

Outline

- 1) Hierarchical Bayes
- 2) Markov Chain Monte Carlo
- 3) Gibbs Sampler

Hierarchical Bayes

[Full power of Bayes is realized in large, complex problems with repeat structure, allowing us to pool information across many observations.]

Ex Predict a batter's "true" batting average from n at-bats. $X = \# \text{ of hits} \sim \text{Binom}(n, \theta)$

Prior info: Most batting avg.s are between 0.1 and 0.3, $\theta = 0.8$ very unlikely. Can represent using a Beta dist., but how to pick α, β ?

Solution: Pool info across players $i=1, \dots, m$ via hierarchical model

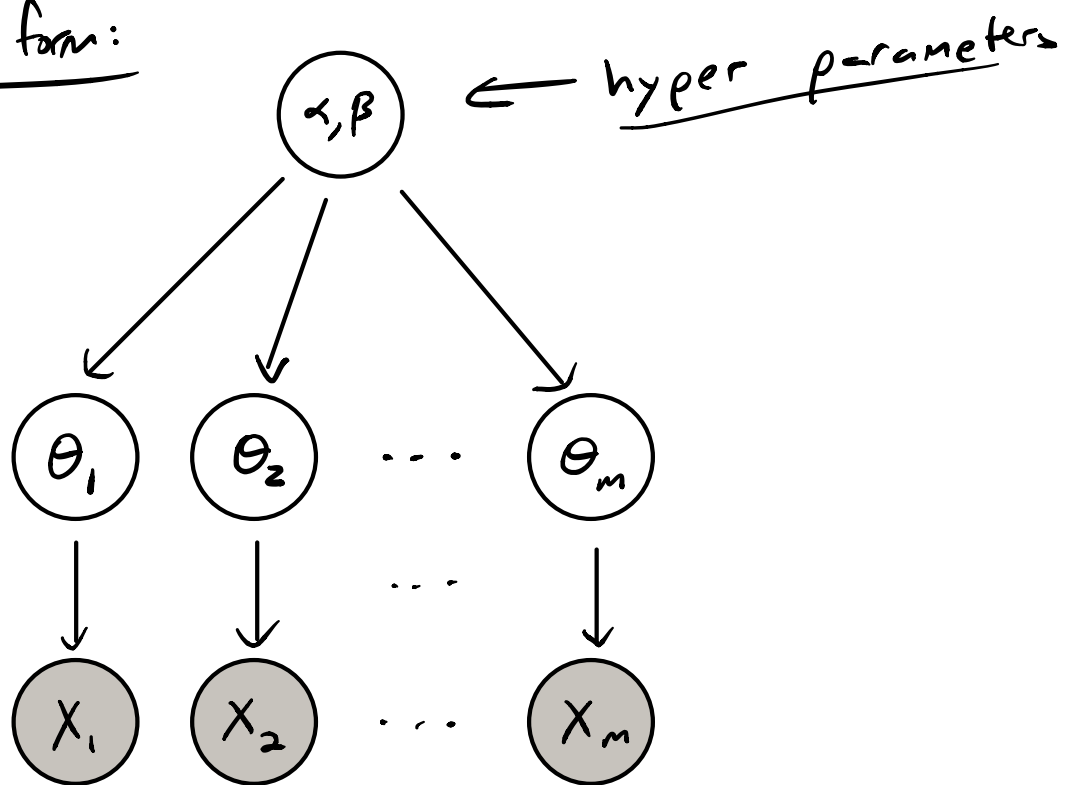
$$\alpha, \beta \sim \lambda_{\alpha, \beta} \quad (\text{say, indep. Exp}(1))$$

$$\theta_i | \alpha, \beta \stackrel{\text{iid}}{\sim} \text{Beta}(\alpha, \beta) \quad i \leq m$$

$$X_i | \theta_i \stackrel{\text{indep}}{\sim} \text{Binom}(n_i, \theta_i) \quad i \leq m$$

[Note: there is always an equivalent model where we marginalize over α, β and just write a more complicated prior on θ . Hierarchical version often gives better intuition or suggests computational strategies]

Graphical form:



This is a directed graphical model. Implies the distribution may be factorized with one factor for each vertex in a DAG (V, E)

$$p(z_1, \dots, z_{|V|}) = \prod_{i=1}^{|V|} p_i(z_i \mid \text{Pa}(z_i))$$

↑ parents

For this model,

$$\begin{aligned} p(\alpha, \beta, \theta_1, \dots, \theta_m, X_1, \dots, X_m) \\ = p(\alpha, \beta) \cdot \prod_i p(\theta_i \mid \alpha, \beta) \cdot \prod_i p(X_i \mid \theta_i) \end{aligned}$$

Practical implication:

X_2, \dots, X_m indirectly influence the estimate of X_1 , by teaching us what values of θ are plausible.

