

Lecture 19 - scribbles

Tuesday, November 13, 2018 14:30

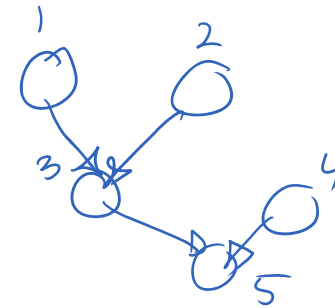
today: more sampling
MCMC

sampling for DGM is easy: ancestral sampling

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

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

suppose WLOG, $1, \dots, p$ be a top-sort for G



ancestral sampling:

for $i=1, \dots, p$

sample $X_i \sim p(x_i = \cdot | X_{\pi_i})$

these are already sampled
by top. sort.

end

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

aside: formal tool to show that (X_1, X_2) has right joint dist

2 nodes example: $U_1 \sim p(x_1)$

$$U_2 | U_1 \sim p(x_2 | x_1 = U_1)$$

to show that two R. vectors are equal in dist. i.e. $U \stackrel{d}{=} V$

$$\Leftrightarrow E_U[f(U)] = E_V[f(V)] \text{ for all functions in a big enough class}$$

(e.g. cts, \int bounded functions)

here, I want to show that:

$$E[f(x_1, x_2)] \stackrel{\text{want to show this}}{=} E[f(U_1, U_2)]$$

$$\begin{aligned} & \stackrel{\text{tower property of expect.}}{=} E_{U_1} [E_{U_2 | U_1} [f(U_1, U_2) | U_1]] \\ & \stackrel{\text{by sampling mechanism}}{=} \int_{u_1} p(u_1) du_1 \left[\int f(u_1, u_2) p(u_2 | u_1) du_2 \right] \\ & = \int_{u_1} \int_{u_2} f(u_1, u_2) \underbrace{[p(u_2 | u_1) p(u_1)]}_{p(u_1, u_2)} du_1 du_2 \end{aligned}$$

⊛ important sidenote: when sample from joint, you are also sampling from "marginal" by just ignoring the joint aspect...

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

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

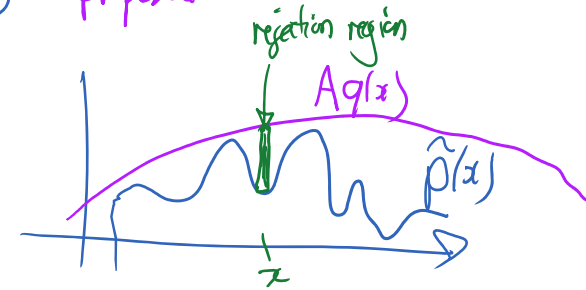
rejection sampling:

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

Let's say we can find a $q(x)$ which easy to sample from

s.t. $Aq(x) \geq \tilde{p}(x) \forall x$
proposal budget (norm.)

"proposal"



rule :

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

Let's show that accept samples have correct dist.

(say X is discrete)

$$P\{ \underbrace{X=x, X \text{ is accepted}}_{\text{Sampling Mechanism}} \} = \underbrace{P\{ X \text{ is accepted} \mid X=x \}}_{\delta(x)} \underbrace{P(X=x)}_{\alpha(x)}$$

$$= \frac{\tilde{p}(x)}{A} \quad \text{Ag(x)}$$

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

(marginal prob. of acceptance
→ want this high)

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

application to DGM:

say we want to sample from $p(x \mid \bar{x}_E)$ complement of E

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

$$Z_p = p(\bar{x}_E)$$

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

$$q(x) = p(x_{E^c}, \bar{x}_E)$$

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

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

alg. : • do ancestral sampling (rejection sampling)

alg. : • do ancestral sampling
 • accept if $Z_E = \bar{x}_E$

(rejection sampling for OBM)

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

15h33

Importance sampling :

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

→ can "weight" sample $x^{(i)}$

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

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

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

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

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

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

↓
"importance weights"

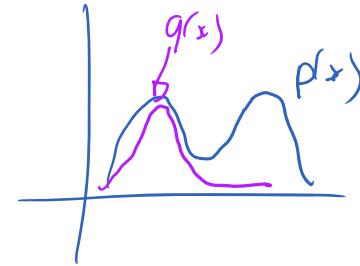
$$w_i = \frac{p(y_i)}{q(y_i)}$$

...proportional weights
(1) (2)

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

(to be checked)

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



issues here when q is small
and p is big

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

extension to un-normalised dist.

