

Left off with "Bayesian Learning"

eg.
$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

but what does that mean? how to compute in practice?

Today: Sample-based approximate bayesian inference

running example: bayesian linear regression: height & weight

model: $w = b_0 + b_1 h + \epsilon$, $\epsilon \sim N(0, \sigma^2)$

$b_0, b_1 \sim ??$

ols: $y = X\beta + \epsilon$

$\hat{\beta} = (X^T X)^{-1} X^T y$

$\hat{\beta} \sim N[\hat{\beta}, \sigma^2 (X^T X)^{-1}]$

$D = \begin{bmatrix} h_1 & w_1 \\ h_2 & w_2 \\ \vdots & \vdots \\ h_n & w_n \end{bmatrix}$

$X = \begin{bmatrix} 1 & h_1 \\ 1 & h_2 \\ \vdots & \vdots \\ 1 & h_n \end{bmatrix}$

purpose: queries of form $P(w_{\text{new}} | h_{\text{new}}, D) = E[f(h) | D] = \int f(h, \theta) P(\theta | D) d\theta$

or more generally: $E[f(x)] = \int f(x) p(x) dx$

1) Ordinary Monte-carlo:

* $E[f(x)] \approx \frac{1}{n} \sum_{i=1}^n f(x_i)$, $x_i \overset{\text{iid}}{\sim} P(x)$ (general) (1)

$E[f(h|D)] \approx \frac{1}{n} \sum_{i=1}^n [b_0 + b_1 h]$, $b_0, b_1 \sim P(b_0, b_1 | D)$

how to get these samples??

~~$E[f(h|D)] = E[E[f(h|D) | \theta]] = \int \int f(h|D, \theta) P(\theta) d\theta dh$~~
 $\frac{1}{n} \sum_{i=1}^n \{ \begin{matrix} b_0 \\ b_1 \end{matrix} \}_{\text{over } \theta}$

ideas??

ideas (cont'd)

→ OLS works in special cases where we have closed-form estimator

→ why hard? Recall: $P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)} = \frac{P(D|b_0, b_1)P(b_0, b_1)}{\iint_{b_0, b_1} P(D|b_0, b_1)P(b_0, b_1)} \leftarrow \text{hard}$

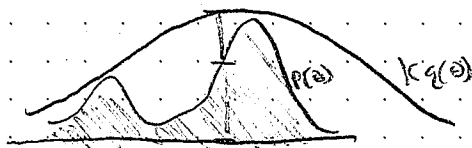
↳ to do analytically, need tractable parametric form for product of both terms

observe: $P(\theta|D) = \frac{1}{Z} P(D|\theta)P(\theta) \propto P(D|\theta)P(\theta)$
proportional to

idea: ignore normalizing constant $\frac{1}{Z}$ somehow?

A) Rejection sampling

* ↳ introduce simple proposal distribution $q(\theta)$, & constant k s.t. $kq(\theta) \geq P(\theta)$



$$\text{accept prob } a(\theta) = \frac{P(\theta)}{kq(\theta)} \quad (2)$$

problems?

↳ obtaining k

↳ what if q is "bad" → rejects most proposal

B) importance sampling

idea: plug (2) directly into OMC (1)

$$\int f(x; \theta) P(\theta) d\theta = \int \frac{f(x; \theta) P(\theta)}{q(\theta)} q(\theta) d\theta$$

problems? if peaks can be dominated by low-weight samples, & be inefficient

ie: sample from easy proposal dist. & then correct them using weights from (2)

3) MCMC* : Markov-chain Monte-Carlo

goal: more efficient method for getting samples from $P(\theta|D)$

↳ same assumptions of (1) & (2): only have proposal dist. & unnormalized prob.

↳ intuitively, we'll get high-prob samples by searching for them

quick review of Markov Chains

↳ MC is an indexed sequence of RV's: defined entirely by transition probs

$$P_{ij} = P(\theta_{t+1} = j | \theta_t = i)$$

[transition dist]

$$(i) \begin{bmatrix} 1 & 2 & \dots & n \\ P_{ij} = P(j|i) \end{bmatrix} \rightarrow Z=1$$

↳ interested in equilibrium cond. & convergence rate

key condition: reversibility ("detailed balance")

notation:

$$\pi_i = P(\theta = i)$$

↳ in equilibrium if $\pi = \pi P \rightarrow \pi_i P_{ji} = \pi_j P_{ij}, \forall i, j$ (4)

π = eq. distribution

↳ π is a distribution over Θ

* (4) lets us build a sampling strategy which defines a MC, & *
therefore lets us guarantee convergence to a stationary (posterior) dist.
using a sequential random walk (ie search) algorithm

this is generic { simple algorithm:
Metropolis-Hastings

↳ define transition rule as $a_{ij} = \min\left[1, \frac{\pi_j q_{ji}}{\pi_i q_{ij}}\right]$

