

- Today:
- Markov Chain Monte Carlo Methods (MCMC)
 - Metropolis - Hastings
 - Gibbs Sampling
 - Particle Filters

1. MCMC* Recap: Metropolis - Hastings

- Setup & Goal: Target distribution $P_x(x) = \frac{P_x^*(x)}{Z}$
- Want: samples $x^{(1)}, x^{(2)}, \dots$ from $P_x(x)$.
- $P_x^*(x)$ ← known
 Z ← not known

• MCMC does work

for continuous variables

→ Assume: $x \in \{1, 2, \dots, |\mathcal{X}| \}^N$

• We'll derive for discrete case for simplicity of presentation.

- Basic Idea:

- Construct a Markov Chain P , s.t. P has unique stationary distribution π & $\pi(x) = P_x(x) \quad \forall x \in |\mathcal{X}|^N$
- The states of P are clearly all possible $x \in \mathcal{X}^N$
 \Rightarrow We need to specify the transition probabilities. $[P_{ij}]$

- Metropolis - Hastings Algorithm:

This is a matrix of size $|\mathcal{X}|^N \times |\mathcal{X}|^N$

(1). Find some transition probability matrix $[K_{ij}]$

- $[K_{ij}]$ can be chosen as you like. e.g. uniform.
- Different choices of $[K_{ij}]$ don't affect correctness, but can affect mixing time
- $K_{ij} > 0 \quad \forall i \in \{1, 2, \dots, N\}$, and any state x can go to any other state x' in finite # of steps.

(2) Define $P_{ij} = \begin{cases} K_{ij} \cdot \min\left\{1, \frac{P_x^*(j) \cdot K_{ji}}{P_x^*(i) \cdot K_{ij}}\right\} & \text{(if } i \neq j) \\ 1 - \sum_{k \neq i} P_{ik} & \text{(if } i = j) \end{cases}$

$R(i,j)$

(3) Pick any initial state $x^{(0)}$

For $t = 1 : \text{Max-Nb-Iter}$

- Given the current state $i \triangleq x^{(t-1)}$, propose a new state j according to $[K_{ij}]$
- With probability $R(i, j)$, set $x^{(t)} = j$
With probability $1 - R(i, j)$, set $x^{(t)} = i$

End

- Comments: (1). The algorithm will spit out a sequence of samples

$x^{(0)}, x^{(1)}, x^{(2)}, \dots, x^{(t)}, x^{(t+1)}, \dots, x^{(t+s)}, x^{(t+s+1)}, \dots, x^{(t+2s)}$

"burn-in" period

choose one sample from every s successive samples to reduce correlation

(2) Intuition about $R(i, j)$

$P_x(j) K_{ji}$: "flow from $j \rightarrow i$ "

$P_x(i) K_{ij}$: "flow from $i \rightarrow j$ "

- if "flow $j \rightarrow i$ " \geq "flow $i \rightarrow j$ ": add flow $i \rightarrow j$, i.e. make the transition from i to j
- if "flow $j \rightarrow i$ " $<$ "flow $i \rightarrow j$ ": control flow $i \rightarrow j$, i.e. only make the transition $i \rightarrow j$ with some probability < 1

(3) In practice, don't compute/store $[P_{ij}]$

$|x| \times |x|$, too big

* Gibbs Sampling

- can be viewed as a special case of Metropolis-Hastings, with $[K_{ij}]$ chosen using the following process: current state = x

(1). uniformly pick a coordinate $k \in \{1, 2, \dots, N\}$

(2). $\forall l \neq k, l \in \{1, 2, \dots, N\}$. set $x'_l = x_l$

x'_k is sampled from $P_{x_k | x_{\setminus k}}(\cdot | x_{\setminus k})$ $\leftarrow x_{\setminus k}$ denotes all other coordinates of x

- Claim: Let us denote the new state as x' . If $R(x, x')$ is defined as before (i.e. $R(x, x') = \min\{1, \frac{P_x(x') \cdot K_{x'x}}{P_x(x) \cdot K_{xx'}}\}$) is always 1.