Markov Chain Monte Carlo

Hierarchical models can get very complex very fast,
creating big computational headaches

$$\lambda(\theta|x) = \frac{\int_{\Omega} p_{\theta}(x) \lambda(\theta) d\theta}{\int_{\Omega} p_{\theta}(x) d\theta} \quad \leftarrow \text{usually nice}$$

$\leftarrow \text{often intractable.}$

Computational strategy: set up a Markov chain
with stationary dist $\propto p_{\theta}(x) \lambda(\theta)$, run it
to get approximate samples from $\lambda(\theta|x)$

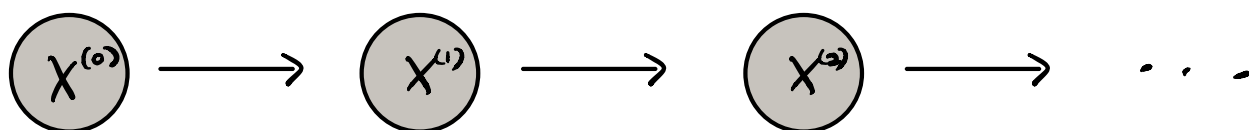
Definition: A (stationary) Markov chain with trans.
kernel $Q(y|x)$ and initial dist. $\pi_0(x)$ is
a sequence of r.v.s $X^{(0)}, X^{(1)}, \dots$ where $X^{(0)} \sim \pi_0$
and $X^{(t+1)} | X^{(0)}, \dots, X^{(t)} \sim Q(\cdot | X^{(t)})$

$$Q(y|x) = \mathbb{P}(X^{(t+1)} = y | X^{(t)} = x)$$

Marginal dist. of $X^{(1)}$:

$$\pi_1(y) = \mathbb{P}(X^{(1)} = y) = \int_{\mathcal{X}} Q(y|x) \pi_0(x) d\mu(x)$$

This is a directed graphical model:



If $\pi(y) = \int_{\mathcal{X}} Q(y|x) \pi(x) d\mu(x)$ we say π is a stationary distribution for Q

Sufficient condition is detailed balance:

$$\pi(x) Q(y|x) = \pi(y) Q(x|y) \quad \forall x, y$$

$$\begin{aligned} \Rightarrow \int_{\mathcal{X}} Q(y|x) \pi(x) d\mu(x) &= \pi(y) \int_{\mathcal{X}} Q(x|y) d\mu(x) \\ &= \pi(y) \end{aligned}$$

A Markov chain with detailed balance is called reversible since $(X^{(0)}, \dots, X^{(t)}) \stackrel{D}{=} (X^{(t)}, \dots, X^{(0)})$

Theorem: If an MC with stationary dist. π is:

- 1) Irreducible: $\forall x, y \exists n: p(X^{(n)} = y | X^{(0)} = x) > 0$ (for cts \mathcal{X})
- 2) Aperiodic: $\forall x, \gcd \{n > 0: p(X^{(n)} = x | X^{(0)} = x) > 0\} = 1$ can be generalized to cts \mathcal{X}

Then $\mathcal{L}(X^{(t)}) \xrightarrow{t \rightarrow \infty} \pi$ (in TV distance),
regardless of π_0 (chain "forgets" π_0)

[Proof beyond scope of our class]

Strategy: Find Q with stationary dist $\lambda(\theta|x)$,
start at any x , run chain for a long time
 $\leadsto X^{(t)} \approx$ sample from posterior, for large t .

