Sampling for graphical models

UNIVERSITY of WASHINGTON

Goal: compute queries like $P(I = i^1 | G = g^2, S = s^0)$ or in words "the probability of high intelligence given medium grade and low SAT?"



P(L, G, S, D, I) = P(L | G)P(G | D, I)P(S | I)P(D)P(I)

Goal: compute queries like $P(I = i^1 | G = g^2, S = s^0)$ or in words "the probability of high intelligence given medium grade and low SAT?"

$$P(I = i^{1} | G = g^{2}, S = s^{0}) = \frac{P(I = i^{1}, G = g^{2}, S = s^{0})}{P(G = g^{2}, S = s^{0})}$$

Approaches:

• Exact inference by variable elimination:

$$P(L, G, S, D, I) = P(L | G)P(G | D, I)P(S | I)P(D)P(I)$$

$$P(I = i^{1}, G = g^{2}, S = s^{0}) = \sum_{l,d} P(L = l, G = g^{2}, S = s^{0}, D = d, I = i^{1})$$

Goal: compute queries like $P(I = i^1 | G = g^2, S = s^0)$ or in words "the probability of high intelligence given medium grade and low SAT?"

$$P(I = i^{1} | G = g^{2}, S = s^{0}) = \frac{P(I = i^{1}, G = g^{2}, S = s^{0})}{P(G = g^{2}, S = s^{0})}$$

Approaches:

• Exact inference by variable elimination:

$$P(L, G, S, D, I) = P(L \mid G)P(G \mid D, I)P(S \mid I)P(D)P(I)$$

$$P(I = i^{1}, G = g^{2}, S = s^{0}) = \sum_{l,d} P(L = l, G = g^{2}, S = s^{0}, D = d, I = i^{1})$$
$$= \sum_{l,d} P(L = l | G = g^{2}) P(G = g^{2} | D = d, I = i^{1}) P(S = s^{0} | I = i^{1}) P(D = d) P(I = i^{1})$$

Goal: compute queries like $P(I = i^1 | G = g^2, S = s^0)$ or in words "the probability of high intelligence given medium grade and low SAT?"

$$P(I = i^{1} | G = g^{2}, S = s^{0}) = \frac{P(I = i^{1}, G = g^{2}, S = s^{0})}{P(G = g^{2}, S = s^{0})}$$

Approaches:

• Exact inference by variable elimination:

$$P(L, G, S, D, I) = P(L | G)P(G | D, I)P(S | I)P(D)P(I)$$

$$P(I = i^{1}, G = g^{2}, S = s^{0}) = \sum_{l,d} P(L = l, G = g^{2}, S = s^{0}, D = d, I = i^{1})$$

$$= \sum_{l,d} P(L = l | G = g^{2}) P(G = g^{2} | D = d, I = i^{1}) P(S = s^{0} | I = i^{1}) P(D = d) P(I = i^{1})$$

$$= \left(\sum_{d} P(G = g^{2} | D = d, I = i^{1}) P(D = d)\right) P(S = s^{0} | I = i^{1}) P(I = i^{1})$$

$$P(G = g^2, S = s^0) = \sum_{i} \left(\sum_{d} P(G = g^2 | D = d, I = i) P(D = d) \right) P(S = s^0 | I = i) P(I = i)$$

Goal: compute queries like $P(I = i^1 | G = g^2, S = s^0)$ or in words "the probability of high intelligence given medium grade and low SAT?"

$$P(I = i^{1} | G = g^{2}, S = s^{0}) = \frac{P(I = i^{1}, G = g^{2}, S = s^{0})}{P(G = g^{2}, S = s^{0})}$$

Approaches:

- Exact inference by variable elimination
 - Straightforward, but can be computationally prohibitive
- Variational inference
 - Approximate, biased
- Sampling strategies
 - Rejection sampling
 - Importance weighted sampling
 - MCMC

Goal: compute queries like $P(I = i^1 | G = g^2, S = s^0)$ or in words "the probability of high intelligence given medium grade and low SAT?"

$$P(I = i^{1} | G = g^{2}, S = s^{0}) = \frac{P(I = i^{1}, G = g^{2}, S = s^{0})}{P(G = g^{2}, S = s^{0})}$$

