

(65)

9/27/19

Physics 8805: Learning from Data: Lecture 11

On board: Reminder: all handwritten lecture notes are available on Carmen in scanned pdf form.

Notebooks for today:

- MCMC - diagnostics, ipynb

Great reference: "The Matrix Cookbook". Full of useful identities.

Quick follow-up to last time:

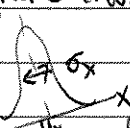
Suppose you have sampled the posterior $p(\vec{\theta} | D)$ where $\vec{\theta}$ is a vector of 8 parameters. You have N samples, each with 8 numbers. Now you want the joint posterior for θ_0 and θ_1 , with $\theta_2 - \theta_7$ marginalized. What do you do?

answer: Take θ_0 and θ_1 only from the N samples.

If you want a 2-d posterior, histogram them.

If you want $\langle g(\theta_0, \theta_1) \rangle \approx \frac{1}{N} \sum_{i=1}^N g(\theta_0^{(i)}, \theta_1^{(i)})$, ignoring $\theta_2 - \theta_7$!

• The confusion might be that you know $p(\theta_0, \theta_1) = \int d\theta_2 \dots d\theta_7 p(\theta_0, \theta_1, \dots, \theta_7)$ but don't realize that to do that integral in this context is just to ignore those variables (because projecting = marginalizing).

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma_x} e^{-\frac{(x-\mu_x)^2}{2\sigma_x^2}}$$


• Recap of combining random variables:

$$X \sim N(\mu_x, \sigma_x^2) \text{ and } Y \sim N(\mu_y, \sigma_y^2) \Rightarrow aX + bY \sim N(a\mu_x + b\mu_y, a^2\sigma_x^2 + b^2\sigma_y^2)$$

• generalize to $X_1 + X_2 + X_3 + \dots + X_m \sim N(\mu_1 + \dots + \mu_m, \sigma_1^2 + \dots + \sigma_m^2)$

• Now continue for case of correlated $X, Y \Rightarrow (63), (64)$

9/27/19

MCMC Sampling II: Assessing Convergence

• We've seen that using MCMC with the Metropolis-Hastings algorithm can lead to a Markov chain — a set of configurations of the parameters we are sampling — that enables inference \Rightarrow they are samples of the posterior of interest.

• But how do we know the chain is converged — that is, that it is providing samples of the stationary distribution?
 \Rightarrow we need diagnostics of convergence.

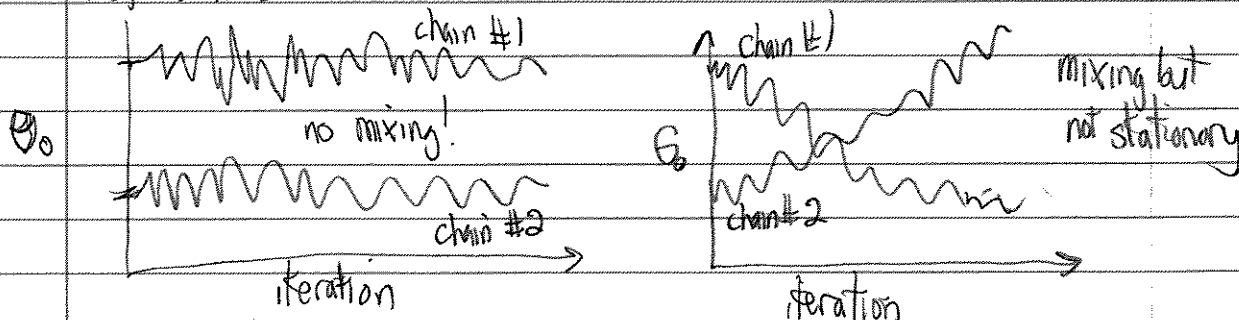
• What can go wrong?

- Could fail to have converged
- Could exhibit apparent convergence but full posterior has not been sampled.
- Convergence could be only approximate because the samples are correlated (need to run longer).

• Strategies to monitor convergence:

1. Run multiple chains with distributed starting points
2. Compute variation between and within chains
 \Rightarrow look for mixing and stationarity
3. Make sure acceptance rate for MC steps is not too low or too high.

Figure 11.3 From ODA-3:



• Future: parallel tempering

9/17/19

Step through MCMC-diagnostics.ipynb notebook

[Note: This notebook needs to be rewritten! Here: laundry list of diagnostics \Rightarrow return when we do pymc3.]

• BDA-3 Figure 11.1

a) not converged

b) 1000 iterations \rightarrow possibly converged

c) (correlated) draws from target distribution

• Straight line fitting again...

• emcee sampler - but with Metropolis-Hasting \Rightarrow stepsize specified

• get chains: sampler, flatchain flattens the chain (all walkers)

• How do we know this chain has converged to the posterior?

• Standard error of the mean

• How does mean of $\vec{\theta}$ deviate in the chain \Rightarrow simulation error of the mean, not the underlying uncertainty of $\vec{\theta}$.

$$SE(\bar{\theta}) = \frac{\text{Posterior Standard Deviation}}{\sqrt{N}}$$

← cf. notebook on CLT and expected distribution of means.

• Visualization with moving average \Rightarrow check for stability

• Autocorrelation: do you recognize the formula

• Acceptance rate. Usually autotuned in packaged MCMC software.

• Assessing mixing: Gelman Rubin diagnostic (page 1)

• Basic idea: multiple chains from different walkers (after warmup) are split up and one looks at the variance within a chain and between chains.

• See BDA-3 pp 284-5 for details. We'll come back to this.

