# CSE 473: Artificial Intelligence

## Bayesian Networks: Inference

## Luke Zettlemoyer

Many slides over the course adapted from either Dan Klein,
Stuart Russell or Andrew Moore

# Outline

- **Bayesian Networks Inference**
  - Exact Inference: Variable Elimination
  - Approximate Inference: Sampling

# Probabilistic Inference

- *Probabilistic inference*: compute a desired probability from other known probabilities (e.g. conditional from joint)

- We generally compute conditional probabilities
  - P(on time | no reported accidents) = 0.90
  - These represent the agent's *beliefs* given the evidence

- Probabilities change with new evidence:
  - P(on time | no accidents, 5 a.m.) = 0.95
  - P(on time | no accidents, 5 a.m., raining) = 0.80
  - Observing new evidence causes *beliefs to be updated*

# Inference by Enumeration

- General case:
  - Evidence variables: $E_1 \ldots E_k = e_1 \ldots e_k$
  - Query* variable: $Q$
  - Hidden variables: $H_1 \ldots H_r$

  $\left.\vphantom{\begin{array}{c}a\\b\\c\end{array}}\right\}$ $X_1, X_2, \ldots X_n$

  *All variables*

- We want: $P(Q|e_1 \ldots e_k)$
- First, select the entries consistent with the evidence
- Second, sum out H to get joint of Query and evidence:

$$P(Q, e_1 \ldots e_k) = \sum_{h_1 \ldots h_r} \underbrace{P(Q, h_1 \ldots h_r, e_1 \ldots e_k)}_{X_1, X_2, \ldots X_n}$$

- Finally, normalize the remaining entries to conditionalize

- Obvious problems:
  - Worst-case time complexity $O(d^n)$
  - Space complexity $O(d^n)$ to store the joint distribution

# Variable Elimination

- **Why is inference by enumeration so slow?**
  - You join up the whole joint distribution before you sum out the hidden variables
  - You end up repeating a lot of work!

- **Idea: interleave joining and marginalizing!**
  - Called "Variable Elimination"
  - Still NP-hard, but usually much faster than inference by enumeration

- **We'll need some new notation to define VE**

# Review: Factor Zoo I

$P(T, W)$

| T | W | P |
|------|------|-----|
| hot | sun | 0.4 |
| hot | rain | 0.1 |
| cold | sun | 0.2 |
| cold | rain | 0.3 |

- Joint distribution: P(X,Y)
  - Entries P(x,y) for all x, y
  - Sums to 1

$P(cold, W)$

| T | W | P |
|------|------|-----|
| cold | sun | 0.2 |
| cold | rain | 0.3 |

- Selected joint: P(x,Y)
  - A slice of the joint distribution
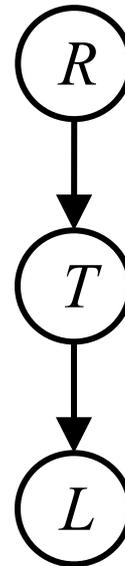  - Entries P(x,y) for fixed x, all y
  - Sums to P(x)

# Example: Traffic Domain

- **Random Variables**
  - R: Raining
  - T: Traffic
  - L: Late for class!

- **First query: P(L)**

$$P(l) = \sum_t \sum_r P(l|t)P(t|r)P(r)$$

$P(R)$

| | |
|---|---|
| +r | 0.1 |
| -r | 0.9 |

$P(T|R)$

| | | |
|---|---|---|
| +r | +t | 0.8 |
| +r | -t | 0.2 |
| -r | +t | 0.1 |
| -r | -t | 0.9 |

$P(L|T)$

| | | |
|---|---|---|
| +t | +l | 0.3 |
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

# Variable Elimination Outline

- Maintain a set of tables called factors

- Initial factors are local CPTs (one per node)

