

# Lecture 18 - scribbles

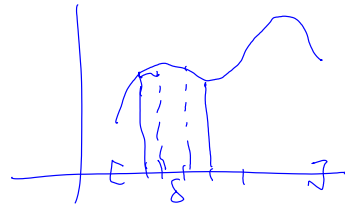
Friday, November 10, 2017

13:18

today : • Sampling  
• MCMC

aside on numerical computing: 1d "is easy"

- numerical integration in 1d  
for function  $L$ -Lipschitz



error  $\propto L\delta$

$$\text{error} \leq \epsilon \Rightarrow \text{complexity } O\left(\frac{1}{\epsilon}\right)$$

use  $\delta \approx \epsilon$

but grid in dimension  $d$   $O\left(\frac{1}{\epsilon^d}\right)$

- global optimization in 1d



error  $\propto L\delta$

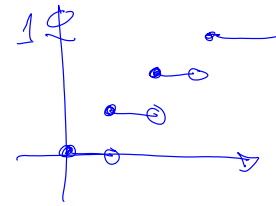
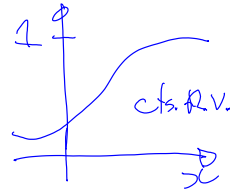
complexity (# of evaluations)  $O\left(\frac{1}{\epsilon}\right)$

How to sample?

- 1)  $X \sim \text{Unif}([0,1])$   $\rightarrow$  pseudo-random generator "rand"
- 2)  $X \sim \text{Bernoulli}(p)$   $X = \mathbb{1}\{U \leq p\}$  where  $U \sim \text{Unif}([0,1])$
- 3) universe transform sampling :

let  $F$  be target cdf of distribution  $p$  for  $X$   $F(x) \triangleq P\{X \leq x\}$  ( $x \in \mathbb{R}$ )  
 (first, suppose  $F$  is invertible)

let  $X \triangleq F^{-1}(U)$  with  $U \sim \text{Unif}(0,1)$



claim that  $X$  has cdf  $F(x)$

$F$  is invertible

proof:  $P\{X \leq y\} = P\{F^{-1}(U) \leq y\} \stackrel{b}{=} P\{U \leq F(y)\} = F(y)$

[if  $F$  is not invertible, define  $X \triangleq \min\{x : F(x) \geq U\}$ ]

(recall  $F$  is cfs from the right)

example:

want  $X \sim \text{Exp}(1)$

density  $p(x) = \lambda e^{-\lambda x} \mathbb{1}_{\mathbb{R}^+}(x)$

cdf  $F(x) = 1 - e^{-\lambda x}$

inverse  $F^{-1}(u) = -\frac{1}{\lambda} \ln(1-u)$

Multivariate distribution?

generalize above trick using "chain rule"

$X_{1:p}$  (dim  $p$ ) cdf  $F(x_{1:p}) \triangleq P\{X_1 \leq x_1, \dots, X_p \leq x_p\}$

$F_{X_{1:p}}(x_{1:p}) = F_{X_1}(x_1) F_{X_2|X_1}(x_2|x_1) \cdots F_{X_p|X_{1:p-1}}(x_p|x_1, \dots, x_{p-1})$

$F_{X_2|X_1}(x_2|x_1) \triangleq P\{X_2 \leq x_2 | X_1 \leq x_1\}$

recall also  $U_1, \dots, U_n \stackrel{iid}{\sim} \text{Unif}$ .

could use  $u_1, \dots, u_p \stackrel{iid}{\sim} \text{Unif}$

$$x_1 = F_{x_1}^{-1}(u_1)$$

$\vdots$

$$x_p = F_{x_p | x_{1:p-1}}^{-1}(u_p | x_{1:p-1})$$

is very complicated function

(curse of dimensionality)

