

#### 37262 Mathematical Statistics

Lecture 11



UTS CRICOS 00099F

#### What Does a "Good" MCMC Simulation Look Like?

- The Markov Chain Monte Carlo method allows us to generate samples from a distribution π(x) by constructing a Markov Chain whose equilibrium distribution is π(x) and sampling from its long-run distribution instead.
- In a perfect world, we would have a Markov Chain which was aperiodic and ergodic and we would run the chain for a very long time before sampling each individual point, restarting the chain after each observation.
- Clearly this is hugely wasteful in terms of time and computer resources, so we prefer to be able to use all (or almost all) positions of the chain once we are content that it has reached equilibrium (i.e. after a "burn in" period.)

# **MCMC** Diagnostics

- If we want to generate independent samples from an equilibrium, we want to see that, ideally, every state is reachable immediately after every other state.
- A traceplot shows which state the chain is in at each iteration.
- We say a chain is **mixing** well if no obvious patterns are observable in the traceplot.
- If the samples from the equilibrium distribution were truly independent, then there would be perfect mixing with all iterations having the same probability of showing a given state, irrespective of what position the chain was in at the previous iteration.







Good Mixing

# MCMC Diagnostics

- A **running mean** plot can also be used to identify obvious flaws in the mixing.
- Plotting the mean of all iterations up to and including the current step should show no strong signal once the burn in period has passed.



Bad Mixing





## **MCMC** Diagnostics

• Problems with MCMC mixing can also show up in an autocorrelation plot.



## MCMC and Bayes' Theorem

• Recall the Metropolis-Hastings algorithm. Given proposal distribution g(x), a realisation is generated from  $\pi(x)$  by accepting the move proposed by g(x),  $x_p$ , with probability

$$A(x_{p}|x) = \min\left\{1, \frac{\pi(x_{p})g(x|x_{p})}{\pi(x)g(x_{p}|x)}\right\}.$$

• One of the main strengths of MCMC is that we do not need to be able to calculate  $\pi(x)$  as the

acceptance probability only relies on the ratio 
$$\frac{\pi(x_{\rho})}{\pi(x)}$$
, which is often easier.

- As we will see, this allows us to sidestep one of the most difficult parts of Bayesian inference.
- Especially when we are dealing with multivariate probabilities, this property is extremely useful.

🕉 UTS

- Bayesian inference uses Bayes' Theorem to update the belief around a parameter as more evidence is gathered.
- Consider flipping a single coin repeatedly and recording its outcome. We do not know what the probability θ that the coin will lands Heads is, but rather will have an initial belief which we update each time more evidence is observed. Call the outcome of the *i*th flip x<sub>i</sub>.
- We have a **prior distribution**  $\pi(\theta)$  which describes our belief in the parameter before any observations have been gathered.
- After observations  $\mathbf{X} = (x_1, x_2, ..., x_n)$  have been gathered, we update our belief of the value of  $\theta$ . (Assume  $x_i = 1$  if the coin lands Heads and 0 if it lands Tails.)



- We call our initial belief the prior distribution  $\pi(\theta)$ .
- From this, we calculate the **likelihood** of seeing the evidence which we did, assuming  $\theta$ ,  $f(\mathbf{X}|\theta)$
- Bayes' Theorem then gives  $f(\theta | \mathbf{X}) = \frac{f(\mathbf{X} | \theta) \pi(\theta)}{f(\mathbf{X})} = \frac{f(\mathbf{X} | \theta) \pi(\theta)}{\int_{\mathbf{X}} f(\mathbf{X}, \theta) d\theta}$  where  $f(\theta | \mathbf{X})$  is the **posterior**
- The posterior density is equal to the prior multiplied by the likelihood, divided by the integral of this over all possible values of θ.
- Without knowing the denominator, this is sometimes said as "posterior  $\propto$  prior  $\times$  likelihood."

- Consider the problem of describing our uncertainty about the probability  $\theta$  that a coin will land Heads when flipped.
- Initially, the experimenter believes the coin is fair, or close to fair and will update this belief based on future observations.
- The prior belief on the probability is described by  $\Theta \sim N(0.5, \sigma^2)$ where  $\sigma^2$  is some known constant.
- (This is not a very good prior as does put some belief onto probabilities of below 0 or above 1, which is clearly not sensible.)





