### 36-463/663: Hierarchical Linear Models

Taste of MCMC / Bayes for 3 or more "levels" Brian Junker 132E Baker Hall brian@stat.cmu.edu

11/3/2016

### Outline

- Practical Bayes
- Mastery Learning Example
- A brief taste of JAGS and RUBE
- Hierarchical Form of Bayesian Models
- Extending the "Levels" and the "Slogan"
- Mastery Learning Distribution of "mastery"
- (to be continued...)

#### Practical Bayesian Statistics

- (posterior)  $\propto$  (likelihood)imes(prior)
- We typically want quantities like
  - point estimate: Posterior mean, median, mode
  - uncertainty: SE, IQR, or other measure of 'spread'
  - credible interval (CI)
    - $(\hat{\theta} 2SE, \hat{\theta} + 2SE)$
    - ( $\theta_{0.025}$ ,  $\theta_{0.975}$ )
  - Other aspects of the "shape" of the posterior distribution
- Aside: If (prior)  $\propto$  1, then
  - $\square$  (posterior)  $\propto$  (likelihood)
  - □ posterior mode = mle, posterior SE =  $1/I(\theta)^{1/2}$ , etc.

11/3/2016

## Obtaining posterior point estimates, credible intervals

- Easy if we recognize the posterior distribution and we have formulae for means, variances, etc.
- Whether or not we have such formulae, we can get similar information by simulating from the posterior distribution.

• Key idea: 
$$E[g(\theta)|\text{data}] = \int g(\theta)f(\theta|\text{data})d\theta$$
  
 $\approx \frac{1}{M}\sum_{m=1}^{M}g(\theta_m)$ 

where  $\theta_{_1}, \theta_{_2}, ..., \theta_{_M}$  is a sample from f( $\theta$ |data).

#### Example: Mastery Learning

- Some computer-based tutoring systems declare that you have "mastered" a skill if you can perform it successfully r times.
- The number of times x that you erroneously perform the skill before the r<sup>th</sup> success is a measure of how likely you are to perform the skill correctly (how well you know the skill).
- The distribution for the number of failures x before the r<sup>th</sup> success is the negative binomial distribution.

11/3/2016

#### **Negative Binomial Distribution**

Let X = x, the number of failures before the r<sup>th</sup> success. The negative binomial distribution is

$$f(x|p,r) = \binom{x+r-1}{x} p^r (1-p)^x$$

 The conjugate prior distribution is the beta distribution

$$f(p|\alpha,\beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1} (1-p)^{\beta-1}$$

Posterior inference for negative binomial and beta prior

• (posterior)  $\propto$  (likelihood) $\times$ (prior):

$$f(p|X = x) \propto p^r (1-p)^x \times p^{\alpha-1} (1-p)^{\beta-1}$$

Since this is Beta(p|α+r,β+x),

$$E[p|X = x] = \frac{\alpha + r}{\alpha + r + \beta + x}$$
$$Var(p|X = x) = \frac{(\alpha + r)(\beta + x)}{(\alpha + r + \beta + x)^2(\alpha + r + \beta + x + 1)}$$

 and this lets us compute point estimate and approximate 95% Cl's.

11/3/2016

#### Some specific cases...

- Let α = 1, β = 1 (so the prior is Unif(0,1)), and suppose we declare mastery for 3 successes
- If X = 4, then

$$E[p|X = x] = \frac{\alpha + r}{\alpha + r + \beta + x} = \frac{1+3}{1+3+1+4} = \frac{4}{9}$$

$$Var(p|X = x) = \frac{(1+3)(1+4)}{(1+3+1+4)^2(1+3+1+4+1)} = \frac{20}{810}$$
Approx 95% CI: (0.12, 0.76)

#### If X = 10, then

$$E[p|X = x] = \frac{\alpha + r}{\alpha + r + \beta + x} = \frac{1+3}{1+3+1+10} = \frac{4}{15}$$
$$Var(p|X = x) = \frac{(1+3)(1+10)}{(1+3+1+10)^2(1+3+1+10+1)} = \frac{44}{3600}$$
Approx 95% CI: (0.05, 0.49)

11/3/2016



11/3/2016

Using a non-conjugate prior distribution

Instead of a uniform prior, suppose psychological theory says the prior density should increase linearly from 0.5 to 1.5 as  $\theta$  moves from 0 to 1?



Posterior inference for negative binomial and linear prior

• (posterior)  $\propto$  (likelihood) $\times$ (prior):

$$f(p|X = x) \propto p^{r}(1-p)^{x} \times (p+0.5)$$
• r = 3, x = 4:  

$$f(p|X = x) \propto p^{3}(1-p)^{4}(p+0.5) ??$$
• r = 3, x = 10:  

$$f(p|X = x) \propto p^{3}(1-p)^{10}(p+0.5) ??$$

11/3/2016

What to do with 
$$p^r(1-p)^x(p+0.5)$$
 ?