Proof: $P_x(x) \cdot K_{x,x'} = \underline{P_x(x)} \cdot \frac{1}{N} P(x_k' | x_{1:k})$

$$= \frac{1}{N} \underline{P(x_k | x_{1:k})} \cdot P(x_{1:k}) P(x_k' | x_{1:k})$$

$$= \frac{1}{N} \underline{P(x_k | x_{1:k})} \cdot \underline{P(x_{1:k})} P(x_k' | x_{1:k})$$

$$= \frac{1}{N} P(x_k | x_{1:k}) \underline{P_x(x')}$$

$$= K_{x',x} P_x(x')$$

$$\Rightarrow R(x, x') = \min \left\{ 1, \frac{P_x(x') K_{x',x}}{P_x(x) K_{x,x'}} \right\} = 1 \Rightarrow \text{always accept.}$$

- Notice the correlation is even higher in Gibbs sampling, because only 1 coordinate is changed at one time.

Variants exist: e.g. check homework problem 8.2 for **block gibbs sampling**

* Failure Modes of MCMC

- Islands of high-probability states, with low-probability states in between
- All states have very small probability except for one state.

e.g. $P(x_0) = 1/2$ $P(x) = \frac{1}{2(2^{100}-1)}$ $\forall x \neq x_0$ & $x, x_0 \in \{0,1\}^{100}$

will have long sequences of x_0 & long sequences of $x \neq x_0$

2. Particle Filters

* Importance Sampling

- Importance Sampling produces $\left\{ \begin{array}{l} \text{estimate of expectation of a given function } E_p[f(x)] \\ \text{samples from } P_x(\cdot) \end{array} \right. \quad \times$

- Algorithm: (1) Propose some distribution $q_x(\cdot)$ that is easy to sample from.
 "proposal distribution" e.g. uniform distribution

(2) Let $x^{(1)}, x^{(2)}, \dots, x^{(k)}$ are ~~some~~ iid samples from $q_x(\cdot)$

(3) Compute weights: $W^k = W(x^k) = \frac{P^*(x^k)}{q(x^k)} \quad k=1, 2, \dots, K$

(4) Compute expected value of given function f :

$$E_p[f(x)] \approx \frac{\frac{1}{K} \sum_{k=1}^K W^k \cdot f(x^k)}{\frac{1}{K} \sum_{k=1}^K W^k}$$

- Comments: (1). Can prove: $\lim_{K \rightarrow \infty} \frac{\frac{1}{K} \sum_{k=1}^K W^k \cdot f(x^k)}{\frac{1}{K} \sum_{k=1}^K W^k} = E_{P_x} [f(x)]$ (*)

(2). $q_x(\cdot)$ can be any distribution (that is easy to sample from).

But choice of $q_x(\cdot)$ will affect the speed of convergence in (*).

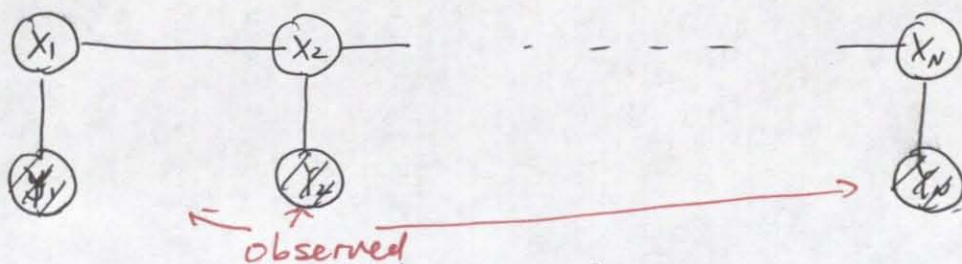
(3) Undesirable situation: $q_x(\cdot)$ has low probability where $P_x(\cdot)$ has high probability.

(4) One way to look at importance sampling is that we're approximating $P_x(\cdot)$ using a bunch of weighted particles.

$\{W^k, x^k\}_{k=1}^K$. Particle filter is built based on this intuition.

