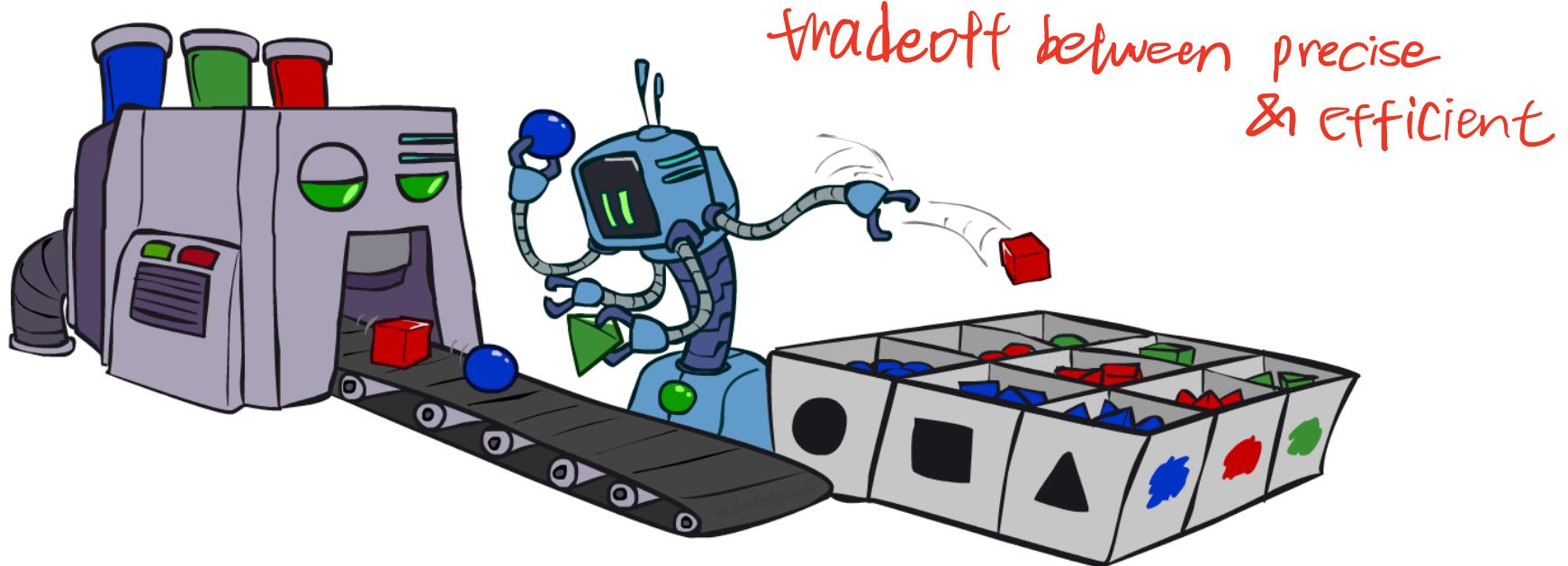


Bayes Nets: Approximate Inference



AIMA Chapter 14.5, PRML Chapter 11

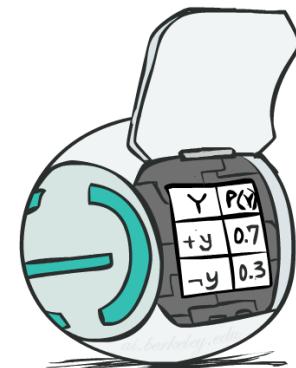
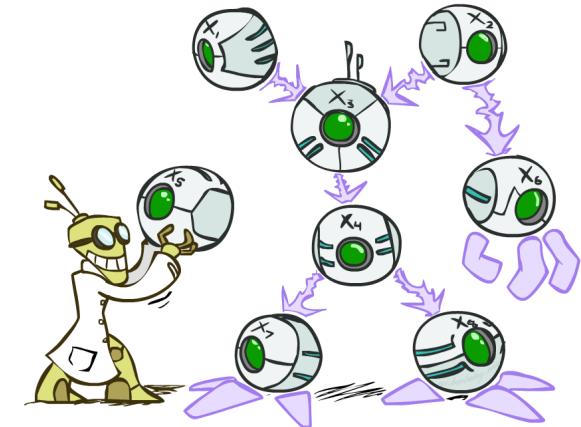
Bayes' Net Representation

- A directed, acyclic graph, one node per random variable
- A conditional probability table (CPT) for each node
 - A collection of distributions over X, one for each combination of parents' values

$$P(X|a_1 \dots a_n)$$

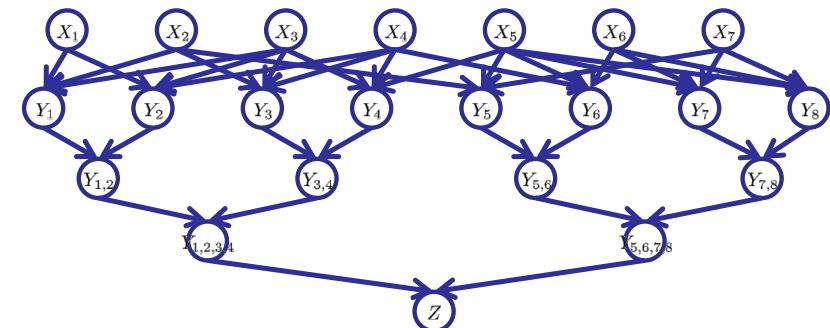
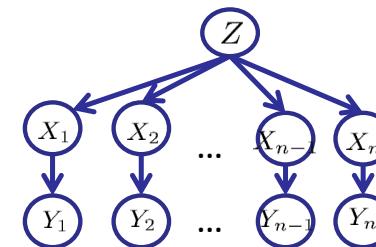
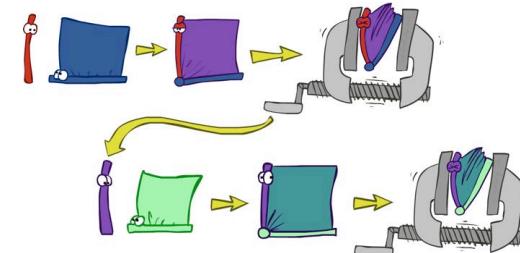
- Bayes' nets implicitly encode joint distributions
 - As a product of local conditional distributions
 - To see what probability a BN gives to a full assignment, multiply all the relevant conditionals together:

