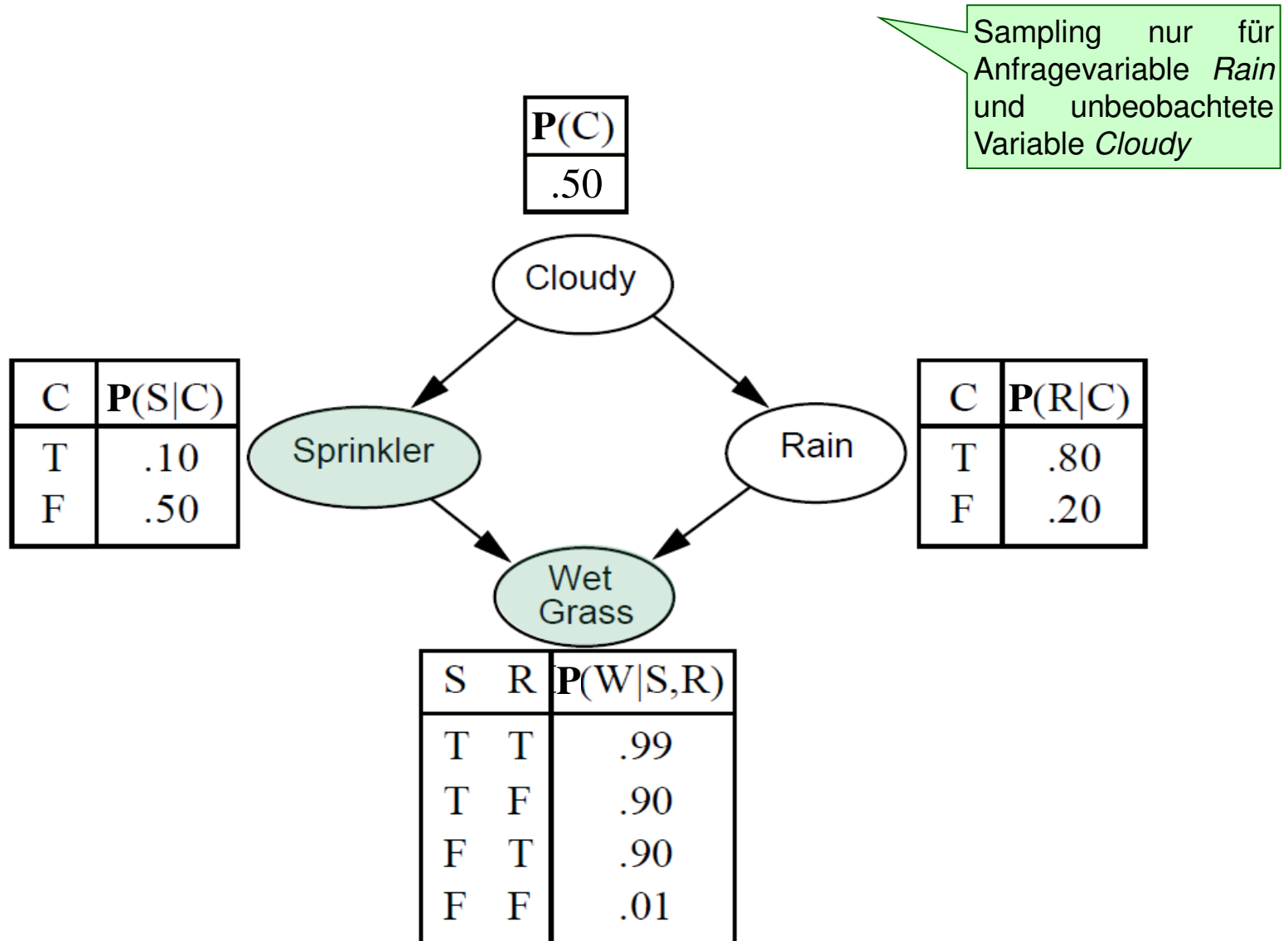
```
function WEIGHTED-SAMPLE( $bn, e$ ) returns an event and a weight
   $w \leftarrow 1$ ;  $x \leftarrow$  an event with  $n$  elements initialized from  $e$ 
  for  $i = 1$  to  $n$  do
    if  $X_i$  has a value  $x_i$  in  $e$ 
      then  $w \leftarrow w \times P(X_i = x_i \mid \text{parents}(X_i))$ 
      else  $x_i \leftarrow$  a random sample from  $P(X_i \mid \text{parents}(X_i))$ 
  return  $x, w$ 
```

