# Inference via Sampling (Contd)

Piyush Rai
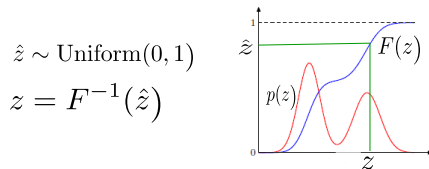
Topics in Probabilistic Modeling and Inference (CS698X)

March 2, 2019

## Recap: Basic Sampling Methods

- Inverse CDF method. Assume $F(z)$ to be CDF of our distribution of interest $p(z)$

$$\hat{z} \sim \text{Uniform}(0, 1)$$
$$z = F^{-1}(\hat{z})$$



- Reparametrization method (also used in VI - pathwise gradient methods), e.g.,

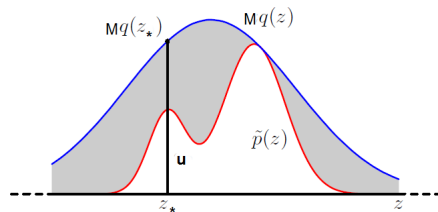$$\hat{z} \sim \mathcal{N}(0, 1) \Rightarrow z = \mu + \sigma\hat{z} \sim \mathcal{N}(\mu, \sigma^2)$$

- Note: The above are examples of the more general idea of transformation of distributions

$$p(\boldsymbol{z}) = q(\hat{\boldsymbol{z}}) \left| \frac{\partial \hat{\boldsymbol{z}}}{\partial \boldsymbol{z}} \right|$$

.. where $\left| \frac{\partial \hat{\boldsymbol{z}}}{\partial \boldsymbol{z}} \right|$ is the determinant of the Jacobian

# Recap: Basic Sampling Methods



- Rejection Sampling: Sample from $p(z) = \frac{\tilde{p}(z)}{Z_p}$ by sampling from $q(z)$, s.t., $Mq(z) \geq \tilde{p}(z)$
    - $z_* \sim q(z)$ and $u \sim \text{Uniform}(0, Mq(z))$
    - If $u \leq \tilde{p}(z_*)$, accept $z_*$ else reject
    - Repeat the above two steps until we have generated the desired number of samples

## Computing Expectations via Monte Carlo Sampling

- Often we are interested in computing expectations of the form

$$\mathbb{E}[f] = \int f(\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

  where $f(\mathbf{z})$ is some function of a random variable $\mathbf{z} \sim p(\mathbf{z})$

- A simple approximation scheme: Monte Carlo integration

- Suppose we can generate $L$ independent samples from $p(\mathbf{z})$: $\{\mathbf{z}^{(\ell)}\}_{\ell=1}^{L} \sim p(\mathbf{z})$

- Monte-Carlo approximation replaces the expectation by an empirical average

$$\hat{f} \approx \frac{1}{L}\sum_{\ell=1}^{L} f(\mathbf{z}^{(\ell)})$$

- Since the samples are independent of each other, can show the following (exercise)

$$\mathbb{E}[\hat{f}] = \mathbb{E}[f] \qquad \text{and} \quad \text{var}[\hat{f}] = \frac{1}{L}\text{var}[f] = \frac{1}{L}\mathbb{E}[(f - \mathbb{E}[f])^2]$$

- Note that the variance in the estimate of expectation decreases as $L$ increases

# Computing Expectations via Importance Sampling

- Monte Carlo assumes we know how to generate samples from $p(z)$. What if we don't know?
- Transformation methods can be one way to handle this situation
- Importance Sampling is another way: Generate from a "proposal" $q(z)$, i.e., $\{z^{(\ell)}\}_{\ell=1}^{L} \sim q(z)$
- Additionally, suppose we can evaluate $p(z)$ at any given $z$
- Importance Sampling then approximates the original expectation as

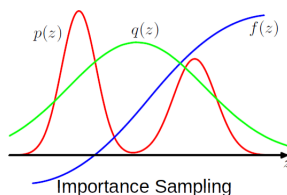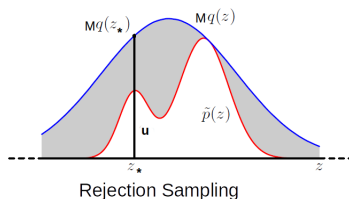$$\mathbb{E}[f] = \int f(z)p(z)dz = \int f(z)\frac{p(z)}{q(z)}q(z)dz \approx \frac{1}{L}\sum_{\ell=1}^{L} f(z^{(\ell)})\frac{p(z^{(\ell)})}{q(z^{(\ell)})}$$