Metropolis 1950

Hastings 1970

proof:

$$\pi_i q_{ij} a_{ij} = \pi_i q_{ij} \cdot \frac{\pi_j q_{ji}}{\pi_i q_{ij}}$$

$$= \pi_j q_{ji}$$

$$= \pi_j q_{ji} a_{ji} \int \dots \text{reversible} \quad \text{www.icra2013.org}$$

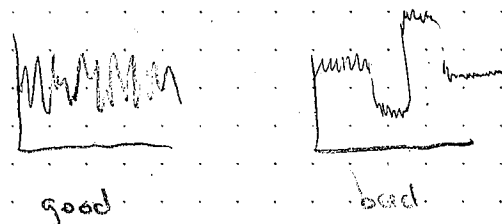
Neighbour
ratio
proposal
ratio

key idea: moves based
on relative probability
of different setting of θ

Metropolis-Hastings MCMC (cont'd)

Basic Convergence ideas

1) Mixing: how well is chain moving around parameter space?

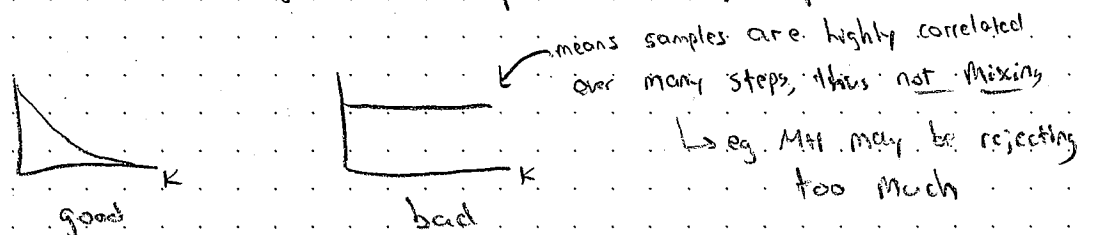


2) Auto-correlation: a statistic for examining mixing.

- lag- k autocorr $\hat{=}$ ρ_k

$$\rho_k = \frac{\sum_{i=1}^{n-k} (x_i - \bar{x})(x_{i+k} - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

- $\hat{=}$ the average corr of points k steps apart



Example: non-linear physics inference

- harping on "allows inference when joint posterior only known to a constant"
- how often does this really come up?? All the time!

eg. physics

$$x_{t+1} = F(x_t, u_t) + \epsilon, \quad F = \langle \text{some black box dynamics} \rangle$$

$$\text{generative model: } \theta, \sigma \sim P(\theta, \sigma | h) \\ \epsilon \sim N(0, \sigma^2)$$

$$\text{likelihood: } L(\theta | \sigma, x, u) = N(F(x, u), \sigma^2)$$

notation

$$\Theta = \{\text{model params}\}$$

$$\pi = P(\Theta)$$

$$\pi_i = P(\Theta = i)$$

$$m = |\Theta|$$

$$i = (i_1, i_2, \dots, i_m)$$

i_c = randomly chosen component

$$i_{\bar{c}} = \text{remaining components} \\ = (i_1, \dots, i_{c-1}, i_{c+1}, \dots, i_m)$$

j is a neighbor of i if:

$$j_{\bar{c}} = i_{\bar{c}}$$

→ IE can move from i to j by changing i_c

Gibbs sampling: Special case of M+I where q is a proper (normalized) conditional distribution

$$q_{ij} = \frac{1}{m} \frac{\pi_j}{\sum_{k: k_{\bar{c}} = i_{\bar{c}}} \pi_k}$$

proposes some neighbor j with prob. normalized by all other possible neighbors

why better?

→ have that $a_{ij} \equiv 1$

requires:

→ deriving conditional expression

also:

→ can only make steps along axes in which conditionals are defined

Proof that $a_{ij} = 1$ for Gibbs proposals:

$$\pi_i q_{ij} = \pi_i q_{ij} \frac{\pi_j}{\pi_j}$$

$$= \pi_j \frac{\pi_i q_{ij}}{\pi_j}$$

$$= \pi_j \left[\frac{\pi_i}{\pi_j} \cdot \frac{\pi_j}{m \sum_{k: k_{\bar{c}} = i_{\bar{c}}} \pi_k} \right]$$

$$= \pi_j \left[\frac{\pi_i}{m \sum_{k: k_{\bar{c}} = i_{\bar{c}}} \pi_k} \right]$$

$$= \pi_j q_{ji}$$

∴ $\pi_i q_{ij} = \pi_j q_{ji} \Rightarrow$ Gibbs sampling is reversible, so no acceptance needed
[IE $a_{ij} = a_{ji} = 1$]

b/c π_i & π_j are neighbors, we have that $i_{\bar{c}} = j_{\bar{c}} = k_{\bar{c}}$

(IE they share the same normalizing constant)

