

today: sampling - approximate inference

Approximate inference - sampling

motivation: NP hard to do exact inference in Ising model \rightarrow need approximation

why sampling? $X = (x_1, \dots, x_p)$

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

b) approximate marginal $p(x_i)$

\rightarrow special case of expectations

consider $f: \mathbb{R}^p \rightarrow \mathbb{R}$

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

special case: if $f(x) \triangleq \mathbb{1}\{x_A = z_A\} \Rightarrow \mathbb{E}_p[f(X)] = p(x_A = z_A)$

Monte Carlo integration / estimation \rightarrow appears in physics, applied math, ML, statistics

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

MC estimation: alg:

- n samples $X^{(i)} \stackrel{iid}{\sim} p$
- estimate $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n f(X^{(i)}) = \mathbb{E}_{\hat{p}_n}[f(X)]$

$\hat{p}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x^{(i)} = x\}$

properties:

1) unbiased estimator $\mathbb{E}[\hat{\mu}] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n f(X^{(i)})\right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[f(X^{(i)})] = \frac{1}{n} \sum_{i=1}^n \mu = \mu$

expectation over $(X^{(i)})_{i=1}^n$

this is still true even if $X^{(i)}$'s are dependent

2) expected error (l2-error) $\mathbb{E}[\|\mu - \hat{\mu}\|_2^2] = \mathbb{E}\left[\left\langle \frac{1}{n} \sum_{i=1}^n f(X^{(i)}) - \mu, \frac{1}{n} \sum_{j=1}^n f(X^{(j)}) - \mu \right\rangle\right]$

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

by independence \Rightarrow off diagonal term are zero
 i.e. $\mathbb{E}\langle f(X^{(i)}) - \mu, f(X^{(j)}) - \mu \rangle$
 $(i \neq j)$ terms $\rightarrow \langle \mathbb{E}[f(X^{(i)}) - \mu], \mathbb{E}[f(X^{(j)}) - \mu] \rangle$

$$\begin{aligned} & \text{for } (i=j) \text{ terms} \\ & \rightarrow \frac{1}{n^2} \sum_{i,j} \mathbb{E} \left[\langle f(x^{(i)}) - \mu, f(x^{(j)}) - \mu \rangle \right] = \frac{1}{n^2} \mathbb{E} \left[\langle f(x^{(i)}) - \mu, f(x^{(i)}) - \mu \rangle \right] = \frac{1}{n^2} \mathbb{E} \left[\|f(x^{(i)}) - \mu\|^2 \right] \triangleq \sigma^2 \\ & \quad \downarrow \\ & \quad \text{tr}(\text{cov}(f(x), f(x))) = \frac{\sigma^2}{n} \end{aligned}$$

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

Note: there is no explicit dimension in rate

(apart σ^2 which could depend implicitly)

e.g. $f(x) = x$

$X_j \sim N(0, \sigma^2)$

$\mathbb{E} \left[\|f(x) - \mu\|^2 \right] = p \sigma^2$

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 trick

Let F be target c.d.f. of dist p for X

$F(x) \triangleq \mathbb{P}\{X \leq x\}$

(first, suppose F is invertible)

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

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

F is monotone

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

\downarrow
 $= \mathbb{P}\{F(F^{-1}(U)) \leq F(y)\}$

F is invertible

$= \mathbb{P}\{U \leq F(y)\} = F(y)$

[if F is not invertible,

define $X \triangleq \min \{x : F(x) \geq U\}$

(recall that F is cts from right)

example:

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

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

$F(x) = 1 - \exp(-\lambda x)$

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

multivariate distribution?

can generalize above trick using "chain rule"

$X_i: p$ (dim p)

$X_{1:p}$ (dim p)

from $p(x_{1:p}) = \prod_{i=1}^p p(x_i | x_{1:i-1})$

use cdf for this conditional

"conditional density sense" for $X_{1:i-1} = x_{1:i-1}$

$$F_{X_i | X_{1:i-1}}(x_i | x_{1:i-1}) \triangleq \mathbb{P}\{X_i \leq x_i \mid X_{1:i-1} = x_{1:i-1}\}$$

$$\triangleq \int_{-\infty}^{x_i} p(x_i' | x_{1:i-1}) dx_i'$$

could use $U_1, \dots, U_p \stackrel{iid.}{\sim} \text{Unif}[0,1]$

$X_1 = F_{X_1}^{-1}(u_1)$

$X_2 = F_{X_2 | X_1}^{-1}(u_2 | x_1)$

\vdots
 $X_p = F_{X_p | X_{1:p-1}}^{-1}(u_p | x_{1:p-1})$

inverse of $F_{X_{i+1}}(\cdot | x_i)$
inverse of this argument

is a very complicated function

(curse of dimensionality)

[aside: "copulas" → model for multivariate data with uniform marginals]

exception is multivariate Gaussian

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

(where $U U^T = I_p$
 Λ is diagonal)

generate $V \sim N(0, I_p)$

(V_p iid $N(0,1)$)