$P(R)$

| | |
|---|---|
| +r | 0.1 |
| -r | 0.9 |

$P(T|R)$

| | | |
|---|---|---|
| +r | +t | 0.8 |
| +r | -t | 0.2 |
| -r | +t | 0.1 |
| -r | -t | 0.9 |

$P(L|T)$

| | | |
|---|---|---|
| +t | +l | 0.3 |
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

- Any known values are selected
  - E.g. if we know $L = +\ell$ , the initial factors are

$P(R)$

| | |
|---|---|
| +r | 0.1 |
| -r | 0.9 |

$P(T|R)$

| | | |
|---|---|---|
| +r | +t | 0.8 |
| +r | -t | 0.2 |
| -r | +t | 0.1 |
| -r | -t | 0.9 |

$P(+\ell|T)$

| | | |
|---|---|---|
| +t | +l | 0.3 |
| -t | +l | 0.1 |

- VE: Alternately join factors and eliminate variables

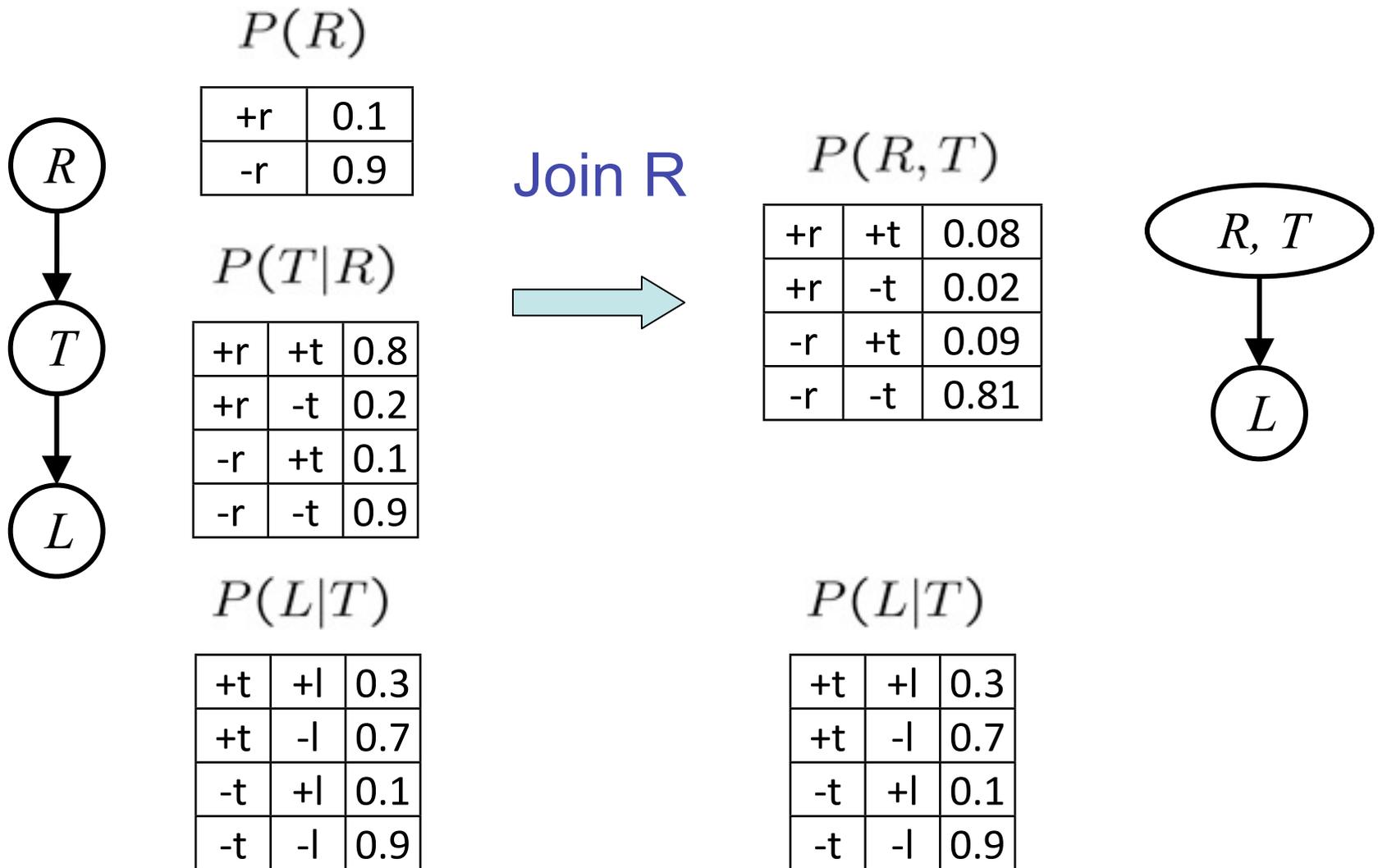# Operation 1: Join Factors

- First basic operation: joining factors
- Combining factors:
  - Just like a database join
  - Get all factors over the joining variable
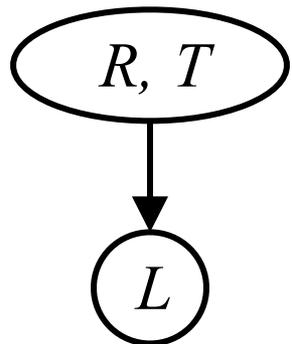  - Build a new factor over the union of the variables involved

- Example: Join on R

$P(R) \quad \times \quad P(T|R) \quad \Longrightarrow \quad P(R,T) \qquad \boxed{R,T}$

| +r | 0.1 |
|----|-----|
| -r | 0.9 |

| +r | +t | 0.8 |
|----|----|-----|
| +r | -t | 0.2 |
| -r | +t | 0.1 |
| -r | -t | 0.9 |

| +r | +t | 0.08 |
|----|----|------|
| +r | -t | 0.02 |
| -r | +t | 0.09 |
| -r | -t | 0.81 |