- No formulae for means, variances!
- No nice function in R for simulation!
- There are many simulation methods that we can build "from scratch"
  - □ B. Ripley (1981) *Stochastic Simulation*, Wiley.
  - L. Devroye (1986) *Nonuniform random variate generation*, Springer.
- We will focus on just one method:
  - Markov Chain Monte Carlo (MCMC)
  - Programs like WinBUGS and JAGS automate the simulation/sampling method for us.



the *inverse probability transform* 

11/3/2016

#### Aside: the Inverse Probability

Transform

- If U ~ Unif(0,1) and F(z) is the CDF of a continuous random variable, then Z = F<sup>-1</sup>(U) will have F(z) as its CDF (exercise!!).
- The prior for p in our model has pdf p + 1/2, so its CDF is F(p)=(p<sup>2</sup> + p)/2.

• 
$$F^{-1}(u) = (2u + 1/4)^{1/2} - 1/2$$

# JAGS: Need to specify likelihood and prior

Likelihood is easy:

 $x \sim dnegbin(p,r)$ 

Prior for p is: p <- (u + 1/4)^{1/2} + 1/2 where u ~ dunif(0,1)

11/3/2016

15

### The JAGS model

```
mastery.learning.model <- "model {
    # specify the number of successes needed for mastery
    r <- 3
    # specify likelihood
    x ~ dnegbin(p,r)
    # specify prior (using inverse probability xform)
    p <- sqrt(2*u + 0.25) - 0.5
    u ~ dunif(0,1)</pre>
```

} "

#### Estimate p when student has 4 failures before 3 successes > library(R2jags) > library(rube) Histogram of p > mastery.data <- list(x=4)</pre> > rube.fit <- rube(mastery.learning.model,</pre> mastery.data, mastery.init, +2.0 parameters.to.save=c("p"), +n.iter=20000, n.thin=3) +9 > # generates about 10,000 samples Den > p <- rube.fit\$sims.list\$p</pre> 0 > hist(p,prob=T) > lines(density(p)) 0.5 > mean(p) + c(-2, 0, 2) \* sd(p)[1] 0.1549460 0.4685989 0.7822518 0:0 > quantile(p,c(0.025, 0.50, 0.975)) 0.0 0.2 0.4 0.6 0.8 2.5% 50% 97.5% p 0.1777994 0.4661076 0.7801111

11/3/2016

# Estimate p when student has 10 failures before 3 successes



#### Hierarchical Form of Bayesian Models

- In the past few lectures, we've seen many models that involve a likelihood and a prior
- A somewhat more compact notation is used to describe the models, and we will start using it now
  - □ Level 1: the likelihood (the distribution of the data)
  - □ Level 2: the prior (the distribution of the parameter(s))
- Eventually there will be other levels as well!
- Examples on the following pages...

11/3/2016

#### Hierarchical Beta-Binomial model

Likelihood is binomial:
 Level 1: x ~ Binom(x|n,p)

 $f(x|n,p) = \binom{n}{x} p^x (1-p)^{n-x}$ 

- Prior is beta distribution:
- $f(p|\alpha,\beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1} (1-p)^{\beta-1}$
- (posterior) ∝
   (likelihood)×(prior)
- Level 2:  $p \sim \text{Beta}(p \mid \alpha, \beta)$
- (posterior) ∝
   (level 1)×(level 2)
   = Beta(p| α+x,β+n-x)

#### Hierarchical Gamma-Exponential Model

Likelihood is Exponential:
$$f(x_1, \dots, x_n) = \prod_{i=1}^n \lambda e^{-\lambda x_i}$$
Prior is Gamma:
Prior is Gamma:
Level 1:  
 $x_i \stackrel{iid}{\sim} Expon(x|\lambda),$   
 $i = 1, \dots, n$ 
Level 2:  
 $f(\lambda|\alpha, \beta) = \frac{\beta^{\alpha}}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\lambda\beta}$ 
 $\lambda \sim Gamma(\lambda|\alpha, \beta)$ 
(posterior)  $\propto$   
(likelihood)×(prior)
(posterior)  $\propto$   
(level 1)×(level 2)  
 $= Gamma(\lambda|\alpha+n,\beta+n\bar{x})$ 

11/3/2016

### Hierarchical Normal-Normal model

#### Hierarchical Beta-Negative Binomial Model

 Likelihood is Negative-Binomial:

$$f(x|p,r) = \binom{x+r-1}{x} p^r (1-p)^x$$

- Prior is Gamma:  $f(p|\alpha,\beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}p^{\alpha-1}(1-p)^{\beta-1}$ 
  - (posterior)∝
     (likelihood)×(prior)

Level 1:

$$x \sim NB(x|r, p)$$

Level 2:

$$p \sim Beta(p|\alpha,\beta)$$

• (posterior) $\propto$ (level 1)×(level 2)  $p \sim Beta(p|\alpha + r, \beta + x)$ 

