

# Spezielle Themen der Künstlichen Intelligenz

10. Termin:

Bayesian Networks Inference Algorithms

## Modeling with Bayesian networks

Using Bayesian networks for real-world problems requires two steps:

- ▶ constructing an appropriate Bayesian network
- ▶ solve the problems by applying one of the possible queries

How to construct a Bayesian network?

1. **define network variables and their values**
  - distinguish between *query*, *evidence*, and *intermediary* variables
2. **define network structure**
  - for each var  $X$ , answer the question: what set of var's are direct causes of  $X$ ?
3. **define network parameters (CPTs)**
  - difficulty and objectivity depend on problem

## Reasoning with Bayesian networks

Four general types of queries one can pose:

- ▶ **probability of evidence**: how likely is a variable instantiation  $\mathbf{e} \rightarrow Pr(\mathbf{e})=?$
- ▶ **prior and posterior marginals**: how probable is an instantiation of a limited set of variables  $\rightarrow Pr(x_1, \dots, x_n)=?$  or  $Pr(x_1, \dots, x_n | \mathbf{e})=?$
- ▶ **most probable explanation (MPE)**: what is the most probable instantiation of all network var's given some evidence  $\mathbf{e} \rightarrow \mathbf{x}$  with  $Pr(x_1, \dots, x_n | \mathbf{e}) = \max$ ?
- ▶ **maximum a posteriori hypothesis (MAP)**: what is the most probable instantiation of a subset of var's given some evidence  $\mathbf{e} \rightarrow \mathbf{x}$  with  $Pr(x_1, \dots, x_m | \mathbf{e}) = \max$ ?

All these queries can be computed from a Bayesian network.

But, how?

## Inference algorithms

- ▶ variable elimination (see last lecture)
- ▶ jointtree algorithm
- ▶ recursive conditioning
- ▶ belief propagation
- ▶ Monte Carlo Markov Chain

# Variable elimination

## Variable elimination:

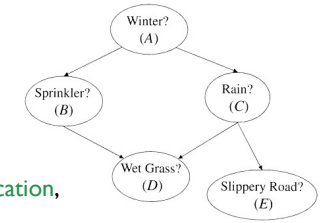
- ▶ given a distribution  $Pr(A,B,C,D,E)$ , variable  $A$  with values  $a_i$  can be summed out by

$$Pr(B, C, D, E) = \sum_{a_i} Pr(a_i, B, C, D, E)$$

**Definition:** factor  $f$  over var's  $\mathbf{X}$  is a function that maps each instantiation  $\mathbf{x}$  of  $\mathbf{X}$  to a number  $f(\mathbf{x}) \geq 0$

- ▶ can represent any marginal or conditional distribution
- ▶ **Summing out a variable from a factor:**  $(\sum_X f)(y) := \sum_x f(x, y)$ 
  - marginalizing  $X$ , projecting on  $Y$
- ▶ **Multiplying factors:**  $(f_1 f_2)(z) := f_1(x) f_2(y)$  with  $\mathbf{x} \sim \mathbf{z}, \mathbf{y} \sim \mathbf{z}$

# Variable elimination



## Use for computing prior marginals:

- ▶ express joint distribution as factor multiplication, viewing CPTs as factors, e.g.  $Pr(a,b,c,d,e)$ :

$$Pr(a, b, c, d, e) = \Theta_{E|C} \Theta_{D|BC} \Theta_{C|A} \Theta_{B|A} \Theta_A$$

- ▶ compute marginal distribution by summing out variables from these factors, e.g.  $Pr(D,E)$ :

$$Pr(D, E) = \sum_{A,B,C} \Theta_{E|C} \Theta_{D|BC} \Theta_{C|A} \Theta_{B|A} \Theta_A$$

Can (and should) be simplified according to:

$$\sum_X f_1 f_2 = f_1 \sum_X f_2 \text{ if } X \text{ appears only in } f_2$$

**Example:** compute prior marginal  $Pr(C)$  by eliminating first  $A$ , then  $B$

$$Pr(C) = \sum_B \Theta_{C|B} \sum_A \Theta_A \Theta_{B|A}$$

Order of elimination irrelevant for result, but not for computational costs !!

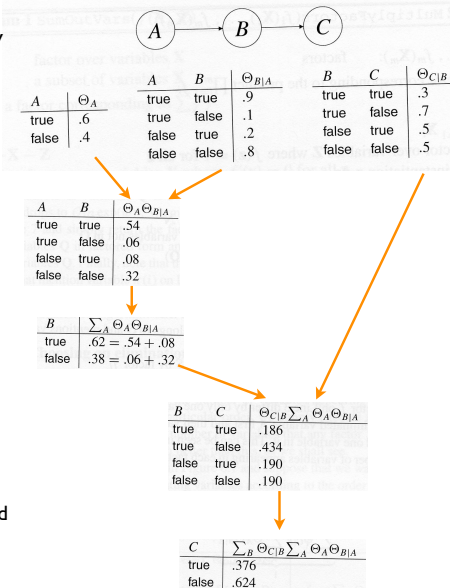
Other possibility: first  $B$  then  $A$

$$Pr(C) = \sum_A \Theta_A \sum_B \Theta_{B|A} \Theta_{C|B}$$

largest factor with 2 variables

Best order: smallest width = number of var's in the largest factor constructed

- ▶ can be computed offline, but NP-hard
- ▶ heuristics



# Variable elimination

## Use for computing posterior marginals:

- ▶ need to compute the factor  $Pr(\mathbf{Q}|\mathbf{e})$

Better to compute joint marginal  $Pr(\mathbf{Q}, \mathbf{e})$  and normalize to get  $Pr(\mathbf{Q}|\mathbf{e})$