- Computation for each entry: pointwise products

$$\forall r, t : \quad P(r,t) = P(r) \cdot P(t|r)$$

# Example: Multiple Joins

$P(R)$

| | |
|---|---|
| +r | 0.1 |
| -r | 0.9 |

Join R

$P(R, T)$

| | | |
|---|---|---|
| +r | +t | 0.08 |
| +r | -t | 0.02 |
| -r | +t | 0.09 |
| -r | -t | 0.81 |

$R, T$

$P(T|R)$

| | | |
|---|---|---|
| +r | +t | 0.8 |
| +r | -t | 0.2 |
| -r | +t | 0.1 |
| -r | -t | 0.9 |

$L$

$P(L|T)$

| | | |
|---|---|---|
| +t | +l | 0.3 |
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

$P(L|T)$

| | | |
|---|---|---|
| +t | +l | 0.3 |
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

# Example: Multiple Joins

$P(R, T)$

| +r | +t | 0.08 |
|----|----|------|
| +r | -t | 0.02 |
| -r | +t | 0.09 |
| -r | -t | 0.81 |

R, T

L

Join T

$R, T, L$

$P(R, T, L)$

| +r | +t | +l | 0.024 |
|----|----|----|-------|
| +r | +t | -l | 0.056 |
| +r | -t | +l | 0.002 |
| +r | -t | -l | 0.018 |
| -r | +t | +l | 0.027 |
| -r | +t | -l | 0.063 |
| -r | -t | +l | 0.081 |
| -r | -t | -l | 0.729 |

$P(L|T)$

| +t | +l | 0.3 |
|----|----|-----|
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

# Operation 2: Eliminate

- **Second basic operation: marginalization**
- **Take a factor and sum out a variable**
  - Shrinks a factor to a smaller one
  - A projection operation
- **Example:**

