

today: sampling - approximate inference

Approximate inference ~ sampling

example: 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 to 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.: • n samples $X^{(i)}$ iid p

$$\cdot \text{estimate } \hat{\mu} = \frac{1}{n} \sum_{i=1}^n f(X^{(i)}) = \mathbb{E}_{p_n}[f(X)]$$

Properties: 1) unbiased estimator $\mathbb{E}_p[\hat{\mu}] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_p[f(X^{(i)})] = \frac{1}{n} \mu = \mu$ expectation over $(X^{(i)})_{i=1}^n$

$$f_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x^{(i)} = x\}$$

this is still true if $X^{(i)}$ were dependent

$$\begin{aligned} 2) \text{ expected error (L2-error)} & \mathbb{E}[\|\mu - \hat{\mu}\|_2^2] = \mathbb{E}\left[\left\langle \frac{1}{n} \sum_{i=1}^n f(X^{(i)}) - \mu, \frac{1}{n} \sum_{j=1}^n f(X^{(j)}) - \mu \right\rangle\right] \\ & \text{tr}(\text{cov}(\hat{\mu}, \hat{\mu})) = \mathbb{E}\left[\frac{1}{n^2} \sum_{i=1, j=1}^n \langle f(X^{(i)}) - \mu, f(X^{(j)}) - \mu \rangle\right] \end{aligned}$$

by independence \Rightarrow off-diagonal term are zero

$$\text{i.e. } \langle \mathbb{E}_p[f(X^{(i)})] - \mu, \mathbb{E}_p[f(X^{(j)})] - \mu \rangle$$

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

$$\mathbb{E}[\|\hat{\mu} - \mu\|_2^2] = \frac{\sigma^2}{n}$$

note: there is no explicit dimension in rate

$$\mathbb{E}[\|f(x^{(i)}) - \mu\|^2] \triangleq \sigma^2 = \text{tr}(\text{cov}(f(x), f(x))) = \frac{\sigma^2}{n}$$

(apart σ^2
which could depend implicitly)

e.g. $f(A) = X$
 $X_j \sim N(0, \sigma^2)$

$$\mathbb{E}[\|f(x) - \mu\|^2] = d\tilde{\sigma}^2$$

How to sample?

- 1) $X \sim \text{Unif}(0, 1)$ → pseudo-random generator "rand"
- 2) $X \sim \text{Bernoulli}(p)$ $X = \mathbf{1}\{U \leq p\}$ where $U \sim \text{Unif}(0, 1)$

3) inverse transform sampling trick:

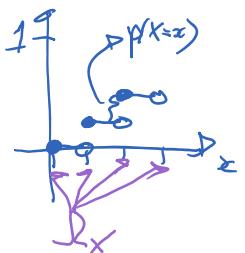
Let F be target c.d.f. of dist. p for X $\xrightarrow{\text{cum. dist. fn}}$ $F(x) \triangleq P\{X \leq x\}$

(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 and monotone

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



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

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

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

multivariate distribution?

can generalize above trick using "chain rule"

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

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

use cdf for this conditional

"conditional density" sense
 $\int_X p(x_{1:i-1}, x_{i:i}) dx_{1:i-1}$

use cdf for this conditional \rightarrow cts.

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

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

$$\left\{ \begin{array}{l} X_1 = F_{X_1}^{-1}(U_1) \\ X_2 = F_{X_2|X_1}^{-1}(U_2|x_1) \end{array} \right. \begin{array}{l} \text{universe of } F_{X_2|X_1}(\cdot|x_1) \\ \text{inverse of this argument} \end{array}$$

$$\left. \begin{array}{l} \vdots \\ X_p = F_{X_p|X_{1:p-1}}^{-1}(U_p|x_{1:p-1}) \end{array} \right. \begin{array}{l} \text{is a very complicated fct.} \\ \text{(curse of dimensionality)} \end{array}$$

[aside: "copulas" \rightarrow 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 $V \sim N(0, I_p)$

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

$$E[X] = \mu$$

$$\text{cov}(X) = U \Lambda U^T \text{cov}(V) U^T = \Sigma$$

Box-Muller transformation to sample 2d Gaussian

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

16h14

sampling from a DGM is (relatively) easy: ancestral sampling

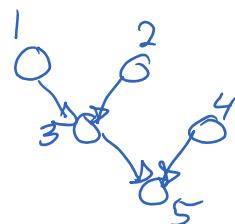
$$(x_1, \dots, x_d) \sim p \in \mathcal{P}(G) \text{ where } G \text{ is a DAG}$$