Algorithm 12.1 Forward Sampling in a Bayesian network

Procedure Forward-Sample (
 \mathcal{B} // Bayesian network over \mathcal{X})1Let X_1, \ldots, X_n be a topological ordering of \mathcal{X} 2for $i = 1, \ldots, n$ 3 $u_i \leftarrow x \langle \operatorname{Pa}_{X_i} \rangle$ // Assignment to Pa_{X_i} in x_1, \ldots, x_{i-1} 4Sample x_i from $P(X_i | u_i)$ 5return (x_1, \ldots, x_n)

Goal: compute queries like $P(I = i^1 | G = g^2, S = s^0)$ or in words "the probability of high intelligence given medium grade and low SAT?"

$$P(I = i^{1} | G = g^{2}, S = s^{0}) = \frac{P(I = i^{1}, G = g^{2}, S = s^{0})}{P(G = g^{2}, S = s^{0})}$$

Draw *M* samples $\{(L_k, G_k, S_k, D_k, I_k)\}_{k=1}^M$ using forward sampling and Set $\chi = \{k : G_k = g^2, S_k = s^0\}$ and output $\widehat{P}(I = i^1 | G = g^2, S = s^0) = \frac{1}{|\chi|} \sum_{k \in \chi} \mathbf{1}\{I_k = i^1, G_k = g^2, S_k = s^0\}$

$$(I = i^1, G = g^2, S = s^0)$$
 $(G = g^2, S = s^0)$ All outcomes

Goal: compute queries like $P(I = i^1 | G = g^2, S = s^0)$ or in words "the probability of high intelligence given medium grade and low SAT?"

$$P(I = i^{1} | G = g^{2}, S = s^{0}) = \frac{P(I = i^{1}, G = g^{2}, S = s^{0})}{P(G = g^{2}, S = s^{0})}$$

 $\mathbb{E}\Big[\widehat{P}(I = i^{1} | G = g^{2}, S = s^{0}) = \mathbb{E}\Big[\frac{1}{|\chi|} \sum_{k \in \chi} \mathbf{1}\{I_{k} = i^{1}, G_{k} = g^{2}, S_{k} = s^{0}\} \\ = \sum_{n=1}^{\infty} \mathbb{E}\Big[\frac{1}{n} \sum_{k \in \chi} \mathbf{1}\{I_{k} = i^{1}, G_{k} = g^{2}, S = s^{0}) \Big] = \mathbb{E}\Big[\frac{1}{|\chi|} \sum_{k \in \chi} \mathbf{1}\{I_{k} = i^{1}, G_{k} = g^{2}, S_{k} = s^{0}\} \Big]$

 $=\sum_{n=1}^{\infty} \frac{nP(I=i^1 | G=g^2, S=s^0)}{n} \mathbb{P}(|\chi|=n) \quad \text{Unbiased estimator of } P(I=i^1 | G=g^2, S=s^0)!$

Goal: compute queries like $P(I = i^1 | G = g^2, S = s^0)$ or in words "the probability of high intelligence given medium grade and low SAT?"

$$P(I = i^{1} | G = g^{2}, S = s^{0}) = \frac{P(I = i^{1}, G = g^{2}, S = s^{0})}{P(G = g^{2}, S = s^{0})}$$

Draw *M* samples $\{(L_k, G_k, S_k, D_k, I_k)\}_{k=1}^M$ using forward sampling and Set $\chi = \{k : G_k = g^2, S_k = s^0\}$ and output $\widehat{P}(I = i^1 | G = g^2, S = s^0) = \frac{1}{|\chi|} \sum_{k \in \chi} \mathbf{1}\{I_k = i^1, G_k = g^2, S_k = s^0\}$

$$(I = i^1, G = g^2, S = s^0)$$
 $(G = g^2, S = s^0)$ All outcomes

How big goes M need to be to get an accurate estimate?