- ▶ also gives  $Pr(\mathbf{e})$  for free as  $Pr(\mathbf{e}) = \sum_q Pr(\mathbf{q}, \mathbf{e})$

**Example:**

D	E	$Pr(\mathbf{Q} \mathbf{e})$
true	true	.448
true	false	.192
false	true	.112
false	false	.248

D	E	$Pr(\mathbf{Q}, \mathbf{e})$
true	true	.21504
true	false	.09216
false	true	.05376
false	false	.11904

$.11904 / .48 \rightarrow \sum = .48 \ Pr(\mathbf{e})$

# Variable elimination

Use elimination for computing joint marginals:

- ▶ zero out all rows that are not compatible with evidence  $\mathbf{e}$

Definition: reduction of factor  $f(\mathbf{X})$  given evidence  $\mathbf{e}$  is another factor over  $\mathbf{X}$  denoted by  $f^e$ , defined by

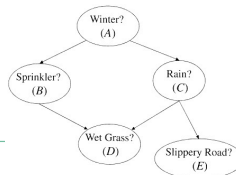
$$f^e(\mathbf{x}) := \begin{cases} f(\mathbf{x}) & \text{if } \mathbf{x} \sim \mathbf{e} \\ 0 & \text{otherwise} \end{cases}$$

- ▶ distributivity with factor multiplication:  $(f_1 f_2)^e = f_1^e f_2^e$

Joint marginal  $Pr(\mathbf{Q}, \mathbf{e})$  can hence be computed as follows:

- ▶ Example:  $\mathbf{Q}=\{D,E\}$

$$Pr(\mathbf{Q}, \mathbf{e}) = \sum_{A,B,C} \Theta_{E|C}^e \Theta_{D|BC}^e \Theta_{C|A}^e \Theta_{B|A}^e \Theta_A^e$$



# Variable elimination

Summary

- ▶ combines factor multiplication, reduction & projecting
- ▶ useful for computing prior and posterior marginals and probability of evidence (joint marginals)
- ▶ structure-based algorithm, performance depends on network structure (e.g. no. parents per node, loops, paths between nodes)

Usually combined with pruning for a given query  $Pr(\mathbf{Q}, \mathbf{e})$

- ▶ iteratively remove any leaf node not in  $\mathbf{Q}$  or  $\mathbf{E}$
- ▶ remove edge  $U \rightarrow X$  from any node  $U$  in  $\mathbf{E}$  and  $\Theta_{X|U} \leftarrow \sum_U \Theta_{X|U}^u$

Example: compute posterior marginal  $Pr(\mathbf{Q}=\{C\}, \mathbf{e}:A=true)$  by eliminating first  $A$ , then  $B$

$$Pr(\mathbf{Q}, \mathbf{e}) = \sum_B \sum_A \Theta_A^e \Theta_{B|A}^e \Theta_{C|B}^e = \sum_B \Theta_{C|B}^e \Theta_A^e \Theta_{B|A}^e$$

Therefore:

- ▶  $Pr(C=true, A=true) = .192$
- ▶  $Pr(C=false, A=true) = .408$
- ▶  $Pr(A=true) = .6$
- ▶  $Pr(C=true|A=true) = .192/.6 = .32$

Bayesian network: A → B → C

Initial joint distribution table:

A	B	C	$\Theta_{C B}$
true	true	true	.3
true	true	false	.7
true	false	true	.5
true	false	false	.5
false	true	true	.9
false	true	false	.1
false	false	true	.2
false	false	false	.8

Step 1: Eliminate A. Marginalize over A.

B	C	$\Theta_{C B}^e$
true	true	.3
true	false	.7
false	true	.5
false	false	.5

Step 2: Eliminate B. Marginalize over B.