Bewertung der Kompatibilität von Evidenz mit Ereignis

Sampling von Anfrage- und unbeobachteten Variablen

Beispiel: Likelihood-Gewichtung

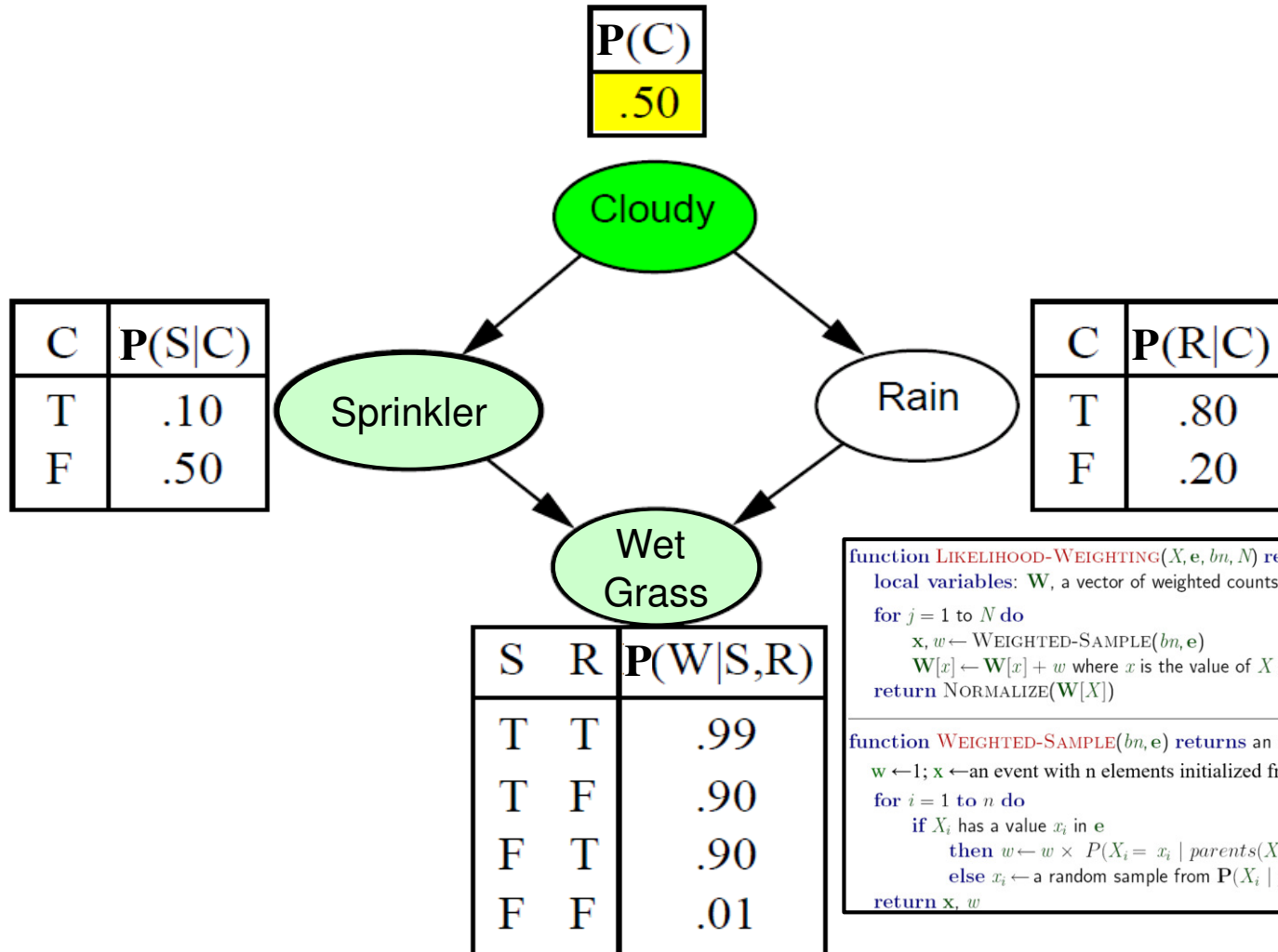
Beispiel: Anfrage $P(Rain \mid \textit{Sprinkler} = \textit{true}, \textit{WetGrass} = \textit{true})$