- Bayes' Theorem gives  $f(\theta | \mathbf{X}) = \frac{f(\mathbf{X} | \theta) f(\theta)}{f(\mathbf{X})} = \frac{f(\mathbf{X} | \theta) f(\theta)}{\int_{\Theta} f(\mathbf{X}, \theta) d\theta}$
- The probability that, say,  $\mathbf{X} = (x_1, x_2, \dots, x_n) = (0, 1, 1, 0, 0)$  for a known  $\theta$  is  $f(\mathbf{X}|\theta) = (1-\theta)\theta\theta(1-\theta)(1-\theta) = (1-\theta)^3\theta^2$ .
- The numerator is trivial to calculate. If  $\Theta \sim N(0.5, \sigma^2)$ , then  $f(\theta) = \frac{1}{\sqrt{2\pi\sigma^2}}e^{-\frac{(\theta-0.5)^2}{2\sigma^2}}$ .
- The denominator is, however, much more problematic  $f(\mathbf{X}) = \int_{-\infty}^{\infty} (1-\theta)^3 \theta^2 \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\theta-0.5)^2}{2\sigma^2}} d\theta$ .



• If, however, we wished to know what the posterior distribution of the parameter  $\theta$  was, we could sample from this easily by MCMC.



 As more and more evidence becomes available, for a well-defined Bayesian framework, the posterior distribution should converge ever more closely to the information provided by the observations.



- As models (and resulting probability density functions or probability mass functions) increase in complexity, the advantages of using MCMC become ever greater.
- Consider drawing samples from  $X \sim N(\mu, \sigma^2)$  where neither  $\mu$  nor  $\sigma^2$  are definitely known.
- We might, though, know what these parameters should be, plus or minus some uncertainties.
- For example, we might say  $\mu \sim N(10, 0.5)$  and  $\sigma^2 \sim N(0.1, 0.0001)$ .
- How would we update these beliefs as we observe more and more realisations of X?

Les Kirkup and Bob Frenkel

#### An Introduction to Uncertainty in Measurement



• Given observations of X,  $\mathbf{X} = (x_1, x_2, ..., x_n)$ , our posterior belief of the values of  $\mu$  and  $\sigma^2$  is

$$f(\mu,\sigma^{2}|\mathbf{X}) = \frac{f(\mathbf{X}|\mu,\sigma^{2})f(\mu,\sigma^{2})}{f(\mathbf{X})} = \frac{\frac{1}{\sqrt{2\pi\sigma^{2}}}e^{-\frac{(x-\mu)^{2}}{2\sigma^{2}}}\frac{1}{\sqrt{2\pi(0.5)}}e^{-\frac{(\mu-10)^{2}}{2(0.5)}}\frac{1}{\sqrt{2\pi(0.01)}}e^{-\frac{(\sigma^{2}-0.1)^{2}}{2(0.0001)}}}{\int_{-\infty}^{\infty}\int_{-\infty}^{\infty}\frac{1}{\sqrt{2\pi\sigma^{2}}}e^{-\frac{(x-\mu)^{2}}{2\sigma^{2}}}\frac{1}{\sqrt{2\pi(0.5)}}e^{-\frac{(\mu-10)^{2}}{2(0.5)}}\frac{1}{\sqrt{2\pi(0.01)}}e^{-\frac{(\sigma^{2}-0.1)^{2}}{2(0.0001)}}d\mu d\sigma^{2}}$$

- The denominator is hardly trivial to calculate...
- Using Metropolis-Hastings, however, we would only ever need to calculate the numerator, which requires no integration and into which proposed values for  $\mu$  and  $\sigma^2$  could simply be plugged.



- In some cases, we can choose a prior which gives a simple form for the posterior distribution which can sidestep the need for MCMC (or similar) methods.
- Consider again our flipping of a coin which lands Heads with probability  $\theta$  for some unknown  $\theta$ .
- Given  $\theta$ , the likelihood of a given sequence of outcomes  $\mathbf{X} = (x_1, x_2, ..., x_n)$  is  $f(\mathbf{X}|\theta) = \theta^k (1-\theta)^{n-k}$  where k is the number of times Heads was observed.
- If we do not know θ, but have no reason to believe (before any observations) that it is not close to 0.5, we may choose a prior distribution which has expectation 0.5.



• Recall the beta distribution, with density function  $f(y) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} y^{\alpha-1} (1-y)^{\beta-1}$  for  $\alpha, \beta > 0$  which takes values on [0,1].