- This is basically "weighted" Monte Carlo integration
  - $w_\ell = \frac{p(z^{(\ell)})}{q(z^{(\ell)})}$ denotes the importance weight of each sample $z^{(\ell)}$
- Works even when we can evaluate $p(z) = \frac{\tilde{p}(z)}{Z_p}$ only up to a prop. constant (PRML 11.1.4)
- Note: Monte Carlo and Importance Sampling are NOT sampling methods!
  - .. that is, not used for generating samples but only for computing expectations using samples

# Limitations of Basic Sampling Methods

- Transformation based methods: Usually limited to drawing from standard distributions

- Rejection Sampling and Importance Sampling: Require good proposal distributions



Rejection Sampling

Importance Sampling

- Difficult to find good prop. distr. especially when $z$ is high-dim. (e.g., models with many params)

  - In high dimensions, most of the mass of $p(z)$ is concentrated in a tiny region of the $z$ space

  - Difficult to *a priori* know what those regions are, thus difficult to come up with good proposal dist.

- A solution to these: MCMC methods

# Markov Chain Monte Carlo (MCMC)

- Goal: Generate samples from some target distribution $p(z) = \frac{\tilde{p}(z)}{Z}$, where $z$ is high-dimensional

- Assume we can evaluate $p(z)$ at least up to a proportionality constant (i.e., can compute $\tilde{p}(z)$)

- Basic idea: MCMC uses a Markov Chain which, when converged, starts giving samples from $p(z)$

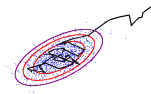$$\underbrace{z^{(1)} \to z^{(2)} \to z^{(3)} \to}_{\text{initial samples typically garbage}} \cdots \to \underbrace{z^{(L-2)} \to z^{(L-1)} \to z^{(L)}}_{\text{after convergence, actual samples from } p(z)}$$

- Given a current sample $z^{(\ell)}$ from the chain, MCMC generates the next sample $z^{(\ell+1)}$ as
  - Use a proposal distribution $q(z|z^{(\ell)})$ to generate a candidate sample $z^*$
  - Accept/reject $z^*$ as the next sample based on an acceptance criterion (will see later)
  - If accepted, $z^{(\ell+1)} = z^*$. If rejected, $z^{(\ell+1)} = z^{(\ell)}$

- Note that in MCMC, the proposal distribution $q(z|z^{(\ell)})$ depends on the previous sample (unlike methods such as rejection sampling)

# MCMC: The Basic Scheme

- MCMC chain run infinitely long (i.e., post-convergence) will give ONE sample from the target $p(z)$



- But we usually require <u>several samples</u> to approximate $p(z)$. How do we get those?

  - Start at an initial $z^{(0)}$. Using a prop. dist. $q(z^{(\ell+1)}|z^{(\ell)})$, run the chain long enough, say $T_1$ steps
  - Discard the first $(T_1 - 1)$ samples (called "burn-in" samples) and take the last sample $z^{(T_1)}$
  - Continue from $z^{(T_1)}$ up to $T_2$ steps, discard intermediate samples, take the last sample $z^{(T_2)}$
    - This helps ensure that $z^{(T_1)}$ and $z^{(T_2)}$ are uncorrelated
  - Repeat the same for a total of $S$ times
  - In the end, we have $S$ i.i.d. samples from $p(z)$, i.e., $z^{(T_1)}, z^{(T_2)}, \ldots, z^{(T_S)} \sim p(z)$
  - Note: Good choices for $T_1$ and $T_i - T_{i-1}$ are usually based on heuristics
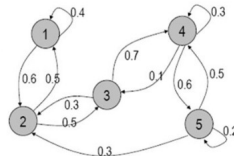  - Note: MCMC is an approximate method because we don't usually know what $T_1$ is "long enough"

## MCMC: Some Basic Theory

- A first order Markov Chain assumes $p(z^{(\ell+1)}|z^{(1)}, \ldots, z^{(\ell)}) = p(z^{(\ell+1)}|z^{(\ell)})$

- A 1st order Markov Chain $z^{(0)}, z^{(1)}, \ldots, z^{(L)}$ is a sequence of r.v.'s and is defined by
  - An initial state distribution $p(z^{(0)})$
  - A Transition Function (TF): $T_\ell(z^{(\ell)} \to z^{(\ell+1)}) = p(z^{(\ell+1)}|z^{(\ell)})$.

