

SPR Day 20

Sampling Methods

Reading Bishop PRML 11.1.1 Standard distrib.

Skim for broader knowledge

11.1.2-5 rejection importance

11.2.1 Markov chains

Topics

Monte Carlo estimation

easy way to approx. hard expectations

Graphical Models + Ancestral Sampling

Method: Uniform Sample + Inverse CDF

Transformations of Sampled Variables

Markov Chain Monte Carlo

Why samples?

- Represent distributions that have no simple analytical form

- Easily estimate expectations

Task: You have model for r.v. z (continuous or discrete) given by $p(z)$

You want to know expected value of function $f(z)$

Ideal:
$$\bar{f} = \mathbb{E}_{z \sim p(z)} [f(z)]$$
$$= \int p(z) f(z) dz$$

integral might be hard!

Monte Carlo Estimate:

$$\hat{f} = \frac{1}{S} \sum_{s=1}^S f(z^s)$$

where $z^s \stackrel{iid}{\sim} p(z)$

Uses S samples, z^1, z^2, \dots, z^S each iid from $p(z)$

What is expected value of MC estimate?

$$E[\hat{f}] = E_{z^1, \dots, z^S \sim p(z)} \left[\frac{1}{S} \sum_{s=1}^S f(z^s) \right]$$

$$= \frac{1}{S} \sum_{s=1}^S E_{z^s \sim p(z)} [f(z^s)]$$

by linearity of expectations and iid assumption

$$= \frac{1}{S} \sum_{s=1}^S \int p(z) f(z) dz$$

$$= \frac{1}{S} \sum_{s=1}^S \bar{f}$$

by definition of \bar{f}

$$= \bar{f}$$

MC estimate is UNBIASED.

What is variance of MC estimate?

$$\text{Var}[\hat{f}] = \frac{1}{S} \text{Var}[f(z)]$$

variance of function f under original $p(z)$

So, for large S , MC will be "close" to ideal \bar{f} !

Variance of MC estimate: Derivation

$$\begin{aligned} \text{Var}[\hat{f}] &= \mathbb{E}_{z^1, \dots, z^S} [(\hat{f} - \bar{f})^2] \\ &= \mathbb{E}_{z^1, \dots, z^S} [\hat{f}^2 - 2\hat{f}\bar{f} + \bar{f}^2] \end{aligned}$$

$$= \mathbb{E}_{z^1, \dots, z^S} [\hat{f}^2] - \bar{f}^2$$

$$= \mathbb{E} \left[\frac{1}{S} \sum_{s=1}^S f(z^s) \right]^2 - \bar{f}^2$$

$$= \frac{1}{S^2} \mathbb{E} \left[\left(\underbrace{f(z^1) + \dots + f(z^S)}_{\text{self product}} \right) \left(\underbrace{f(z^1) + \dots + f(z^S)}_{\text{cross product}} \right) \right]$$

Simplify cross terms

$$\begin{aligned} &\text{bc } z^1 \text{ is indep of } z^S \\ &\mathbb{E}[f(z^1)f(z^S)] \\ &= \mathbb{E}[f(z^1)]\mathbb{E}[f(z^S)] = \bar{f}^2 \end{aligned}$$

$$= \frac{1}{S^2} \left[\underbrace{S \mathbb{E}[f(z)^2]}_{\text{self terms } f(z^i)f(z^i)} + \underbrace{(S^2 - S) \bar{f}^2}_{\text{cross terms } f(z^i)f(z^j) \text{ } i \neq j} \right] - \bar{f}^2$$

expand + group into self + cross terms

$$= \frac{1}{S} \mathbb{E}[f(z)^2] + \cancel{\frac{S^2 - S}{S^2} \bar{f}^2} - \cancel{\frac{S^2}{S^2} \bar{f}^2} \quad (\text{cancel these})$$

$$= \frac{1}{S} \mathbb{E}_{p(z)} [f(z)^2] - \frac{1}{S} \bar{f}^2 = \frac{1}{S} \text{Var}_{p(z)} [f(z)]$$