$P(R,T)$

| +r | +t | 0.08 |
|----|----|------|
| +r | -t | 0.02 |
| -r | +t | 0.09 |
| -r | -t | 0.81 |

sum $R$ →

$P(T)$

| +t | 0.17 |
|----|------|
| -t | 0.83 |

# Multiple Elimination

$R, T, L$          $T, L$          $L$

$P(R, T, L)$

| | | | |
|---|---|---|---|
| +r | +t | +l | 0.024 |
| +r | +t | -l | 0.056 |
| +r | -t | +l | 0.002 |
| +r | -t | -l | 0.018 |
| -r | +t | +l | 0.027 |
| -r | +t | -l | 0.063 |
| -r | -t | +l | 0.081 |
| -r | -t | -l | 0.729 |

**Sum out R**

$P(T, L)$

| | | |
|---|---|---|
| +t | +l | 0.051 |
| +t | -l | 0.119 |
| -t | +l | 0.083 |
| -t | -l | 0.747 |

**Sum out T**

$P(L)$

| | |
|---|---|
| +l | 0.134 |
| -l | 0.886 |

# P(L) : Marginalizing Early!

$P(R)$

| +r | 0.1 |
|----|-----|
| -r | 0.9 |

**Join R**

**Sum out R**

$P(T|R)$



| +r | +t | 0.8 |
|----|----|-----|
| +r | -t | 0.2 |
| -r | +t | 0.1 |
| -r | -t | 0.9 |

$P(R,T)$

| +r | +t | 0.08 |
|----|----|------|
| +r | -t | 0.02 |
| -r | +t | 0.09 |
| -r | -t | 0.81 |

$P(T)$

| +t | 0.17 |
|----|------|
| -t | 0.83 |

$P(L|T)$

| +t | +l | 0.3 |
|----|----|-----|
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

$P(L|T)$

| +t | +l | 0.3 |
|----|----|-----|
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

$P(L|T)$

| +t | +l | 0.3 |
|----|----|-----|
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

# Marginalizing Early (aka VE*)



T → L

Join T    (T, L)    Sum out T    L

P(T)

| +t | 0.17 |
|----|------|
| -t | 0.83 |

P(L|T)

| +t | +l | 0.3 |
|----|----|-----|
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

P(T, L)

| +t | +l | 0.051 |
|----|----|-------|
| +t | -l | 0.119 |
| -t | +l | 0.083 |
| -t | -l | 0.747 |

P(L)

| +l | 0.134 |
|----|-------|
| -l | 0.886 |

\* VE is variable elimination

# Evidence

- **If evidence, start with factors that select that evidence**
  - No evidence uses these initial factors:

$P(R)$

| +r | 0.1 |
|----|-----|
| -r | 0.9 |

$P(T|R)$

| +r | +t | 0.8 |
|----|----|-----|
| +r | -t | 0.2 |
| -r | +t | 0.1 |
| -r | -t | 0.9 |

$P(L|T)$

| +t | +l | 0.3 |
|----|----|-----|
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

  - Computing $P(L| + r)$, the initial factors become:

$P(+r)$

| +r | 0.1 |
|----|-----|

$P(T| + r)$

| +r | +t | 0.8 |
|----|----|-----|
| +r | -t | 0.2 |

$P(L|T)$

| +t | +l | 0.3 |
|----|----|-----|
| +t | -l | 0.7 |
| -t | +l | 0.1 |
| -t | -l | 0.9 |

- **We eliminate all vars other than query + evidence**

# Evidence II

- Result will be a selected joint of query and evidence
    - E.g. for P(L | +r), we'd end up with:

$P(+r, L)$

| +r | +l | 0.026 |
|----|----|-------|
| +r | -l | 0.074 |

Normalize →

$P(L| + r)$

| +l | 0.26 |
|----|------|
| -l | 0.74 |

- To get our answer, just normalize this!

- That's it!