Goal: compute queries like $P(I = i^1 | G = g^2, S = s^0)$ or in words "the probability of high intelligence given medium grade and low SAT?"

$$P(I = i^{1} | G = g^{2}, S = s^{0}) = \frac{P(I = i^{1}, G = g^{2}, S = s^{0})}{P(G = g^{2}, S = s^{0})}$$

Draw *M* samples $\{(L_k, G_k, S_k, D_k, I_k)\}_{k=1}^M$ using forward sampling and Set $\chi = \{k : G_k = g^2, S_k = s^0\}$ and output $\widehat{P}(I = i^1 | G = g^2, S = s^0) = \frac{1}{|\chi|} \sum_{k \in \chi} \mathbf{1}\{I_k = i^1, G_k = g^2, S_k = s^0\}$

Rejection sampling takeaways:

- Very simple to implement
- May require an enormous amount of samples if the conditional statement is rare.
 - Consider P(disease | symptoms). Any precise set of symptoms is going to be rare.

Goal: compute queries like $P(I = i^1 | G = g^2, S = s^0)$ or in words "the probability of high intelligence given medium grade and low SAT?"

$$P(I = i^{1} | G = g^{2}, S = s^{0}) = \frac{P(I = i^{1}, G = g^{2}, S = s^{0})}{P(G = g^{2}, S = s^{0})}$$

Approaches:

- Exact inference by variable elimination
 - Straightforward, but can be computationally prohibitive
- Variational inference
 - Approximate, biased
- Sampling strategies
 - Rejection sampling 💙
 - Importance weighted sampling
 - MCMC

Fix any function $f : \mathscr{X} \to [0,1]$ and suppose we wish to estimate $\mu = \mathbb{E}_P[f(X)]$ If I draw $X_1, \dots, X_M \sim P$ and define $\hat{\mu}_P := \frac{1}{M} \sum_{i=1}^M f(X_i)$ then $\mathbb{E}_P[\hat{\mu}_P] = \mathbb{E}_P[f(X)]$

Fix any function $f : \mathscr{X} \to [0,1]$ and suppose we wish to estimate $\mu = \mathbb{E}_P[f(X)]$ If I draw $X_1, \dots, X_M \sim P$ and define $\hat{\mu}_P := \frac{1}{M} \sum_{i=1}^M f(X_i)$ then $\mathbb{E}_P[\hat{\mu}_P] = \mathbb{E}_P[f(X)]$

If I draw
$$Y_1, \dots, Y_M \sim Q$$
 and define $\hat{\mu}_Q := \frac{1}{M} \sum_{i=1}^M f(Y_i) \frac{P(Y_i)}{Q(Y_i)}$ then $\mathbb{E}_Q[\hat{\mu}_Q] = \mathbb{E}_P[f(X)]$

$$\mathbb{E}_{Q}\left[\frac{1}{M}\sum_{i=1}^{M}f(Y_{i})\frac{P(Y_{i})}{Q(Y_{i})}\right] =$$

Fix any function $f : \mathscr{X} \to [0,1]$ and suppose we wish to estimate $\mu = \mathbb{E}_P[f(X)]$ If I draw $X_1, \dots, X_M \sim P$ and define $\hat{\mu}_P := \frac{1}{M} \sum_{i=1}^M f(X_i)$ then $\mathbb{E}_P[\hat{\mu}_P] = \mathbb{E}_P[f(X)]$

If I draw
$$Y_1, ..., Y_M \sim Q$$
 and define $\hat{\mu}_Q := \frac{1}{M} \sum_{i=1}^M f(Y_i) \frac{P(Y_i)}{Q(Y_i)}$ then $\mathbb{E}_Q \Big[\hat{\mu}_Q \Big] = \mathbb{E}_P [f(X)]$
 $\mathbb{E}_Q \Big[\frac{1}{M} \sum_{i=1}^M f(Y_i) \frac{P(Y_i)}{Q(Y_i)} \Big] = \mathbb{E}_Q \Big[f(Y_1) \frac{P(Y_1)}{Q(Y_1)} \Big]$
 $= \sum_x Q(x) \cdot f(x) \frac{P(x)}{Q(x)}$
 $= \sum_x P(x) \cdot f(x)$
 $= \mathbb{E}_P [f(X)]$

Fix any function $f : \mathscr{X} \to [0,1]$ and suppose we wish to estimate $\mu = \mathbb{E}_P[f(X)]$ If I draw $X_1, \dots, X_M \sim P$ and define $\hat{\mu}_P := \frac{1}{M} \sum_{i=1}^M f(X_i)$ then $\mathbb{E}_P[\hat{\mu}_P] = \mathbb{E}_P[f(X)]$