Gibbs sampler

Parameter vector $\theta = (\theta_1, \dots, \theta_d)$

Algorithm:

Initialize $\theta = \theta^{(0)}$

For $t = 1, \dots, T$:

For $j = 1, \dots, d$:

Sample $\theta_j \sim \lambda(\theta_j | \theta_{-j}, x)$ } (*)

Record $\theta^{(t)} = \theta$

Variations on (*) :

- Update one random coordinate $J^{(t)} \sim \text{Unif}\{0, \dots, d\}$
- Update coordinates in random order

Advantage for hier-archical priors: only need to sample low-dimensional conditional dists:

$$\lambda(\theta_j | \theta_{-j}, x) \propto p(\theta_j | \theta_{P_2(j)}) \cdot \prod_{i: j \in P_2(i)} p(\theta_i | \theta_{P_2(i)})$$

Especially easy if using conjugate priors at all levels, often can be parallelized.

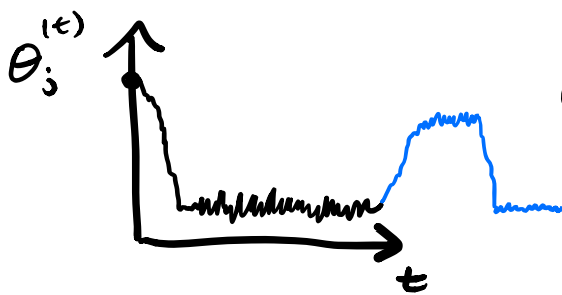
MCMC in Practice

In theory: Pick any initialization $\theta^{(0)}$ and valid kernel Q , sample long enough $\rightarrow \theta^{(t)} \approx \lambda(\theta | x)$

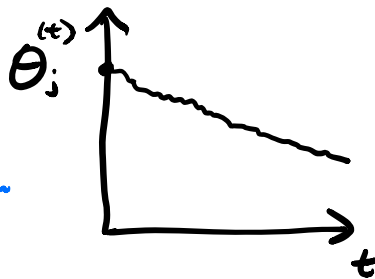
Do it again N more times $\rightarrow N$ samples from $\lambda(\theta | x)$

In practice, how do we know we've sampled long enough?

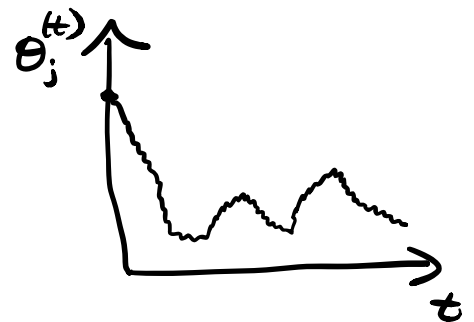
Trace plots: Show how fast the MC mixes



GOOD (?)



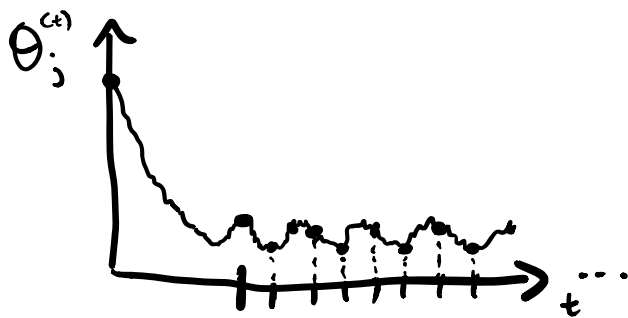
BAD



NOT GREAT

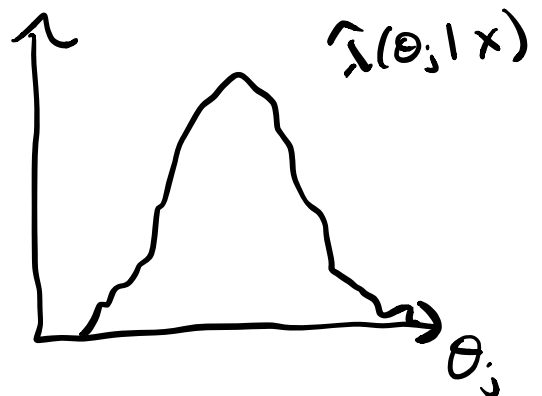
Can be deceived!