- TF defines a distribution over the values of next state given the value of the current state

- Assuming a discrete state-space, the TF is defined by a $K \times K$ probability table

Transition probabilities can be defined using a $K$x$K$ table if **z** is a discrete r.v. with $K$ possible values

$$
\begin{array}{c}
\quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \\
\begin{array}{c}1\\2\\3\\4\\5\end{array}
\begin{bmatrix}
0.4 & 0.6 & 0.0 & 0.0 & 0.0 \\
0.5 & 0.0 & 0.5 & 0.0 & 0.0 \\
0.0 & 0.3 & 0.0 & 0.7 & 0.0 \\
0.0 & 0.0 & 0.1 & 0.3 & 0.6 \\
0.0 & 0.3 & 0.0 & 0.5 & 0.2
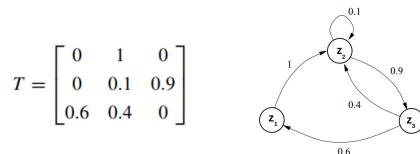\end{bmatrix}
\end{array}
$$



- Homogeneous Markov Chain: The TF is the same for all $\ell$, i.e., $T_\ell = T$

## MCMC: Some Basic Theory

- Consider the following simple TF with $K = 3$ (want to sample from a multinoulli)

$$T = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0.1 & 0.9 \\ 0.6 & 0.4 & 0 \end{bmatrix}$$



- Consider the initial state distribution $p(\mathbf{z}^{(0)}) = p(z_1^{(0)}, z_2^{(0)}, z_3^{(0)}) = [0.5, 0.2, 0.3]$
- Easy to see that $p(\mathbf{z}^{(0)}) \times T = [0.2, 0.6, 0.2] \Rightarrow$ distribution of $\mathbf{z}^{(1)}$
- Also easy to see that, after a few (say $m$) iterations, $p(\mathbf{z}^{(0)}) \times T^m = [0.2, 0.4, 0.4] = p(\mathbf{z})$ (say)
- For the above $T$, <u>any choice</u> of $p(\mathbf{z}^{(0)})$ leads to multinoulli $p(\mathbf{z})$i with $\pi = [0.2, 0.4, 0.4]$
  - Such a $p(\mathbf{z})$ is called the stationary/invariant distribution of this Markov Chain
- A Markov Chain has a stationary distribution if $T$ has the following properties
  - Irreducibility: $T$'s graph is connected (ensures reachability from anywhere to anywhere)
  - Aperiodicity: $T$'s graph has no cycles (ensures that the chain isn't trapped in cycles)

## MCMC: Some Basic Theory

- A sufficient (but not necessary) condition: A Markov Chain with transition function $T$ has stationary distribution $p(z)$ if $T$ satisfies Detailed Balance

- For any two states $z$ and $z'$, the Detailed Balanced condition is

$$p(z)T(z \rightarrow z') = p(z')T(z' \rightarrow z)$$

- Integrating out (or summing over) both sides w.r.t. $z'$ gives

$$p(z) = \int p(z')T(z' \rightarrow z)dz'$$

- Therefore $p(z)$ is a stationary distribution of this chain

- Thus a Markov Chain with detailed balance will always converge to a stationary distribution

# Some MCMC Algorithms

# Metropolis-Hastings (MH) Sampling (Hastings, 1970)

- Suppose we wish to generate samples from a distribution $p(\mathbf{z}) = \frac{\tilde{p}(\mathbf{z})}{Z_p}$

- Assume a proposal distribution $q(\mathbf{z}|\mathbf{z}^{(\tau)})$, e.g., $\mathcal{N}(\mathbf{z}|\mathbf{z}^{(\tau)}, \sigma^2\mathbf{I}_D)$

- In each step, draw $\mathbf{z}^* \sim q(\mathbf{z}|\mathbf{z}^{(\tau)})$ and accept $\mathbf{z}^*$ with probability

$$A(\mathbf{z}^*, \mathbf{z}^{(\tau)}) = \min\left(1, \frac{\tilde{p}(\mathbf{z}^*)q(\mathbf{z}^{(\tau)}|\mathbf{z}^*)}{\tilde{p}(\mathbf{z}^{(\tau)})q(\mathbf{z}^*|\mathbf{z}^{(\tau)})}\right)$$

