

today: Sampling - approximate inference

## Approximate inference

Sampling example: NP hard to do exact inference in Ising model  
 $\rightarrow$  need approximate

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

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

b) approximate  $p(x_i)$

$\rightarrow$  special case of expectations

Consider:  $f: \mathbb{R}^P \rightarrow \mathbb{R}$

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

e.g. 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 def: • n samples  $X^{(i)} \stackrel{iid}{\sim} P$

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

properties: 1) unbiased  $\mathbb{E}_P[\hat{\mu}] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_P[f(X^{(i)})] = \underbrace{\frac{1}{n} \mu}_{\text{empirical dist}} = \hat{\mu}$   
 $\mathbb{E}_P[\hat{\mu}] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_P[\mathbb{1}\{X^{(i)} = x^{(i)}\}]$   
 expectation over  $(X^{(i)})_{i=1}^n$

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

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

( $L_2$ -error)  $\text{tr}(\text{cov}(\hat{\mu}, \hat{\mu})) = \mathbb{E}\left[\left\langle \frac{1}{n} \sum_{i,j} f(X^{(i)}) - \mu, f(X^{(j)}) - \mu \right\rangle\right]$

by independence  $\Rightarrow$  off diagonal terms are zero

$$\begin{aligned}
 &= \frac{1}{n^2} \sum_{i=1}^n \underbrace{\mathbb{E}[f(x^{(i)}) - \mu, f(x^{(i)}) - \mu]}_{\mathbb{E}[\|f(x^{(i)}) - \mu\|^2] \triangleq \sigma^2} \\
 &= \frac{1}{n^2} \sigma^2
 \end{aligned}$$

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

note : that no dimension in rate ( $\propto \sigma^2$   
which could depend  
implicitly)

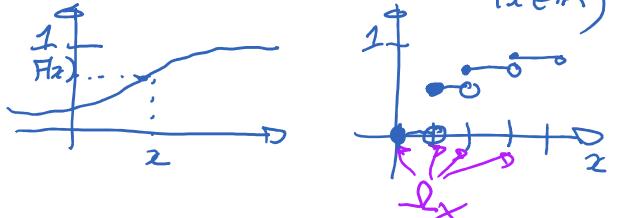
## How to sample?

- 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  $\xrightarrow{\text{cum. dist. fct.}}$

let  $F$  be target c.d.f. 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 and monotone

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

[ if  $F$  is not invertible , define  $X \triangleq \min\{x : F(x) \geq u\}$  ] (recall  $F$  is cts from right)

example :

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

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

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

multivariate distribution?

can generalize above trick using "chain rule"

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

change from class scribbles

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

use cdf for this conditional:

$$F_{x_p | x_{1:p-1}}(x_p | x_{1:p-1}) \triangleq P\{X_p \leq x_p | X_{1:p-1} = x_{1:p-1}\}$$

(not  $P\{X_p \leq x_p | X_{1:p-1} = x_{1:p-1}\}$ )

as I had written in class

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

$$X_1 = F_{X_1}^{-1}(U_1) \quad \text{inverse of } F_{X_1 | x_{1:p-1}}(\cdot | X_{1:p-1})$$

$$X_2 = F_{X_2 | X_1}^{-1}(U_2 | X_1) \quad \text{is a very complicated function}$$

$$\vdots \quad X_p = F_{X_p | X_{1:p-1}}^{-1}(U_p | X_{1:p-1}) \quad (\text{curse of dimensionality?})$$

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

14h3)

exception is multivariate Gaussian:

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

(where  $U U^T = I_p$   
 $U$  is diagonal)

(Cholesky decomposition  
 $L = U \sqrt{\Sigma}$ )

$$\Sigma = L L^T$$

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

( $v_p \stackrel{iid}{\sim} N(0, 1)$ )

$$X \triangleq \sum_L U_L V + \mu$$

$$\begin{aligned} \mathbb{E}X &= \mu \\ \text{cov}(X) &= U_L L \underbrace{V \text{ cov}(V) V^T}_{I_p} L^T U_L^T \\ &= \Sigma \end{aligned}$$

Box-Muller transformation to sample (2d) Gaussian

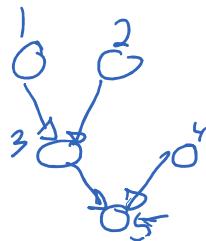
$$\begin{aligned} R^2 &\sim \text{Exp}(1) \\ \theta &\sim \text{Unif}(0, 2\pi) \end{aligned} \Rightarrow \begin{aligned} X &\triangleq R \cos \theta \\ Y &\triangleq R \sin \theta \end{aligned} \quad (X, Y) \sim N(0, I_2)$$

