

Lecture 19 - scribbles

Tuesday, November 8, 2016  
14:28

today: approximate inference - sampling

Why sampling?  $X = (X_1, \dots, X_p)$

a) simulation:  $X^{(i)} \sim p$

b) approximate  $p(x_j)$ ; more generally,

$f: \mathbb{R}^p \rightarrow \mathbb{R}^d$ , approximate  $\mu = \mathbb{E}_p[f(X)]$

e.g.  $f(X) = \mathbb{1}\{X_A = x_A\}$   $\mathbb{E}_p[f(X)] = P\{X_A = x_A\}$

Monte-Carlo integration/estimation:  $\rightarrow$  appears in physics, applied math, ML

to approximate  $\mu = \mathbb{E}_p[f(X)]$

alg:  $n$  samples  $X^{(i)} \stackrel{iid}{\sim} p$   
 estimate  $\hat{\mu} \triangleq \frac{1}{n} \sum_{i=1}^n f(X^{(i)})$

this even true if  $X^{(i)}$  are dependent

properties: 1) unbiased  $\mathbb{E}[\hat{\mu}] = \frac{1}{n} \sum_i \mathbb{E}[f(X^{(i)})] = \mu$

2) expected error:  $\mathbb{E}[\|\hat{\mu} - \mu\|^2] = \mathbb{E}[\frac{1}{n^2} \sum_{i,j} \langle f(X^{(i)}) - \mu, f(X^{(j)}) - \mu \rangle]$   
 (variance in 1d)  
 $\text{tr}(\text{cov}(\hat{\mu}, \hat{\mu}))$  (key independence)

$$\mathbb{E}[\|\hat{\mu} - \mu\|^2] = \frac{\sigma^2}{n}$$

$$= \frac{1}{n^2} \sum_{i,j} \mathbb{E}[\langle f(X^{(i)}) - \mu, f(X^{(j)}) - \mu \rangle]$$

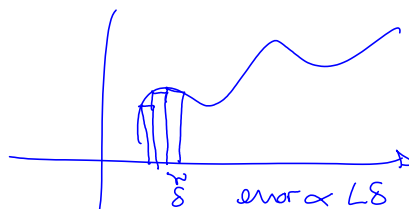
$\underbrace{\mathbb{E}[\|f(X) - \mu\|^2]}_{\triangleq \sigma^2}$   
 $\text{tr}(\text{cov}(f(X), f(X)))$

[aside: asymptotic analysis LLN  $\hat{\mu} \xrightarrow{as} \mu$

CLT ( $\mu \in \mathbb{R}$ )  $\sqrt{n}(\hat{\mu} - \mu) \xrightarrow{d} N(0, \sigma^2)$   
 $\sigma^2 \ll n$

aside on numerical computing: 1d is "easy"

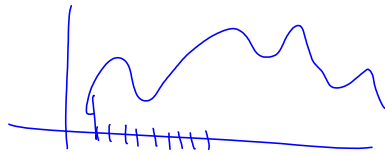
- numerical integration in 1d  
 $L$ -Lipschitz functions



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

but grid in high d,  $O\left(\frac{1}{\epsilon^d}\right) \rightarrow$  curse of dimensionality

• global optimization in 1d



$$\text{error} \propto \text{LS} \Rightarrow \text{complexity } O\left(\frac{1}{\epsilon}\right) \text{ (# of evaluations)}$$

How to sample?

1)  $X \sim \text{Unif}([0,1]) \rightarrow$  pseudo-random generator "rand"

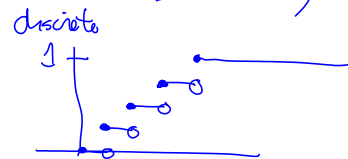
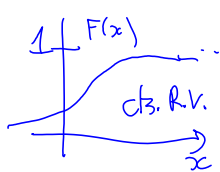
2)  $X \sim \text{Bernoulli}(p) \quad X = \mathbb{1}\{U \leq p\}$  where  $U \sim \text{Unif}([0,1])$

3) inverse transform sampling:

let  $F$  be target cdf of distribution  $F(x) \triangleq \mathbb{P}\{X \leq x\} \quad (x \in \mathbb{R})$

(say  $F$  is invertible)

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



$\Rightarrow$   $X$  has cdf  $F(x)$

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

(aside if  $F$  not invertible, define  $X = \min\{x: F(x) \geq U\}$ )

example:  $X \sim \text{Exp}(\lambda)$  density  $p(x) = \lambda e^{-\lambda x} \mathbb{1}_{\mathbb{R}^+}(x)$

$$\begin{aligned} \text{cdf } F(x) &= 1 - e^{-\lambda x} \\ \text{inverse } F^{-1}(u) &= \frac{-1}{\lambda} \ln(1-u) \end{aligned}$$

Multivariate distribution?

use "chain rule"

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

could use  $U_1, \dots, U_p \stackrel{\text{i.i.d.}}{\sim} \text{Unif}$

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

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

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

(cause of dim.)

is very complicated fct.

