

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 $f(x) \triangleq \mathbb{1}\{x_A = x_A\}$ $\mathbb{E}_p[f(x)] = p(x_A = x_A)$

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

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

MC estimation alg. \dots n samples $X^{(i)} \stackrel{i.i.d.}{\sim} p$

\cdot estimate $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n f(x^{(i)}) = \mathbb{E}_{\hat{p}_n}[f(x)]$

\hookrightarrow emp. dist.

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

Properties!

1) unbiased estimator $\mathbb{E}[\hat{\mu}] = \frac{1}{n} \sum_{i=1}^n \underbrace{\mathbb{E}[f(x^{(i)})]}_{\mu} = \frac{1}{n} \cdot n \mu = \mu$

\uparrow expectation over $(x^{(i)})_{i=1}^n$

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

2) expected error (L2-error) $\mathbb{E}[\|\mu - \hat{\mu}\|_2^2] = \mathbb{E}\left[\left\langle \frac{1}{n} \sum_{i=1}^n f(x^{(i)}) - \mu, \frac{1}{n} \sum_{j=1}^n f(x^{(j)}) - \mu \right\rangle\right]$

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

by independence \Rightarrow off-diagonal terms are zero
 i.e. $\mathbb{E} \langle f(x^{(i)}) - \mu, f(x^{(j)}) - \mu \rangle$
 $i \neq j$
 $\langle \mathbb{E} \frac{f(x^{(i)})}{\sqrt{n}} - \mu, \mathbb{E} \frac{f(x^{(j)})}{\sqrt{n}} - \mu \rangle$
 $= \frac{1}{n^2} \sum_{i,j} \mathbb{E} \langle f(x^{(i)}) - \mu, f(x^{(j)}) - \mu \rangle = \frac{p\sigma^2}{n^2}$
 $\mathbb{E} [\|f(x^{(i)}) - \mu\|^2] \triangleq \sigma^2$
 $\text{tr}(\text{cov}(f(x), f(x))) \quad \frac{\sigma^2}{n}$

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

note: there is no explicit dimension in rate

(constant σ^2 which could depend implicitly)

eg $f(x) = X$
 $X_j \sim N(0, \sigma^2)$
 $\mathbb{E} [\|f(x) - \mu\|^2] = p\sigma^2$

How to sample?

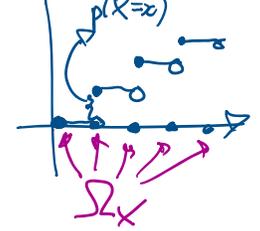
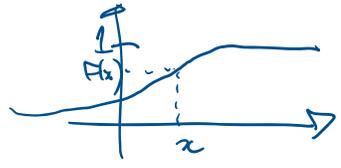
- 1) $X \sim \text{Unif}([0,1]) \rightarrow$ pseudo-random generator "rand"
- 2) $X \sim \text{Bernoulli}(p) \quad X = \mathbb{1}\{U \leq p\}$ where $U \sim \text{unif}([0,1])$
- 3) inverse transform sampling trick

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)

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



claim that X has cdf $F(x)$

F is invertible (and monotone)

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

[if F is not invertible, define

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

(recall that F is cdf from right)

n.s.m.b.

F is cdf 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"

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

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

use cdf for this conditional

$F_{X_i | X_{1:i-1}}(x_i | x_{1:i-1}) \stackrel{\Delta}{=} \mathbb{P}\{X_i \leq x_i | X_{1:i-1} = x_{1:i-1}\}$
 $\stackrel{\Delta}{=} \int_{-\infty}^{x_i} p(x | x_{1:i-1}) dx$

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

could use (u_1, \dots, u_p) i.i.d. Unif(0,1)

$X_1 = F_{X_1}^{-1}(u_1)$ *inverse of $F_{X_1|X_1}(\cdot | X_1)$*

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

\vdots

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

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)

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

generate $V \sim N(0, I_p)$
(V_p i.i.d. $N(0,1)$)

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

$$\mathbb{E}X = \mu$$

$$\text{cov}(X) = U \Omega^{1/2} \overset{I_P}{\text{cov}(U)} \Omega^{1/2} U^T = \Sigma$$