Gibbs Sampling (cont'd)

Example model: Simple linear regression

$$y = \beta_0 + \beta_1 x + e, \quad e \sim N(0, \sigma^2)$$

Bayesian version:

$$P(y_i | x_i, \beta_0, \beta_1, \tau) \sim N(\beta_0 + \beta_1 x_i, 1/\tau)$$

$$\beta_0 \sim N(\mu_0, 1/\tau_0)$$

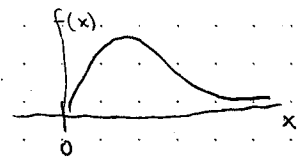
$$\beta_1 \sim N(\mu_1, 1/\tau_1)$$

$$\tau \sim \text{Gamma}(a, b)$$

note: $\tau = \frac{1}{\sigma^2}$

$$X \sim \text{Gamma}(a, b) \Rightarrow f(x) = \frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx}$$

recall: $\Gamma(a) = (a-1)!$



inference

Assume $[\mu_0, \mu_1, \tau_0, \tau_1, a, b]$ known constants (for now)

model params $\theta = [\beta_0, \beta_1, \tau]$

want: $P(\theta | x, y) = P(\beta_0, \beta_1, \tau | x, y) = P(y | \beta_0, \beta_1, \tau, x) \underbrace{P(\beta_0)P(\beta_1)P(\tau)}_{\text{why independent? why not? model says so}} \underbrace{P(x)}_{\text{always observed}}$

joint posterior = $P(y | \beta_0, \beta_1, \tau, x) P(\beta_0) P(\beta_1) P(\tau)$

$$= N(\beta_0 + \beta_1 x, 1/\tau) \cdot N(\mu_0, 1/\tau_0) N(\mu_1, 1/\tau_1) \text{Gamma}(\tau, a, b)$$

$$= \left[\prod_{i=1}^n \frac{\tau}{\sqrt{2\pi}} e^{-\frac{(y_i - \beta_0 - \beta_1 x_i)^2}{2(1/\tau)}} \right] \left[\frac{\tau_0}{\sqrt{2\pi}} e^{-\frac{(\beta_0 - \mu_0)^2}{2(1/\tau_0)}} \right] \left[\frac{\tau_1}{\sqrt{2\pi}} e^{-\frac{(\beta_1 - \mu_1)^2}{2(1/\tau_1)}} \right] \left[\frac{b^a}{\Gamma(a)} \tau^{a-1} e^{-b\tau} \right] = \underline{\underline{\text{ugh!}}}$$

Gibbs for the save: consider each θ_i in turn

$$P(\beta_0 | x, y, \beta_1, \tau) \sim N\left(\frac{\tau_0 \mu_0 + \tau \sum_i (y_i - \beta_1 x_i)}{\tau_0 + n\tau}, \frac{1}{\tau_0 + n\tau}\right)$$

$$P(\beta_1 | x, y, \beta_0, \tau) \sim N\left(\frac{\tau_1 \mu_1 + \tau \sum_i (y_i - \beta_0)}{\tau_1 + \tau \sum_i x_i^2}, \frac{1}{\tau_1 + \tau \sum_i x_i^2}\right)$$

$$P(\tau | x, y, \beta_0, \beta_1) \sim \text{Gamma}(a + n/2, b + \sum_i (y_i - \beta_0 - \beta_1 x_i)^2 / 2)$$

Gibbs sampling (cont'd)

So: have conditional distributions for each θ_i :

$$B_0 | x, y, B_1, \tau$$

$$B_1 | x, y, B_0, \tau$$

$$\tau | x, y, B_0, B_1$$

inference:
 $i=0$

while (not bored):

$$B_0^{(i)} \sim N\left(\frac{\tau_0 \mu_0 + \tau \sum_i (y_i - B_1 x_i)}{\tau_0 + n\tau}, \frac{1}{\tau_0 + n\tau}\right)$$

$$B_1^{(i)} \sim N\left(\frac{\tau_1 \mu_1 + \tau \sum_i (y_i - B_0 x_i)}{\tau_1 + \tau \sum_i x_i^2}, \frac{1}{\tau_1 + \tau \sum_i x_i^2}\right)$$

$$\tau^{(i)} \sim \text{Gamma}\left(a + n/2, b + \sum_i (y_i - B_0 - B_1 x_i)^2/2\right)$$

$i++$

< inspect $\vec{B}_0, \vec{B}_1, \vec{\tau}$, eg $\hat{B}_0 = \vec{B}_0$ >

take home:

- in general, MCMC does approximate bayesian inference for arbitrarily complex models
- standard Metropolis-Hastings is most general, but must be tuned (& can be slow)
- Gibbs sampling is special case of MH when we have conditional distributions for some θ_i

↳ Strategy: write down joint posterior, isolate a single parameter, & plow through the horrible algebra