If I draw
$$Y_1, \dots, Y_M \sim Q$$
 and define $\hat{\mu}_Q := \frac{1}{M} \sum_{i=1}^M f(Y_i) \frac{P(Y_i)}{Q(Y_i)}$ then $\mathbb{E}_Q[\hat{\mu}_Q] = \mathbb{E}_P[f(X)]$

Moreover, as
$$M \to \infty$$
 we have $\frac{1}{M} \sum_{i=1}^{M} f(Y_i) \frac{P(Y_i)}{Q(Y_i)} \sim \mathcal{N}\left(\mathbb{E}_P[f(X)], \sigma_Q^2/M\right)$ where

$$\sigma_Q^2 = \mathbb{E}_Q\left[\left(f(Y) \frac{P(Y)}{Q(Y)}\right)^2\right] - \mathbb{E}_Q[f(Y) \frac{P(Y)}{Q(Y)}]^2 \qquad \sigma_Q^2 \text{ is minimizes when}$$

$$= \mathbb{E}_P[f(X)^2 \frac{P(X)}{Q(X)}] - \mathbb{E}_P[f(X)]^2 \qquad Q(x) \propto |f(x)| P(x)$$

Importance sampling with candidate distribution Q:

If I draw
$$Y_1, \ldots, Y_M \sim Q$$
 and define $\hat{\mu}_Q := \frac{1}{M} \sum_{i=1}^M f(Y_i) \frac{P(Y_i)}{Q(Y_i)}$ then $\mathbb{E}_Q \left[\hat{\mu}_Q \right] = \mathbb{E}_P[f(X)]$



Example: rare event sampling

If I draw
$$Y_1, \ldots, Y_M \sim Q$$
 and define $\hat{\mu}_Q := \frac{1}{M} \sum_{i=1}^M f(Y_i) \frac{P(Y_i)}{Q(Y_i)}$ then $\mathbb{E}_Q[\hat{\mu}_Q] = \mathbb{E}_P[f(X)]$



Self-normalized Importance sampling

Fix any function $f : \mathscr{X} \to [0,1]$ and suppose we wish to estimate $\mu = \mathbb{E}_P[f(X)]$ If I draw $X_1, \ldots, X_M \sim P$ and define $\hat{\mu}_P := \frac{1}{M} \sum_{i=1}^M f(X_i)$ then $\mathbb{E}_P[\hat{\mu}_P] = \mathbb{E}_P[f(X)]$

Importance sampling with candidate distribution Q:

If I draw
$$Y_1, \dots, Y_M \sim Q$$
 and define $\hat{\mu}_Q := \frac{1}{M} \sum_{i=1}^M f(Y_i) \frac{P(Y_i)}{Q(Y_i)}$ then $\mathbb{E}_Q[\hat{\mu}_Q] = \mathbb{E}_P[f(X)]$

Self-normalized Importance sampling with candidate distribution Q:

If I draw
$$Y_1, \dots, Y_M \sim Q$$
 and define $\hat{\mu}_Q^{sn} := \frac{\sum_{i=1}^M f(Y_i) \frac{P(Y_i)}{Q(Y_i)}}{\sum_{i=1}^M \frac{P(Y_i)}{Q(Y_i)}}$ then $\mathbb{E}_Q\left[\hat{\mu}_Q\right] = \mathbb{E}_P[f(X)]$

Biased $\mathbb{E}_{Q}[\hat{\mu}_{Q}^{sn}] \neq \mu$ but is asymptotically consistent since $\mathbb{E}_{Q}[\frac{P(Y_{i})}{Q(Y_{i})}] = 1$

Self-normalized Importance sampling

Fix any function $f : \mathscr{X} \to [0,1]$ and suppose we wish to estimate $\mu = \mathbb{E}_P[f(X)]$ If I draw $X_1, \ldots, X_M \sim P$ and define $\hat{\mu}_P := \frac{1}{M} \sum_{i=1}^M f(X_i)$ then $\mathbb{E}_P[\hat{\mu}_P] = \mathbb{E}_P[f(X)]$