(Note  $\Gamma(z)$  is the gamma function  $\Gamma(z) = \int_{0}^{\infty} t^{z-1}e^{-t}dt$ . For integer values of z,  $\Gamma(z) = (z-1)!$  and for all z,  $\Gamma(z+1) = z\Gamma(z)$ )

• The expectation of this is easy to calculate

$$E(Y) = \int_{0}^{1} y \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} y^{\alpha - 1} (1 - y)^{\beta - 1} dy = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \int_{0}^{1} y^{\alpha} (1 - y)^{\beta - 1} dy = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(\alpha + 1)\Gamma(\beta)}{\Gamma(\alpha + \beta + 1)}$$

$$=\frac{\Gamma(\alpha+1)}{\Gamma(\alpha)}\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha+\beta+1)} = \frac{\alpha}{\alpha+\beta}$$

- When examining Bernoulli trials (and hence Binomial distributions), it is often desirable to use a beta prior on the probability parameter.
- Since the likelihood is  $f(\mathbf{X}|\theta) = \theta^{k}(1-\theta)^{n-k}$ , then setting a prior of  $f(\theta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}\theta^{\alpha-1}(1-\theta)^{\beta-1}$ gives a posterior proportional to  $\theta^{k}(1-\theta)^{n-k} \times \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}\theta^{\alpha-1}(1-\theta)^{\beta-1} \propto \theta^{k+\alpha-1}(1-\theta)^{n-k+\beta-1}$
- This therefore gives a posterior which is also beta distributed,  $\theta | \mathbf{X} \sim Beta(\alpha + k, \beta + n k)$ .
- This has expectation  $\frac{\alpha + k}{\alpha + \beta + n}$ .
- If we set a large  $\alpha$  and  $\beta$ , this is  $\approx \frac{\alpha}{\alpha + \beta}$ , however as *n* and *k* get bigger, it is  $\approx \frac{k}{n}$ .

- We say that a distribution is a **conjugate prior** for a given likelihood function if it results in a posterior distribution which is of the same form (usually with updated parameters.)
- We have already seen that if  $X \sim Bern(\theta)$  then  $\theta \sim Beta(\alpha, \beta)$  is a conjugate prior since we get a posterior distribution  $\theta | \mathbf{X} \sim Beta(\alpha + k, \beta + n k)$  (where  $\mathbf{X}$  contains k 1s and n k 0s.)
- Other conjugate priors include, for  $X \sim N(\mu, \sigma^2)$  where  $\sigma^2$  is known but  $\mu$  is uncertain,  $\cdot \mu \sim N(m, s^2)$ .
- Although not simple to derive, this gives a posterior of

$$\mu \left| \boldsymbol{X} \sim N\left( \left( \left( \frac{m}{s^2} + \frac{1}{\sigma^2} \sum_{i=1}^n \boldsymbol{X}_i \right) \left( \frac{1}{\sigma^2} + \frac{n}{s^2} \right)^{-1} \right), \left( \frac{1}{\sigma^2} + \frac{n}{s^2} \right)^{-1} \right) \right|$$

- Suppose that it is known that faults arise in a system at the instants of a Poisson process where the rate of this process is not known.
- This means that the time between successive arrivals are independent realisations of an exp(λ) variable for some unknown λ > 0.



 Initially, the prior belief is that λ ~ Gamma(1,2) i.e. it is thought that λ is probably around 0.5, but that there is some uncertainty around this value. As observations become available, this belief is updated.

- Suppose now that observations of the interarrival times  $W_1, W_2, ..., W_n$  are recorded.
- If λ were known, we could work out the likelihood of obtaining the observations we did.

$$L(w_1, w_2, ..., w_n | \lambda) = \left[ \lambda e^{-\lambda w_1} \right] \left[ \lambda e^{-\lambda w_2} \right] ... \left[ \lambda e^{-\lambda w_n} \right] 0.0$$

$$= \lambda^n \left[ e^{-\lambda (w_1 + w_2 + ... + w_n)} \right]$$

$$0.0$$

0.6

0.4

0.2

• Bayes' theorem gives 
$$f(\lambda | w_1, w_2, \dots, w_n) = \frac{L(w_1, w_2, \dots, w_n | \lambda) f(\lambda)}{f(w_1, w_2, \dots, w_n)}$$
.