[aside: "copulas"  $\rightarrow$  model for multivariate dependence with uniform marginals]

exception is multivariate Gaussian:

$$N(\mu, \Sigma) \quad \Sigma = U \Lambda U^T$$

(where  $U U^T = I_p$   
and  $\Lambda$  is diagonal)

(Cholesky decomposition)

$$L = U \Lambda^{1/2}$$

$$\Sigma = L L^T$$

generate  $V \sim N(0, I_p)$   
(ie.  $v_p \stackrel{iid}{\sim} N(0,1)$ )

$$X \triangleq U \Lambda^{1/2} V + \mu$$

$$\begin{aligned} \mathbb{E}X &= \mu \\ \text{Cov}(X) &= \sum_{\substack{X \sim N(\mu, \Sigma) \\ \Rightarrow}} \end{aligned}$$

Box-Muller transform to sample Gaussian: (2d)

$$\begin{aligned} x &= r \cos \theta \\ y &= r \sin \theta \end{aligned}$$

$$\begin{pmatrix} x \\ y \end{pmatrix} \sim N(0, I)$$

$$\leadsto r^2 \sim \text{Exp}(1)$$

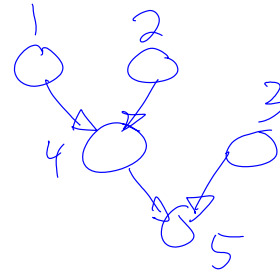
$\theta$  is  $\text{Unif}([0, 2\pi])$

Sampling for DBM is easy: ancestral sampling

$$(x_1, \dots, x_p) \sim p \in \mathcal{F}(\mathcal{G}) \quad \text{where } \mathcal{G} \text{ is a DAG}$$

$$p(x_1, \dots, x_p) = \prod_{i=1}^p p(x_i | x_{\pi_i})$$

suppose WLOG,  $1, \dots, p$  is a top sort of  $G$



ancestral sampling:

for  $i=1, \dots, p$  do

sample  $x_i \sim p(x_i = \cdot | x_{\pi_i})$

end

(can show (by induction) that  $(x_1, \dots, x_p)$  has distribution  $p$ )

aside: formal tool to show that  $(x_1, x_2)$  has right distribution

2 node example:

$$x_1 \sim p(x_1)$$

$$x_2 | x_1 \sim p(x_2 | x_1)$$

can show that two R.V. are equal in dist. i.e.  $U \stackrel{d}{=} V$

$$\Leftrightarrow \mathbb{E}_U[f(U)] = \mathbb{E}_V[f(V)] \text{ for all functions in a big enough class (eg. Cb. \& bounded functions)}$$

here, I want to show that

$$\mathbb{E}[f(x_1, x_2)] = \mathbb{E}_{x_1}[\mathbb{E}_{x_2|x_1}[f(x_1, x_2) | x_1]]$$

$$= \int_{\mathcal{X}_1} p_{x_1}(x_1) dx_1 \left[ \int f(x_1, x_2) p(x_2 | x_1) dx_2 \right]$$

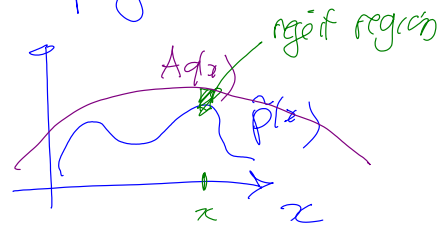
$$\begin{aligned}
 &= \int_{x_1} p_{x_1}(x_1) dx_1 \left[ \int_{x_2} f(x_1, x_2) p(x_2|x_1) dx_2 \right] \\
 &= \int_{x_1} \int_{x_2} f(x_1, x_2) \underbrace{[p_{x_1}(x_1) p(x_2|x_1)]}_{p(x_1, x_2)} dx_1 dx_2
 \end{aligned}$$

⊗ important side note: when sample from joint, you are also sampling from marginal  
 i.e.  $(X, Y) \sim p(x, y)$   
 then looking at  $X$  by itself, you have  $X \sim p(x)$

rejection sampling:

say  $p(x) = \frac{\tilde{p}(x)}{Z_p}$ ; let's say we can find  $q(x)$ , a dist. we can easily sample from

stb.  $A q(x) \geq \tilde{p}(x)$   
 proposal target



rule:

- sample  $X \sim q(x)$
- Accept with probability  $\frac{\tilde{p}(x)}{Aq(x)} \in [0, 1]$   
 (reject o.w.)

let's show accepted samples have correct dist.

