

CS786

Lecture 13: May 14, 2012

Sampling techniques
[KF Chapter 12]

CS786 P. Poupart 2012

1

Sampling Techniques

- Direct sampling
- Rejection sampling
- Likelihood weighting
- Importance sampling
- Markov chain Monte Carlo (MCMC)
 - Gibbs Sampling
 - Metropolis-Hastings
- Sequential Monte Carlo sampling (a.k.a. particle filtering)

CS786 P. Poupart 2012

2

Approximate Inference by Sampling

- Expectation: $E_P[f(x)] = \int_x P(x)f(x)dx$
 - Approximate integral by sampling:
 $E_P[f(x)] \approx \frac{1}{n} \sum_{i=1}^n f(x_i)$ where $x_i \sim P(x)$
- Inference query: $\Pr(\mathbf{X}|e) = \sum_Y \Pr(\mathbf{X}, Y|e)$
 - Approximate exponentially large sum by sampling:
 $\Pr(\mathbf{X}|e) = \frac{1}{n} \sum_{i=1}^n \Pr(\mathbf{X}|\mathbf{y}_i, e)$ where $\mathbf{y}_i \sim P(\mathbf{Y}|e)$

CS786 P. Poupart 2012

3

Direct Sampling (a.k.a. forward sampling)

- Unconditional inference queries (i.e., $\Pr(V = t)$)
- Bayesian networks only
 - Idea: sample each variable given the values of its parents according to the topological order of the graph.

CS786 P. Poupart 2012

4

Direct Sampling Algorithm

Sort the variables by topological order

For $i = 1$ to n do (sample n particles)

For each variable V_j do

Sample $v_j^{(i)} \sim \Pr(V_j | \mathbf{pa}_V)$

- Approximation: $\Pr(V_k = t) \approx \frac{1}{n} \sum_{i=1}^n \delta(v_k^{(i)} = t)$

CS786 P. Poupart 2012

5

Example

CS786 P. Poupart 2012

6

Analysis

- Complexity: $O(n|V|)$ where $|V| = \text{\#variables}$
- Accuracy
 - Absolute error ϵ : $P(|\hat{P}(V) - P(V)| > \epsilon) \leq \delta = 2e^{-2n\epsilon^2}$
 - Sample size $n \geq \frac{\ln(\frac{2}{\delta})}{2\epsilon^2}$
 - Relative error ϵ : $P\left(\frac{\hat{P}(V)}{P(V)} \notin [1 - \epsilon, 1 + \epsilon]\right) \leq \delta = 2e^{-\frac{nP(V)\epsilon^2}{3}}$
 - Sample size $n \geq \frac{3 \ln(\frac{2}{\delta})}{2P(V)\epsilon^2}$

CS786 P. Poupart 2012

7

Rejection Sampling

- Conditional inference queries (i.e., $\Pr(V = t|e)$)
- Bayesian networks only
 - Idea: sample each variable given the values of its parents according to the topological order of the graph, however reject samples that do not agree with evidence

CS786 P. Poupart 2012

8

Rejection Sampling Algorithm

Sort the variables by topological order

For $i = 1$ to n do (sample n particles)

For each variable V_j do

Sample $v_j^{(i)} \sim \Pr(V|pa_V)$

Reject $\mathbf{v}^{(i)}$ if $\mathbf{v}^{(i)}$ is inconsistent with \mathbf{e} (i.e., $\mathbf{v}_E^{(i)} \neq \mathbf{e}$)

- Approximation: $\Pr(V_k = t|\mathbf{e}) \approx \frac{\sum_{i=1}^n \delta(v_k^{(i)}=t \wedge \mathbf{v}_E^{(i)}=\mathbf{e})}{\sum_{i=1}^n \delta(\mathbf{v}_E^{(i)}=\mathbf{e})}$

CS786 P. Poupart 2012

9

Example

CS786 P. Poupart 2012

10

Analysis

- Complexity: $O(n|V|)$ where $|V| = \text{\#variables}$
- Expected # samples that are accepted: $O(n \Pr(\mathbf{e}))$
 - Since $\Pr(\mathbf{e})$ often decreases exponentially with the number of evidence variables, the number of samples also decreases exponentially.
 - For good accuracy: exponential # of samples often needed in practice.