# General Variable Elimination

- Query: $P(Q|E_1 = e_1, \ldots E_k = e_k)$

- Start with initial factors:
  - Local CPTs (but instantiated by evidence)

- While there are still hidden variables (not Q or evidence):
  - Pick a hidden variable H
  - Join all factors mentioning H
  - Eliminate (sum out) H

- Join all remaining factors and normalize

# Variable Elimination Bayes Rule

## Start / Select

$P(B)$

| B | P |
|---|---|
| +b | 0.1 |
| ¬b | 0.9 |



$P(A|B) \rightarrow P(a|B)$

| B | A | P |
|---|---|---|
| +b | +a | 0.8 |
| ~~b~~ | ~~¬a~~ | ~~0.2~~ |
| ¬b | +a | 0.1 |
| ~~¬b~~ | ~~¬a~~ | ~~0.9~~ |

## Join on B

$a, B$

$P(a, B)$

| A | B | P |
|---|---|---|
| +a | +b | 0.08 |
| +a | ¬b | 0.09 |

## Normalize

$P(B|a)$

| A | B | P |
|---|---|---|
| +a | +b | 8/17 |
| +a | ¬b | 9/17 |

# Example

Query: $P(B|j, m)$

$$P(B) \qquad P(E) \qquad P(A|B, E) \qquad P(j|A) \qquad P(m|A)$$

## Choose A

$$P(A|B, E)$$
$$P(j|A) \qquad \times \qquad P(j, m, A|B, E) \qquad \Sigma \qquad P(j, m|B, E)$$
$$P(m|A)$$

$$P(B) \qquad P(E) \qquad P(j, m|B, E)$$

# Example

$$P(B) \qquad P(E) \qquad P(j, m | B, E)$$

## Choose E

$$P(E)$$
$$P(j, m | B, E)$$
$\xrightarrow{\times}$
$$P(j, m, E | B)$$
$\xrightarrow{\Sigma}$
$$P(j, m | B)$$

$$P(B) \qquad P(j, m | B)$$

## Finish with B

$$P(B)$$
$$P(j, m | B)$$
$\xrightarrow{\times}$
$$P(j, m, B)$$
$\xrightarrow{\text{Normalize}}$
$$P(B | j, m)$$

# Exact Inference: Variable Elimination

- **Remaining Issues:**
  - Complexity: exponential in tree width (size of the largest factor created)
  - Best elimination ordering? NP-hard problem

- **What you need to know:**
  - Should be able to run it on small examples, understand the factor creation / reduction flow
  - Better than enumeration: saves time by marginalizing variables as soon as possible rather than at the end

- **We have seen a special case of VE already**
  - HMM Forward Inference

# Approximate Inference

- **Simulation has a name: sampling**

- **Sampling is a hot topic in machine learning, and it's really simple**

- **Basic idea:**
  - Draw N samples from a sampling distribution S
  - Compute an approximate posterior probability
  - Show this converges to the true probability P

- **Why sample?**
  - Learning: get samples from a distribution you don't know
  - Inference: getting a sample is faster than computing the right answer (e.g. with variable elimination)

$F$

$S$

$A$

# Prior Sampling

$P(C)$

| +c | 0.5 |
|----|-----|
| -c | 0.5 |

Cloudy

$P(S|C)$

| | +s | 0.1 |
|----|----|-----|
| +c | -s | 0.9 |
| | +s | 0.5 |
| -c | -s | 0.5 |

Sprinkler

Rain

$P(R|C)$

| | +r | 0.8 |
|----|----|-----|
| +c | -r | 0.2 |
| | +r | 0.2 |
| -c | -r | 0.8 |

$P(W|S,R)$

| | | +w | 0.99 |
|----|----|-----|------|
| | +r | -w | 0.01 |
| | | +w | 0.90 |
| +s | -r | -w | 0.10 |
| | | +w | 0.90 |
| | +r | -w | 0.10 |
| | | +w | 0.01 |
| -s | -r | -w | 0.99 |

WetGrass

Samples:

+c, -s, +r, +w

-c, +s, -r, +w

…

# Prior Sampling

- This process generates samples with probability:

$$S_{PS}(x_1 \ldots x_n) = \prod_{i=1}^{n} P(x_i | \text{Parents}(X_i)) = P(x_1 \ldots x_n)$$

…i.e. the BN's joint probability

- Let the number of samples of an event be $N_{PS}(x_1 \ldots x_n)$

- Then
$$\begin{aligned}
\lim_{N \to \infty} \hat{P}(x_1, \ldots, x_n) &= \lim_{N \to \infty} N_{PS}(x_1, \ldots, x_n)/N \\
&= S_{PS}(x_1, \ldots, x_n) \\
&= P(x_1 \ldots x_n)
\end{aligned}$$

- I.e., the sampling procedure is consistent

# Example

- We'll get a bunch of samples from the BN:

    +c, -s, +r, +w

    +c, +s, +r, +w

    -c, +s, +r,  -w

    +c, -s, +r, +w

    -c,  -s,  -r, +w

- If we want to know P(W)
  - We have counts <+w:4, -w:1>
  - Normalize to get P(W) = <+w:0.8, -w:0.2>
  - This will get closer to the true distribution with more samples
  - Can estimate anything else, too
  - What about P(C| +w)?   P(C| +r, +w)?  P(C| -r, -w)?
  - Fast: can use fewer samples if less time (what's the drawback?)

# Rejection Sampling

- ## Let's say we want P(C)

    - No point keeping all samples around
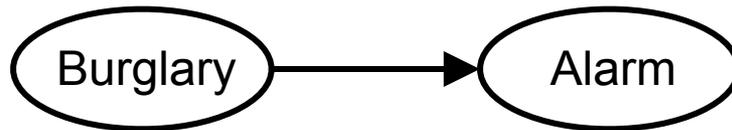    - Just tally counts of C as we go

- ## Let's say we want P(C| +s)

    - Same thing: tally C outcomes, but ignore (reject) samples which don't have S=+s
    - This is called rejection sampling
    - It is also consistent for conditional probabilities (i.e., correct in the limit)

+c, -s, +r, +w
+c, +s, +r, +w
-c, +s, +r, -w
+c, -s, +r, +w
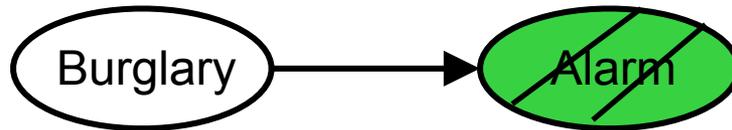-c, -s, -r, +w