$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$



Variable Elimination

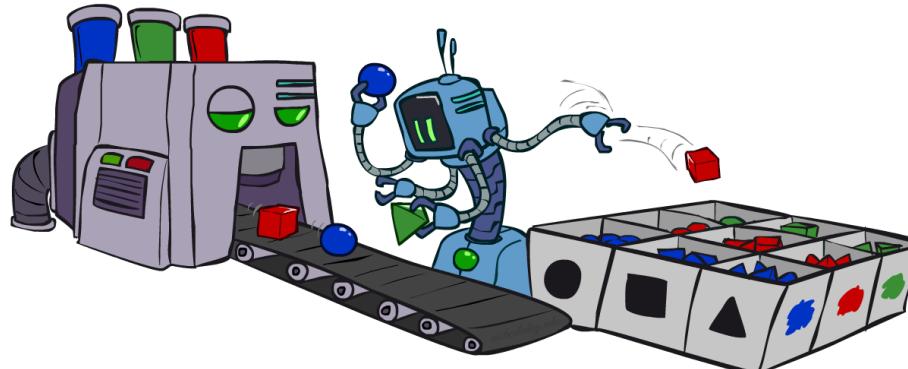
- Interleave joining and marginalizing
- d^k entries computed for a factor over k variables with domain sizes d
- Ordering of elimination of hidden variables can affect size of factors generated
- Worst case: running time exponential in the size of the Bayes' net



Sampling

频率 & 概率

- Goal: probability P
- Basic idea
 - Draw N samples from a sampling distribution S
 - Compute some quantity from the samples
 - Show this converges to the true probability P
- Why sample?
 - Often very fast to get a decent approximate answer
 - The algorithms are very simple and general (easy to apply to fancy models)
 - They require very little memory ($O(n)$)



Sampling from a discrete distribution

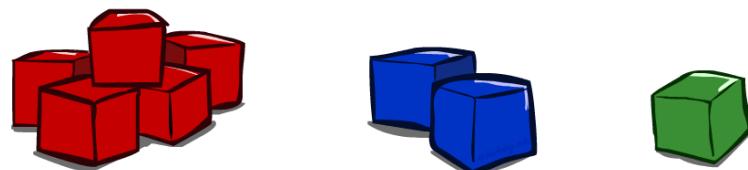
- Sampling from given distribution
 - Step 1: Get sample u from uniform distribution over $[0, 1)$
 - `Random()` in many programming languages
 - Step 2: Convert this sample u into an outcome for the given distribution by associating each outcome x with a $P(x)$ -sized sub-interval of $[0,1)$

- Example

C	P(C)
red	0.6
green	0.1
blue	0.3

$0 \leq u < 0.6, \rightarrow C = \text{red}$
 $0.6 \leq u < 0.7, \rightarrow C = \text{green}$
 $0.7 \leq u < 1, \rightarrow C = \text{blue}$

- If `random()` returns $u = 0.83$, then our sample is $C = \text{blue}$
- E.g., after sampling 8 times:



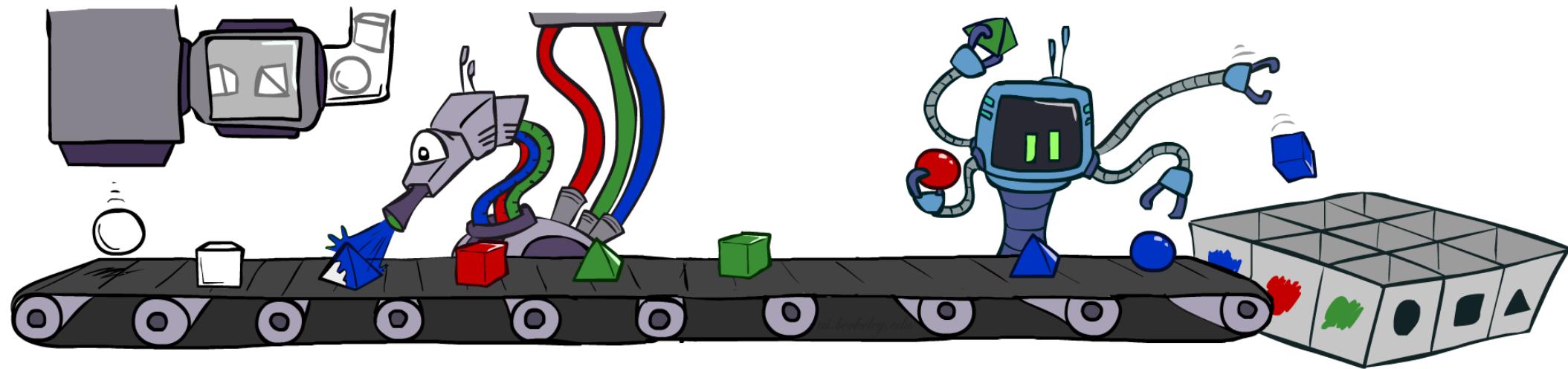
Sampling in Bayes Nets

- Prior Sampling
- Rejection Sampling
- Likelihood Weighting
- Gibbs Sampling

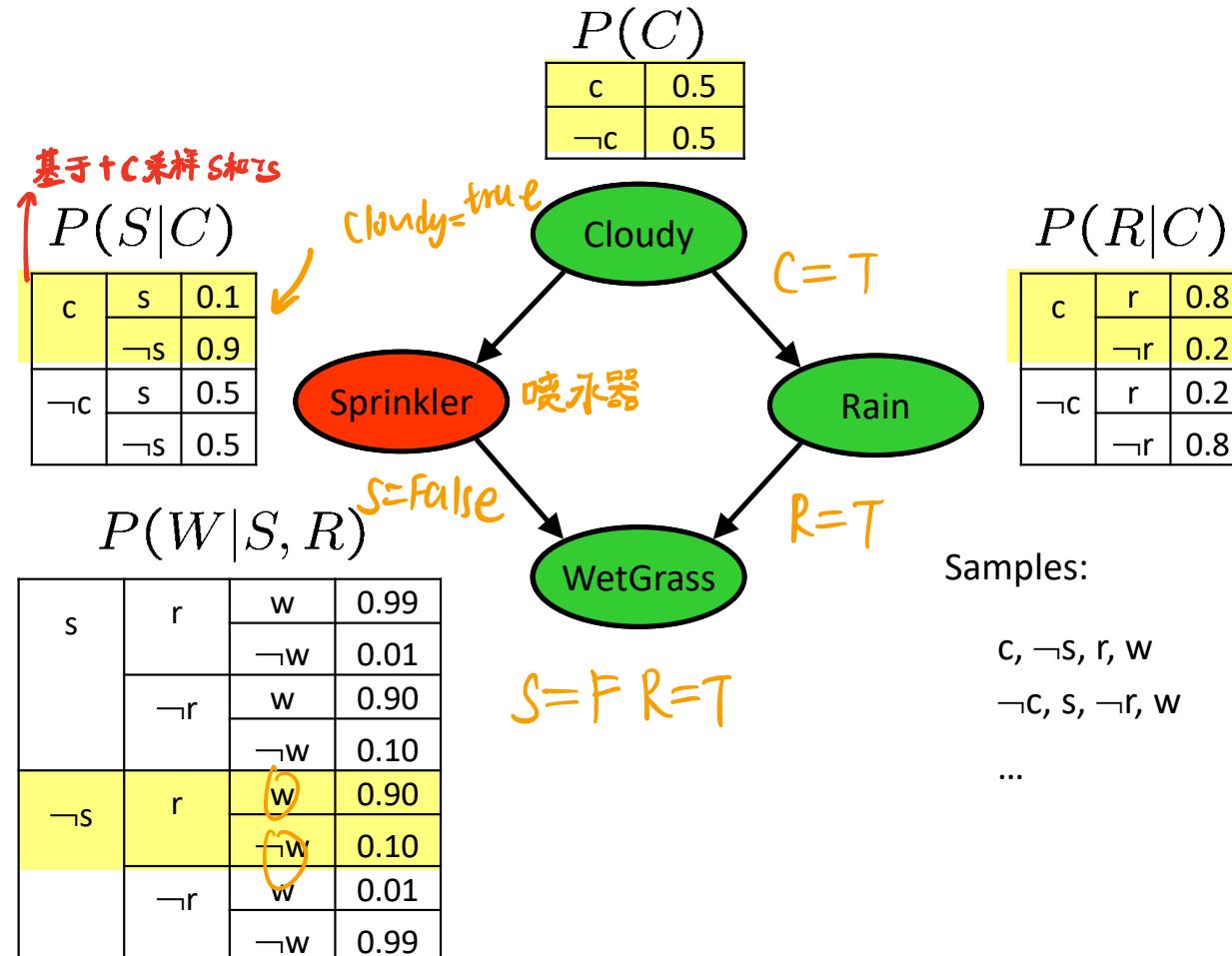
提高



Prior Sampling

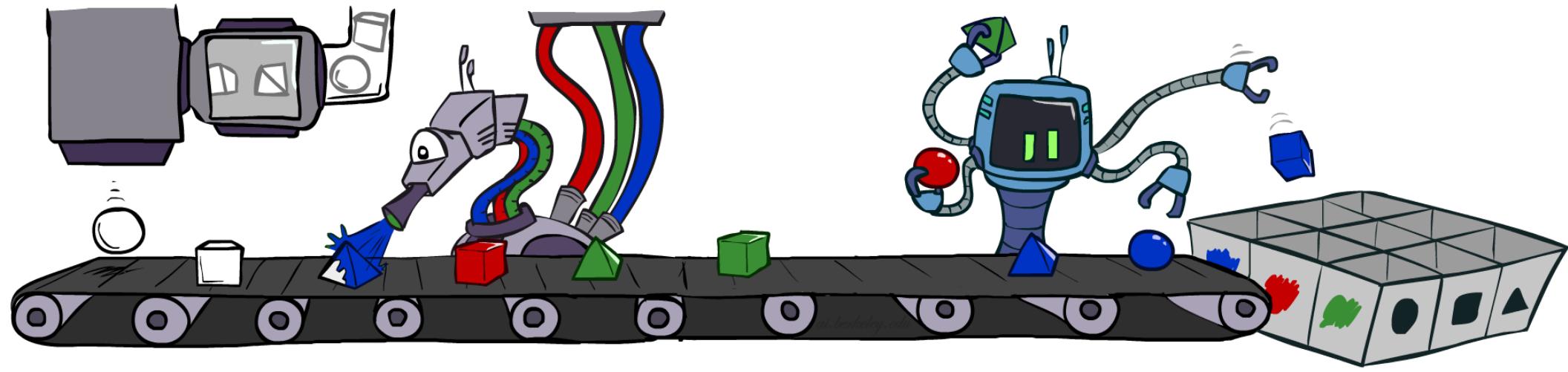


Prior Sampling



Prior Sampling

- For $i=1, 2, \dots, n$ (in topological order)
 - Sample X_i from $P(X_i | parents(X_i))$
- Return (x_1, x_2, \dots, x_n)



Using samples

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

$C, \neg S, R, W$

$\neg C, S, R, W$

$\neg C, S, R, \neg W$

$C, \neg S, R, W$

$\neg C, \neg S, \neg R, W$

- If we want to know $P(W)$
 - We have counts $\langle W:4, \neg W:1 \rangle$
 - Normalize to get $P(W) = \langle W:0.8, \neg W:0.2 \rangle$
 - This will get closer to the true distribution with more samples
- If we want to know $P(C| r, w)$
 - Count (C, r, w) and $(\neg C, r, w)$
 - Normalize to get $P(C| r, w) = \langle C:0.67, \neg C:0.33 \rangle$

$$\begin{array}{c} \nearrow \frac{2}{3} \\ \searrow \frac{1}{3} \end{array}$$

Prior Sampling

- This process generates samples with probability:

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

采样概率 *生成特定事件的概率* ...i.e. the BN's joint probability *特定事件在样本中出现次数*

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

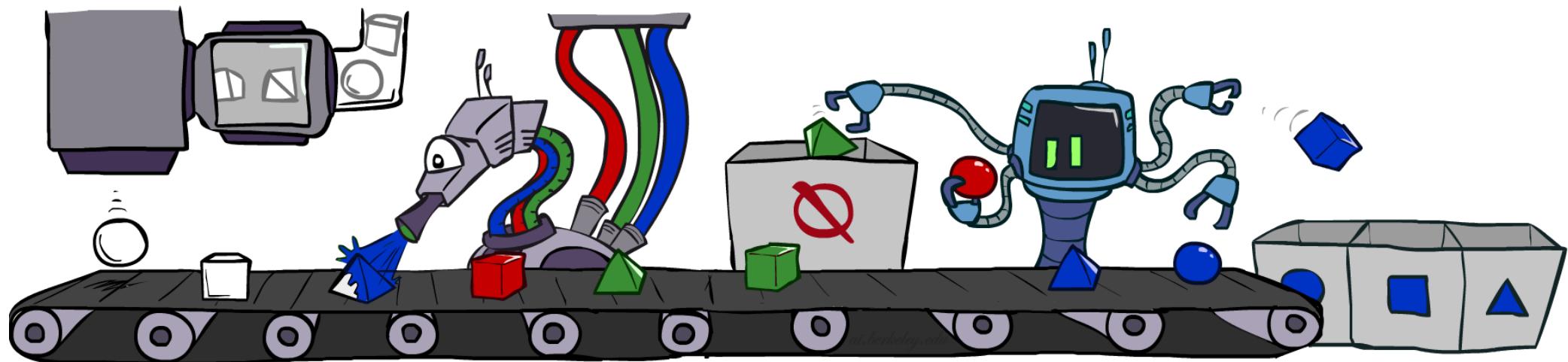
- So $\hat{P}(x_1, \dots, x_n) = N_{PS}(x_1, \dots, x_n)/N$

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

- I.e., the sampling procedure is **consistent**

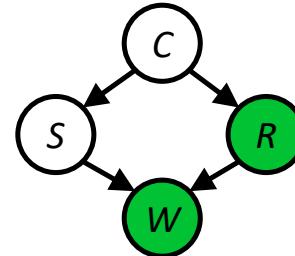
Rejection Sampling 拒绝采样

给定一个易于采样的分布,为一个难于采样的分布生成采样样本的通法



Rejection Sampling

- A simple modification of prior sampling for conditional probabilities
- Let's say we want $P(C \mid r, w)$
- When generating a sample, reject it immediately if not $R=\text{true}$, $W=\text{true}$
- It is consistent for conditional probabilities (i.e., correct in the limit)