CS786 P. Poupart 2012

11

Likelihood Weighting

- Conditional inference queries (i.e., $\Pr(V = t|\mathbf{e})$)
- Bayesian networks only
 - Idea: sample each non-evidence variable given the values of its parents in topological order. Assign weights to samples based on the probability of the evidence.

CS786 P. Poupart 2012

12

Likelihood Weighting Algorithm

Sort the variables by topological order

For $i = 1$ to n do (sample n particles)

$w_i \leftarrow 1$

For each variable V_j do

If V_j is not an evidence variable do

Sample $v_j^{(i)} \sim \Pr(V_j | \mathbf{pa}_{V_j})$

else

$w_i \leftarrow w_i * \Pr(v_j | \mathbf{pa}_{V_j})$

- Approximation: $\Pr(V_k = t | \mathbf{e}) \approx \frac{\sum_{i=1}^n w_i \delta(v_k^{(i)} = t)}{\sum_{i=1}^n w_i}$

CS786 P. Poupart 2012

13

Example

CS786 P. Poupart 2012

14

Analysis

- Complexity: $O(n|V|)$ where $|V| = \text{\#variables}$
- Effective sample size: $O(n \Pr(e))$
 - Even though all samples are accepted, their importance is reweighted to a fraction equal to $\Pr(e)$
 - For good accuracy: the # of samples will be the same as for rejection sampling (hence exponential with the number of evidence variables).

CS786 P. Poupart 2012

15

Importance Sampling

- Likelihood weighting is a special case of importance sampling
- General approach to estimate $E_P[f(x)]$ by sampling from Q instead of P
 - Works for Bayes nets and probability densities
- Idea: generate samples x from Q and assign weights $P(x)/Q(x)$

CS786 P. Poupart 2012

16

Importance Sampling Algorithm

For $i = 1$ to n do (sample n particles)

Sample $x^{(i)}$ from Q

Assign weight: $w_i \leftarrow P(x^{(i)})/Q(x^{(i)})$

- Approximation: $E_P[f(x)] \approx \frac{1}{n} \sum_{i=1}^n w_i f(x^{(i)})$
 - Unbiased estimator
 - Variance of estimator decreases linearly with sample size

CS786 P. Poupart 2012

17

Normalized Importance Sampling

- Often the reason why we are sampling from Q instead of P is that we don't know P .
- But, we may know \tilde{P} an unnormalized version of P
 - Markov nets: $P(\mathbf{X}) = k \prod_i f_i(\mathbf{X})$ while $\tilde{P}(\mathbf{X}) = \prod_i f_i(\mathbf{X})$
 - Bayes nets: $P(\mathbf{X}|\mathbf{e})$ while $\tilde{P}(\mathbf{X}, \mathbf{e})$
- Idea: generate samples x from Q and assign weights $\tilde{P}(x)/Q(x)$. Normalize the estimator.

CS786 P. Poupart 2012

18

Normalized Importance Sampling Algorithm

For $i = 1$ to n do (sample n particles)

Sample $x^{(i)}$ from Q

Assign weight: $w_i \leftarrow \tilde{P}(x^{(i)})/Q(x^{(i)})$

- Approximation: $E_P[f(x)] \approx \frac{\sum_{i=1}^n w_i f(x^{(i)})}{\sum_{i=1}^n w_i}$
 - Biased estimator for finite n (unbiased for $n = \infty$)
 - Variance of estimator decreases linearly with sample size

CS786 P. Poupart 2012

19

Markov Chain Monte Carlo

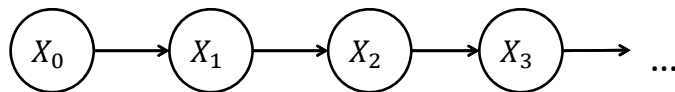
- Iterative sampling technique that converges to the desired distribution in the limit
- Idea: set up a Markov chain such that its stationary distribution is the desired distribution