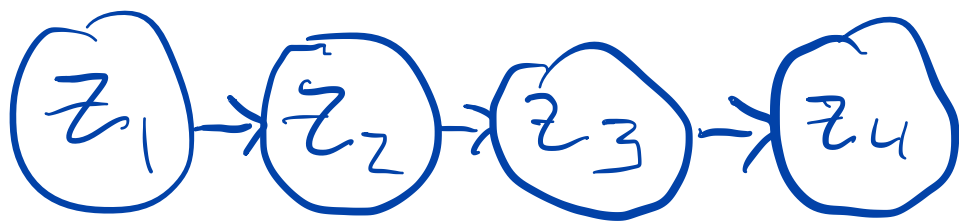
How to Sample from Models with 5 Multiple Variables

Suppose we have model of many r.v., we wish to sample from joint distribution.

$p(z_1, z_2, \dots, z_T)$ e.g. draw from Markov model.

Suppose further we know conditional independence assumptions, and can use these to define a graph (must be directed, acyclic)

where vertices/nodes are random vars
directed edges represent conditional dependence assumptions



Markov model
with $T=4$

Can use graph to re-write joint as:

$$p(z_{1:4}) = \prod_{i \in V} p(z_i | z_{pa(i, E)})$$

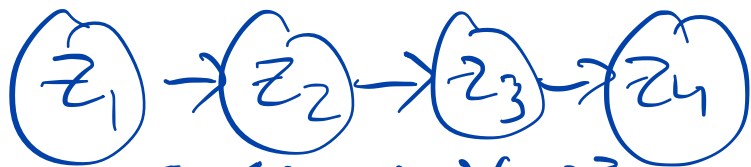
V is vertex list

E is edge list

$pa(i, E)$: "parents" of node i

$\{j : (j, i) \in E\}$

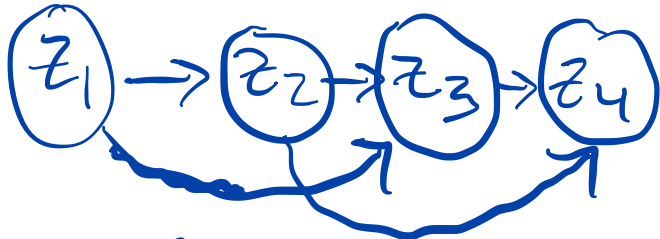
So for our 1st order Markov model 6



$$E = \{(1,2), (2,3), (3,4)\}$$

$$P(z_{1:4}) = P(z_1)P(z_2|z_1)P(z_3|z_2)P(z_4|z_3)$$

2nd order Markov



$$E = \{(1,2), (2,3), (3,4), (1,3), (2,4)\}$$

$$P(z_{1:4}) = P(z_1)P(z_2|z_1)$$

- $P(z_3|z_2, z_1)$

- $P(z_4|z_3, z_2)$

We call this representation a directed graphical model

Useful for sampling and many other things

(computing marginals $P(z_t)$
joints $P(z_1, z_2)$
conditionals $P(z_2|z_1)$)

How to sample?

Assume we have way to sample from
the simple conditional $P(z_i | z_{pa(i, E)})$
for all nodes i

Then, we can sample $z_{1:T} \sim p(z_{1:T})$
using ancestral sampling.

Arrange node indices in order $1, 2, \dots, T$
s.t. for any index $j \in \{1, 2, \dots, T\}$
if i is a parent of j , then $i < j$

This order is always achievable if graph has
no cycles (remember, our directed graph is acyclic)

Ancestral Sampling: \leftarrow assume
TOPO SORT order!

for i in $1, 2, \dots, T$:

$$z_i \sim p(z_i | z_{\text{pa}(i, E)})$$

return $[z_1, z_2, \dots, z_T]$

Guaranteed to sample from joint $p(z_{1:T})$!

What can we do with samples
from a joint distribution?