Beispiel: Likelihood-Gewichtung (Forts. 1)

1) Stichprobe aus $\mathbf{P}(\text{Cloudy}) = \langle 0.5, 0.5 \rangle$ liefere *true* :

$w \leftarrow 1.0$



```
function LIKELIHOOD-WEIGHTING( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
  local variables:  $W$ , a vector of weighted counts over  $X$ , initially zero
  for  $j = 1$  to  $N$  do
     $x, w \leftarrow \text{WEIGHTED-SAMPLE}(bn, e)$ 
     $W[x] \leftarrow W[x] + w$  where  $x$  is the value of  $X$  in  $x$ 
  return NORMALIZE( $W[X]$ )
```

```
function WEIGHTED-SAMPLE( $bn, e$ ) returns an event and a weight
   $w \leftarrow 1$ ;  $x \leftarrow$  an event with  $n$  elements initialized from  $e$ 
  for  $i = 1$  to  $n$  do
    if  $X_i$  has a value  $x_i$  in  $e$ 
      then  $w \leftarrow w \times P(X_i = x_i \mid \text{parents}(X_i))$ 
    else  $x_i \leftarrow$  a random sample from  $P(X_i \mid \text{parents}(X_i))$ 
  return  $x, w$ 
```

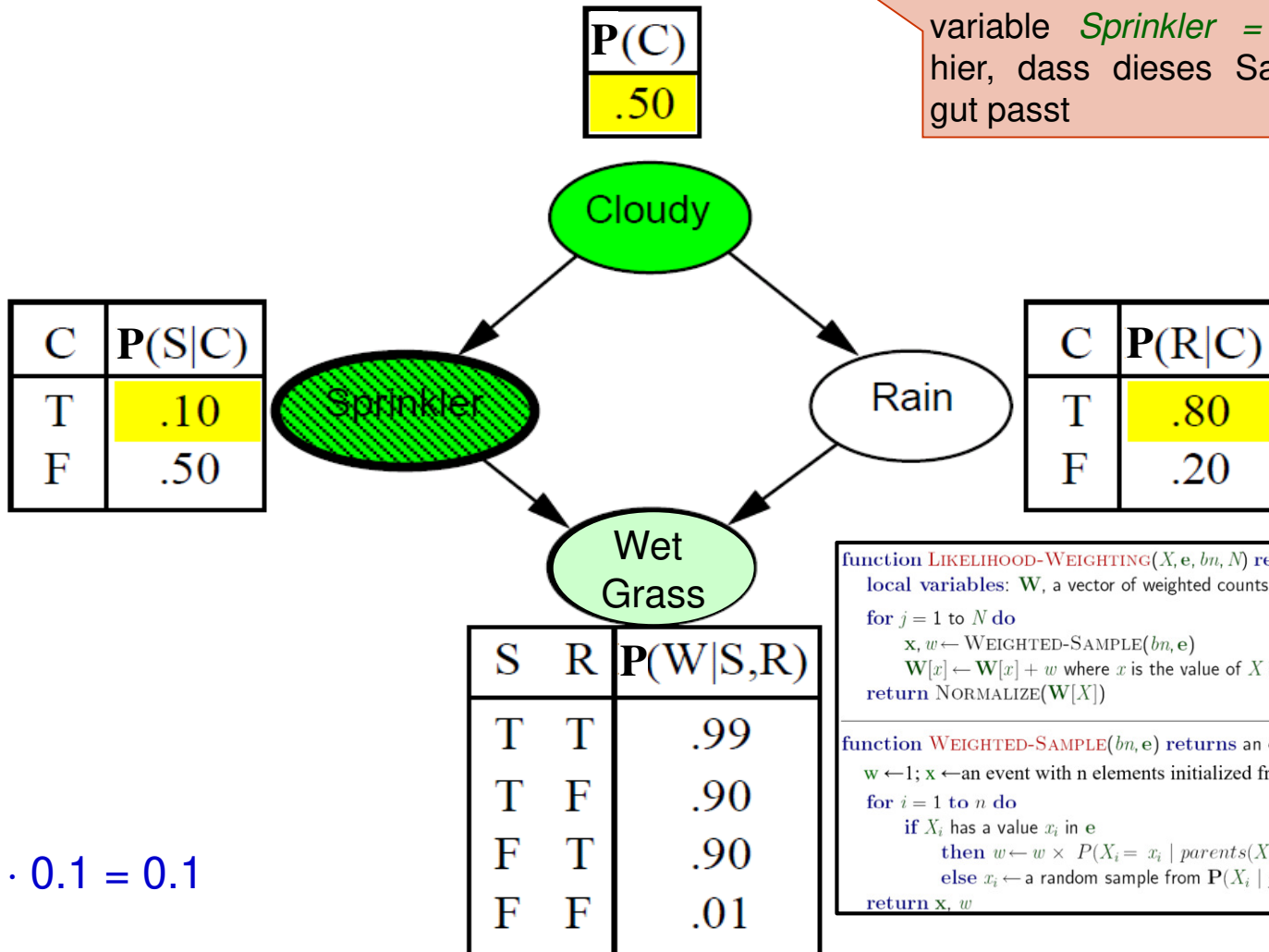
$w = 1.0$

Beispiel: Likelihood-Gewichtung (Forts. 2)

2) *Sprinkler* ist Evidenzvariable mit Wert *true* (was bei *cloudy* unwahrscheinlich ist):

$$w \leftarrow w \cdot P(\text{sprinkler} \mid \text{cloudy}) = 1.0 \cdot 0.1 = 0.1$$

Gewichtsreduktion: Die Evidenzvariable *Sprinkler* = *true* zeigt hier, dass dieses Sample nicht gut passt