$c, \neg s, r, w$

~~$c, s, \neg r$~~

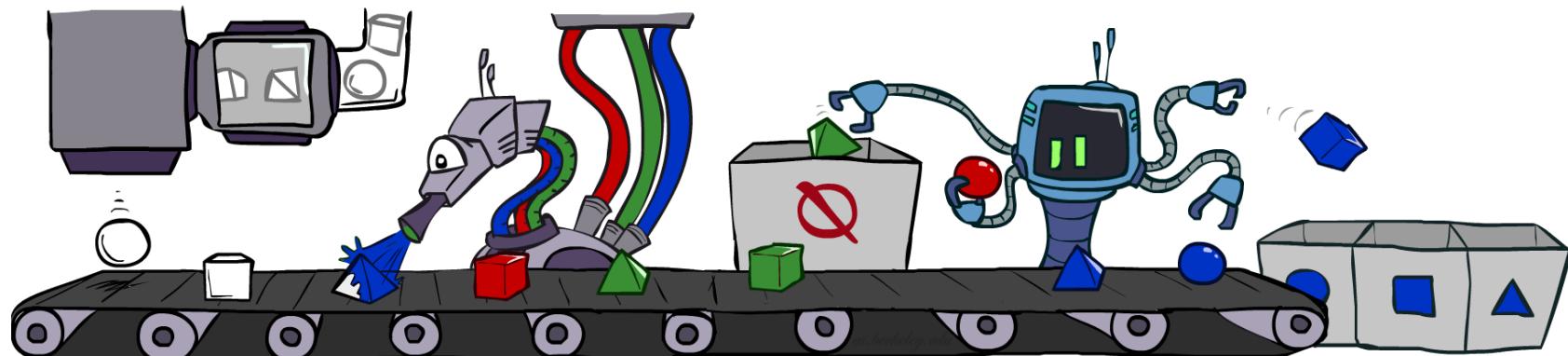
~~$\neg c, s, r, \neg w$~~

~~$c, \neg s, \neg r$~~

$\neg c, \neg s, r, w$

Rejection Sampling

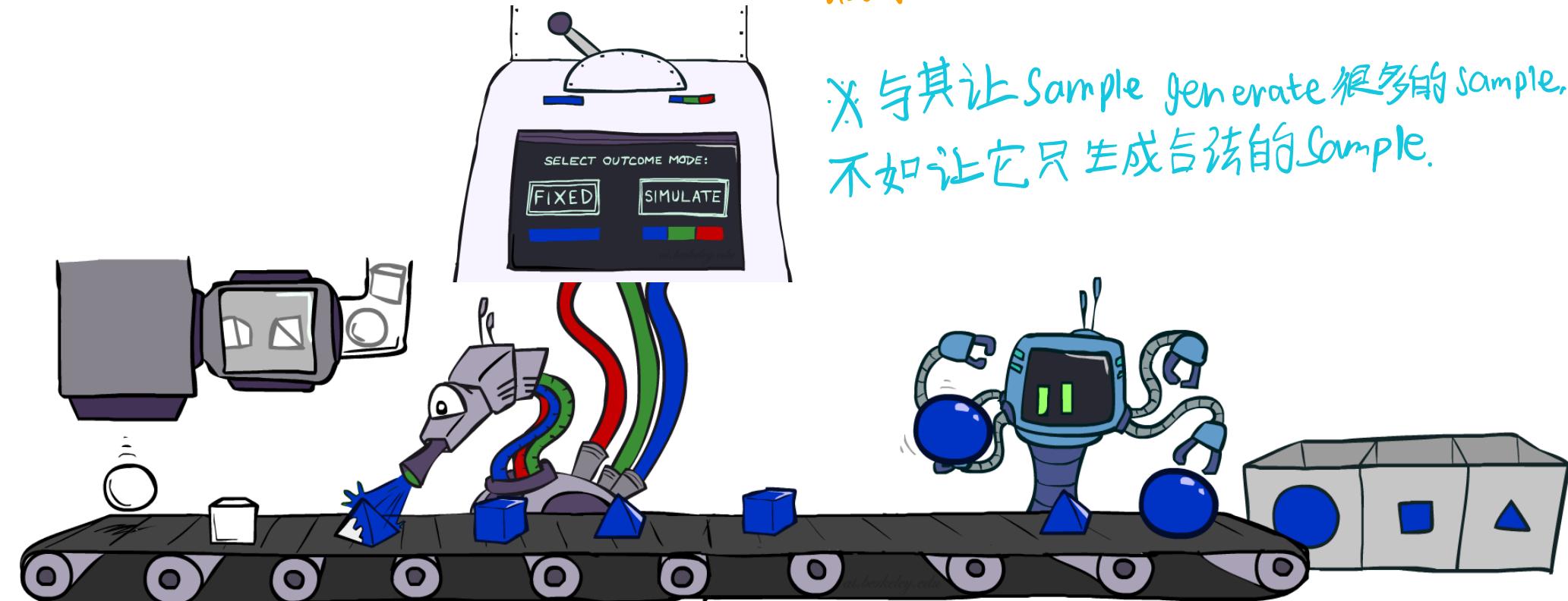
- Input: evidence e_1, \dots, e_k
- For $i=1, 2, \dots, n$
 - Sample X_i from $P(X_i | \text{parents}(X_i))$
 - If x_i not consistent with evidence
 - Reject: Return, and no sample is generated in this cycle
- Return (x_1, x_2, \dots, x_n)



Likelihood Weighting 倏然如风

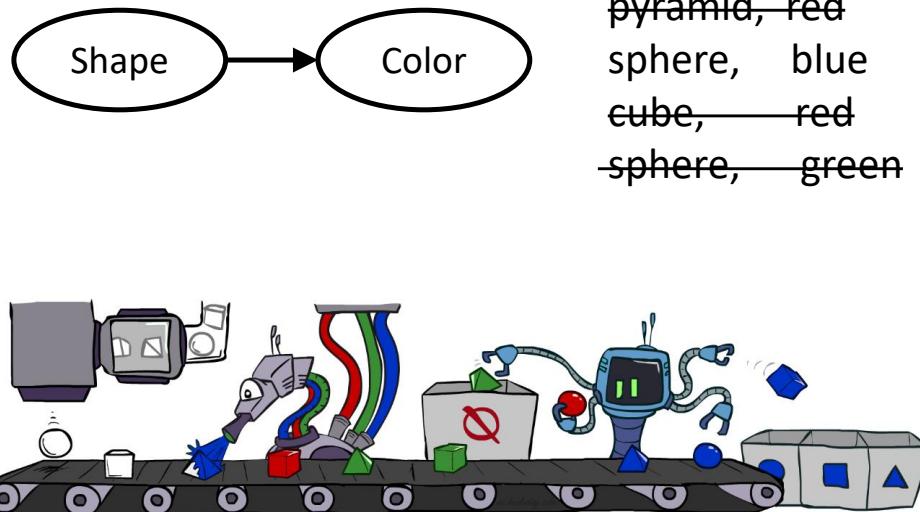
属于重要性采样的一部分

与其让 Sample generate 很多的 sample,
不如让它只生成合法的 sample.