(0) Compute a Monte Carlo (MC) expectation

(1) Sample from a marginal

Given S samples $z_{1:T}^{(1)}, \dots, z_{1:T}^{(S)}$

we can get samples of z_t

by just keeping $z_t^{(1)}, z_t^{(2)}, \dots, z_t^{(S)}$

(2) Sample from a conditional

Given S samples $z_{1:T}^{(1)}, z_{1:T}^{(2)}, \dots, z_{1:T}^{(S)}$

we can get samples from $p(z_t | z_u = k)$

by keeping $\{z_t^{(s)} : z_u^{(s)} = k\}$

and discarding others

Sampling via inverse CDF

9

Consider real valued random var x

with pdf: $p(x)$

and cumulative distribution function

$$\text{cdf}(a) = p(x \leq a) = \int_{-\infty}^a p(x) dx$$

Properties of cdf: $0 \leq \text{cdf}(a) \leq 1$
for all $a \in \mathbb{R}$

If cdf function F is invertible analytically,
we can sample using the simple transformation

$$\textcircled{1} u \sim \text{Unit}([0, 1]) \quad \textcircled{2} x \leftarrow F^{-1}(u)$$

Why?

$$p(x \leq a) = \int_0^{F(a)} 1 du = F(a)$$

thus, x by construction must have cdf F

Sampling via transformations ^{"change of variables"} 10

Suppose we have a "target" rand var X
and a "source" rand. var. U

U has known pdf $f: \mathbb{R} \rightarrow [0, +\infty]$
that is easy to evaluate

U has known sampling procedure DRAW_U
so we can easily generate
 u^1, u^2, \dots, u^S

We have an invertible transformation function $T: \mathbb{R} \rightarrow \mathbb{R}$

We generate samples x via:
(1) $u \sim \text{DRAW}_U$
(2) $x \leftarrow T(u)$

What is pdf of rand. var. x ?

Remember CHANGE of VARIABLES from calculus

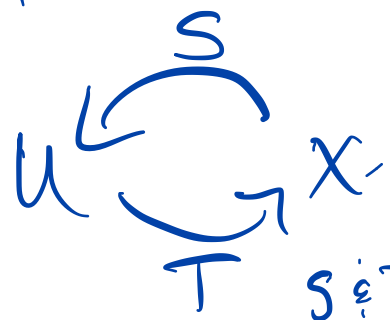
For any smooth function g

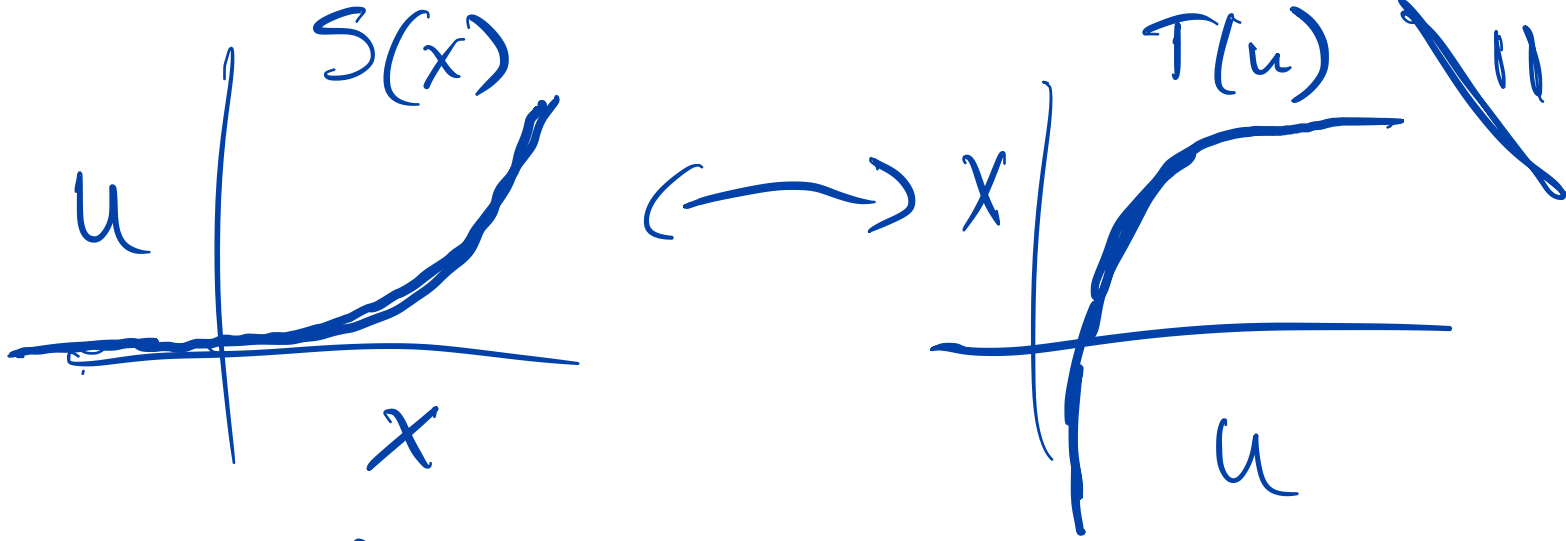
$$\int_a^b g(u) du = \int_{T(a)}^{T(b)} g(S(x)) S'(x) dx$$