* Particle Filters

- Introduced on HMM ← can generalize to trees
 Can deal with distributions over ^{any} continuous variables ← not just Gaussians



- View 1: (as introduced in lecture)

• Target Distribution: $P_{x_1^N | y_1^N} (x_1^N | y_1^N) = \frac{P_{x_1^N, y_1^N} (x_1^N, y_1^N)}{P_{y_1^N} (y_1^N)}$
← P^x known
← \bar{z} not known

• use $P_{x_1^N} (x_1^N)$ as the proposal distribution $q(\cdot)$

$(P_{x_1^N} (x_1^N) = P_{x_1} (x_1) \prod_{i=1}^{N-1} P_{x_{i+1} | x_i} (x_{i+1} | x_i))$ can be sampled easily)

• Let $x^{(1)}, x^{(2)}, \dots, x^{(K)}$ be iid samples from $P_{x_1^N}(\cdot)$

• $W^{(k)} = \frac{P^*(x^{(k)})}{q(x^{(k)})} = \frac{P_{x_1^N, y_1^N} (x_1^{(k)}, y_1^N)}{P_{x_1^N} (x_1^{(k)})} = P_{y_1^N | x_1^N} (y_1^N | x_1^{(k)}) = \prod_{i=1}^N P_{y_i | x_i} (y_i | x_i^{(k)})$

• $E_{P_{x_1^N | y_1^N}} [f(x)] \approx \frac{\frac{1}{K} \sum_{i=1}^K W^{(k)} f(x^{(k)})}{\frac{1}{K} \sum_{i=1}^K W^{(k)}}$

- View 2. (

• Target distribution: $P_{X_n | Y_1^n}(\cdot | y_1^n)$ $n=1, 2, \dots, N$

• initialization: $P_{X_1}(\cdot)$: $\{W_1^{(k)}, x_1^{(k)}\}_{k=1}^K$ $W_1^{(k)} = \frac{1}{K}$.

• For $n=1: N-1$

~~$x_{n+1}^{(k)} \sim P_{X_{n+1} | X_n}(\cdot | x_n^{(k)})$~~ $x_{n+1}^{(k)} \sim P_{X_{n+1} | X_n}(\cdot | x_n^{(k)})$ $k=1, 2, \dots, K$

$\tilde{W}_{n+1}^{(k)} = \tilde{W}_n^{(k)}$ $k=1, 2, \dots, K$.

$W_{n+1}^{(k)} = \tilde{W}_{n+1}^{(k)} \cdot P_{Y_{n+1} | X_{n+1}}(y_{n+1} | x_{n+1}^{(k)})$ $k=1, 2, \dots, K$. "update step"

"prediction step."

$\{W_n^{(k)}, x_n^{(k)}\}_{k=1}^K$
 $\Rightarrow \{\tilde{W}_{n+1}^{(k)}, x_{n+1}^{(k)}\}_{k=1}^K$

$\{\tilde{W}_{n+1}^{(k)}, x_{n+1}^{(k)}\}_{k=1}^K$

$\Rightarrow \{W_{n+1}^{(k)}, x_{n+1}^{(k)}\}_{k=1}^K$

End

• resampling: if $\hat{N}_{eff} = \frac{1}{\sum_{k=1}^K (W_n^{(k)})^2}$ too small.
 (heuristic)

• Related to forward pass of sum-product:

$P_{n+1|n}(x_{n+1} | y_1^n) = \int_{x_n} P_{X_{n+1} | X_n}(x_{n+1} | x_n) P_{n|n}(x_n | y_1^n) dx_n$

$P_{n+1|n+1}(x_{n+1} | y_1^{n+1}) = \frac{1}{Z} \int_{x_n} P_{Y_{n+1} | X_{n+1}}(y_{n+1} | x_{n+1}) P_{n+1|n}(x_{n+1} | y_1^n) dx_n$

* Beyond HMM: Particle Filters on Trees

- BP equations on trees:

$M_{i \rightarrow j}(x_j) = \int_{x_i} \psi_{ij}(x_i, x_j) \phi_i(x_i) \prod_{l \in \text{N}(i) \setminus j} M_{l \rightarrow i}(x_i) dx_i$

separate into 2 steps:

(1): $\phi_{ij}(x_i) \triangleq \phi_i(x_i) \prod_{l \in \text{N}(i) \setminus j} M_{l \rightarrow i}(x_i)$

(2): $M_{i \rightarrow j}(x_j) = \int_{x_i} \psi_{ij}(x_i, x_j) \phi_{ij}(x_i) dx_i$