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

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

ancestral sampling:

$\left\{ \begin{array}{l} \text{for } i=1, \dots, d \\ \text{sample } x_i \sim p(x_i = \cdot | x_{\pi_i}) \\ \text{end} \end{array} \right. \begin{array}{l} \text{those are already observed} \\ \text{by top. sort. property} \end{array}$



end

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

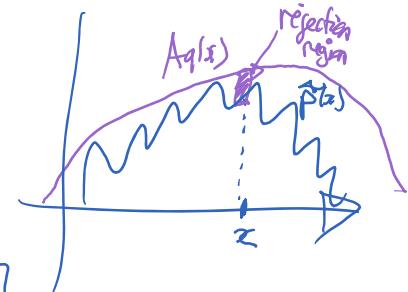
Important sidenote: when you sample from joint
you are also sampling from "marginals" by just
ignoring joint aspect

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

rejection sampling:

say $p(x) = \frac{\tilde{p}(x)}{Z_p}$; Let's say can find a $q(x)$ "proposal"
which is easy to sample from

$$Z_p \triangleq \int_x \tilde{p}(x) dx \quad \text{s.t. } Aq(x) \geq \tilde{p}(x) \forall x$$



rule

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

Let's show that accepted samples has correct dist.

(say X is discrete)

$$\underset{\text{Sampling mechanism}}{P_q^S \{X=x, X \text{ is accepted}\}} = P_q^S \{X \text{ is accepted} | X=x\} P_q^S \{X=x\} = \frac{\tilde{p}(x)}{A}$$

$$P_q^S \{X \text{ is accepted}\} = \sum_x \frac{\tilde{p}(x)}{A} = \frac{Z_p}{A} \quad (\text{marginal prob. of acceptance})$$

$$P_q^S \{X=x | X \text{ is accepted}\} = \frac{\tilde{p}(x)/A}{Z_p/A} = \frac{\tilde{p}(x)}{Z_p} = p(x) \quad (\text{want this to be high})$$

application to conditioning in a DGM

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

here, could use $\tilde{p}(x) = p(x_{\bar{x}_C}, x_E) \delta(x_E, \bar{x}_E) \Rightarrow p(x) = p(x_{\bar{x}_C} | \bar{x}_E) \cdot \delta(x_E - \bar{x}_E)$

$$\Delta p = p^{(1-t)} - p^t$$

Let $q(x)$ be original joint in DGM [sample using ancestral sampling]

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

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

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

alg.

- do ancestral sampling
- accept if $x_E = \tilde{x}_E$
- o.w. reject

(rejection sampling
for DGM Sampling)

$$P\{\text{accept}\} = \frac{1}{A} = p(\tilde{x}_E)$$

Importance sampling:

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

→ can "weight" sample $X^{(i)}$

$$\begin{aligned} \mathbb{E}_p[f(x)] &= \sum_x f(x)p(x) = \sum_x f(x) \frac{p(x)}{q(x)} \cdot q(x) \quad \text{for some dist } q \\ &= \mathbb{E}_q\left[\frac{f(y)p(y)}{q(y)}\right] \quad \text{st. supp}(q) \supseteq \text{supp}(p) \\ &\approx \frac{1}{n} \sum_{i=1}^n g(y^{(i)}) \quad \text{where } Y^{(i)} \sim q \end{aligned}$$

where $Y^{(i)} \sim q$
and $g(y) \triangleq f(y)w(y)$

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

Indeed $\text{var}(\hat{\mu})$ can be ∞ sometimes!

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

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

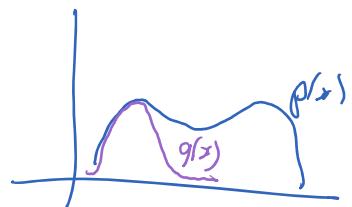
("importance weights")

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

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

issues when q is small
and p is big

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



extension to un-normalized dist.

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

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

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

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

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

- note:
- $\hat{\mu}_{UIS}$ is (slightly biased), but asymptotically unbiased $\xrightarrow{n \rightarrow \infty}$
 - this estimator has often lower variance than $\hat{\mu}_{IS}$ even when $Z_p = Z_q = 1$
 (normalize "stabilizes" estimator now 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