$$w = 1.0 \cdot 0.1 = 0.1$$

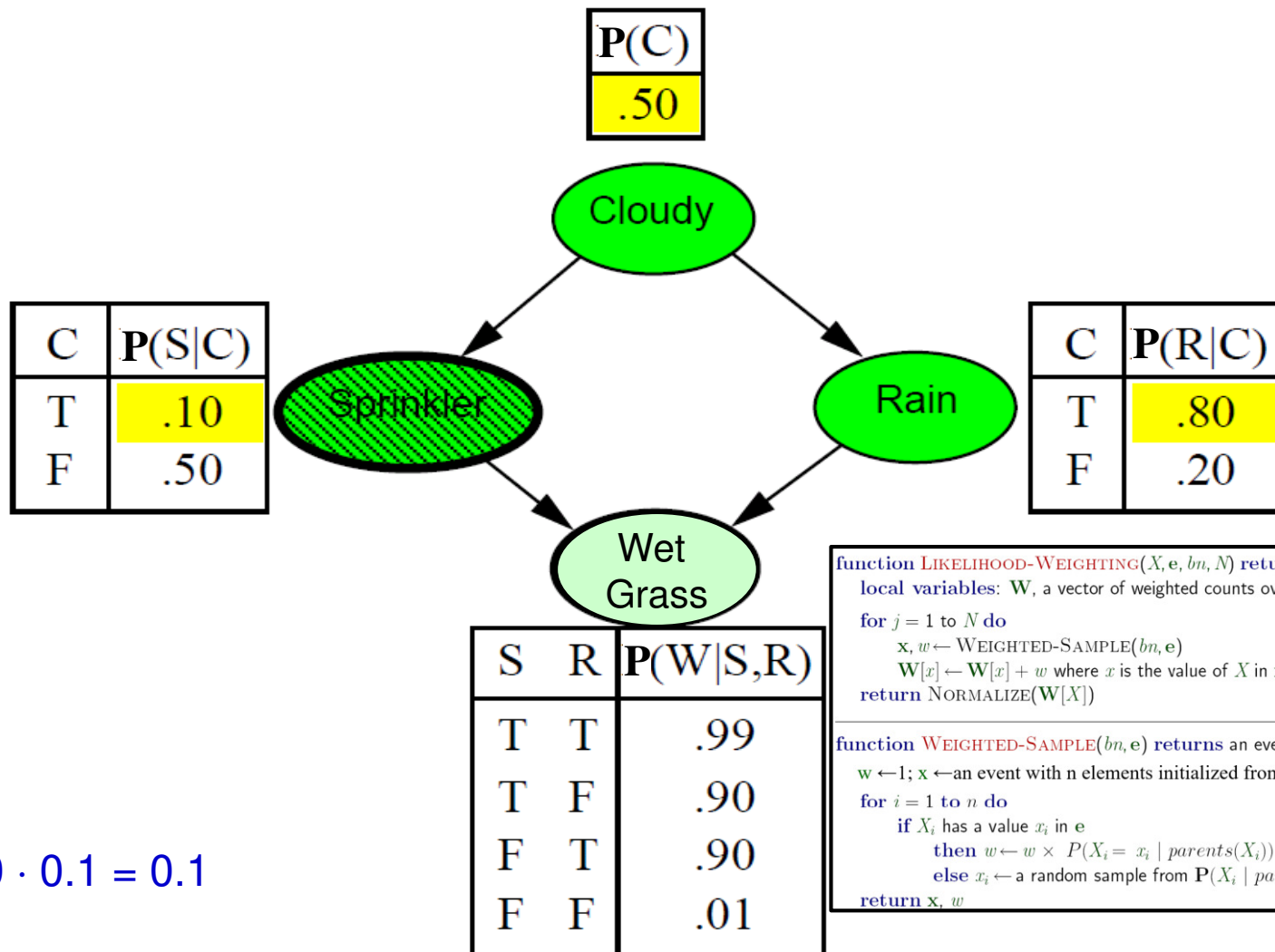
```
function LIKELIHOOD-WEIGHTING( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
  local variables:  $W$ , a vector of weighted counts over  $X$ , initially zero
  for  $j = 1$  to  $N$  do
     $x, w \leftarrow \text{WEIGHTED-SAMPLE}(bn, e)$ 
     $W[x] \leftarrow W[x] + w$  where  $x$  is the value of  $X$  in  $x$ 
  return NORMALIZE( $W[X]$ )

function WEIGHTED-SAMPLE( $bn, e$ ) returns an event and a weight
   $w \leftarrow 1$ ;  $x \leftarrow$  an event with  $n$  elements initialized from  $e$ 
  for  $i = 1$  to  $n$  do
    if  $X_i$  has a value  $x_i$  in  $e$ 
      then  $w \leftarrow w \times P(X_i = x_i \mid \text{parents}(X_i))$ 
    else  $x_i \leftarrow$  a random sample from  $P(X_i \mid \text{parents}(X_i))$ 
  return  $x, w$ 
```

Beispiel: Likelihood-Gewichtung (Forts. 3)

3) Stichprobe aus $\mathbf{P}(\text{Rain}) = \langle 0.8, 0.2 \rangle$ liefere *true* :

