

Approximate Inference 2: Importance Sampling

1

(Unnormalized) Importance Sampling

(Unnormalized) Importance Sampling

- Likelihood weighting is a special case of a general approach called **importance sampling**
- Let \mathbf{X} be a set of variables that takes on values in some space $Val(\mathbf{X})$
- Importance sampling is a way to estimate $E_{P(\mathbf{x})}[f(\mathbf{x})]$ ie. the expectation of a function $f(\mathbf{x})$ relative to some distribution $P(\mathbf{X})$, typically called the **target distribution**

(Unnormalized) Importance Sampling

- Generate samples $\mathbf{x}[1], \dots, \mathbf{x}[M]$ from P
- Then estimate:

$$E_P[f] \approx \frac{1}{M} \sum_{m=1}^M f(\mathbf{x}[m])$$

(Unnormalized) Importance Sampling

- Sometimes you might want to generate samples from a different distribution Q (called a **proposal distribution** or **sampling distribution**)
- Why?
 - Might be impossible or computationally expensive to sample from P
- Proposal distribution can be arbitrary
 - Require that $Q(\mathbf{x}) > 0$ whenever $P(\mathbf{x}) > 0$
 - But computational performance of importance sampling depends strongly on how similar Q is to P

(Unnormalized) Importance Sampling

How to use the proposal distribution:

$$\begin{aligned} E_{Q(\mathbf{x})} \left[f(\mathbf{X}) \frac{P(\mathbf{X})}{Q(\mathbf{X})} \right] &= \sum_{\mathbf{x}} Q(\mathbf{x}) f(\mathbf{x}) \frac{P(\mathbf{x})}{Q(\mathbf{x})} \\ &= \sum_{\mathbf{x}} f(\mathbf{x}) P(\mathbf{x}) = E_{P(\mathbf{x})} [f(\mathbf{X})] \end{aligned}$$

Generate a set of samples $\mathcal{D} = \{\mathbf{x}[1], \dots, \mathbf{x}[M]\}$ from Q then estimate:

$$\hat{E}_D(f) = \frac{1}{M} \sum_{m=1}^M f(\mathbf{x}[m]) \frac{P(\mathbf{x}[m])}{Q(\mathbf{x}[m])}$$

Unnormalized
importance sampling
estimator

(Unnormalized) Importance Sampling

This estimator is **unbiased**:

$$\begin{aligned} E_D[\hat{E}_D(f)] &= E_{Q(\mathbf{x})} \left[f(\mathbf{X}) \frac{P(\mathbf{X})}{Q(\mathbf{X})} \right] \\ &= E_{Q(\mathbf{x})} [f(\mathbf{X}) w(\mathbf{X})] = E_{P(\mathbf{x})} [f(\mathbf{X})] \end{aligned}$$

Normalized Importance Sampling

Normalized Importance Sampling

- Frequently, P is known only up to a normalizing constant Z ie.

$$\tilde{P}(\mathbf{X}) = ZP(\mathbf{X})$$

- Happens when:
 - We know $P(\mathbf{X}, \mathbf{e})$ but need $P(\mathbf{X}|\mathbf{e})$
 - We have the unnormalized product of clique potentials for a Markov network

Normalized Importance Sampling

- Define

$$w(\mathbf{X}) = \frac{\tilde{P}(\mathbf{X})}{Q(\mathbf{X})}$$

- The expected value of the $w(\mathbf{X})$ under $Q(\mathbf{X})$ is

$$E_{Q(\mathbf{X})}[w(\mathbf{X})] = \sum_{\mathbf{x}} Q(\mathbf{x}) \frac{\tilde{P}(\mathbf{x})}{Q(\mathbf{x})} = \sum_{\mathbf{x}} \tilde{P}(\mathbf{x}) = Z$$

Normalized Importance Sampling

$$\begin{aligned} E_{P(\mathbf{X})}[f(\mathbf{X})] &= \sum_{\mathbf{x}} P(\mathbf{x}) f(\mathbf{x}) \\ &= \sum_{\mathbf{x}} Q(\mathbf{x}) f(\mathbf{x}) \frac{P(\mathbf{x})}{Q(\mathbf{x})} \\ &= \frac{1}{Z} \sum_{\mathbf{x}} Q(\mathbf{x}) f(\mathbf{x}) \frac{\tilde{P}(\mathbf{x})}{Q(\mathbf{x})} \quad (\text{since } P(\mathbf{x}) = \frac{1}{Z} \tilde{P}(\mathbf{x})) \\ &= \frac{1}{Z} E_{Q(\mathbf{X})}[f(\mathbf{X}) w(\mathbf{X})] \\ &= \frac{E_{Q(\mathbf{X})}[f(\mathbf{X}) w(\mathbf{X})]}{E_{Q(\mathbf{X})}[w(\mathbf{X})]} \end{aligned}$$

