

Lecture 08 - MCMC Convergence

Recap: Helpful convergence visualizations

- Traceplot of the sampling chain
- Histograms over subsets of the chain
- Traceplots of multiple chains with different random start points.

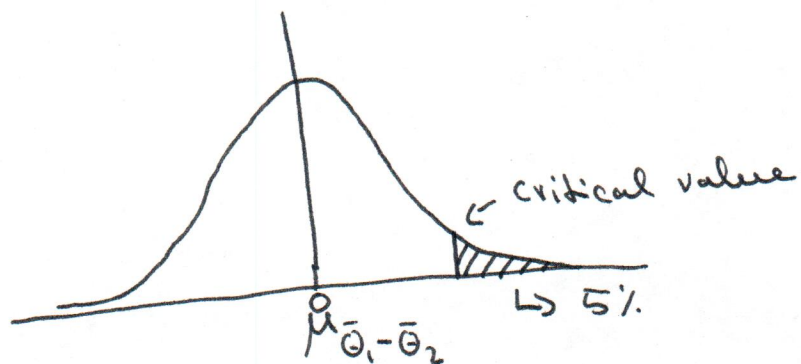
A bit more formal testing for convergence:

Geweke: Take two non-overlapping samples and compare the means.

⇒ Hypothesis test for difference of means

$$H_0: \mu_{\theta_1} - \mu_{\theta_2} = 0 \quad \Rightarrow \quad \mu_{\bar{\theta}_1 - \bar{\theta}_2} = 0$$

↳ mean of the distribution of the differences



what is the standard deviation of this distribution?

$$\sigma_{\bar{\theta}_1 - \bar{\theta}_2} = \sqrt{\frac{\text{var}(\theta_1)}{n_1} + \frac{\text{var}(\theta_2)}{n_2}}$$

Pavlos rule of thumb:

$$-2 < \frac{\bar{\theta}_1 - \bar{\theta}_2}{\sigma_{\bar{\theta}_1 - \bar{\theta}_2}} < 2 \quad \Rightarrow \quad 2 \cdot \sigma_{\bar{\theta}_1 - \bar{\theta}_2} > \bar{\theta}_1 - \bar{\theta}_2$$

↳ this is the 68, 95, 99, 7 rule

⇒ If the null hypothesis is correct there is only ~ 5% chance of the mean difference being larger than $2 \cdot \sigma_{\bar{\theta}_1 - \bar{\theta}_2}$.

Gelman-Rubin test

- uses multiple chains
- Compares between chain variance and within chain variance.
- Large deviation ⇒ have not converged yet.

Assuming m chains of length n :

within chain variance:

$$W = \frac{1}{m} \sum_{j=1}^m s_j^2$$

↳ mean of the chain variances

$$s_j^2 = \frac{1}{n-1} \sum_{i=1}^n (\theta_{ij} - \bar{\theta}_j)^2$$

↳ variance of the j th chain.

Between chain variance:

$$B = \frac{n}{m-1} \sum_{j=1}^m (\bar{\theta}_j - \bar{\bar{\theta}})^2 \quad \text{with} \quad \bar{\bar{\theta}} = \frac{1}{m} \sum_{j=1}^m \bar{\theta}_j$$

↳ mean of the chain means

⇒ variance of the chain means multiplied by the number of samples in each chain.

$$\text{var}(\theta|Y) = (1 - \frac{1}{n})W + \frac{1}{n}B \quad \Rightarrow \text{just a weighted average}$$

potential scale reducing factor:

$$\hat{R} = \sqrt{\frac{\text{var}(\theta|Y)}{W}} \quad \Rightarrow \text{should be close to one when we converge.}$$

If $\hat{R} \gg 1$ then between chain variance is greater than within chain variance.

Note: The starting points should be overdispersed to make this analysis useful.

Maize example:

Data:

Length: L

Diameter: D

Weight: W

Goal:

Estimate a simple linear relationship between these variables.

How do we find our model?