- The acceptance probability makes intuitive sense
  - It favors accepting $\mathbf{z}^*$ if $\tilde{p}(\mathbf{z}^*)$ has a higher value than $\tilde{p}(\mathbf{z}^{(\tau)})$
  - Unfavors $\mathbf{z}^*$ if the proposal distribution $q$ unduly favors its generation (i.e., if $q(\mathbf{z}^*|\mathbf{z}^{(\tau)})$ is large)
  - Favors $\mathbf{z}^*$ if we can "reverse" to $\mathbf{z}^{(\tau)}$ from $\mathbf{z}^*$ (i.e., if $q(\mathbf{z}^{(\tau)}|\mathbf{z}^*)$ is large). Needed for good "mixing"

- Transition function of this Markov Chain: $T(\mathbf{z}^{(\tau)} \to \mathbf{z}^*) = A(\mathbf{z}^*, \mathbf{z}^{(\tau)})q(\mathbf{z}^*|\mathbf{z}^{(\tau)})$

- Exercise: Show that $T(\mathbf{z} \to \mathbf{z}^{(\tau)})$ satisfies the detailed balance property

$$T(\mathbf{z} \to \mathbf{z}^{(\tau)})p(\mathbf{z}) = T(\mathbf{z}^{(\tau)} \to \mathbf{z})p(\mathbf{z}^{(\tau)})$$

# The MH Sampling Algorithm

- Initialize $z^{(0)}$ randomly

- For $\ell = 0, \ldots, L-1$

  - Sample $z^* \sim q(z^*|z^{(\ell)})$ and $u \sim \text{Unif}(0,1)$

  - If $u < A(z^*, z^{(\ell)}) = \min\left(1, \frac{\tilde{p}(z^*)q(z^{(\ell)}|z^*)}{\tilde{p}(z^{(\ell)})q(z^*|z^{(\ell)})}\right)$

$$z^{(\ell+1)} = z^* \qquad \text{(meaning: accepting with probability } A(z^*, z^{(\ell)}))$$

    else

$$z^{(\ell+1)} = z^{(\ell)}$$

# MH Sampling in Action: A Toy Example..

Target $p(\boldsymbol{z}) = \mathcal{N}\left(\begin{bmatrix} 4 \\ 4 \end{bmatrix}, \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}\right)$, Proposal $q(\boldsymbol{z}^{(t)}|\boldsymbol{z}^{(t-1)}) = \mathcal{N}\left(\boldsymbol{z}^{(t-1)}, \begin{bmatrix} 0.01 & 0 \\ 0 & 0.01 \end{bmatrix}\right)$
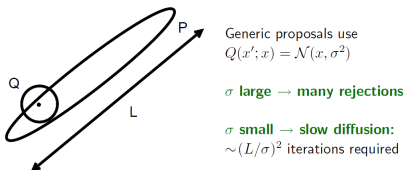
# MH Sampling: Some Comments

- If proposal distrib. is symmetric, we get Metropolis Sampling algorithm (Metropolis, 1953) with

$$A(\mathbf{z}^*, \mathbf{z}^{(\tau)}) = \min\left(1, \frac{\tilde{p}(\mathbf{z}^*)}{\tilde{p}(\mathbf{z}^{(\tau)})}\right)$$

- Some limitations of MH sampling

  - MH can have a very slow convergence. Figure below: $P$ is the target dist., $Q$ is the proposal



Generic proposals use
$Q(x'; x) = \mathcal{N}(x, \sigma^2)$

$\sigma$ large $\rightarrow$ **many rejections**

$\sigma$ small $\rightarrow$ **slow diffusion:**
$\sim (L/\sigma)^2$ iterations required

  - Computing acceptance probability can be expensive: When $p(\mathbf{z}) = \frac{\tilde{p}(\mathbf{z})}{Z_p}$ represents a posterior distribution of some model, $\tilde{p}$ is the unnormalized posterior that depends on all the data (note: a lot of recent work on speeding up this step using subsets of data[*])

---

[*] Austerity in MCMC Land: Cutting the Metropolis-Hastings Budget (Korattikara et al, 2014)