- step (1): $M_{l \rightarrow i}(x_i)$ is now represented by a weighted particle set

$\{W_{l \rightarrow i}^{(k)}, x_{l \rightarrow i}^{(k)}\}_{k=1}^K$ i.e. $M_{l \rightarrow i}(x_i) \approx \sum_{k=1}^K W_{l \rightarrow i}^{(k)} \delta(x_i - x_{l \rightarrow i}^{(k)})$

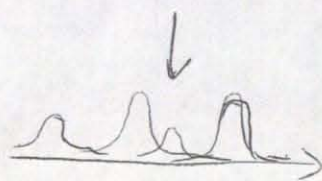
But we have problems multiplying $M_{l_1 \rightarrow i}(x_i)$ & $M_{l_2 \rightarrow i}(x_i)$ in this form

(continuous $x_i \Rightarrow P(x_{l_1 \rightarrow i}^{(k)} = x_{l_2 \rightarrow i}^{(k)}) = 0$)

Instead, use small Gaussians.

$$M_{l \rightarrow i}(x_i) \approx \sum_{k=1}^K W_{l \rightarrow i}^{(k)} \mathcal{N}(x_i; x_{l \rightarrow i}^{(k)}, \Sigma_{l \rightarrow i}^{-1})$$

mixture of Gaussians



$\Rightarrow M_{l_1 \rightarrow i}(x_i) \cdot M_{l_2 \rightarrow i}(x_i)$ also a mixture of Gaussians.

For further details: "Efficient Multiscale Sampling from Products of Gaussian Mixtures" by Ihler et al 2003.

- step (2). Just a summation over the samples.

Rosencrantz & Guildenstern Are Dead

Film Discussion Questions

1. What do you make of the opening credits? As the play opens in a place “without visible character,” the film begins with credits on a black screen. However, the film introduces “audible character” in the form of a western-themed soundtrack. How does this contrast with your expectation of “two Elizabethans”?
2. Unlike a play, film can cut from scene to scene instantaneously. It can completely change the setting and mise-en-scene without limitation. Does the film take advantage of this fact, or does it try to respect the story’s original medium by forcing characters around in circles so that they repeatedly end up in the same room of the castle?
3. How are sound effects used in the film? Extradiegetic chimes are added at crucial moments – first after Guildenstern brings up the question of suspense, then again when the coin is finally tails. The film inserts dramatic echoes to some of Guildenstern’s most excited moments that momentarily halt the flow of dialogue and narrative. Ambient animal noises can be heard throughout the film and are at one point revealed to be coming from Rosencrantz. Do these elements add anything to the narrative other than perhaps to emphasize moments of significance to the less sophisticated audience of cinema?
4. As discussed in Thursday’s class, two individuals on a stage or screen are far more distinguishable from one another than words on a page. It was mentioned that being able to visually tell them apart would make it easier to tell which character was which. After watching the film, was it clear who was Rosencrantz and who was Guildenstern? Or did the fact that you could differentiate them as Gary Oldman and Tim Roth (or the long-haired one and short-haired one) eliminate the necessity of thinking of each one specifically as Rosencrantz or Guildenstern?
5. What were the effects of having the question game played on an actual tennis court?
6. What is the significance of wind in the film? Hamlet is said to be “at the mercy of the elements,” but Rosencrantz and Guildenstern are also plagued by their quest to find out the direction of the wind. As an example, the windmill toy spins after Rosencrantz insists that “there isn’t any wind.”
7. Let’s discuss the puppet show. The tragedians use puppets in addition to live action to enact the deaths from Hamlet. Why does the film version include this story-telling mode when the original play does not? Is it simply a way to visually emphasize the complex literary layers of *Rosencrantz and Guildenstern are Dead*?
8. The film ends with the Tragedians driving away in their caravan. What is the effect of ending with life rather than ending on a stage filled with death (as both Hamlet and the play of R&G do)?

MIT OpenCourseWare
<http://ocw.mit.edu>

6.438 Algorithms for Inference
Fall 2014

For information about citing these materials or our Terms of Use, visit: <http://ocw.mit.edu/terms>.