Esp. for bimodal posterior



Burn-in: "Forget" initialization

thinning: makes samples more independent



Estimate posterior based on $\{\theta_j^{(B)}, \theta_j^{(B+s)}, \dots, \theta_j^{(B+Ns)}\}$

Posterior mean: $\frac{1}{N+1} \sum_{k=0}^N \theta_j^{(B+ks)} \xrightarrow{N \rightarrow \infty} \mathbb{E}[\theta_j | x]$

Implementation details matter!

$$\theta_1, \theta_2 \stackrel{\text{iid.}}{\sim} N(0, 1)$$

$$X_i | \theta \stackrel{\text{iid.}}{\sim} N(\theta_1 + \theta_2, 1) \quad i=1, \dots, n$$

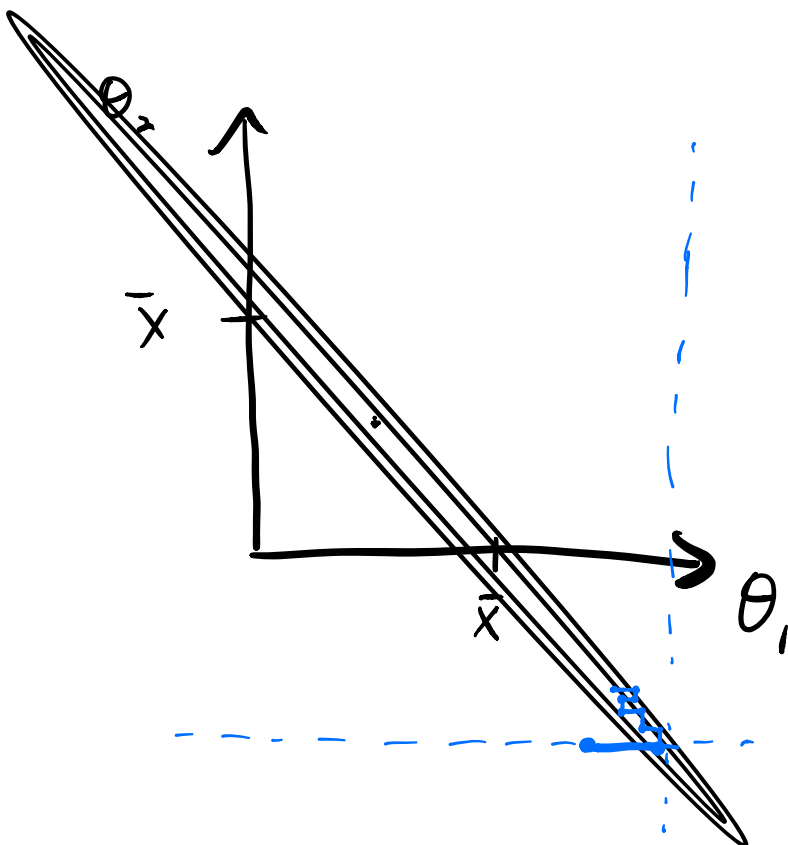
$$\Rightarrow \begin{pmatrix} \theta_1 \\ \theta_2 \\ \bar{x} \end{pmatrix} \sim N_3 \left(0, \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 1 & 2 + \frac{1}{n} \end{pmatrix} \right)$$

$$\theta | \bar{x} \sim N_2 \left(m(\bar{x}), \Sigma(\bar{x}) \right)$$

$$m(\bar{x}) = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \left(2 + \frac{1}{n} \right)^{-1} \bar{x} = \frac{n\bar{x}}{2n+1} \cdot \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\Sigma(\bar{x}) = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} - \begin{pmatrix} 1 \\ 1 \end{pmatrix} \left(2 + \frac{1}{n} \right)^{-1} \begin{pmatrix} 1 & 1 \end{pmatrix}$$

$$= \frac{n+1}{2n+1} \begin{pmatrix} 1 & -\frac{n}{n+1} \\ \frac{n}{n+1} & 1 \end{pmatrix}$$



Gibbs takes a long time to mix

Better parameterization:

$$\beta_1 = \theta_1 + \theta_2$$

$$\beta_2 = \theta_1 - \theta_2$$

$$\beta_1 \perp \beta_2 | x$$

Gibbs \Leftrightarrow Directly sampling from posterior.