\Rightarrow start simple and make it more complex if you need to

\Rightarrow All models are wrong, but some models are useful. (George P. Box)

\Rightarrow Maize kind of looks like a cylinder:

$$W \propto L \cdot D^2$$

$$W = \beta_0 + \beta_1 \cdot L \cdot D^2 + \epsilon \quad \epsilon \sim N(0, \sigma_\epsilon)$$

\hookrightarrow error

Now we need our likelihood:

$$P(\text{Data} | \theta) = P(W_i, L_i, D_i | \beta_0, \beta_1, \sigma_\epsilon)$$

$\underbrace{L_i, D_i}_{\text{fixed input, no stochasticity}}$

\hookrightarrow stochasticity caused by measurement error.

$$P(w_i, L_i, D_i | \beta_0, \beta_1, \sigma_\epsilon)$$

$$\propto \prod_i P(w_i | \beta_0, \beta_1, \sigma_\epsilon, x_i)$$

\hookrightarrow measurement input
 L_i and D_i

$$\propto \prod_i e^{-\frac{(w_i - (\beta_0 + \beta_1 x_i))^2}{2\sigma_\epsilon^2}}$$

$$\propto e^{-\sum_i \frac{(w_i - (\beta_0 + \beta_1 x_i))^2}{2\sigma_\epsilon^2}}$$

Now we need to think about priors.

\Rightarrow We don't really know anything about β_0, β_1

\Rightarrow Make priors flat and uninformative

\Rightarrow Normal priors are nice, because they are the conjugate prior to our likelihood.

$$P(\beta_0) = N(\mu_{\beta_0}, \sigma_{\beta_0})$$

$$P(\beta_1) = N(\mu_{\beta_1}, \sigma_{\beta_1})$$

What about $P(\sigma_\epsilon^2)$?

Although our likelihood is a Gaussian, it is not normal with respect to σ_ϵ^2 !

The conjugate prior for the variance of a normal distribution is an Inverse gamma.

Instead we can also use the precision $\tau_G = \frac{1}{\sigma_G^2}$ and the conjugate prior is a gamma distribution.

In all we now have our posterior as:

$$p(\beta_0, \beta_1, \tau_G | w_i, x_i) \propto P(w_i | \beta_0, \beta_1, \tau_G, x_i) \\ \cdot P(\beta_0) \cdot P(\beta_1) \cdot P(\tau_G)$$

We can sample from this using Metropolis's-Hastings.

⇒ We need our proposal distribution.

$$q(\theta^* | \theta^{(i)}) \quad \text{for } \theta: \beta_0, \beta_1, \tau_G$$

normal or uniform are generally good

choices ⇒ $N(\theta^{(i)}, \lambda)$

↳ we need to tune

λ (corresponds to our step size).

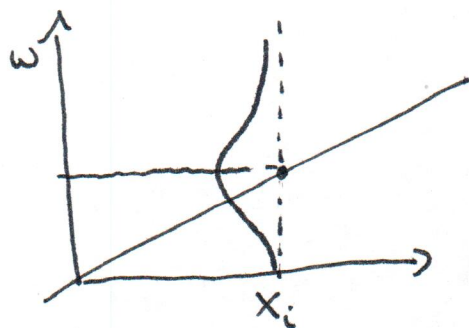
Predictive Probability:

How do we check if our model works?

- Sanity test: visual inspection
- Better: Evaluate model by making data predictions.
 - ↳ Simulated data should have the same distribution as actual data.

In the Bayesian view we have probability distributions over our parameters, not value estimates.

⇒ for each input we have a predictive distribution



⇒ we need to average over all possible parameter values weighted by their posterior probability.

predictive probability distribution:

$$p(w^* | x^*, x, w) = \underbrace{\int p(\theta | x, w)}_{\text{posterior}} \cdot \underbrace{p(w^* | x^*, \theta)}_{\text{likelihood}} d\theta$$

⇒ likelihood gives the likelihood of the new unseen data point for a given choice of parameter.

\Rightarrow posterior says how likely this choice of parameter is.

How do we do this in praxis?

- We have an analytic form for the likelihood (because we designed it)
- We have (or can get) samples from our posterior.

From the homework we remember:

$$\int h(x) \cdot g(x) dx \propto \sum_{x \sim g(x)} h(x)$$

$$\propto \sum_{i \sim \pi} h(x_i) \cdot g(x_i)$$

\Rightarrow take the samples from the posterior, plug them into the likelihood and we are good.

$$\Rightarrow w_{MAP} = \underset{w}{\operatorname{argmax}} P(w^* | x^*, w, x)$$