Importance sampling with candidate distribution Q:

If I draw
$$Y_1, \dots, Y_M \sim Q$$
 and define $\hat{\mu}_Q := \frac{1}{M} \sum_{i=1}^M f(Y_i) \frac{P(Y_i)}{Q(Y_i)}$ then $\mathbb{E}_Q[\hat{\mu}_Q] = \mathbb{E}_P[f(X)]$

Self-normalized Importance sampling with candidate distribution Q:

If I draw
$$Y_1, \dots, Y_M \sim Q$$
 and define $\widehat{\mu}_Q^{sn} := \frac{\sum_{i=1}^M f(Y_i) \frac{P(Y_i)}{Q(Y_i)}}{\sum_{i=1}^M \frac{P(Y_i)}{Q(Y_i)}}$ then $\mathbb{E}_Q[\widehat{\mu}_Q] = \mathbb{E}_P[f(X)]$

If we only know P up to a normalizing constant such that we have $\widetilde{P}(x) = ZP(x)$, $\widehat{\mu}_{Q}^{sn} = \frac{\sum_{i=1}^{M} f(Y_i) \frac{\widetilde{P}(Y_i)}{Q(Y_i)}}{\sum_{i=1}^{M} \frac{\widetilde{P}(Y_i)}{Q(Y_i)}} = \frac{\sum_{i=1}^{M} f(Y_i) \frac{P(Y_i)}{Q(Y_i)}}{\sum_{i=1}^{M} \frac{\widetilde{P}(Y_i)}{Q(Y_i)}}$



$$\begin{split} P(X) &:= P(L, G, S, D, I) \\ &= P(L \mid G) P(G \mid D, I) P(S \mid I) P(D) P(I) \end{split}$$

If X := (L, G, S, D, I) and P(X) denotes a Bayesian network then as we saw earlier, sampling from conditionals of P directly is awkward and inefficient. Can we define a convenient Q and use important sampling?

Original network \mathscr{B}





 $P(X) := P(L, G, S, D, I) \qquad Q_{I=i^1, G=g^2}(Y) = P(L \mid G = g^2)P(S \mid I = i^1)P(D)$ = $P(L \mid G)P(G \mid D, I)P(S \mid I)P(D)P(I)$

If X := (L, G, S, D, I) and P(X) denotes a Bayesian network then as we saw earlier, sampling from conditionals of P directly is awkward and inefficient. Can we define a convenient Q and use important sampling?

Original network \mathscr{B}



"Mutilated" network $\mathscr{B}_{I=i^1,G=g^2}$



P(X) := P(L, G, S, D, I) = P(L|G)P(G|D, I)P(S|I)P(D)P(I) $Q_{I=i^{1},G=g^{2}}(Y) = P(L|G=g^{2})P(S|I=i^{1})P(D)$ Note that $P(I=i^{1},G=g^{2}) = \mathbb{E}_{-}[1\{I=i^{1},G=g^{2}\}]$

Note that
$$P(I = i^1, G = g^2) = \mathbb{E}_P[\mathbf{1}\{I = i^1, G = g^2\}]$$

= $\mathbb{E}_Q[\frac{1}{M}\sum_{k=1}^M \mathbf{1}\{I_k = i^1, G_k = g^2\}\frac{P(Y_k)}{Q(Y_k)}]$

If $\{I = i^1, G = g^2\}$ is a rare event, this could require substantial fewer samples!

How do we compute
$$\frac{P(Y)}{Q(Y)}$$
? Let $Y = (L = l, g = g^2, S = s, D = d, I = i^1)$

Then
$$\frac{P(Y)}{Q(Y)} = \frac{P(L=l, g=g^2, S=s, D=d, I=l^2)}{Q_{I=i^1, G=g^2}(L=l, g=g^2, S=s, D=d, I=i^1)}$$

 $P(X) := P(L, G, S, D, I) \qquad Q_{I=i^{1}, G=g^{2}}(Y) = P(L | G = g^{2})P(S | I = i^{1})P(D)$ = P(L | G)P(G | D, I)P(S | I)P(D)P(I)

Note that
$$P(I = i^1, G = g^2) = \mathbb{E}_P[\mathbf{1}\{I = i^1, G = g^2\}]$$

= $\mathbb{E}_Q[\frac{1}{M}\sum_{k=1}^M \mathbf{1}\{I_k = i^1, G_k = g^2\}\frac{P(Y_k)}{Q(Y_k)}]$

If $\{I = i^1, G = g^2\}$ is a rare event, this could require substantial fewer samples!