C	$\Theta_{C B}^e$
true	.192
false	.408

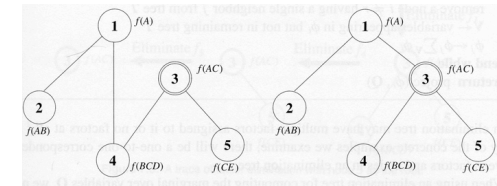
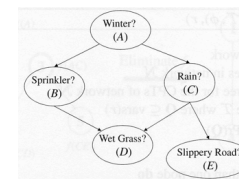
# Factor elimination

Generalization of variable elimination to factor elimination, i.e. elimination of sets of variables (Lauritzen & Spiegelhalter 1988)

- ▶ elimination order → elimination trees

Definition: An elimination tree  $(T, \Theta)$  for a set of factors  $\mathbf{S}$  is a tree  $T$ , in which each node  $\Theta_i$  is assigned exactly one factor in  $\mathbf{S}$

- ▶ factors are the CPTs in the Bayesian network
- ▶ different tree structures are possible:



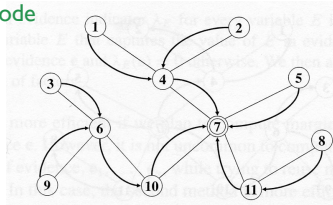
## Factor elimination

### Factor elimination in elimination trees:

- ▶ eliminate a node (factor) if all its neighbors, except the one closer to the root, have been eliminated
- ▶ when a node  $i$  is about to be eliminated, it will have a single neighbor  $j$  and  $i$ 's factor is projected and multiplied into factor of  $j$ 
  - viewed as passing a „message“ from  $i$  to  $j$

### Using factor elimination for computing marginal over $\mathcal{Q}$

- ▶ pick one node  $r$  with  $\mathcal{Q} \subseteq \text{vars}(r)$  as **root node**
- ▶ push messages toward the root
- ▶ when all messages are available in root, multiply with factor  $r$  and project to  $\mathcal{Q}$



## Jointtree algorithm

**Definition:** A **jointtree**  $(T, \mathcal{C})$  for a DAG  $G$  is a tree  $T$  in which each node has a label  $C_i$  (called **cluster**), satisfying the properties:

- ▶ each cluster is a set of nodes from  $G$
- ▶ each **family\*** in  $G$  appears in some cluster
- ▶ if a node appears in two clusters  $C_i, C_j$ , it must appear in every cluster on the path connecting nodes  $i$  and  $j$  in the jointtree

the **separator** of edge  $i-j$  is defined as  $S_{ij} := C_i \cap C_j$

\*family = a node along with its parents

Also known as **junction trees, clique trees, Markov trees, hypertrees**

An **evidence indicator** is a factor over variable  $X$  that captures the value of  $X$  in evidence  $\mathbf{e}$ :  $\lambda_X(x) = 1$  if  $x$  consistent with  $\mathbf{e}$ , 0 otherwise

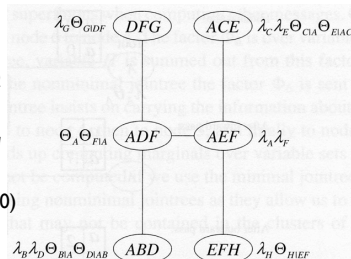
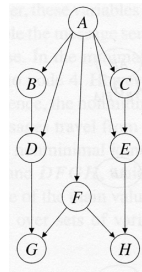
## Jointtree algorithm

### Algorithm:

- ▶ construct jointtree for a given Bayesian network
- ▶ assign each CPT  $\Theta_{X|U}$  to a cluster that contains  $X$  and  $U$
- ▶ assign each evidence indicator  $\lambda_X$  to a cluster that contains  $X$
- ▶ select a root node that contains the query  $\mathcal{Q}$
- ▶ start eliminating factors (using projecting and multiplying) inwards/outwards\*
- ▶ finally project cluster in the root node onto  $\mathcal{Q}$

\*different propagation strategies with different space and time complexities

- ▶ **Shenoy-Shafer architecture** (Shenoy & Shafer 1990)
- ▶ **Hugin architecture** (Jensen et al. 1990)



## Recursive conditioning

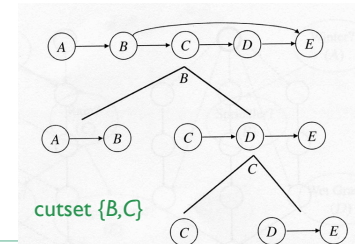
**Idea:** simplify a problem by solving a number of cases and combining the results to a solution to the original problem (**case analysis**)

$$Pr(x) = \sum_c Pr(x, c)$$

**Approach:** reduce query on a network into a queries on simpler networks

- ▶ if var  $E$  given as evidence, the network can be **pruned** (see above)
- ▶ in general: any query  $Pr(\mathbf{q}, \mathbf{e})$  leads to decomposition into networks  $N_e^r$  and  $N_e^l$  such that

$$\begin{aligned} \text{„cutset“ } \mathbf{e} \quad Pr(q) &= \sum_e Pr(q, e) \\ &= \sum_e Pr_e^l(q^l, e^l) Pr_e^r(q^r, e^r) \end{aligned}$$



## Recursive conditioning

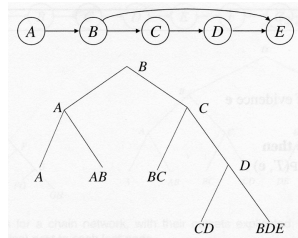
### Recursive condition algorithm:

- ▶ decompose network in a divide-and-conquer fashion, following an appropriate cutset
- ▶ when at leaf node, look up the conditioned CPT
- ▶ propagate value back according to  $\sum_e Pr_e^l(q^l, e^l) Pr_e^r(q^r, e^r)$

Question (again): what is an appropriate cutset (order)?