 $f(\boldsymbol{x}|\boldsymbol{\alpha},\boldsymbol{\beta}) = \frac{\boldsymbol{\beta}^{\alpha} \boldsymbol{x}^{(\alpha-1)} \boldsymbol{e}^{-\boldsymbol{\beta}\boldsymbol{x}}}{\boldsymbol{\Gamma}}$ 

 $\lambda \sim Gamma(1,2)$ 



• Since the posterior distribution

$$f(\lambda|w_1, w_2, \dots, w_n) = \frac{\lambda^{(n+\alpha)-1} \left[ e^{-\lambda(\beta+w_1+w_2+\dots+w_n)} \right] \left[ \beta^{\alpha} \right]}{f(w_1, w_2, \dots, w_n) \Gamma(\alpha)}$$

is also of a form proportional to  $\lambda^{(\alpha-1)}e^{-\beta\lambda}$ we note that our posterior distribution after observing the data is

$$\lambda \sim Gamma\left(n+\alpha, \sum_{j=1}^{n} w_j + \beta\right).$$



- For example, if we had observations
  - $W_1 = 0.5$  $W_2 = 0.3$  $W_3 = 0.3$

 $W_4 = 0.2$ 

**ÖUTS** 

our posterior belief about  $\lambda$  would be

$$\lambda \sim Gamma\left(n + \alpha, \sum_{j=1}^{n} w_j + \beta\right)$$
  
i.e.  $\lambda \sim Gamma(5, 3.3)$ 



• As the expectation of a  $Gamma(\alpha, \beta)$  variable is  $\frac{\alpha}{\beta}$ , our smaller than expected observations of *w* led to our belief in the rate parameter increasing.

- A negative binomial variable with range {r, r + 1, r + 2,...}can be considered as the sum of r independent identically distributed geometric random variables. It describes how many independent identically distributed Bernoulli variables must be observed until r successes (1s) have been observed.
- For *N* ~ *NegBin*(*r*,*p*), the probability mass function is

$$f(n) = P(N = n) = \begin{cases} \frac{(n-1)!}{(r-1)!(n-r)!} p^r (1-p)^{n-r} & n \in \{r, r+1, r+2, ...\} \\ 0 & \text{otherwise} \end{cases}$$



- Consider a series of observations of  $N \sim NegBin(r, p)$ ,  $\mathbf{N} = (n_1, n_2, ..., n_k)$ .
- The likelihood of this sample (assuming *r* is known) would then be

$$L(\boldsymbol{N} \mid \boldsymbol{p}) = \prod_{i=1}^{m} \frac{(n_i - 1)!}{(r - 1)!(n_i - r)!} \boldsymbol{p}^r (1 - \boldsymbol{p})^{n_i - r}. \quad (\text{assuming each } n_i \in \{r, r + 1, r + 2, ...\})$$

We now consider quantifying our uncertainty about the value of *p*. Initially, we place a prior belief on *p*, *P* ~ *Beta*(*α*, *β*) since this only takes values on [0,1].

• 
$$f(p) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha - 1} (1 - p)^{\beta - 1}$$
 for  $\alpha, \beta > 0$ 

• We can show that this is, in fact, a conjugate prior.

•  $f(p \mid \mathbf{N}) \propto L(\mathbf{N} \mid p) f(p)$  hence  $f(p \mid \mathbf{N}) \propto \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha - 1} (1 - p)^{\beta - 1} \prod_{i=1}^{m} \frac{(n_i - 1)!}{(r - 1)!(n_i - r)!} p^r (1 - p)^{n_i - r}$ 

• 
$$f(p \mid \mathbf{N}) \propto \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \prod_{i=1}^{m} \frac{(n_i - 1)!}{(r - 1)!(n_i - r)!} p^{\alpha - 1} p^{\sum_{i=1}^{m} r} (1 - p)^{\beta - 1} (1 - p)^{\sum_{i=1}^{m} n_i - r} \propto p^{\alpha + mr - 1} (1 - p)^{\beta + \sum_{i=1}^{m} n_i - mr - 1}$$

- This gives a posterior density of  $Beta\left(\alpha + mr, \beta + \sum_{i=1}^{m} n_i mr\right)$
- As both the prior and posterior distribution are beta distributed, the prior is a conjugate prior for negative binomial outcomes.