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

(Cholesky decomposition)
 $\Sigma = L L^T$

$\mathbb{E}X = \mu$
 $\text{cov}(X) = U \Lambda^{1/2} \underbrace{\text{cov}(V)}_{I_p} \Lambda^{1/2} U^T = \Sigma$

Box-Muller transformation to sample (2d) Gaussian

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

$\Rightarrow X \stackrel{\Delta}{=} R \cos \theta$

$Y = R \sin \theta$

$\begin{pmatrix} X \\ Y \end{pmatrix} \sim N(0, I)$

$\theta \sim \text{Unif}[0, 2\pi)$

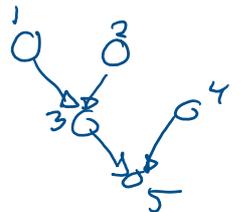
13h34

sampling from a DGM is (relatively) easy e.g. ancestral sampling

$(x_1, \dots, x_d) \sim p \in \mathcal{S}(G)$ where G is a DAG

$p(x_1, \dots, x_d) = \prod_{i=1}^d p(x_i | x_{pa(i)})$

suppose WHOG $1, \dots, d$ is a top-sort of G



ancestral sampling:

for $1, \dots, d$ sample $x_i \sim p(x_i = \cdot | x_{pa(i)})$

these are already sampled by top-sort process

for $1, \dots, d$ sample $x_i \sim p(x_i = \cdot | \bar{x}_{-i})$ these are already sampled by pp. int property
 end

can show by induction (x_1, \dots, x_d) has dist. p

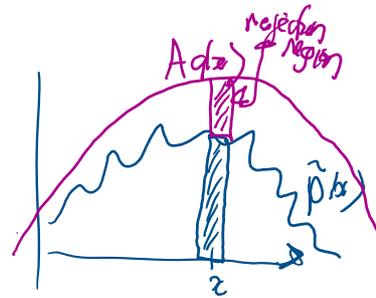
important side note: when you sample from a joint

you are also sampling from "marginals" by just ignoring joint aspect

i.e. $(x, y) \sim p(x, y)$ then look at x by itself
 $X \sim p(x)$

rejection sampling:

say $p(x) = \frac{\tilde{p}(x)}{Z_p}$; let's say can form a $q(x)$ "proposal" that is easy to sample from
 s.t. $Aq(x) \geq \tilde{p}(x) \forall x$



rule:

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

let's show that accepted samples have correct dist.

(say X discrete)

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

↑
sampling mechanism

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

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

(marginal prob. of acceptance \rightarrow want this to be high)

application to conditioning in a DGM

say want to sample $p(x | \bar{x}_E)$

here could use $\check{p}(x) = p(x_E, x_{E^c}) \delta(x_E, \bar{x}_E) \Rightarrow p(x) = p(x_E | \bar{x}_E) \delta(x_E, \bar{x}_E)$

let $q(x)$ be original joint in DGM (i.e. $p(x) = q(x)$)

↳ can sample using ancestral sampling

$$q(x) = p(x_E, x_{E^c})$$

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

$$\text{acceptance prob. : } \frac{\check{p}(x)}{Aq(x)} = \delta(x_E, \bar{x}_E)$$

alg.:

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

(rejection sampling for DGM conditionals)

$$P\{\text{accept}\} = \frac{\sum p}{A} = p(\bar{x}_E)$$

Importance sampling

in context of computing $\mathbb{E}_p[f(x)] = \mu$ $x \sim p$

↳ can "weight" sample $x^{(i)}$

$$\mathbb{E}_p[f(x)] = \sum_x f(x)p(x) = \sum_x f(x) \frac{p(x)}{q(x)} q(x) \quad \text{for some dist } q$$

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

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

$$\text{and } g(y) \triangleq f(y) w(y)$$

$$\text{where } w(y) \triangleq \frac{p(y)}{q(y)} \quad \text{"weights"}$$

$$\hat{\mu}_{\text{I.S.}} = \frac{1}{n} \sum_{i=1}^n f(y_i) w_i \quad y_i \stackrel{i.i.d.}{\sim} q$$

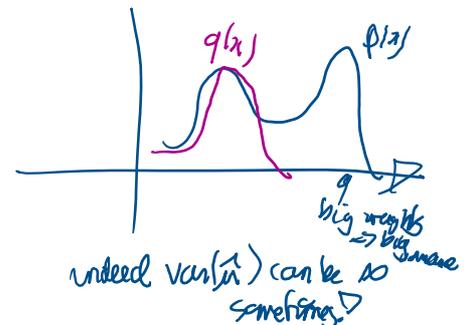
$w_i \triangleq \frac{p(y_i)}{q(y_i)}$
"importance weights"

$$\mathbb{E}[\hat{\mu}_{\text{I.S.}}] = \mu$$

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

issues when q is small and p is large

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



extension to un-normalized dist:

...

extension to un-normalized dist:

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

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

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

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

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

$$= \frac{1}{n} \sum_{i=1}^n w_i$$

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

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

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

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

see 2017 notes for
 • variance reduction (link with SAGA)
 • Rao-Blackwellization

Good reference on sampling:
 Monte Carlo Statistical Methods, Robert & Casella, 2004