Answer: all are valid, some lead to less work

- ▶ need to minimize total number of considered cases
- ▶ use **decomposition trees**: full binary trees, leaves are CPTs in the network
- ▶ useful to employ caching techniques (Darwiche, chapt. 8)



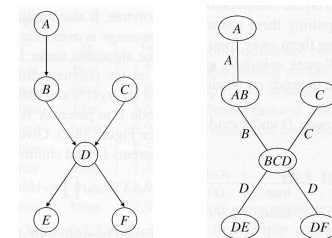
## Belief propagation

### Approximative inference algorithms

- ▶ originally for exact inference in polytree\* networks, then generalized to approximative solution for arbitrary networks
- ▶ spectrum of approximations: trade-off quality with computational costs

Belief propagation algorithm for computing joint marginals  $Pr(X, \mathbf{e})$ :

- ▶ identical to jointtree algorithm for jointtrees that coincide with the polytree network structure
- ▶ Example:



\*polytree = network with only one path between any two nodes

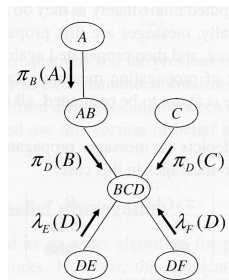
## Belief propagation

Computing joint marginals  $Pr(X, \mathbf{e})$ :

- ▶ node  $i$  in jointtree has cluster  $C_i = X\mathbf{U}$  (with  $\mathbf{U}$  parents of  $X$ )
- ▶ edge  $i-j$  in jointtree corresp. to edge  $X-Y$  in network has „separator“  $S_{ij}=X$
- ▶ „messages“ to eliminate factors:
  - from  $\mathbf{U}$  to  $X$ : **causal support**  $\pi_X(\mathbf{U})$
  - from  $Y$  to parent  $X$ : **diagnostic support**  $\lambda_Y(X)$
- ▶ messages are sent by a node, when it has received messages from all other nodes
  - start with those that do not depend on others

Example:

- ▶ Belief propagation toward node  $D$ , evidence  $E=true$   
 $Pr(BCD, \mathbf{e}) = \Theta_{D|BC} \pi_D(B) \pi_D(C) \lambda_E(D) \lambda_F(D)$



## Belief propagation

problem: leads to „deadlocks“ in some network, when each message is dependent on any other

solution: **iterative belief propagation**

- ▶ assume initial values to each message in the network
- ▶ propagate beliefs and re-iterate
- ▶ converge to a „fixed point“ solution
  - may generally have multiple fixed points on a given network
  - may oscillate on some networks, loop forever

## Stochastic sampling

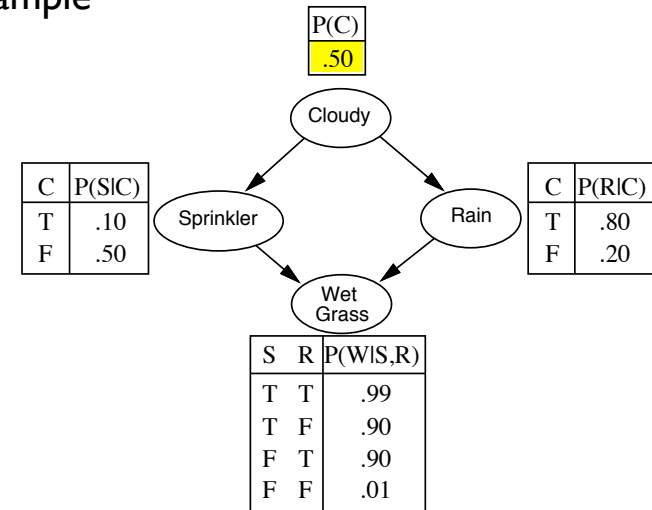
*Idea:* simulate an event according to some probability of occurrence, estimate the prob. of this event from its frequency in these simulations

Simulating a Bayesian network:

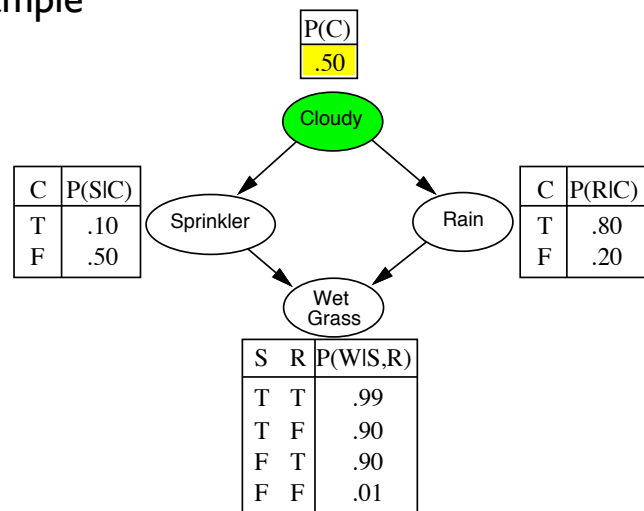
A Bayesian network induces a distribution  $Pr(\mathbf{X})$