Normalized Importance Sampling

With M samples $\mathcal{D} = \{\mathbf{x}[1], \dots, \mathbf{x}[M]\}$ from Q , we can estimate:

$$\hat{E}_{\mathcal{D}}(f) = \frac{\sum_{m=1}^M f(\mathbf{x}[m]) w(\mathbf{x}[m])}{\sum_{m=1}^M w(\mathbf{x}[m])}$$

This is called the **normalized importance sampling estimator** or **weighted importance sampling estimator**

Normalized Importance Sampling

- Normalized importance sampling estimator is **biased**
- But has smaller variance than the unnormalized estimator
- Normalized estimator often used instead of unnormalized estimator, even when P is known and can be sampled from effectively

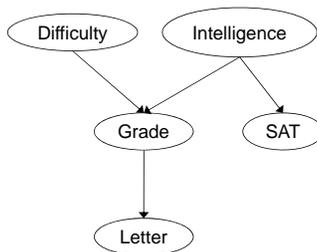
Importance Sampling for Bayesian Networks

Importance Sampling for Bayesian Networks

Student Example

D	P(D)
low	0.6
high	0.4

D	I	G	P(G D,I)
low	low	C	0.3
low	low	B	0.4
low	low	A	0.3
low	high	C	0.02
low	high	B	0.08
low	high	A	0.9
high	low	C	0.7
high	low	B	0.25
high	low	A	0.05
high	high	C	0.2
high	high	B	0.3
high	high	A	0.5



I	P(I)
low	0.7
high	0.3

I	S	P(S I)
low	low	0.95
low	high	0.05
high	low	0.2
high	high	0.8

G	L	P(L G)
C	weak	0.99
C	strong	0.01
B	weak	0.4
B	strong	0.6
A	weak	0.1
A	strong	0.9

Importance Sampling for Bayesian Networks

- What proposal distribution do we use?
- Suppose we want an event $Grade=B$ either as a query or as evidence
 - Easy to sample $P(Letter | Grade = B)$
 - Difficult to account for $Grade=B$'s influence on $Difficulty$, $Intelligence$ and SAT
- In general:
 - Want to account for effect of the event on the descendants
 - But avoid accounting for its effects on the nondescendants

Importance Sampling for Bayesian Networks

- Let B be a network, and $Z_1 = z_1, \dots, Z_k = z_k$, abbreviated $\mathbf{Z}=\mathbf{z}$, an instantiation of variables.
- We define the **mutilated network** $B_{\mathbf{Z}=\mathbf{z}}$ as follows:
 - Each node $Z_i \in \mathbf{Z}$ has no parents in $B_{\mathbf{Z}=\mathbf{z}}$
 - The CPD of Z_i in $B_{\mathbf{Z}=\mathbf{z}}$ gives probability 1 to $Z_i = z_i$ and probability 0 to all other values $z_i' \in \text{Val}(Z_i)$
 - The parents and CPDs of all other nodes $X \notin \mathbf{Z}$ are unchanged

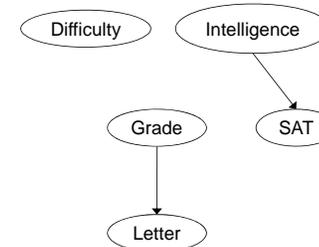
Importance Sampling for Bayesian Networks

Mutilated Network:

$\beta_{\text{Intelligence=high, Grade=B}}^{\text{student}}$

D	P(D)
low	0.6
high	0.4

G	P(G D,I)
C	0
B	1
A	0



I	P(I)
low	0
high	1

I	S	P(S I)
low	low	0.95
low	high	0.05
high	low	0.2
high	high	0.8

G	L	P(L G)
C	weak	0.99
C	strong	0.01
B	weak	0.4
B	strong	0.6
A	weak	0.1
A	strong	0.9

Importance Sampling for Bayesian Networks

- Proposition 12.2:** Let ξ be a sample generated by Likelihood Weighting and w be its weight. Then the distribution over ξ is as defined by the network $B_{\mathbf{Z}=\mathbf{z}}$, and

$$w(\xi) = \frac{P_B(\xi)}{P_{B_{\mathbf{Z}=\mathbf{z}}}(\xi)}$$

- (Informally) Importance sampling using a mutilated network as a proposal distribution is equivalent to Likelihood Weighting with $P_B(X, \mathbf{z})$ and proposal distribution Q induced by the mutilated network $B_{\mathbf{E}=\mathbf{e}}$.