CS786 P. Poupart 2012

20

Markov Chain

- Definition: A Markov chain is a linear chain Bayesian network with a stationary conditional distribution known as the transition function



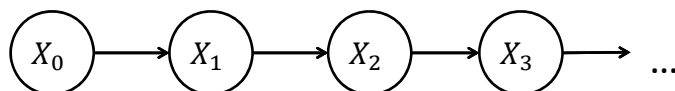
- Initial distribution: $\Pr(X_0)$
- Transition distribution: $\Pr(X_t|X_{t-1})$

CS786 P. Poupart 2012

21

Markov Chain

- Definition: A Markov chain is a linear chain Bayesian network with a stationary conditional distribution known as the transition function



- Initial distribution: $\Pr(X_0)$
- Transition distribution: $\Pr(X_t|X_{t-1})$

CS786 P. Poupart 2012

22

Asymptotic Behaviour

- Let $\Pr(X_t)$ be the distribution at time step t

$$\begin{aligned}\Pr(X_t) &= \sum_{X_{0..t-1}} \Pr(X_{0..t}) \\ &= \sum_{X_{t-1}} \Pr(X_{t-1}) \Pr(X_t | X_{t-1})\end{aligned}$$

- In the limit (i.e., when $t \rightarrow \infty$), the Markov chain may converge to stationary distribution $\pi(x) = \Pr(X_\infty = x)$

$$\begin{aligned}\pi(x) &= \Pr(X_\infty = x) \\ &= \sum_{X_{\infty-1}} \Pr(X_{\infty-1} = x') \Pr(X_\infty = x | X_{\infty-1} = x') \\ &= \sum_{x'} \pi(x') \Pr(x | x')\end{aligned}$$

CS786 P. Poupart 2012

23

Stationary distribution

- Let $T_{x|x'} = \Pr(x|x')$ be a matrix that represents the transition function
- If we think of π as a column vector, then π is an eigenvector of T with eigenvalue 1

$$T\pi = \pi$$

CS786 P. Poupart 2012

24

Ergodic Markov Chain

- Definition: A Markov chain is ergodic when there is a non-zero probability of reaching any state from any state in a finite number of steps
- When the Markov chain is ergodic, there is a unique stationary distribution

- Sufficient condition: detailed balance

$$\pi(x)\Pr(x'|x) = \pi(x')\Pr(x|x')$$

Detailed balance \rightarrow ergodicity \rightarrow unique stationary dist.

CS786 P. Poupart 2012

25

Markov Chain Monte Carlo

- Idea: set up an ergodic Markov chain such that the unique stationary distribution is the desired distribution
- Since the Markov chain is a linear chain Bayes net, we can use direct sampling (forward sampling) to obtain a sample of the stationary distribution

CS786 P. Poupart 2012

26

Generic MCMC Algorithm

Sample $x_0 \sim \Pr(X_0)$

For $i = 1$ to n do (sample n particles)

Sample $x_t \sim \Pr(X_t | x_{t-1})$

- Approximation: $\pi(x) \approx \frac{1}{n} \sum_{t=1}^n \delta(x_t = x)$
- In practice, ignore the first k samples for a better estimate (burn-in period):

$$\pi(x) \approx \frac{1}{n-k} \sum_{t=k}^n \delta(x_t = x)$$

CS786 P. Poupart 2012

27

Choosing a Markov Chain

- Different Markov chains lead to different algorithms
 - Gibbs sampling
 - Metropolis Hastings

CS786 P. Poupart 2012

28

Gibbs Sampling

- Suppose $\Pr(\mathbf{X})$ defined by a graphical model (Bayes net or Markov net)
- Inference query: $\Pr(\mathbf{Y}|\mathbf{e})$? Where $\mathbf{Y} \subseteq \mathbf{X}$
- Idea: randomly assign values to all non-evidence variables, then repeatedly sample each non-evidence variable given the assigned values for all other variables

CS786 P. Poupart 2012

29

Gibbs Sampling Algorithm

Randomly assign $v_j^{(0)}$ to all non-evidence variables V_j

For $i = 1$ to n do (sample n particles)