- ▶ visit each node in topological order
- ▶ generate value for each node according to  $Pr(x|\mathbf{u})$
- ▶ end with a sample  $\{\mathbf{x}^1, \dots, \mathbf{x}^n\}$  of  $n$  events
- ▶ estimate probability  $^{\wedge}Pr(\mathbf{x})$  of value  $\mathbf{x}$  from its frequency in this sample
- ▶ show that  $^{\wedge}Pr(\mathbf{x})$  converges against  $Pr(\mathbf{x})$  with increasing  $n$

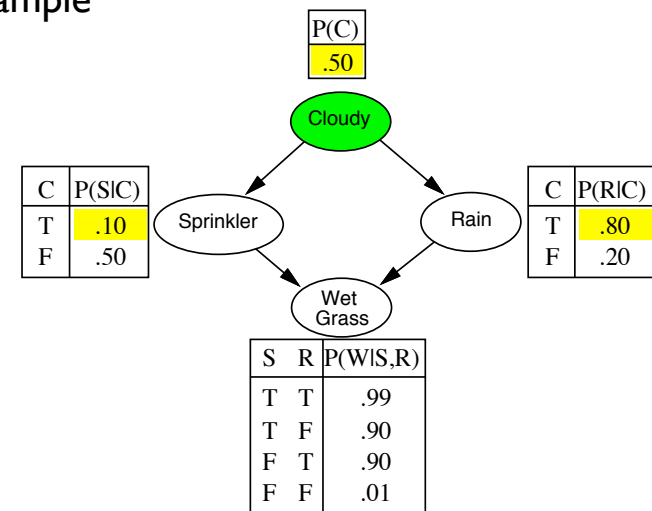
## Example



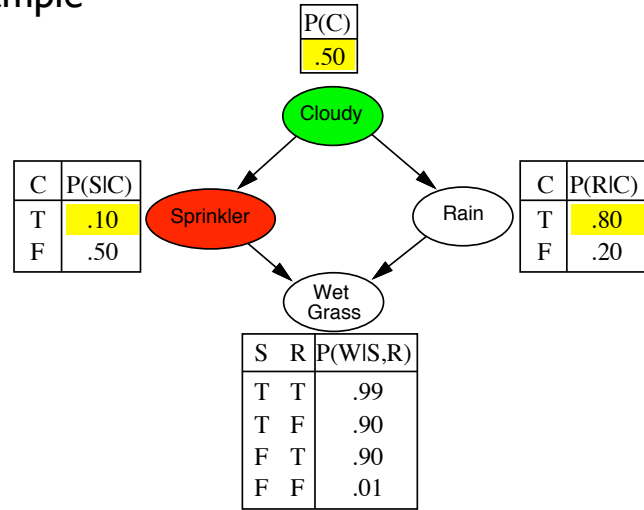
## Example



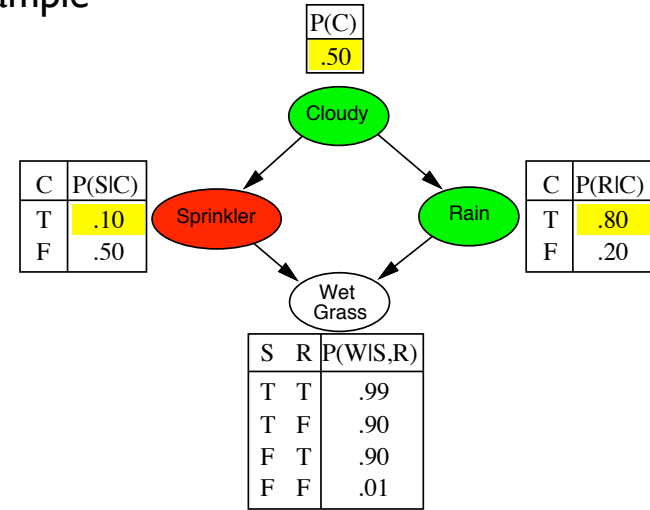
## Example



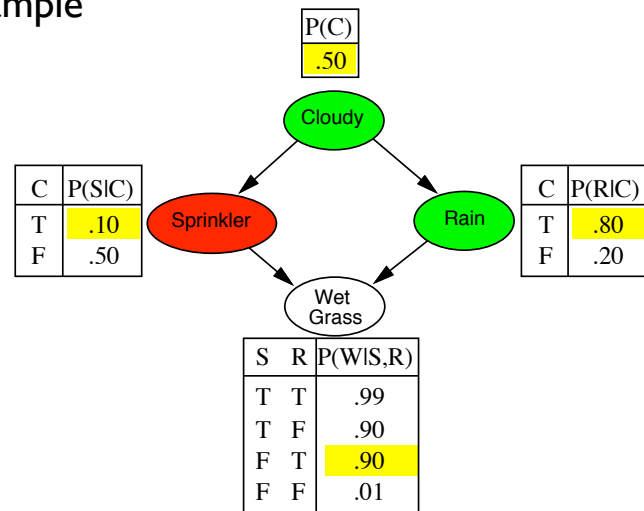
### Example



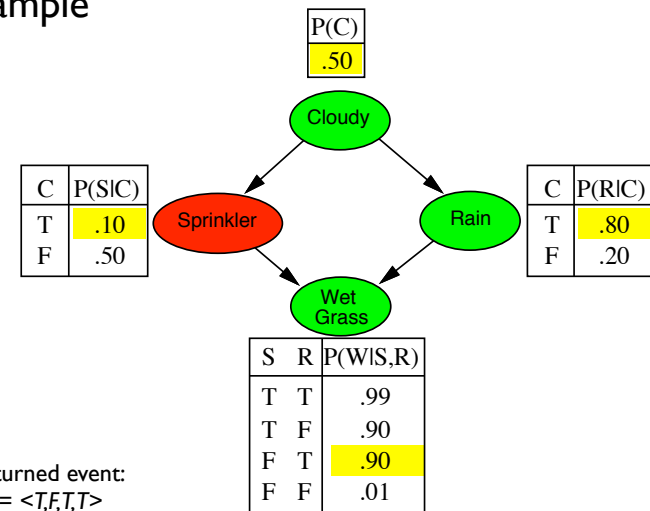
### Example



### Example



### Example



## Stochastic sampling

Sampling relies on taking probability as expectation about a function

- ▶ expectation value of a function  $f(\mathbf{X})$ :  $Ex(f) := \sum_x f(x) \cdot Pr(x)$   
 $\mu$
- ▶ variance of a function  $f(\mathbf{X})$ :  $Var(f) := \sum_x (f(x) - Ex(f))^2 \cdot Pr(x)$   
 $\sigma^2$

Direct sampling function:

- ▶ let  $\hat{\alpha}(x) := 1$  if  $\alpha$  true at  $x$ , 0 otherwise
- ▶ then:  $Ex(\hat{\alpha}) = Pr(\alpha)$   
 $Var(\hat{\alpha}) = Pr(\alpha)Pr(\neg\alpha) = Pr(\alpha) - Pr(\alpha)^2$

That is, approximating  $Pr$  boils down to **estimating the expectation**

**How?**

## Monte Carlo simulation

Principle:

- ▶ simulate random sample  $\mathbf{x}^1, \dots, \mathbf{x}^n$  from **sampling distribution**  $Pr(\mathbf{X})$
- ▶ evaluate function at each instantiation  $f(\mathbf{x}^1), \dots, f(\mathbf{x}^n)$
- ▶ compute arithmetic average of attained values: **sample mean**  
 $Av_n(f) := \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}^i)$