Likelihood Weighting Revisited

Likelihood Weighting Revisited

Two versions of likelihood weighting

1. Ratio Likelihood Weighting
2. Normalized Likelihood Weighting

Likelihood Weighting Revisited

Ratio Likelihood Weighting

$$P(\mathbf{y} | \mathbf{e}) = \frac{P(\mathbf{y}, \mathbf{e})}{P(\mathbf{e})}$$

Use unnormalized importance sampling:

1. For numerator – use LW to generate M samples with $\mathbf{Y}=\mathbf{y}$, $\mathbf{E}=\mathbf{e}$ as the event
2. For denominator – use LW to generate M' samples with $\mathbf{E}=\mathbf{e}$ as the event

$$\hat{P}_D(\mathbf{y} | \mathbf{e}) = \frac{\hat{P}_D(\mathbf{y}, \mathbf{e})}{\hat{P}_D(\mathbf{e})} = \frac{\frac{1}{M} \sum_{m=1}^M w[m]}{\frac{1}{M'} \sum_{m=1}^{M'} w'[m]}$$

Likelihood Weighting Revisited

Normalized Likelihood Weighting

- Ratio Likelihood Weighting estimates a single query $P(\mathbf{y}|\mathbf{e})$ from a set of samples (ie. it sets $\mathbf{Y}=\mathbf{y}$ when sampling)
- Sometimes we want to evaluate a set of queries $P(\mathbf{y}|\mathbf{e})$
- Use normalized likelihood weighting with
$$\tilde{P}(\mathbf{X}) = P_B(\mathbf{X}, \mathbf{e})$$
- Estimate the expectation of a function f:
$$f(\xi) = \mathbf{I}\{\xi(\mathbf{Y}) = \mathbf{y}\}$$

Likelihood Weighting Revisited

Quality of importance sampling depends on how close the proposal distribution Q is to the target distribution P.

Consider the two extremes:

1. All evidence at the roots:
 - Proposal distribution is the posterior
 - Evidence affects samples all along the way and all samples have the same weight P(e)

Likelihood Weighting Revisited

2. All evidence at the leaves:

- Proposal distribution is the prior distribution $P_B(X)$
- Evidence doesn't affect samples, weights have to compensate. LW will only work well if prior is similar to the posterior

Likelihood Weighting Revisited

$$P(X) = \sum_e \underbrace{P(e)}_{\text{Prior}} \underbrace{P(X|e)}_{\text{Posterior}}$$

- If $P(e)$ is high, then the posterior $P(X|e)$ plays a large role and is close to the prior $P(X)$
- If $P(e)$ is low, then the posterior $P(X|e)$ plays a small role and the prior $P(X)$ will likely look very different

Likelihood Weighting Revisited

Summary

Ratio Likelihood Weighting

- Computes $P(\mathbf{y}|\mathbf{e})$ for a specific \mathbf{y} (ie. values for \mathbf{y} are set)
- Uses unnormalized importance sampling for both numerator and denominator in $P(\mathbf{y},\mathbf{e})/P(\mathbf{e})$
- Needs a new set of samples for each query \mathbf{y}
- Lower variance in estimator
- Can bound # of samples required for a good estimate (but under strong conditions)

Likelihood Weighting Revisited

Summary

Normalized Likelihood Weighting

- Samples an assignment for \mathbf{Y} , which introduces additional variance
- Allows multiple queries \mathbf{y} using the same set of samples (conditioned on evidence \mathbf{e})

Likelihood Weighting Revisted

Problems with Likelihood Weighting:

- If there are a lot of evidence variables $P(\mathbf{Y} | \mathbf{E}_1 = \mathbf{e}_1, \dots, \mathbf{E}_k = \mathbf{e}_k)$:
 - Many samples will have ε weight
 - Weighted estimate dominated by a small fraction of samples that have $> \varepsilon$
- If evidence variables occur in the leaves, the samples drawn will not be affected much by the evidence