Gaussian Hierarchical Model:

$$\tau^2 \sim \lambda(\tau) \quad \text{e.g.} \quad \frac{1}{\tau^2} \sim \text{Gamma}(k, s)$$

$$\theta_i | \tau^2 \stackrel{\text{iid}}{\sim} N(0, \tau^2) \quad i \leq n$$

$$X_i | \tau^2, \theta \stackrel{\text{ind.}}{\sim} N(\theta_i, 1)$$

Posterior mean:

$$\begin{aligned} \delta(x_i) &= \mathbb{E}[\theta_i | x] \\ &= \mathbb{E}\{\mathbb{E}[\theta_i | x, \tau^2] | x\} \\ &= \mathbb{E}\left[\frac{\tau^2}{1+\tau^2} x_i | x\right] \\ &= \underbrace{\left(\mathbb{E}\left[\frac{\tau^2}{1+\tau^2} | x\right]\right)}_{\text{shrinkage factor}} \cdot x_i \end{aligned}$$

Bayes estimate of optimal "shrinkage" factor

Define $\zeta = \frac{1}{1+\tau^2}$ ($\zeta = 0 \Leftrightarrow$ no shrinkage)

$$X_i | \tau^2 \stackrel{\text{iid}}{\sim} N(0, 1+\tau^2)$$

$$\Rightarrow X | \zeta \sim N_n(0, \zeta^{-1} \mathbf{I}_n)$$

$$= \left(\frac{\zeta}{2\pi}\right)^{n/2} e^{-\frac{\zeta}{2} \|X\|^2}$$

$$\propto_{\zeta} \text{Gamma}\left(1 + \frac{n}{2}, \frac{2}{\|X\|^2}\right)$$

shape scale

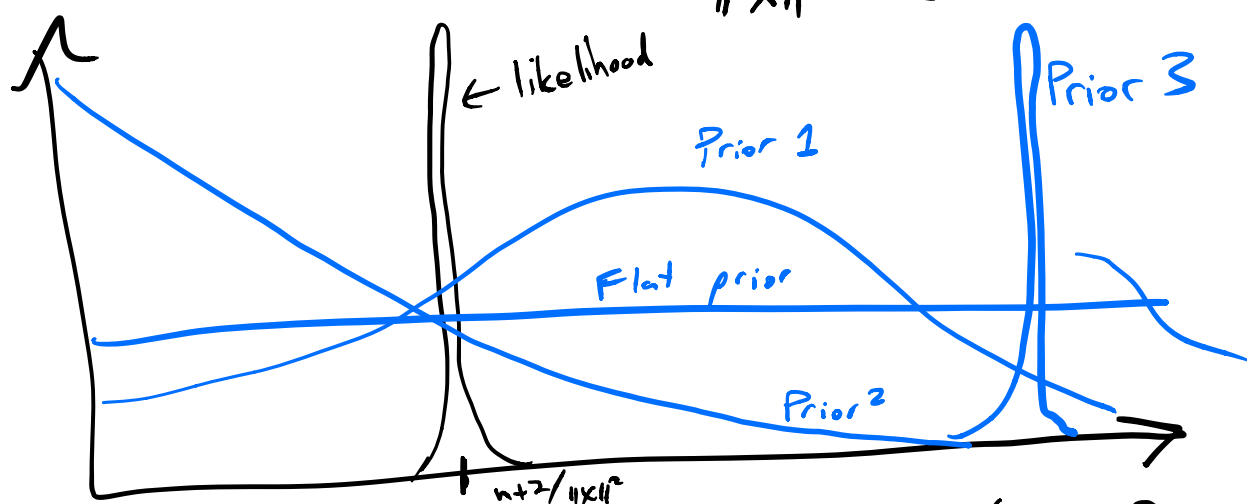
✓ note this is just the likelihood

This likelihood has a sharp peak at $\frac{n+2}{\|X\|^2} \approx \xi$

Why? $Z \sim \text{Gamma}(1 + \frac{n}{2}, 2/\|X\|^2)$

has mean $\frac{n+2}{\|X\|^2}$ ($\xrightarrow{n \rightarrow \infty} \xi$)

variance $\frac{n+4}{\|X\|^4}$ ($\xrightarrow{n \rightarrow \infty} 0$)



For any reasonably "open-minded" prior (not Prior 3),

$$\mathbb{E}[\xi|X] \approx \xi \Rightarrow \hat{\theta}_i \approx (1 - \xi) X_i$$

If prior doesn't matter much, why use one?

Could just estimate ξ from data

however we want, "plug it in"

Called "Empirical Bayes" a hybrid approach
in which hyper parameters treated as fixed,
others treated as random.