- ▶ works because of **law of large numbers**: for function  $f$  with expectation  $\mu$  and every  $\epsilon > 0$ :  $\lim_{n \rightarrow \infty} P(|Av_n(f) - \mu| \leq \epsilon) = 1$

Monte Carlo simulation using  $\hat{\alpha}(x)$  gives **direct sampling**:

- ▶ simulate sample  $\mathbf{x}^1, \dots, \mathbf{x}^n$  from Bayesian network
- ▶ compute values  $\hat{\alpha}(x^1), \dots, \hat{\alpha}(x^n)$
- ▶ estimate  $Pr(\alpha)$  using sample mean  $Av_n(\hat{\alpha})$

## Rejection sampling

Goal: Calculate conditional prob.  $Pr(a|b)$  with  $Pr(\cdot)$  induced by network

Approach:

- ▶ calculate estimate for  $Pr(a \wedge b)$  and  $Pr(b)$ :  $Av_n(\hat{\gamma}), Av_n(\hat{\beta})$  with  $\gamma = a \wedge b$
- ▶ take ratio as estimate for  $Pr(a|b)$ :  $Av_n(\hat{\gamma}) / Av_n(\hat{\beta})$   
-  $c_1 = \# \text{samples with } a \wedge b = \text{true}, c_2 = \# \text{samples with } b = \text{true} \rightarrow (c_1/n) / (c_2/n) = c_1/c_2$
- ▶ reject all samples in which  $b$  is false: **rejection sampling**

Example: estimate  $P(\text{Rain} | \text{Sprinkler} = \text{true})$  from 100 samples; 27 have  $\text{Sprinkler} = \text{true}$ , of these 8 have  $\text{Rain} = \text{true}$ , 19 have  $\text{Rain} = \text{false}$

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

True answer:  $\langle 0.3, 0.7 \rangle$

## Importance sampling

Idea: reduce variance due to rare events by sampling from an **importance distribution**  $Pr'$  emphasizing instantiations consistent with rare event

Monte Carlo simulation using the **importance sampling function**:

$$\tilde{\alpha}(x) = Pr(x) / Pr'(x) \text{ if } \alpha \text{ true at instantiation } x, 0 \text{ otherwise}$$

Improves on direct sampling only when  $Pr'$  emphasizes important events no less than  $Pr$