How do we compute
$$\frac{P(Y)}{Q(Y)}$$
? Let $Y = (L = l, g = g^2, S = s, D = d, I = i^1)$
Then $\frac{P(Y)}{Q(Y)} = \frac{P(L = l, g = g^2, S = s, D = d, I = i^1)}{Q_{I=i^1,G=g^2}(L = l, g = g^2, S = s, D = d, I = i^1)}$
 $= \frac{P(G = g^2 | D = d, I = i^1)P(I = i^1)}{1}$ Note: all but the "mutilated" terms cancel.

Note that
$$P(I = i^1, G = g^2) = \mathbb{E}_P[\mathbf{1}\{I = i^1, G = g^2\}]$$

= $\mathbb{E}_Q[\frac{1}{M}\sum_{k=1}^M \mathbf{1}\{I_k = i^1, G_k = g^2\}\frac{P(Y_k)}{Q(Y_k)}]$

If $\{I = i^1, G = g^2\}$ is a rare event, this could require substantial fewer samples!

Algorithm 12.2 Likelihood-weighted particle generation

Procedure LW-Sample (\mathcal{B} , // Bayesian network over \mathcal{X} Z = z // Event in the network Let X_1, \ldots, X_n be a topological ordering of \mathcal{X} 1 2 $w \leftarrow 1$ for i = 1, ..., n3 $oldsymbol{u}_i \leftarrow oldsymbol{x} \langle \mathrm{Pa}_{X_i}
angle$ // Assignment to Pa_{X_i} in x_1, \ldots, x_{i-1} 4 if $X_i \notin \mathbf{Z}$ then 5 Sample x_i from $P(X_i \mid \boldsymbol{u}_i)$ 6 else 7 $x_i \leftarrow \boldsymbol{z} \langle X_i \rangle$ // Assignment to X_i in \boldsymbol{z} 8 $w \leftarrow w \cdot P(x_i \mid u_i)$ // Multiply weight by probability of desired value 9 return $(x_1,\ldots,x_n), w$ 10

Theorem: If the above algorithm is run on the mutilated network wrt to a set of set variables **z** and Q(Y) represents its probability, then $w = \frac{P(Y)}{Q(Y)}$.

Algorithm to estimate $P(\chi(X) = z)$:

1. Execute algorithm on previous slide M times with event $\{\chi(X) = z\}$ 2. Get $(Y_1, w_1), \dots, (Y_M, w_M)$ back 3. Set $\widehat{P}(\chi(X) = z) = \frac{1}{M} \sum_{k=1}^M \mathbf{1}\{\chi(Y_k) = z\} w_k$

Takeaway: $\widehat{P}(\chi(X) = z)$ is an unbiased estimator of $P(\chi(X) = z)$. By the properties of importance sampling estimators of the previous slides, if Q is close to P, the critical size of M may not even depend on $P(\chi(X) = z)$, even if it is very rare. But if Q is very far, this could be worse than rejection sampling.

Theorem: If the above algorithm is run on the mutilated network wrt to a set of set variables **z** and Q(Y) represents its probability, then $w = \frac{P(Y)}{Q(Y)}$.

How do we use estimators like $\widehat{P}(\chi(X) = z) = \frac{1}{M} \sum_{k=1}^{M} \mathbf{1}\{\chi(Y_k) = z\} w_k$ to compute conditional queries like $P(I = i^1 | G = g^2, S = s^0)$?

Ratio method

$$P(I = i^{1} | G = g^{2}, S = s^{0}) = \frac{P(I = i^{1}, G = g^{2}, S = s^{0})}{P(G = g^{2}, S = s^{0})}$$

Use M samples to compute

 $\widehat{P}(I=i^1, G=g^2, S=s^0)$

Use M' samples to compute

 $\widehat{P}(G=g^2,S=s^0)$

and take the ratio!