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

exception: multivariate Gaussian

$$N(\mu, \Sigma) \quad \Sigma = LL^T \quad \text{generate } V \sim N(0, I)$$

$\downarrow$  define  $X = LL^T V + \mu$

$$X \sim N(\mu, \Sigma)$$

$\Sigma = LL^T$   
 $L = LL^T V$

Box-Muller transform 2d Gaussian  $x = r \cos \theta$   $\begin{pmatrix} x \\ y \end{pmatrix} \sim N(0, I)$   
 $y = r \sin \theta$

$\rightarrow r^2$  is  $\text{Exp}(1)$   
 $\theta$  is  $\text{Unif}([0, 2\pi])$

sampling for DBM is easy: ancestral sampling

$$(x_1, \dots, x_p) \sim p \in \mathcal{J}(G)$$

where  $G$  is DAG

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

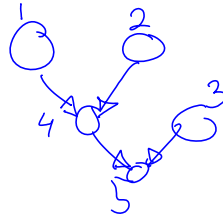
suppose WLOG,  $1, \dots, p$  is a top sort of  $G$

ancestral sampling:

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

draw  $x_i \sim p(x_i | x_{\pi_i})$

end



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

2 node example:

$$(X, Y) \quad X \sim p(x)$$

$$Y|X \sim \tilde{p}(y|x)$$

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

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

here

$$\begin{aligned} \mathbb{E}[f(\tilde{X}, \tilde{Y})] &= \mathbb{E}[\mathbb{E}[f(\tilde{X}, \tilde{Y}) | \tilde{X}]] \\ &= \int_{\tilde{x}} \left[ \int_{\tilde{y}} f(\tilde{x}, \tilde{y}) p_{Y|X}(\tilde{y} | \tilde{x}) d\tilde{y} \right] p_X(\tilde{x}) d\tilde{x} \\ &= \int_{\tilde{x}, \tilde{y}} f(\tilde{x}, \tilde{y}) \underbrace{[p_{Y|X}(\tilde{y} | \tilde{x}) p_X(\tilde{x})]}_{p_{X,Y}(\tilde{x}, \tilde{y})} d\tilde{y} d\tilde{x} \end{aligned}$$

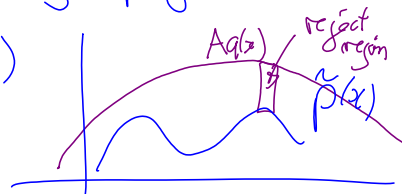
thus  $(\tilde{X}, \tilde{Y})$  has correct joint //

rejection sampling:

say  $p(x) = \frac{\tilde{p}(x)}{Z_p}$ ; say we can find  $q(x)$  a distribution we can easily sample from

$$\text{s.t. } Aq(x) \geq \tilde{p}(x)$$

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



Let's show that accepted draws from rejection sampling have distribution  $p$ :  
(say  $X$  is discrete)

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

$$P\{X \text{ is accepted}\} = \sum_x \frac{\tilde{p}(x)}{A} = \frac{Z_p}{A} \rightarrow \text{marginal probs of acceptance [want this high]}$$

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

Importance sampling:

Suppose  $X \sim p$ , we want to compute  $\mathbb{E}_p[f(X)]$

$$\mathbb{E}_p[f(X)] = \int f(x) p(x) dx = \int \frac{f(x)p(x)}{q(x)} q(x) dx \text{ for some distribution } q$$

$$= \mathbb{E}_q \left[ f(Y) \frac{p(Y)}{q(Y)} \right] \text{ where } Y \sim q$$

$$\approx \frac{1}{n} \sum_{i=1}^n g(Y_i) \text{ where } Y_i \stackrel{iid}{\sim} q \text{ and } g(y) \triangleq \frac{f(y)p(y)}{q(y)}$$

$$= \frac{1}{n} \sum_{i=1}^n f(Y_i) w_i \text{ where } w_i \triangleq \frac{p(Y_i)}{q(Y_i)}$$

↑  
"importance weights"

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n f(Y_i) w_i$$

$$\mathbb{E}[\hat{\mu}] = \mu = \mathbb{E}_p[f(X)]$$

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

intuitively, you want  $q(x) \propto f(x)p(x)$

extension to un-normalized distributions:

$$p(x) = \frac{\tilde{p}(x)}{Z_p} \quad q(x) = \frac{\tilde{q}(x)}{Z_q}$$

$$\mathbb{E}_q \left[ f(Y) \frac{\tilde{p}(Y)}{\tilde{q}(Y)} \right]$$

↑ weights  $w_i$

$$= \mathbb{E}_q \left[ f(Y) \frac{p(Y)}{q(Y)} \frac{Z_p}{Z_q} \right] = \mu \frac{Z_p}{Z_q}$$

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

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n f(Y_i) w_i \text{ where } Y_i \sim q$$

↑  $\frac{1}{n} \sum_{i=1}^n w_i$   $\hat{\frac{Z_p}{Z_q}}$

ideas behind Markov chain Monte-Carlo (MCMC)

idea is to relax independence assumptions;

i.e. we'll run a chain  $X_t | X_{t-1}$  s.t.  $X_t \xrightarrow{t \rightarrow \infty}$  correct distribution

"stationary dist." of the chain

then, we can approximate  $\int$

$$E_p[f(x)] \text{ as } \frac{1}{T-T_0} \sum_{t=T_0+\Delta t} f(x_t)$$

$T_0$  is "burn-in" period  $\rightarrow$  depends on "mixing time"

⊛ no need to thin [use  $\Delta t$  between samples to get more independence]  
as this yields higher variance

$\rightarrow$  better to use all samples after  $T_0$  (unless too expensive)