For each non-evidence variable V_j do

Sample $v_j^{(i)} \sim \Pr(V_j | \mathbf{v}_{\sim j}^{(i-1)}, \mathbf{e})$

- Approximation: $\Pr(V_k = t | \mathbf{e}) \approx \frac{1}{n} \sum_{i=1}^n \delta(v_k^{(i)} = t)$

CS786 P. Poupart 2012

30

Example

CS786 P. Poupart 2012

31

Practical Consideration

- Burn-in period: ignore first k samples:

$$\Pr(V_k = t | \mathbf{e}) \approx \frac{1}{n - k} \sum_{i > k}^n \delta(v_k^{(i)} = t)$$

- Use most recent values to sample $V_j^{(i)}$

$$v_j^{(i)} \sim \Pr(V_j^{(i)} | \mathbf{v}_{1 \dots j-1}^{(i)}, \mathbf{v}_{j+1 \dots |V|}^{(i-1)})$$

- Use conditional independence to restrict parent variables to the Markov blanket

$$v_j^{(i)} \sim \Pr(V_j^{(i)} | \mathbf{v}_{\forall k < j, k \in mb(j)}^{(i)}, \mathbf{v}_{\forall k > j, k \in mb(j)}^{(i-1)})$$

CS786 P. Poupart 2012

32

Convergence

- Let $\Pr(\mathbf{V}^{(i)}|\mathbf{V}^{(i-1)}, \mathbf{e})$ be the transition function of the Markov chain associated with Gibbs sampling
- **Theorem:** Gibbs sampling converges to $\Pr(\mathbf{V}|\mathbf{e})$ when all potentials are strictly positive.
- Proof: $\Pr(\mathbf{V}^{(i)}|\mathbf{V}^{(i-1)}, \mathbf{e})$ satisfies detailed balance i.e. $\Pr(\mathbf{V}|\mathbf{e}) \Pr(\mathbf{V}'|\mathbf{V}, \mathbf{e}) = \Pr(\mathbf{V}'|\mathbf{e}) \Pr(\mathbf{V}|\mathbf{V}', \mathbf{e})$

CS786 P. Poupart 2012

33

Metropolis-Hastings

- Suppose we can compute $\pi(x)$ for a given x but we can't sample from $\pi(x)$ easily.
- Idea: use an arbitrary transition distribution $Q(x'|x)$ and use rejection sampling to correct for the choice of Q .
- Advantage: since Q can be anything, we can always obtain an MCMC algorithm by Metropolis-Hastings
 - It is particularly useful to approximate continuous distributions

CS786 P. Poupart 2012

34

Metropolis-Hastings Algorithm

Randomly select $x^{(0)}$

For $i = 1$ to n do (sample n particles)

Sample $x^{(i)} \sim Q(X|x^{(i-1)})$

Accept $x^{(i)}$ with probability $\min \left[1, \frac{\pi(x^{(i)})Q(x^{(i)}|x^{(i-1)})}{\pi(x^{(i-1)})Q(x^{(i-1)}|x^{(i)})} \right]$

Otherwise reject $x^{(i)}$ (i.e., $x^{(i)} \leftarrow x^{(i-1)}$)

- Approximation: $\pi(x) \approx \frac{1}{n} \sum_{i=1}^n \delta(x^{(i)} = x)$

CS786 P. Poupart 2012

35

Convergence

- The transition distribution in Metropolis-Hastings is

$$\Pr(x'|x) = \begin{cases} Q(x'|x)A(x \rightarrow x') & x \neq x' \\ Q(x|x) + \sum_{x' \neq x} Q(x'|x)(1 - A(x \rightarrow x')) & x = x' \end{cases}$$

$$\text{where } A(x \rightarrow x') = \min \left[1, \frac{\pi(x')Q(x'|x)}{\pi(x)Q(x|x')} \right]$$

- **Theorem:** Metropolis-Hastings converges to $\pi(x)$.
- Proof: $\Pr(x'|x)$ satisfies detailed balance

CS786 P. Poupart 2012

36