Box-Muller transformation to sample (2d) Gaussian

$$\begin{aligned} R^2 &\sim \text{Exp}(1) \\ \Theta &\sim \text{unif}([0, 2\pi]) \end{aligned} \Rightarrow \begin{aligned} X &\triangleq R \cos \Theta \\ Y &\triangleq R \sin \Theta \end{aligned} \quad \begin{pmatrix} X \\ Y \end{pmatrix} \sim N(0, I)$$

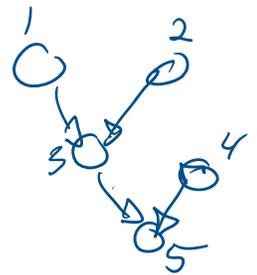
15h32

sampling from a DAG is (relatively) easy & ancestral sampling

$(X_1, \dots, X_d) \sim p \in \mathcal{P}(G)$ where G is a DAG

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

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



ancestral sampling:

for $i=1, \dots, d$
end

sample $X_i \sim p(X_i = \cdot | X_{\pi_i})$ *these are already sampled by top-sort property*

can show by induction (X_1, \dots, X_d) has dist. p

⊗ important side note! when you sample from 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
 $X \sim p(x)$

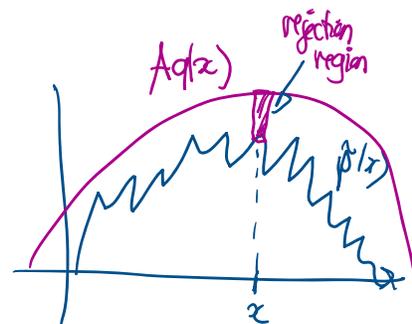
rejection sampling:

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

$$Z_p \triangleq \int_x \tilde{p}(x) dx$$

let's say can form a (d/o) "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 $\frac{\tilde{p}(x)}{Aq(x)} \in [0,1]$
 - reject \rightarrow start again

Let's show that accepted samples have correct dist.

(say X is discrete)

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

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

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

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

application to conditioning in a DGM

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

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

Let $q(x)$ be original joint in DGM (sample using ancestral sampling) $\delta(x_E, \bar{x}_E)$

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

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

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

- alg.:
- do ancestral sampling
 - accept if $x_E = \bar{x}_E$
 - ow. reject

(rejection sampling for DGM sampling)

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

Importance sampling:

in context of computing $\mathbb{E}_p[f(X)] = \mu$ $X \sim p$
 \leadsto can "weight" sample $X^{(i)}$

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

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

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

and $g(Y) \triangleq f(Y)w(Y)$
 where $w(y) \triangleq \frac{p(y)}{q(y)}$ "weights"

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

"importance weight"

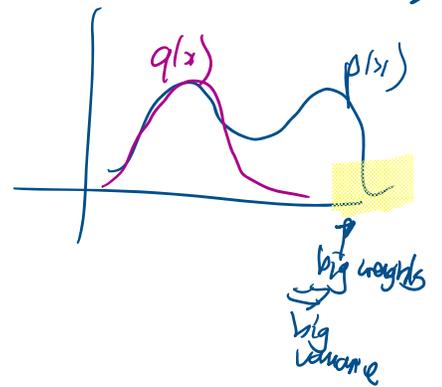
$$w_i \triangleq \frac{p(Y_i)}{q(Y_i)}$$

indeed $\text{Var}(\hat{\mu})$ can be so *sometimes?*

$$\mathbb{E}[\hat{\mu}_{IS}] = \mu$$

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

issues when q is small and p "big"



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

extension to unnormalized dist.

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

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

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

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

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

$$\hat{\mu}_{MIS} = \frac{\frac{1}{n} \sum_{i=1}^n f(x_i) w_i}{\frac{1}{n} \sum_{i=1}^n w_i} \quad \begin{array}{l} y_i \sim q \\ w_i = \frac{p(y_i)}{q(y_i)} \end{array}$$

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

- This estimator has often lower variance than $\hat{\mu}_{IS}$ even when $Z_p = Z_q = \mathbb{I}$ (normalized "stabilizes" estimator new weights $\tilde{w}_i = \frac{w_i}{\sum_{j=1}^n w_j} \in [0, 1]$)

see 2017 notes for

- variance reduction (link with SAGA)
- Rao-Blackwellization

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