where

$$u = S(x)$$

$$\frac{\partial u}{\partial x} = S'(x) = \frac{\partial}{\partial x} S(x)$$





Invertible functions are monotonic bc they need to be one-to-one.

So interval (a, b) in u domain is either from $T(a), T(b)$ if $S'(x) \geq 0$ ^{pos slope}
 or $T(b), T(a)$ if $S'(x) \leq 0$ ^{neg slope}

So, using change of vars, we have:

$$P(a \leq u \leq b) = \int_a^b f(u) du \quad f \text{ known pdf for } u$$

$$\text{Use abs value to fix sign issue} = \int_{T(a)}^{T(b)} f(S(x)) S'(x) dx$$

Thus

$$P(\min \leq x \leq \max) = \int_{\min}^{\max} f(S(x)) |S'(x)| dx$$

$$\text{pdf}(x) = f(S(x)) |S'(x)| \quad \text{uses known simpler pdf } f \text{ \& transform } S$$

Markov Chain Monte Carlo ¹²

Goal :

Want sample z from a target distribution $p^*(z)$

but we may only know the pdf up to a constant c (does not depend on z)

$$p^*(z) = c \tilde{p}(z)$$

unknown

known function, easy to evaluate

$$\log p^*(z) = \log c + \log \tilde{p}(z)$$

Insight :

We can sample a sequence

by

$z_1 \leftarrow$ reasonable guess



$z_t | z_{t-1} \sim T$

Markov proposal distribution

if we are careful about choosing proposal distrib T , then marginal $p(z_s) = p^*(z)$

Markov Chain Monte Carlo

uses Markov distrib.

T to propose

$$z_t / z_{t-1}$$

purpose is to draw

samples from a

target distrib. $p^*(z)$

We want $p^*(z)$ to be the stationary distribution of the Markov chain.

Stationary distribution

$p(z)$ is stationary distrib. of Markov operator T if

Discrete case
 $z \in \{1, 2, \dots, K\}$

$$p(z_{t+1} = k) = \sum_{j=1}^K p(z_t = j) T(z_{t+1} = k | z_t = j)$$

joint prob of being in state j & transition to state k

Continuous case

$$p(z_{t+1}) = \int_{z_t \in \Omega} p(z_t) T(z_{t+1} / z_t) dz_t$$

When does a unique stationary distribution exist? When the Markov chain is

ergodic, which means

if we start in state i at $t=0$, then for some time $T_0 > 0$,

we have for every state k

$$P(Z_\tau = k | Z_0 = i) > 0$$

for all $\tau > T_0$

Intuition: Need to be able to get from any state to any other state

Key Technical Conditions:

- T must be irreducible: path exists between any two states

- T must be aperiodic: no cycles



without shortcuts or "longcuts"
need to transition $A \rightarrow B$ in diff. number of steps between any states