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

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

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

$$\hat{\mu}_{IS0} = \frac{1}{n} \sum_{i=1}^n \frac{f(y_i) w_i}{w_i} \quad y_i \sim q$$

$$\hat{\mu}_{IS0} = \frac{\frac{1}{n} \sum_{i=1}^n f(y_i) w_i}{\frac{1}{n} \sum_{j=1}^n w_j} \quad y_i \sim q$$

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

note: • $\hat{\mu}_{IS0}$ is (slightly) biased, but asymptotically unbiased (as $n \rightarrow \infty$)

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

See 2017 notes for

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

MCMC

Markov Chain Monte-Carlo :

Idea: is to relax indep. assumption between samples
to allow adaptive proposal dist.

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

"stationary dist. of chain"

then, we can approximate

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

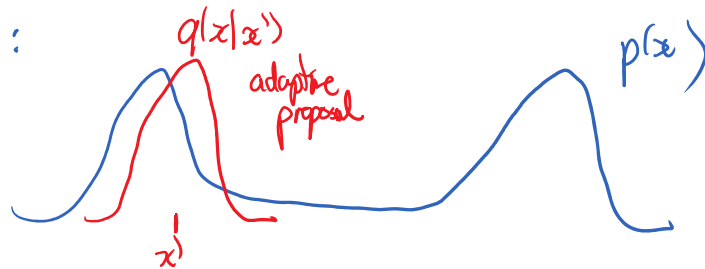
T_0 is called "burn in" period \rightarrow depends on "mixing time" of Markov chain

(*) no need to thin the samples [i.e. use Δt between samples to get more independence]

as this yield higher variance

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

Motivation:



before: samples were $x^{(i)}$ i.i.d

$$\text{MCMC: } x^{(t)} | x^{(t-1)} \sim q(\cdot | x^{(t-1)})$$

\uparrow Markov transition probs

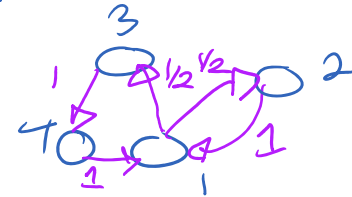
review of (finite state space) Markov chains [$|X| = k$]



- as a DGM, $x^{(0)} \rightarrow x^{(1)} \rightarrow \dots \rightarrow x^{(t-1)} \rightarrow x^{(t)} \rightarrow \dots$

- there is also a transition prob. point of view: use one node per state (probabilistic FSA)

4 states example



[homogeneous M.C.]

\rightarrow i.e. $P\{X_t = i \mid X_{t-1} = j\} = A_{ij}$ (no time dep.)

A is a $K \times K$ matrix st. $\mathbb{1}_K^T A = \mathbb{1}_K$

vector of K ones

"left-stochastic matrix"

⊗ (as in HMM), suppose $P\{X_{t+1} = j\} = \pi(j)$

$$P\{X_t = i\} = \sum_j \underbrace{P\{X_t = i \mid X_{t-1} = j\}}_{A_{ij}} \underbrace{P\{X_{t-1} = j\}}_{\pi_j}$$

$$\pi_{t+1} = A \pi_t$$

$$\Rightarrow \boxed{\pi_t = A^t \pi_0}$$

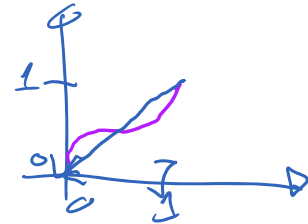
$$\Rightarrow \boxed{\pi_t = A^t \pi_0}$$

Stationary dist π of A is a dist. π s.t. $A\pi = \pi$

[note that π is a right e-vector of A with e-value of 1]

fact: every stochastic matrix has at least 1 stat. dist.

(by Brouwer's fixed pt. thm.)



16h70