(say  $X$  is discrete)

$$\begin{aligned} P\{X=x, X \text{ is accepted}\} &= P\{X \text{ is accepted} | X=x\} P(X=x) \\ &= \frac{\tilde{p}(x)}{Aq(x)} \cdot q(x) \\ &= \frac{\tilde{p}(x)}{A} \end{aligned}$$

$$P\{X \text{ is accepted}\} = \sum_x \frac{\tilde{p}(x)}{A} = \frac{Z_{\tilde{p}}}{A} \quad \left( \begin{array}{l} \rightarrow \text{marginal probs. of acceptance} \\ \text{[want this to be high]} \end{array} \right)$$

$$P\{X=x | X \text{ is accepted}\} = \frac{\tilde{p}(x)}{Z_{\tilde{p}}} = p(x)$$

application to DSM:

say we want to sample from  $p(x | \bar{x}_E)$  and  $p \in \mathcal{S}(G)$

here  $\tilde{p}(x) = p(x_E, \bar{x}_E) \delta(x_E, \bar{x}_E)$

here  $Z_p = p(\bar{x}_E)$

let  $q(x)$  be from original joint in DSM using ancestral sampling

$$q(x) \geq \tilde{p}(x) \quad \forall x \quad [\text{take } A=1]$$

alg: 

- do ancestral sampling
- accept if  $x_E = \bar{x}_E$

$$\text{acceptance prob} = \frac{\tilde{p}(x)}{q(x)} = \begin{cases} 1 & \text{if } x_E = \bar{x}_E \\ 0 & \text{o.w.} \end{cases}$$

[i.e. reject when  $x_E \neq \bar{x}_E$ ]

PS accept } marginally  $= \sum_A = p(\bar{x}_E)$

Importance sampling:

in context of computing  $E_p[f(x)] = \mu \quad X \sim p$

$\sim$  can "weight" sample  $X^{(i)}$

$$E_p[f(x)] = \int f(x)p(x)dx = \int f(x) \frac{p(x)}{q(x)} q(x) dx \quad \text{for some distribution } q$$

$\Rightarrow \text{supp}(q) \supseteq \text{supp}(p)$

$$= E_q \left[ f(y) \frac{p(y)}{q(y)} \right] \quad \text{where } Y \sim q$$

$$\approx \frac{1}{n} \sum_{i=1}^n g(y_i) \quad \text{where } y_i \overset{i.i.d.}{\sim} q$$

and  $g(y) \triangleq f(y) w(y)$

where  $w(y) \triangleq \frac{p(y)}{q(y)}$  "weight"

$$\hat{\mu}_{IS} = \frac{1}{n} \sum_{i=1}^n f(y_i) w_i \quad Y_i \sim q$$

$\downarrow$   
"importance weights"

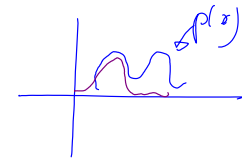
$w_i = \frac{p(y_i)}{q(y_i)}$

$$E[\hat{\mu}] = \mu$$

$$\text{Var}(\hat{\mu}) = \frac{1}{n} \text{Var}_{q(x)} \left[ f(x) \frac{p(x)}{q(x)} \right] = \frac{1}{n} \left[ E_p \left[ f(x)^2 \frac{p(x)}{q(x)} \right] - \mu^2 \right]$$

intuitively, you want  $q(x)$  or  $f(x)p(x)$

issues here when  $q$  small and  $p$  big



extension to un-normalized distributions:

Extension to un-normalized distributions:

$$\begin{aligned}
 p(x) &= \frac{\tilde{p}(x)}{Z_p} & q(x) &= \frac{\tilde{q}(x)}{Z_q} & \mathbb{E}_q \left[ f(y) \frac{p(y)}{q(y)} \right] \\
 & & & & = \mathbb{E}_q \left[ f(y) \frac{\tilde{p}(y)}{\tilde{q}(y)} \right] \frac{Z_q}{Z_p} \\
 & & & & = \mu \cdot \frac{Z_q}{Z_p}
 \end{aligned}$$

estimate  $\frac{Z_q}{Z_p}$  with  $\hat{\frac{Z_q}{Z_p}} \triangleq \frac{1}{n} \sum_{i=1}^n \frac{\tilde{p}(y_i)}{\tilde{q}(y_i)} = \frac{1}{n} \sum_{i=1}^n w_i$

$$\hat{\mu}_{IS0} = \frac{\frac{1}{n} \sum_{i=1}^n f(y_i) w_i}{\frac{1}{n} \sum_{j=1}^n w_j} \hat{\frac{Z_q}{Z_p}} \quad y_i \sim q \quad w_i \triangleq \frac{\tilde{p}(y_i)}{\tilde{q}(y_i)}$$

note:  $\hat{\mu}$  here is (slightly) biased, but asymptotically unbiased (as  $n \rightarrow \infty$ )

- this estimator often have lower variance than  $\hat{\mu}_{IS}$  even when  $Z_p = Z_q = 1$   
(normalization "stabilizes"  $\hat{\mu}_{IS}$ ) new weights  $\hat{w}_i = \frac{w_i}{\sum_j w_j} \in [0, 1]$

Variance reduction:

control variate trick

want  $\hat{\mu}$  for  $\mu = \mathbb{E}[X]$ ; say have  $Y$  s.t.  $\mathbb{E}Y$  is cheap to compute and  $Y$  correlated to  $X$

control variate estimator  $\hat{\mu}_c \triangleq \alpha(X - Y) + \mathbb{E}Y$  (1-sample version)



consider estimator  $\delta_\alpha \triangleq \alpha(X-Y) + \mathbb{E}Y$  (1-sample version)

$$\alpha \in [0,1]$$

$$\mathbb{E}\delta_\alpha = \alpha \mathbb{E}X + (1-\alpha) \mathbb{E}Y \quad (\text{convex combo})$$

if  $\alpha=1$  then  $\delta_\alpha$  is unbiased

$$\delta_\alpha = X + \underbrace{(\mathbb{E}Y - Y)}_{\text{correction}}$$

$$\text{Var}(\delta_\alpha) = \alpha^2 [\text{Var}X + \text{Var}Y - 2\text{Cov}(X,Y)]$$

↓  
smaller  $\alpha$ , reduce the variance

if big enough, then  $\text{var}(\delta_\alpha) \leq \text{var}(X)$

$$\text{recall } \mathbb{E}(\delta_\alpha - \mu)^2 = \text{var}(\delta_\alpha) + \underbrace{(\mathbb{E}[\delta_\alpha] - \mu)^2}_{\text{bias}}$$

$$\alpha=1 \rightsquigarrow \text{SAGA}$$

$$\alpha = \frac{1}{n} \rightsquigarrow \text{SAG}$$

(see lecture 8)

## Rao-Blackwellization

motivation: Rao-Blackwell thm:  $\text{var}(\mathbb{E}[\delta(x)|Z]) \leq \text{var}(\delta(x))$

so if can compute  $\mathbb{E}[\delta(x)|Z]$ ; get better estimator

e.g. say  $X = (Y, Z)$   $Y|Z$  is simple (e.g. Gaussian)

but  $Z$  is complicated

moral of story: if can integrate out analytically some parts, do it!