# Likelihood Weighting

- **Problem with rejection sampling:**
    - If evidence is unlikely, you reject a lot of samples
    - You don't exploit your evidence as you sample
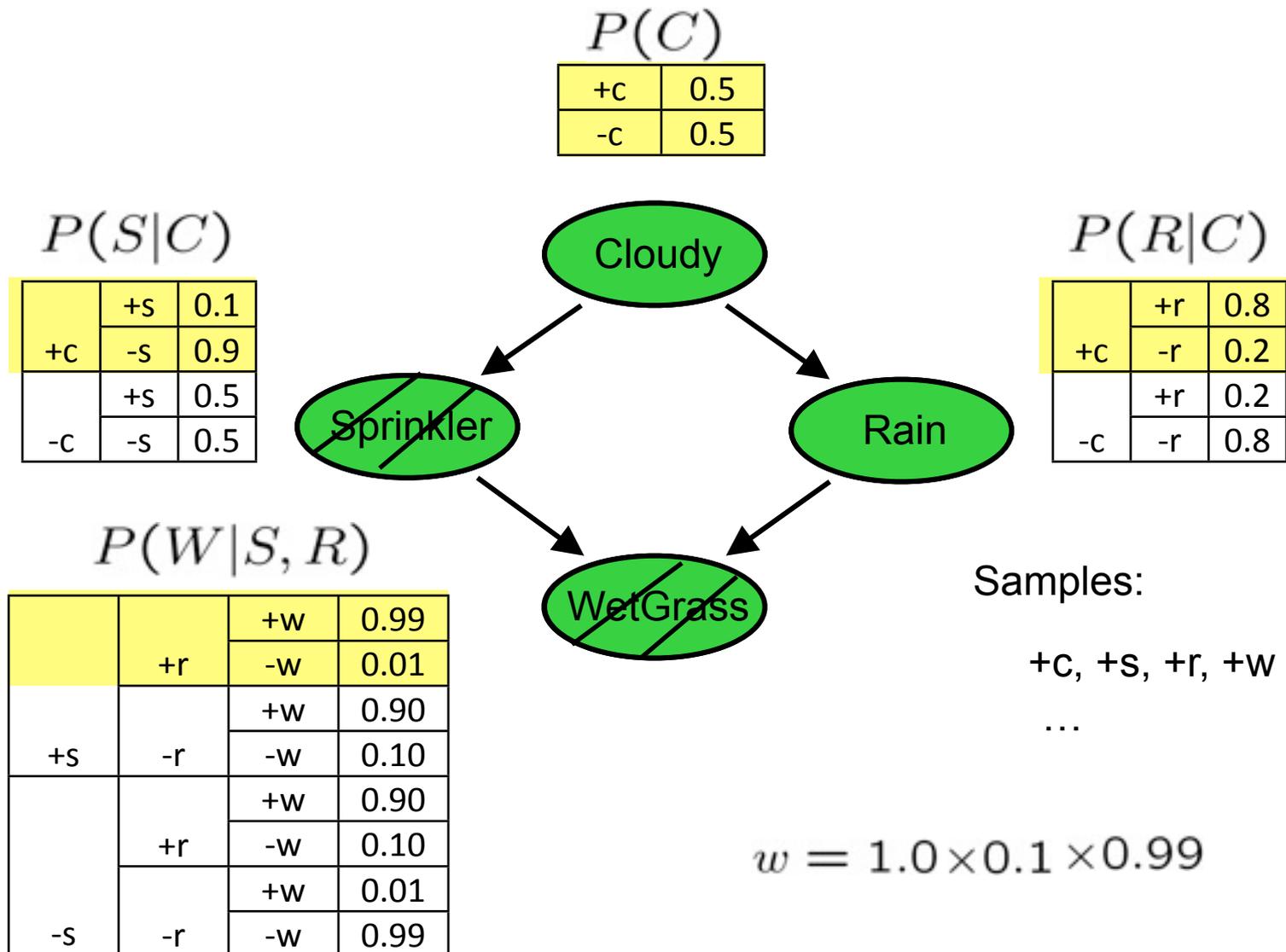    - Consider P(B|+a)



-b,  -a
-b,  -a
-b,  -a
-b,  -a
+b, +a

- **Idea: fix evidence variables and sample the rest**



-b  +a
-b, +a
-b, +a
-b, +a
+b, +a

- **Problem: sample distribution not consistent!**
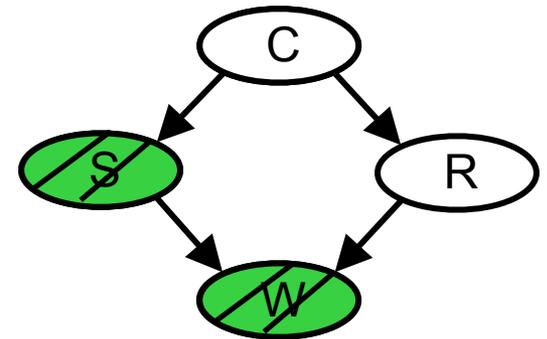- **Solution: weight by probability of evidence given parents**