$w = 0.1$



$$w = 1.0 \cdot 0.1 = 0.1$$

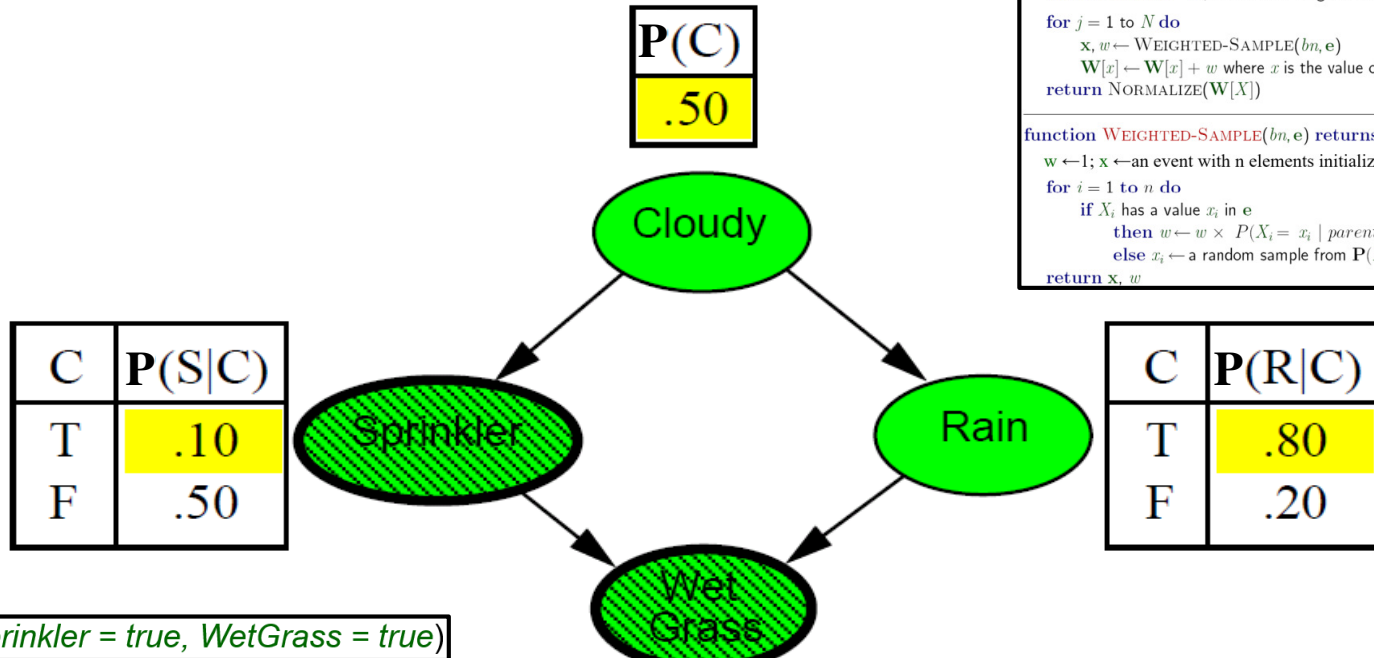
```
function LIKELIHOOD-WEIGHTING( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
  local variables:  $W$ , a vector of weighted counts over  $X$ , initially zero
  for  $j = 1$  to  $N$  do
     $x, w \leftarrow \text{WEIGHTED-SAMPLE}(bn, e)$ 
     $W[x] \leftarrow W[x] + w$  where  $x$  is the value of  $X$  in  $x$ 
  return NORMALIZE( $W[X]$ )

function WEIGHTED-SAMPLE( $bn, e$ ) returns an event and a weight
   $w \leftarrow 1$ ;  $x \leftarrow$  an event with  $n$  elements initialized from  $e$ 
  for  $i = 1$  to  $n$  do
    if  $X_i$  has a value  $x_i$  in  $e$ 
      then  $w \leftarrow w \times P(X_i = x_i \mid \text{parents}(X_i))$ 
    else  $x_i \leftarrow$  a random sample from  $P(X_i \mid \text{parents}(X_i))$ 
  return  $x, w$ 
```


Beispiel: Likelihood-Gewichtung (Forts. 4)

4) *WetGrass* ist Evidenzvariable mit Wert *true* :

$$w \leftarrow w \cdot P(\text{wetgrass} \mid \text{sprinkler}, \text{rain}) = 0.1 \cdot 0.99 = 0.099$$



```
function LIKELIHOOD-WEIGHTING( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
  local variables:  $W$ , a vector of weighted counts over  $X$ , initially zero
  for  $j = 1$  to  $N$  do
     $x, w \leftarrow \text{WEIGHTED-SAMPLE}(bn, e)$ 
     $W[x] \leftarrow W[x] + w$  where  $x$  is the value of  $X$  in  $x$ 
  return NORMALIZE( $W[X]$ )

function WEIGHTED-SAMPLE( $bn, e$ ) returns an event and a weight
   $w \leftarrow 1$ ;  $x \leftarrow$  an event with  $n$  elements initialized from  $e$ 
  for  $i = 1$  to  $n$  do
    if  $X_i$  has a value  $x_i$  in  $e$ 
      then  $w \leftarrow w \times P(X_i = x_i \mid \text{parents}(X_i))$ 
      else  $x_i \leftarrow$  a random sample from  $P(X_i \mid \text{parents}(X_i))$ 
  return  $x, w$ 
```

Weighted-Sample liefert Ereignis (cloudy, sprinkler, rain, wetgrass), was als *Rain = true* gezählt wird, mit Gewicht $w = 0.099$.