Finding ideal distribution generally not feasible, but some other weaker conditions can be ensured easier and still improve on variance



# Importance sampling

**Likelihood weighting:** given evidence  $e$ , what is  $Pr(x|e)$ ?

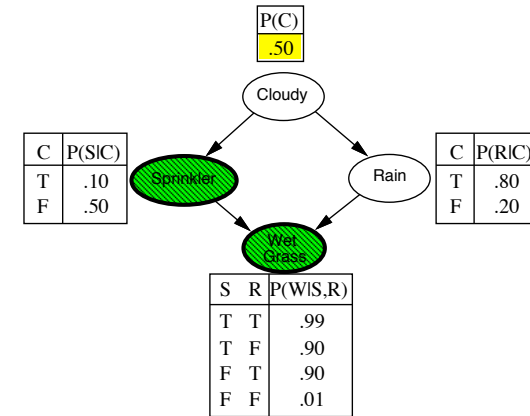
- ▶ generate only samples that are consistent with  $e$
- ▶ fix evidence variables, sample non-evidence var's and **weight sample by likelihood** it accords the evidence
- ▶ consistent estimate, but performance gets worse with growing evidence because few samples have  $\sim$ all total weight

```

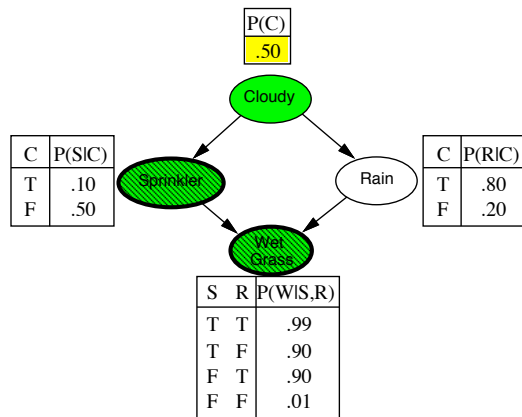
function WEIGHTED-SAMPLE( $bn, e$ ) returns an event and a weight
   $x \leftarrow$  an event with  $n$  elements;  $w \leftarrow 1$ 
  for  $i = 1$  to  $n$  do
    if  $X_i$  has a value  $x_i$  in  $e$ 
      then  $w \leftarrow w \times P(X_i = x_i | parents(X_i))$ 
      else  $x_i \leftarrow$ 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

# Example

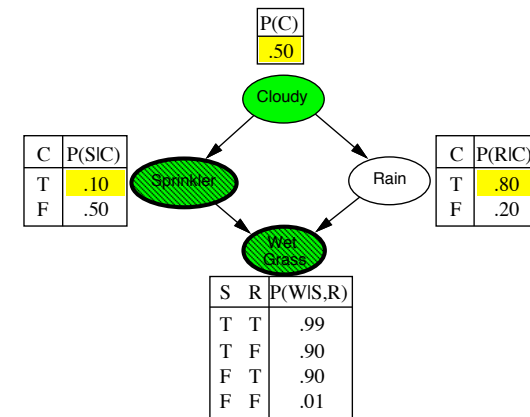
Query:  $P(\text{Rain} | \text{Sprinkler}=\text{true}, \text{WetGrass}=\text{true}) = ??$



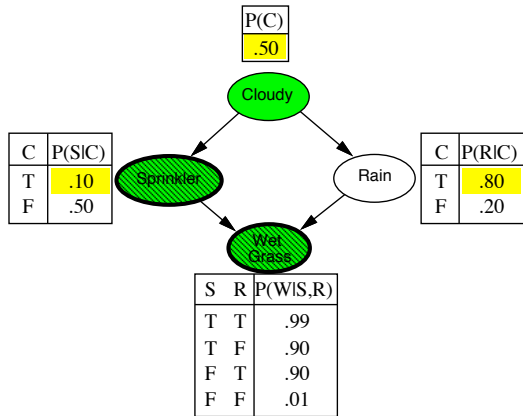
$w = 1.0$



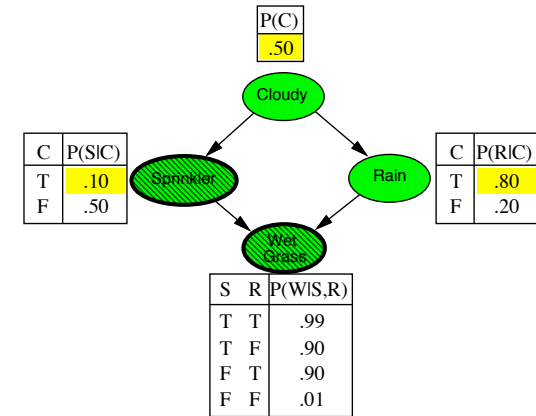
$w = 1.0$



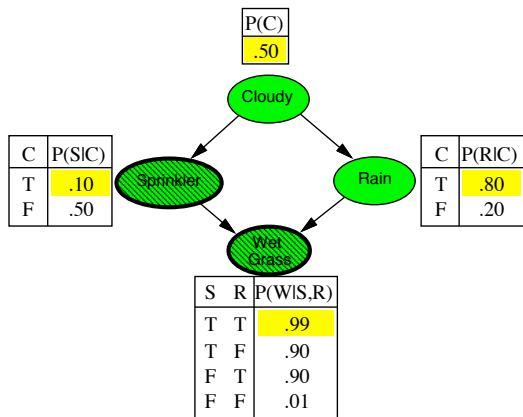
$w = 1.0$



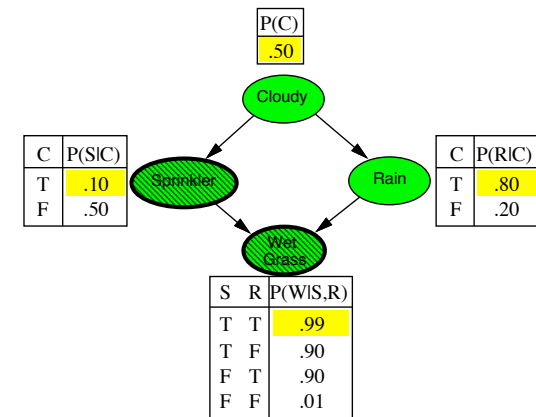
$w = 1.0 \times 0.1$



$w = 1.0 \times 0.1$



$w = 1.0 \times 0.1$

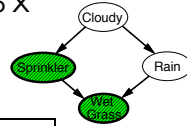


$w = 1.0 \times 0.1 \times 0.99 = 0.099 = \text{weight for event } \langle t, t, t \rangle$

## Markov Chain Monte Carlo (MCMC)

Network is in a **state** = current assignment to variables  
 next state: sample **non-evidence variable**  $X$  given its **Markov blanket**  
 = variables that, when known, make other variables irrelevant to  $X$

- ▶ Markov blanket of *Cloudy* is *Sprinkler* and *Rain*
- ▶ Markov blanket of *Rain* is *Sprinkler*, *Cloudy*, *WetGrass*

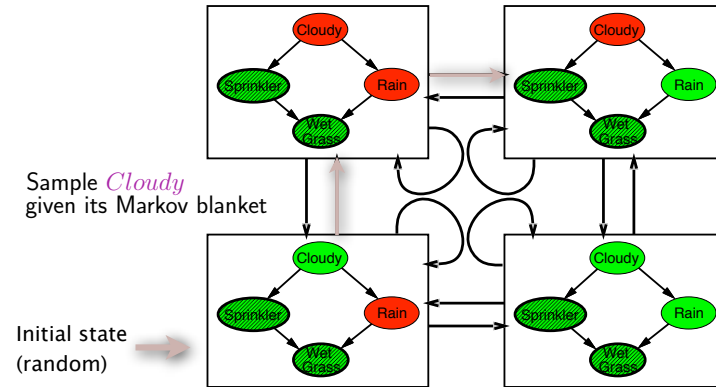