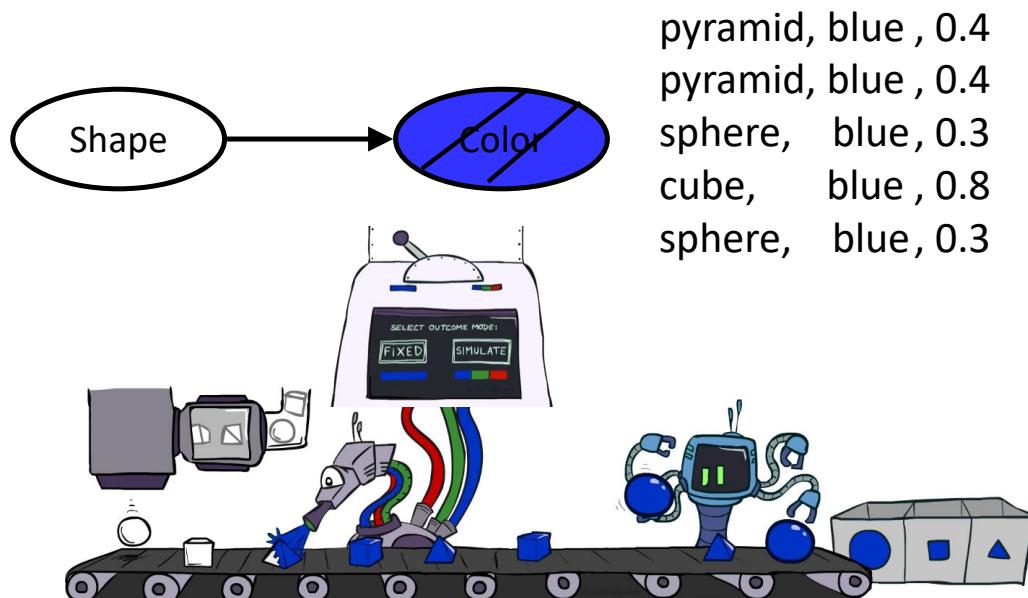


Likelihood Weighting

- Problem with rejection sampling:
 - If evidence is unlikely, rejects lots of samples
 - Evidence not exploited as you sample
 - Consider $P(\text{Shape} | \text{Color}=\text{blue})$



- Idea: fix evidence variables, sample the rest
 - Problem: sample distribution not consistent!
 - Solution: **weight** each sample by probability of evidence variables given parents



Likelihood Weighting

$$P(C|ts, tw)$$

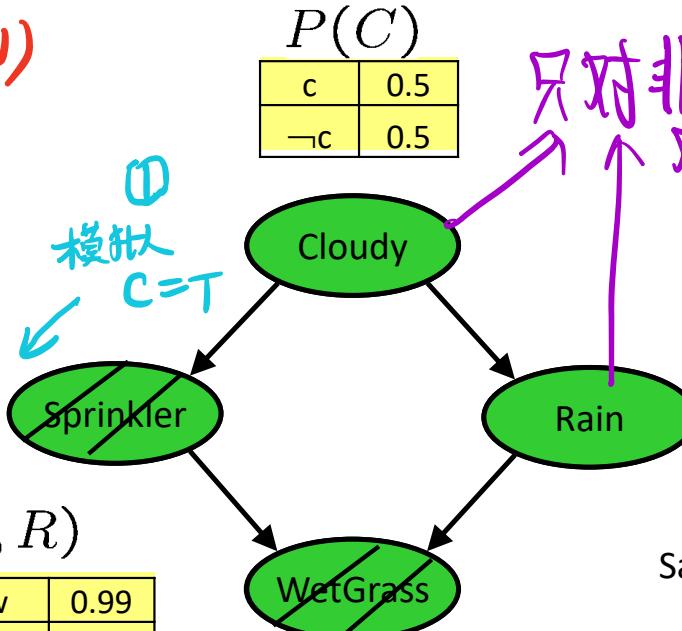
这里强制让 S 变成真的。但由于在现实生活中，我们大概率会取到 $\neg S$ 。所以强制让结果成为 S ，需要使用权重来规避掉这个问题。

$$P(S|C)$$

	s	0.1
c	s	0.9
$\neg c$	s	0.5
	$\neg s$	0.5

$$P(W|S, R)$$

s	r	w	0.99
		$\neg w$	0.01
	$\neg r$	w	0.90
	$\neg r$	$\neg w$	0.10
$\neg s$	r	w	0.90
$\neg s$	r	$\neg w$	0.10
$\neg s$	$\neg r$	w	0.01
$\neg s$	$\neg r$	$\neg w$	0.99



	P(C)
c	0.5
$\neg c$	0.5

$$P(R|C)$$

	r	0.8
c	$\neg r$	0.2
$\neg c$	r	0.2
	$\neg r$	0.8

Samples:

c, s, r, w

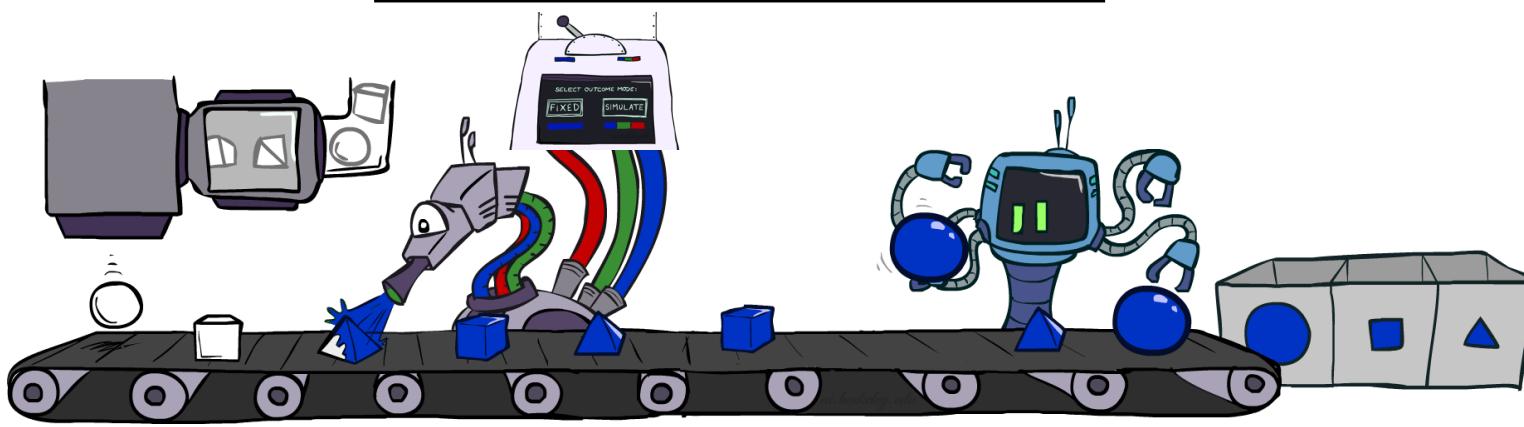
①采样了 $C:T$ $R:T$ ②加权

$$w = 1.0 \times 0.1 \times 0.99$$

$$\begin{array}{c} S:T \\ | \\ \begin{array}{c} C:T \\ | \\ \begin{array}{c} W:T \\ | \\ \begin{array}{c} R:T \\ | \\ \begin{array}{c} S:T \\ | \\ 0.1 \end{array} \end{array} \end{array} \end{array} \end{array}$$

Likelihood Weighting

- Input: evidence e_1, \dots, e_k
- $w = 1.0$
- for $i=1, 2, \dots, n$
 - if X_i is an evidence variable *对 evidence 加权*
 - x_i = observed value_i for X_i
 - Set $w = w * P(x_i | \text{Parents}(X_i))$
 - else *对非 evidence 样*
 - Sample x_i from $P(X_i | \text{Parents}(X_i))$
- return $(x_1, x_2, \dots, x_n), w$



Using samples

0-8 X

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

$c, \neg s, r, w$

0.1

$w: P(r | c) \cdot P(w | \neg s, r)$

c, s, r, w

0.2

$\neg c, s, r, w$

0.3

$c, \neg s, r, w$

0.1

$\neg c, \neg s, r, w$

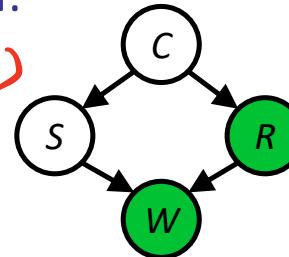
0.5

- If we want to know $P(C | r, w)$

- We have weight sums $\langle c, r, w \rangle: 0.4, (\neg c, r, w): 0.8 \rangle$

- Normalize to get $P(C | r, w) = \langle c: 0.33, \neg c: 0.67 \rangle$

- This will get closer to the true distribution with more samples



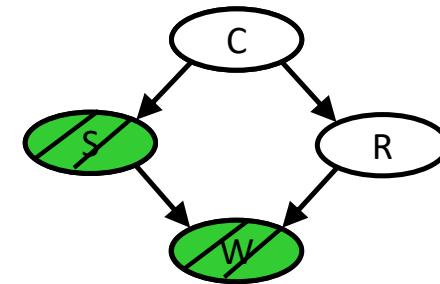
Likelihood Weighting

- Sampling distribution (\mathbf{z} is sampled and \mathbf{e} is fixed evidence)

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

- Now, samples have weights

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

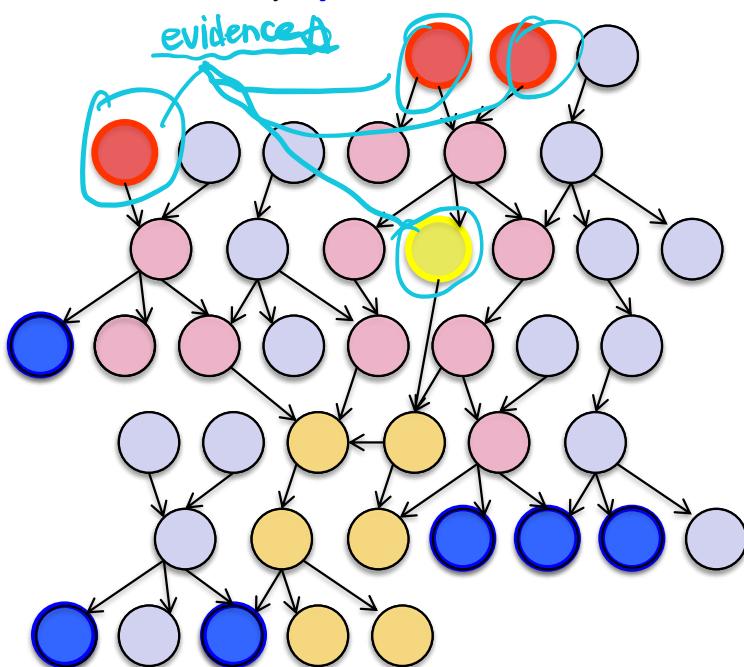


- Together, weighted sampling distribution is consistent

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

Likelihood Weighting

- Likelihood weighting is good
 - All samples are used
 - The values of *downstream* variables are influenced by *upstream* evidence

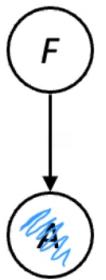


- Likelihood weighting still has weaknesses
 - The values of *upstream* variables are unaffected by *downstream* evidence 下游证据不影响上游变量
 - With many downstream evidence, we may
 - mostly get samples that are inconsistent with the evidence and thus have very small weights
 - get a few lucky samples with very large weights, which dominate the result
- We would like each variable to “see” *all* the evidence!

问题: ① 大多数上游变量采出来之后, 满足下游所有证据的样本太少。
② 即使采出符合证据的样本, 这些样本一家独大, 影响结果

Example: Fire Alarm

- In likelihood weighting, evidence influences the choice of downstream variables, but not upstream ones
- Example: $P(F|+a)$