Das Gewicht ist gering, da das Ereignis einen wolkigen Tag zeigt, bei dem der Sprinkler wahrscheinlich nicht läuft.

Betrachtung der Likelihood-Gewichtung (1)

Seien Evidenzvariable \mathbf{E} mit Werten \mathbf{e} feststehend, Anfragevariable sei X und \mathbf{Z} seien die unbeobachteten Variablen \mathbf{Y} und Anfrage X , also: $\mathbf{Z} = \{X\} \cup \mathbf{Y}$

Teil 1: *Weighted-Sample* sampelt jede Variable aus \mathbf{Z} für bekannte Elternwerte:

$$S_{ws}(\mathbf{z}, \mathbf{e}) = \prod_i P(z_i \mid \text{parents}(Z_i)).$$

$\text{parents}(Z_i)$ kann sowohl verborgene Variable als auch Evidenzvariable enthalten. Anders als die unbedingte Verteilung $\mathbf{P}(\mathbf{z})$ berücksichtigt die Verteilung S_{ws} die Evidenz der Vorfahren von Z_i

Teil 2: Die Gewichtung für ein gegebenes Sample (\mathbf{z}, \mathbf{e}) ist das Produkt der Wahrheiten für jede Evidenzvariable mit bekannten Eltern:

$$w(\mathbf{z}, \mathbf{e}) = \prod_j P(e_j \mid \text{parents}(E_j)).$$

Betrachtung der Likelihood-Gewichtung (2)

Die gewichtete Sampling-Likelihood ist dann das Produkt aus (1) und (2)

$$S_{WS}(\mathbf{z}, \mathbf{e}) \cdot w(\mathbf{z}, \mathbf{e}) = \prod_i P(z_i \mid \text{parents}(Z_i)) \cdot \prod_i P(e_i \mid \text{parents}(E_i)).$$

Das Produkt deckt also alle Variablen des Netzes ab.

Es lässt sich zeigen, dass $\hat{P}(\mathbf{x}|\mathbf{e})$ für große N über die Sampling-Verteilung $\Sigma_{\mathbf{y}}$ $S_{WS}(\mathbf{x}, \mathbf{y}, \mathbf{e}) \cdot w(\mathbf{x}, \mathbf{y}, \mathbf{e})$ gegen $P(\mathbf{x}|\mathbf{e})$ konvergiert: $\hat{P}(\mathbf{x}|\mathbf{e}) \approx P(\mathbf{x}|\mathbf{e})$.

Die Likelihood-Gewichtung erzeugt also konsistente Schätzungen. Sie kann deutlich effizienter sein als Rejection Sampling.

Zusammenfassung Bayessche Netze

- Bayessche Netze erlauben eine **kompakte Repräsentation** der Verbund-w'keit.
- Dies wird erreicht durch **Unabhängigkeitsannahmen**.
- Sie unterstützen **verschiedene Formen des Schließens** bei gegebenen Evidenzen: kausal, diagnostisch, interkausal, gemischt.
- Inferenz bedeutet dabei die **Berechnung der Wahr'keitsverteilung für die Belegungen einer Menge von Variablen** bei gegebenen Evidenzen.
- Die **Komplexität der Inferenz** in Bayesschen Netzen hängt von der **Struktur des Netzwerkes** ab.
- Für **Polybäume** ist die Komplexität **polynomiell** in der Größe des Netzwerks.
- **Im Allgemeinen** ist die **Inferenz** in Bayesschen Netzen **NP-hart**.
- Es gibt **Approximationstechniken**, um diesem Problem zu begegnen.