Numerator and denominator unbiased.

Self-normalized method

Collect data $(Y_1, w_1), \dots, (Y_M, w_M)$ using event $\{G = g^2, S = s^0\}$

Output $\frac{\sum_{k=1}^{M} \mathbf{1} \{I_k = i^1, G_k = g^2, S_k = s^0\} w_k}{\sum_{k=1}^{M} w_k}$

Numerator and denominator unbiased: $\mathbb{E}_{Q}[w_{k}] = \mathbb{E}_{Q}[w_{k}\mathbf{1}\{G_{k} = g^{2}, S_{k} = s^{0}\}]$ $= P(G = g^{2}, S = s^{0})$

Goal: compute queries like $P(I = i^1 | G = g^2, S = s^0)$ or in words "the probability of high intelligence given medium grade and low SAT?"

$$P(I = i^{1} | G = g^{2}, S = s^{0}) = \frac{P(I = i^{1}, G = g^{2}, S = s^{0})}{P(G = g^{2}, S = s^{0})}$$

Approaches:

- Exact inference by variable elimination
 - Straightforward, but can be computationally prohibitive
- Variational inference
 - Approximate, biased
- Sampling strategies
 - Rejection sampling 💙
 - Importance weighted sampling V
 - MCMC

Markov Chain Monte Carlo (MCMC)

Goal: compute queries like $P(I = i^1 | G = g^2, S = s^0)$ or in words "the probability of high intelligence given medium grade and low SAT?"

$$P(I = i^{1} | G = g^{2}, S = s^{0}) = \frac{P(I = i^{1}, G = g^{2}, S = s^{0})}{P(G = g^{2}, S = s^{0})}$$

Sampling from the "mutilated" network will give you sample such that $\{G = g^2, s = s^0\}$. The problem is that this sample is drawn from $Q_{I=i^1,G=g^2}(Y)$ and not $P(I = i^1 | G = g^2, S = s^0)$. Importance sampling tries to "correct" the misalignment by weighting the sample appropriately.

MCMC starts with the candidate distribution $\pi^{(0)} := Q$ and then *evolves* the distribution $\pi^{(t)}$ slowly over time until it converges $\pi^{(t)} \to \pi_* := P$. Sampling from this distribution results in a sample from P

Gibbs sampling

Algorithm 12.4 Generating a Gibbs chain trajectory **Procedure** Gibbs-Sample (X // Set of variables to be sampled Φ // Set of factors defining P_{Φ} $P^{(0)}(\boldsymbol{X})$, // Initial state distribution T // Number of time steps $P(X_i | x_{-i}) = P(X_i | \text{MarkovBlanket}(i))$ Sample $\boldsymbol{x}^{(0)}$ from $P^{(0)}(\boldsymbol{X})$ 1 2 for t = 1, ..., T $= P(X_i | parents(i))$ $P(x_i | parents(j))$ $\boldsymbol{x}^{(t)} \leftarrow \boldsymbol{x}^{(t-1)}$ 3 $i \in child(i)$ for each $X_i \in \mathbf{X}$ 4 Sample $x_i^{(t)}$ from $P_{\Phi}(X_i \mid \boldsymbol{x}_{-i})$ 5 // Change X_i in $\boldsymbol{x}^{(t)}$ 6 return $\boldsymbol{x}^{(0)},\ldots,\boldsymbol{x}^{(T)}$ 7

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



Gibbs sampling



Markov Chain Monte Carlo (MCMC)

Gibb's sampling is just one example of a MCMC algorithm.

In general, the approach:

- Construct a Markov chain transition kernel $q:\mathcal{X}\to\mathcal{X}$ whose stationary distribution is equal to P
- Start with a realization $x^{(0)} \sim \pi^{(0)}$ drawn from an arbitrary starting distribution $\pi^{(0)}$
- Run Markov chain to evolve $x^{(t)} \mapsto x^{(t+1)}$ (equiv. $\pi^{(t)} \to \pi^{(t+1)}$) and return sample once convergence $\pi^{(t)} \approx P$

Easy to see Gives sampling defines a Markov chain with stationary distribution P. Is it guaranteed to converge? How fast does it converge? What are other chains?