$P(+f)$	0.001
$P(-f)$	0.999

$P(+a +f)$	0.9
$P(-a +f)$	0.1

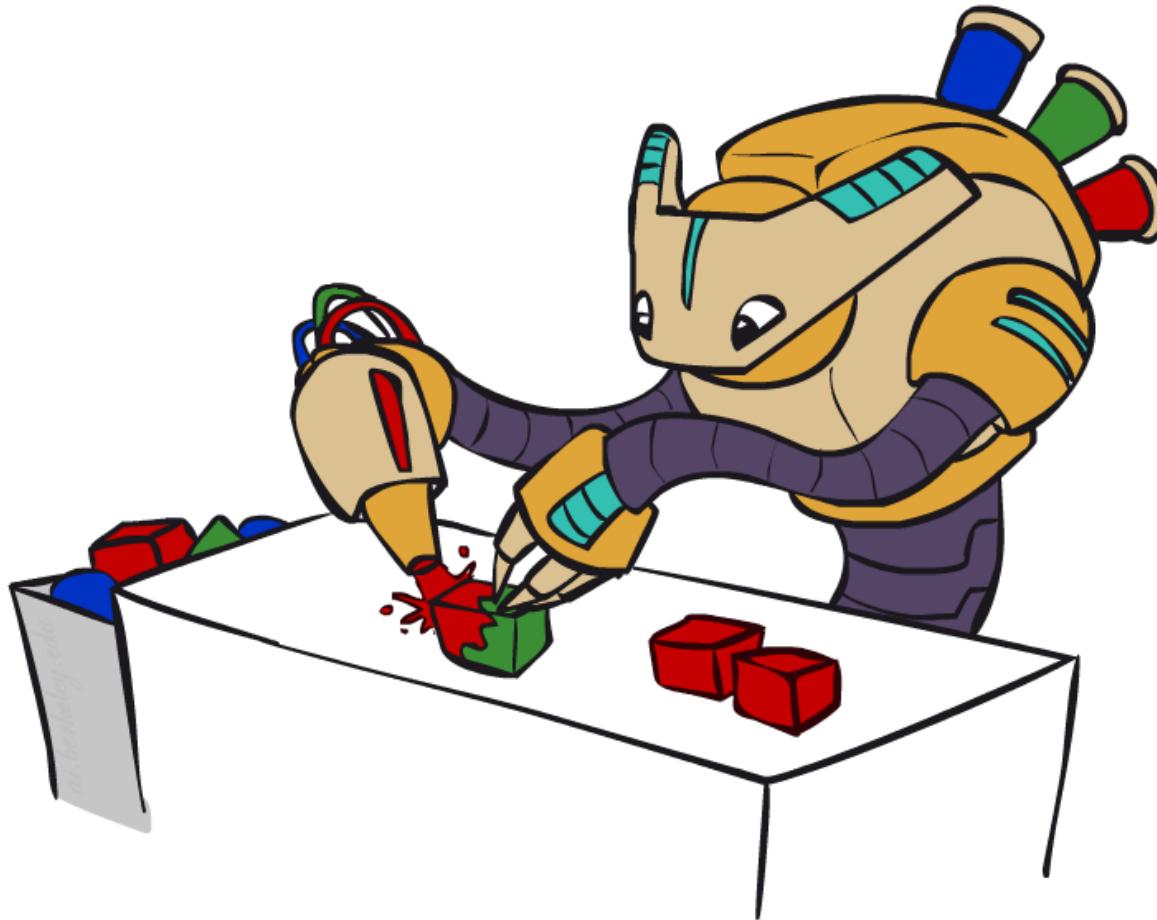
$P(+a -f)$	0.01
$P(-a -f)$	0.99

$$\begin{array}{lll} -F & +a & 0.01 \\ -F & +a & 0.01 \\ -F & +a & 0.01 \end{array} \left. \begin{array}{l} \\ \\ \end{array} \right\} 20 \text{ times}$$
$$+F & +a & 0.9 \left. \begin{array}{l} \\ \end{array} \right\} 1 \text{ time}$$

**you'll get the likelihood weighting
你会得到似然加权。**

现在有一个问题，如果该样本出现的机会特别小，那么需要等很长很长的时间。尤其是如果我的evidence在拓扑序靠后。

Gibbs Sampling



之前的采样，总是在父结点的基础上进行采样，如果父结点很深，downstream

Gibbs Sampling

① 生成一个完整样本 (down stream & up stream)

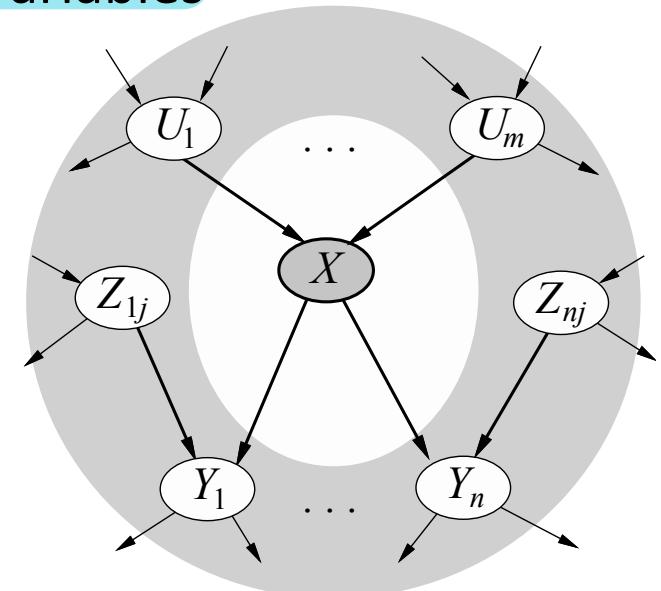
- *Procedure:* keep track of a full instantiation x_1, x_2, \dots, x_n . Start with an arbitrary instantiation consistent with the ~~evidence~~. Sample one variable at a time, conditioned on all the rest, but keep evidence fixed. Keep repeating this for a long time.
- *Property:* in the limit of repeating this infinitely many times the resulting samples come from the correct distribution (i.e. conditioned on evidence).
- *Rationale:* both upstream and downstream variables condition on evidence.
- In contrast: likelihood weighting only conditions on upstream evidence, and hence weights obtained in likelihood weighting can sometimes be very small. Sum of weights over all samples is indicative of how many “effective” samples were obtained, so we want high weight.

you pick up a single value, and you resample that value

您选取一个数值，然后重新对该数值进行抽样

Gibbs Sampling

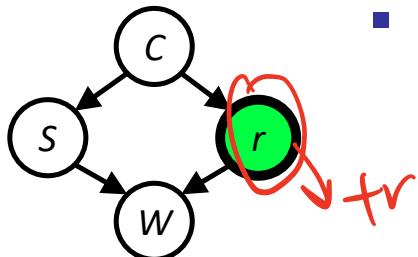
- Generate each sample by making a random change to the preceding sample
 - Evidence variables remain fixed. For each of the non-evidence variable, sample its value conditioned on all the other variables
 - $X'_i \sim P(X_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$
 - In a Bayes net
$$\begin{aligned}P(X_i | x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n) \\= P(X_i | \text{markov_blanket}(X_i)) \\= \alpha P(X_i | u_1, \dots, u_m) \prod_j P(y_j | \text{parents}(Y_j))\end{aligned}$$