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

$(x_1, \dots, x_p) \sim p \in \mathcal{Y}(G)$  where  $G$  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, ..., p
    sample  $x_i \sim p(x_i = \cdot | X_{\pi_i})$ 
end

```

*these are already observed  
by top sort.*

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

⊕ important side note: when you sample from joint,

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

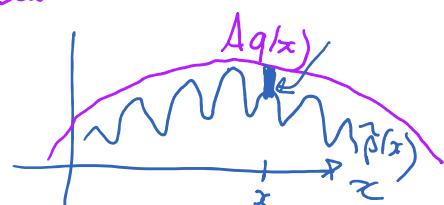
$$\text{i.e. } (X, Y) \sim p(x, y)$$

then looking at  $X$  itself, we have  $X \sim p(x)$

### 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

$$\text{s.t. } A q(x) \geq \tilde{p}(x) \forall x$$



Rule:

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

let's show that accepted sample has correct dist,

(say  $X$  is discrete)

$$P_{\text{sampling mechanism}}^S \{X=x, X \text{ is accepted}\} = P_{\text{sampling mechanism}}^S \{X \text{ is accepted} | X=x\} P_{\text{sampling mechanism}}^S \{X=x\}$$

$$P\{X=x, X \text{ is accepted}\} = P\{X \text{ is accepted } | X=x\} P\{X=x\}$$

sampling mechanism proposal  
dist  
 $\frac{\tilde{p}(x)}{A(x)}$

$$P\{X \text{ is accepted}\} = \sum_x \frac{\tilde{p}(x)}{A} = \frac{z_p}{A} \quad (\text{marginal prob of acceptance})$$

$$P\{X=x | X \text{ is accepted}\} = \left(\frac{\tilde{p}(x)}{A}\right) / \frac{z_p}{A} = p(x) \quad \begin{matrix} \rightarrow \text{want this} \\ \text{to be high} \end{matrix}$$

application to conditioning in a DGM:

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

$$\text{hence, could use } \tilde{p}(x) = p(x_E^c, \bar{x}_E) \delta(x_E, \bar{x}_E)$$

$$z_p = p(\bar{x}_E) \quad p(x) = p(x_E^c | \bar{x}_E) \delta(x_E, \bar{x}_E)$$

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

$$q(x) = p(x_E^c, x_E)$$

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

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

alg:  $\begin{cases} \cdot \text{ do ancestral sampling} \\ \cdot \text{ accept if } x_E = \bar{x}_E \end{cases}$  [rejection sampling for DGM conditioning]

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

Importance sampling:

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

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

$$\begin{aligned} E_p[f(x)] &= \sum_x f(x) p(x) = \sum_x f(x) \frac{p(x)}{q(x)} q(x) \quad \text{for some dist } q \\ &\quad \text{s.t. } \text{supp}(q) \supseteq \text{supp}(p) \\ &= E_q[f(x) \frac{p(x)}{q(x)}] \quad \text{where } X \sim q \\ &\quad \rightsquigarrow \dots \dots \dots \text{ i.e.} \end{aligned}$$

$$= \frac{1}{n} \sum_{i=1}^n f(x_i) \frac{p(x_i)}{q(x_i)}$$

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

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

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

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

"importance weights"

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

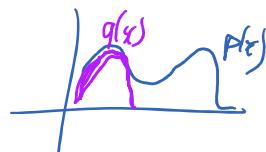
$$\mathbb{E}[\hat{\mu}_{IS}] = \mu$$

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

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

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

indeed  $\text{Var}[\hat{\mu}]$  can sometimes be  $\infty$



extension to un-normalized dist.:

$$p(x) = \frac{\tilde{p}(x)}{\tilde{Z}_p} \quad q(x) = \frac{\tilde{q}(x)}{\tilde{Z}_q}$$

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

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

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

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

note: •  $\hat{\mu}_{ISU}$  is (slightly biased), but asymptotically unbiased

• This estimator has often lower variance than  $\hat{\mu}_{IS}$ , even when  $\tilde{Z}_p = \tilde{Z}_q = 1$

(normalization "stabilizes" estimator) new weights  $\tilde{w}_i = \frac{w_i}{\sqrt{\sum_j w_j}} \in [0, n]$

See 2017 notes for

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

Good reference on sampling:

**Monte Carlo Statistical Methods**, Robert & Casella, 2004