# Gibbs Sampling (Geman & Geman, 1984)

- Suppose we wish to sample from a joint distribution $p(\boldsymbol{z})$ where $\boldsymbol{z} = (z_1, z_2, \ldots, z_M)$
- However, suppose we can't sample from $p(\boldsymbol{z})$ but can sample from each conditional $p(z_i | \boldsymbol{z}_{-i})$
  - Can we done easily if we have a locally conjugate model
- For Gibbs sampling, the proposal is the conditional distribution $p(z_i | \boldsymbol{z}_{-i})$
- Gibbs sampling samples from these conditionals in a cyclic order
- Gibbs sampling is equivalent to Metropolis Hastings sampling with acceptance prob. $= 1$

$$A(\boldsymbol{z}^*, \boldsymbol{z}) = \frac{p(\boldsymbol{z}^*)q(\boldsymbol{z}|\boldsymbol{z}^*)}{p(\boldsymbol{z})q(\boldsymbol{z}^*|\boldsymbol{z})} = \frac{p(z_i^*|\boldsymbol{z}_{-i}^*)p(\boldsymbol{z}_{-i}^*)p(z_i|\boldsymbol{z}_{-i}^*)}{p(z_i|\boldsymbol{z}_{-i})p(\boldsymbol{z}_{-i})p(z_i^*|\boldsymbol{z}_{-i})} = 1$$

where we use the fact that $\boldsymbol{z}_{-i}^* = \boldsymbol{z}_{-i}$

## Gibbs Sampling: Sketch of the Algorithm

$M$: Total number of variables, $T$: number of Gibbs sampling steps

1. Initialize $\{z_i : i = 1, \ldots, M\}$
2. For $\tau = 1, \ldots, T$:
   - Sample $z_1^{(\tau+1)} \sim p(z_1 | z_2^{(\tau)}, z_3^{(\tau)}, \ldots, z_M^{(\tau)})$.
   - Sample $z_2^{(\tau+1)} \sim p(z_2 | z_1^{(\tau+1)}, z_3^{(\tau)}, \ldots, z_M^{(\tau)})$.
     $\vdots$
   - Sample $z_j^{(\tau+1)} \sim p(z_j | z_1^{(\tau+1)}, \ldots, z_{j-1}^{(\tau+1)}, z_{j+1}^{(\tau)}, \ldots, z_M^{(\tau)})$.
     $\vdots$
   - Sample $z_M^{(\tau+1)} \sim p(z_M | z_1^{(\tau+1)}, z_2^{(\tau+1)}, \ldots, z_{M-1}^{(\tau+1)})$.
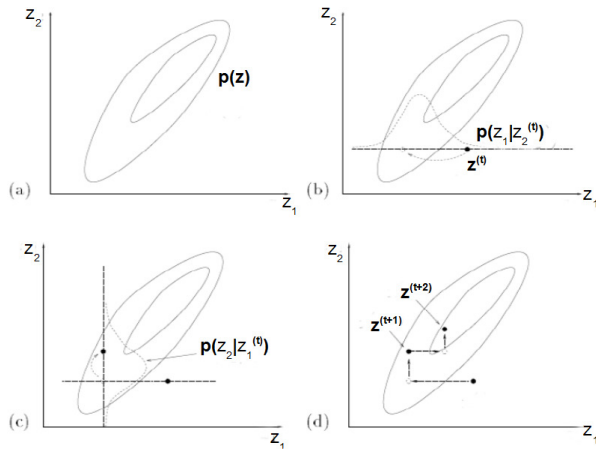
Note: When sampling each variable from its conditional posterior, we use the most recent values of all other variables (this is akin to a co-ordinate ascent like procedure)

Note: Order of updating the variables *usually* doesn't matter (but see "Scan Order in Gibbs Sampling: Models in Which it Matters and Bounds on How Much" from NIPS 2016)

# Gibbs Sampling: A Simple Example

Can sample from a 2-D Gaussian using 1-D Gaussians (recall that if the joint distribution is a 2-D Gaussian, conditionals will simply be 1-D Gaussians)

# Gibbs Sampling: Some Comments

- One of the most popular MCMC algorithm

- Very easy to derive and implement for locally conjugate models

- Many variations exist, e.g.,

  - Blocked Gibbs: sample multiple variables jointly (sometimes possible)

  - Rao-Blackwellized Gibbs: Can collapse (i.e., integrate out) the unneeded variables while sampling. Also called "collapsed" Gibbs sampling

  - MH within Gibbs

- Instead of sampling from the conditionals, an alternative is to use the mode of the conditional.

  - Called the "Iterative Conditional Mode" (ICM) algorithm (doesn't give the posterior though)

## Next Class

- Using posterior's gradient info in sampling algorithms

- Online MCMC algorithms

- Recent advances in MCMC

- Some other practical issues