Gibbs Sampling Example: $P(S | r)$

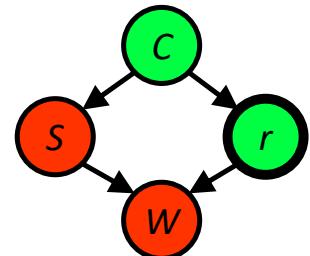
- Step 1: Fix evidence

- $R = \text{true}$



- Step 2: Initialize other variables

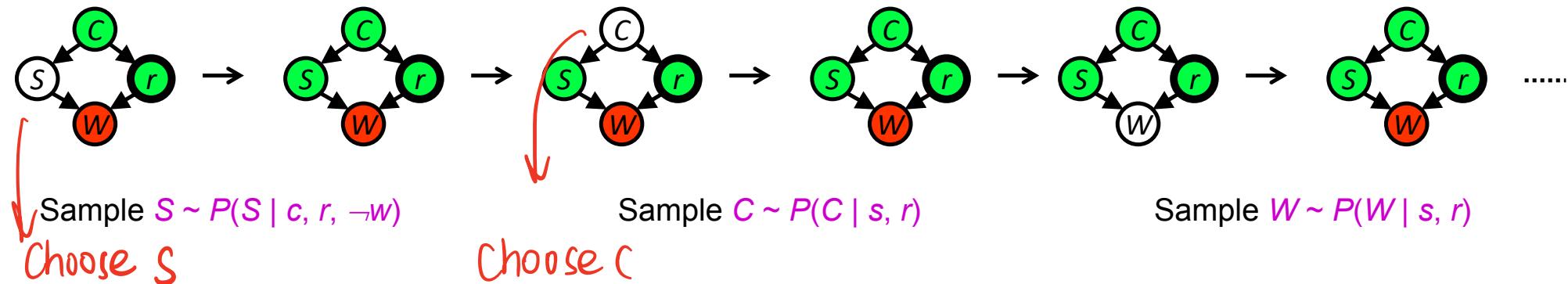
- Randomly



- Step 3: Repeat

- Choose an arbitrary non-evidence variable X
- Resample X from $P(X | \text{markov_blanket}(X))$

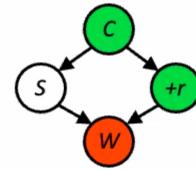
所有的样本都互相关联了。



Efficient Resampling of One Variable

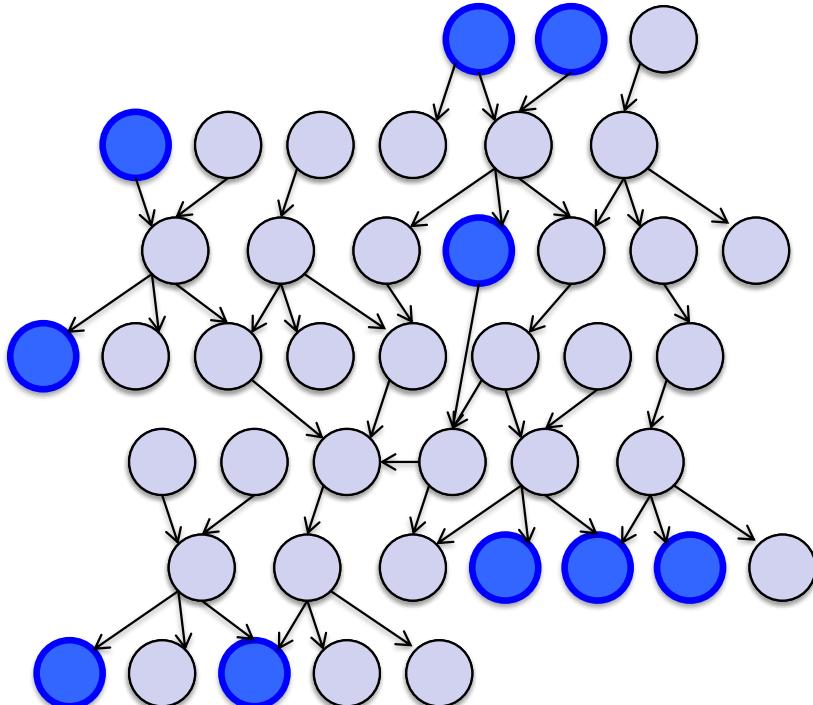
- Sample from $P(S | +c, +r, -w)$

$$\begin{aligned} P(S | +c, +r, -w) &= \frac{P(S, +c, +r, -w)}{P(+c, +r, -w)} && \text{def. of conditional probability} \\ &= \frac{P(S, +c, +r, -w)}{\sum_s P(s, +c, +r, -w)} && \text{introduce summation} \\ &= \frac{P(+c)P(S | +c)P(+r | +c)P(-w | S, +r)}{\sum_s P(+c)P(s | +c)P(+r | +c)P(-w | s, +r)} && \text{def. of Bayes' Nets} \\ &= \frac{P(+c)P(S | +c)P(+r | +c)P(-w | S, +r)}{P(+c)P(+r | +c) \sum_s P(s | +c)P(-w | s, +r)} && \text{move summation term out} \\ &= \frac{P(S | +c)P(-w | S, +r)}{\sum_s P(s | +c)P(-w | s, +r)} && \text{cancel out terms} \end{aligned}$$



So I'll walk you through it.
那我来带你一步步了解。

Why doing this?



- Samples soon begin to reflect all the evidence in the network
- Eventually they are being drawn from the true posterior!
- Theorem: Gibbs sampling is consistent

- MCN
- some
- MCN
- W
- MCN

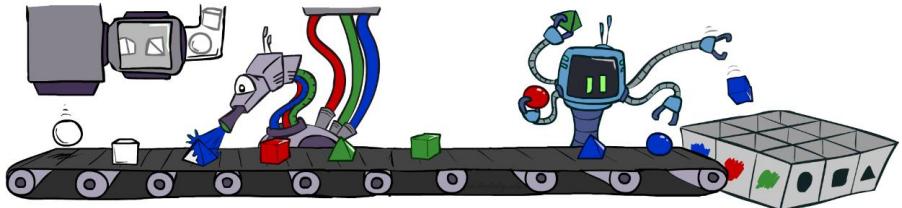


Markov Chain Monte Carlo (MCMC)

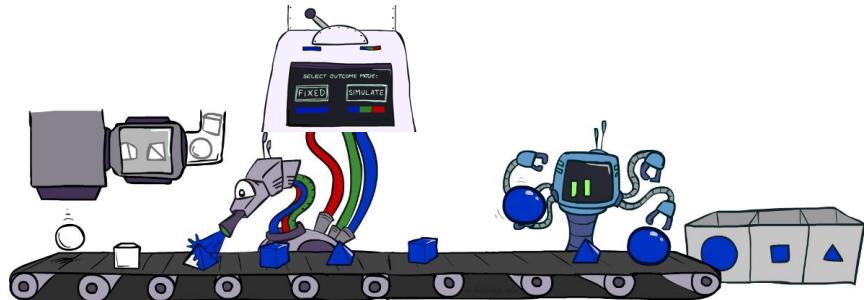
- MCMC is a family of randomized algorithms for approximating some quantity of interest over a very large state space
 - Markov chain = a sequence of randomly chosen states (“random walk”), where each state is chosen conditioned on the previous state
 - ~~Monte Carlo = a very expensive city in Monaco with a famous casino~~
 - Monte Carlo = an algorithm (usually based on sampling) that is likely to find a correct answer
- MCMC = sampling by constructing a Markov chain
- Gibbs, Metropolis-Hastings, Hamiltonian, Slice, etc.

Summary

- Prior Sampling P



- Likelihood Weighting $P(Q | e)$



- Rejection Sampling $P(Q | e)$



- Gibbs Sampling $P(Q | e)$