# Likelihood Weighting

$P(C)$

| +c | 0.5 |
|----|-----|
| -c | 0.5 |

Cloudy

$P(S|C)$

|    | +s | 0.1 |
|----|----|-----|
| +c | -s | 0.9 |
|    | +s | 0.5 |
| -c | -s | 0.5 |

$P(R|C)$

|    | +r | 0.8 |
|----|----|-----|
| +c | -r | 0.2 |
|    | +r | 0.2 |
| -c | -r | 0.8 |

Sprinkler

Rain

$P(W|S,R)$

|    |    | +w | 0.99 |
|----|----|----|------|
|    | +r | -w | 0.01 |
|    |    | +w | 0.90 |
| +s | -r | -w | 0.10 |
|    |    | +w | 0.90 |
|    | +r | -w | 0.10 |
|    |    | +w | 0.01 |
| -s | -r | -w | 0.99 |

WetGrass

Samples:

+c, +s, +r, +w

…

$w = 1.0 \times 0.1 \times 0.99$

# Likelihood Weighting

- Sampling distribution if z sampled and e fixed evidence

$$S_{WS}(\mathbf{z}, \mathbf{e}) = \prod_{i=1}^{l} P(z_i | \text{Parents}(Z_i))$$

- Now, samples have weights

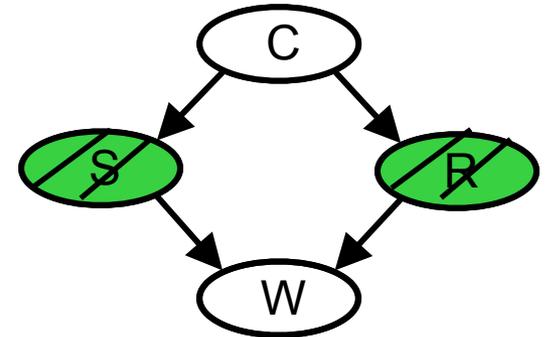$$w(\mathbf{z}, \mathbf{e}) = \prod_{i=1}^{m} P(e_i | \text{Parents}(E_i))$$

- Together, weighted sampling distribution is consistent

$$S_{WS}(z, e) \cdot w(z, e) = \prod_{i=1}^{l} P(z_i | \text{Parents}(z_i)) \prod_{i=1}^{m} P(e_i | \text{Parents}(e_i))$$
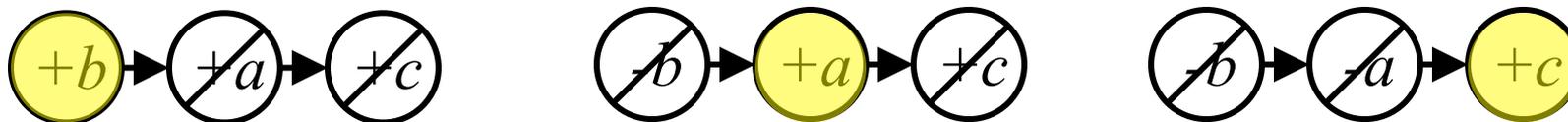
$$= P(\mathbf{z}, \mathbf{e})$$

# Likelihood Weighting

- **Likelihood weighting is good**
    - We have taken evidence into account as we generate the sample
    - E.g. here, W's value will get picked based on the evidence values of S, R
    - More of our samples will reflect the state of the world suggested by the evidence

- **Likelihood weighting doesn't solve all our problems**
    - Evidence influences the choice of downstream variables, but not upstream ones (C isn't more likely to get a value matching the evidence)

- **We would like to consider evidence when we sample every variable**

# Markov Chain Monte Carlo*

- *Idea*: instead of sampling from scratch, create samples that are each like the last one.

- *Gibbs Sampling*: resample one variable at a time, conditioned on the rest, but keep evidence fixed.

- *Properties*: Now samples are not independent (in fact they're nearly identical), but sample averages are still consistent estimators!

- *What's the point*: both upstream and downstream variables condition on evidence.