```
function MCMC-Ask( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
  local variables:  $N[X]$ , a vector of counts over  $X$ , initially zero
                   $Z_i$ , the nonevidence variables in  $bn$ 
                   $x$ , the current state of the network, initially copied from  $e$ 

  initialize  $x$  with random values for the variables in  $Y$ 
  for  $j = 1$  to  $N$  do
    for each  $Z_i$  in  $Z$  do
      sample the value of  $Z_i$  in  $x$  from  $P(Z_i|mb(Z_i))$ 
      given the values of  $MB(Z_i)$  in  $x$ 
       $N[x] \leftarrow N[x] + 1$  where  $x$  is the value of  $X$  in  $x$ 
  return NORMALIZE( $N[X]$ )
```

„transition prob.“ of moving into new state

Estimate  $P(Rain|Sprinkler = true, WetGrass = true)$



Count number of times *Rain* is true and false in the samples.

E.g.: visit 100 states:  $P(Rain|Sprinkler = true, WetGrass = true)$   
 $= \text{NORMALIZE}(\langle 31, 69 \rangle) = \langle 0.31, 0.69 \rangle$

## MCMC - Markov blanket sampling

Because of the transition probabilities, sampling runs into an „equilibrium“ in which **time spent in each state is proportional to its posterior probability**

Transition probability (given the Markov blanket) is:

$$P(x_i|mb(X_i)) = P(x_i|parents(X_i)) \prod_{Z_j \in Children(X_i)} P(z_j|parents(Z_j))$$

- ▶ easily implemented in parallel systems

Main difficulties:

- ▶ difficult to tell if and when convergence has been achieved
- ▶ can be wasteful if Markov blanket large, prob doesn't change much

## Bayes nets inference algorithms - summary

### Exact algorithms

- ▶ Variable Elimination and Factor Elimination
- ▶ Jointtree algorithm
- ▶ Recursive conditioning

### Approximative algorithms

- ▶ Belief propagation
- ▶ Stochastic sampling (Monte Carlo simulation)
  - direct sampling
  - importance sampling, likelihood weighting
- ▶ Monte Carlo Markov Chain