• Now documented internally in notebooks, But we'll see it again.

don't use the code here as an example \Rightarrow use other notebooks

try changing step size from .002 to .005 or higher.

9/27/19

Check the $N=2$ case:

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \end{bmatrix} \begin{bmatrix} b \\ m \end{bmatrix} = \begin{bmatrix} b + mx_1 \\ b + mx_2 \end{bmatrix} \quad \text{so each row is } y_i = b + mx_i \quad \checkmark$$

• Reader: convince yourself that $N=3$ and higher is correct.

Aside: Why don't we just solve the matrix equation?

$$Y = A\theta \stackrel{?}{\Rightarrow} \theta = A^{-1}Y \quad ??$$

Because the equation is overconstrained for $N > 2$ (ok for $N=2$!)Frequentist answer: maximize the likelihood $\propto e^{-\chi^2}$

familiar, uncorrelated case: $\chi^2 = \sum_{i=1}^N \frac{(y_i - f(x_i))^2}{\sigma_{y_i}^2}$

χ^2 is a scalar.

generalized, correlated case: $\chi^2 = \underbrace{[Y - A\theta]^T}_{(1 \times N)} \underbrace{\Sigma^{-1}}_{(N \times N)} \underbrace{[Y - A\theta]}_{(N \times 1)} \rightarrow 1 \times 1$
maximum likelihood estimate MLE

Claim: $\hat{\theta} = [A^T \Sigma^{-1} A]^{-1} [A^T \Sigma^{-1} Y]$

Aside: why can't I say $[A^T \Sigma^{-1} A]^{-1} [A^T \Sigma^{-1} Y] = A^{-1} \Sigma (A^T A)^{-1} \Sigma^{-1} Y = A^{-1} Y$?

\Rightarrow not square, invertible matrices

Plausibility argument for $\hat{\theta}$ result. We need square, invertible matrices

$$\begin{aligned} Y &= A\theta & N \times 1 &= (N \times 2) \cdot (2 \times 1) & \text{for } \theta &= \begin{bmatrix} b \\ m \end{bmatrix} \\ \Sigma^{-1} Y &= \Sigma^{-1} A\theta & (N \times N) \cdot (N \times 1) &= N \times 1 \\ A^T \Sigma^{-1} Y &= A^T \Sigma^{-1} A\theta & (2 \times N) \cdot (N \times N) \cdot (N \times 1) &= 2 \times 1 \\ & & (2 \times N) \cdot (N \times N) \cdot (N \times 2) &= 2 \times 2 \end{aligned}$$

Now invert: $\theta = (A^T \Sigma^{-1} A)^{-1} (A^T \Sigma^{-1} Y)$

9/27/19

Before proving the result more carefully, let's generalize to a higher-order polynomial:

$$\text{If } y_i = qx_i^2 + mx_i + b \Rightarrow \theta = \begin{bmatrix} b \\ m \\ q \end{bmatrix} \Rightarrow A = \begin{bmatrix} 1 & x_1 & x_1^2 \\ \vdots & \vdots & \vdots \\ 1 & x_N & x_N^2 \end{bmatrix}$$

and similarly for higher order.

Prove using indices with summation convention:

$$\chi^2 = (Y - A\theta)^T \Sigma^{-1} (Y - A\theta) = (Y_i - A_{ij}\theta_j) \sum_{i,i'}^{-1} (Y_{i'} - A_{i'j'}\theta_{j'})$$

where i, i' run from 1 to N and j, j' from 1 to p (highest power x^{p-1})

Find MLE from $\frac{\partial \chi^2}{\partial \theta_k} = 0$ for $k=1$ to p

$$\Rightarrow \left. \frac{\partial \chi^2}{\partial \theta_k} \right|_{\theta=\hat{\theta}} = -A_{ij}\delta_{jk} \sum_{i,i'}^{-1} (Y_{i'} - A_{i'j'}\hat{\theta}_{j'}) + (Y_i - A_{ij}\hat{\theta}_j) \sum_{i,i'}^{-1} (-A_{i'j'}\delta_{jk}) = 0$$

• move A terms to the opposite side and show doubled terms are equal (overall multiply by -1)

$$\Rightarrow A_{ik} \sum_{i,i'}^{-1} Y_{i'} + Y_i \sum_{i,i'}^{-1} A_{ik} = A_{ik} \sum_{i,i'}^{-1} A_{i'j'} \hat{\theta}_{j'} + A_{ij} \hat{\theta}_j \sum_{i,i'}^{-1} A_{ik}$$

take $j=k$ symmetric

$$\underbrace{A_{ik} \sum_{i,i'}^{-1} Y_{i'}}_{(A^T \Sigma^{-1} Y)_k} + \underbrace{A_{ik} \sum_{i,i'}^{-1} Y_i}_{i=i' \text{ and } \Sigma_{ii}^{-1} = \Sigma_{ii}^{-1}} = A_{ik} \sum_{i,i'}^{-1} A_{i'j} \hat{\theta}_j + A_{ij} \hat{\theta}_j \sum_{i,i'}^{-1} A_{ik}$$

$$\Rightarrow 2(A^T)_k \sum_{i,i'}^{-1} Y_i = 2(A^T)_k (\Sigma^{-1})_{i,i'} A_{ij} \hat{\theta}_j$$

$$\text{or } (A^T \Sigma^{-1} Y) = (A^T \Sigma^{-1} A) \hat{\theta}$$

$$\Rightarrow \boxed{\hat{\theta} = (A^T \Sigma^{-1} A)^{-1} (A^T \Sigma^{-1} Y)} \quad \leftarrow \text{square, invertible} \quad \text{Q.E.D.}$$