11/3/2016

#### Hierarchical Linear-Negative Binomial Model

 Likelihood is Negative-Binomial:

$$f(x|p,r) = \binom{x+r-1}{x} p^r (1-p)^x$$

Prior is Linear Density:

$$f(p) = p + 0.5, \ 0 \le p \le 1$$

(posterior)∝
 (likelihood)×(prior)

Level 1:

$$x \sim NB(x|r, p)$$

Level 2:

$$p \sim f(p) = p + 0.5$$

• (posterior)
$$\propto$$
  
(level 1)×(level 2)  
 $f(p|data) \propto p^r (1-p)^x (p+0.5)$ 

#### The "levels" and the "slogan"

 The "levels" and the "slogan" are based on the idea that

$$f(x,\theta) = f(x|\theta)f(\theta)$$
  
and so  
$$f(\theta|x) = \frac{f(x,\theta)}{f(x)}$$
$$\propto f(x|\theta)f(\theta)$$
$$= (likelihood) \times (prior)$$
$$= (level 1) \times (level 2)$$

11/3/2016

Extending the "levels" and the "slogan"

Now suppose we have two parameters:

$$f(x, \theta_1, \theta_2) = f(x|\theta_1, \theta_2)f(\theta_1, \theta_2)$$
  
=  $f(x|\theta_1, \theta_2)f(\theta_1|\theta_2)f(\theta_2)$   
This means we can write  
 $f(\theta_1, \theta_2|x) = \frac{f(x, \theta_1, \theta_2)}{f(x)}$   
 $\propto f(x|\theta_1, \theta_2)f(\theta_1|\theta_2)f(\theta_2)$ 

- = (likelihood)  $\times$  (prior)  $\times$  ("hyper-prior")
- = (level 1) × (level 2) × (level 3)
- This idea can be extended with more parameters and more levels...

## Mastery Learning: Distribution of Masteries

- Last time, we considered one student at a time according to
  - Level 1:  $x \sim NB(x|r,p)$
  - Level 2: p  $\sim$  Beta(p| $\alpha$ , $\beta$ )
- Now, suppose we have a sample of n students, and the number of failures before the r<sup>th</sup> success for each student is x<sub>1</sub>, x<sub>2</sub>, ..., x<sub>n</sub>.
  - We want to know the distribution of p, the probability of success, in the population: <u>i.e. we want to estimate  $\alpha \& \beta !</u>$ </u>
- We will model this as
  - Level 1:  $x_i \sim NB(x | r, p_i)$ , i=1, ..., n
  - $\hfill\square$  Level 2:  $\mathbf{p_i}\sim \mathrm{Beta}(\mathbf{p}\,|\,\alpha,\beta)$  , i = 1, ..., n
  - □ Level 3:  $\alpha$  ~ Gamma( $\alpha$ |a<sub>1</sub>,b<sub>1</sub>),  $\beta$  ~ Gamma( $\beta$ |a<sub>2</sub>,b<sub>2</sub>)

11/3/2016

## Mastery Learning: Distribution of Masteries

 We want to know the distribution of p, the probability of success, in the population: <u>i.e. we</u> <u>want to estimate α & β!</u>

• Level 1:  $x_i \sim NB(x|r,p_i)$ , i=1, ..., n n+2 parameters!

- Level 2:  $p_i \sim \text{Beta}(p(\alpha,\beta), i = 1, ..., n$
- □ Level 3:  $\alpha \sim \text{Gamma}(\alpha | a_1, b_1), \beta \sim \text{Gamma}(\alpha | a_2, b_2)$



### Mastery Learning: Distribution of Masteries

### • <u>Level 1</u>: If $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ are the failure counts then $f(x_1, \dots, x_n | p_1, \dots, p_n, r) = \prod_{i=1}^n {\binom{r+x_i-1}{x_i}} p_i^r (1-p_i)^{x_i}$ • <u>Level 2</u>: If $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_n$ are the success probabilties, $f(p_1, \dots, p_n | \alpha, \beta) = \prod_{i=1}^n \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} p_i^{\alpha - 1} (1-p_i)^{\beta - 1}$ • <u>Level 3</u>: $f(\alpha | a_1, b_1) = \frac{b_1^{a_1}}{\Gamma(a_1)} \alpha^{a_1 - 1} e^{-\alpha b_1}$ $f(\beta | a_2, b_2) = \frac{b_2^{a_2}}{\Gamma(a_2)} \beta^{a_2 - 1} e^{-\beta b_2}$

11/3/2016

#### Distribution of Mastery probabilities...

#### Applying the "slogan" for 3 levels:



Distribution of Mastery probabilities...

If we take a<sub>1</sub>=1, b<sub>1</sub>=1, a<sub>2</sub>=1, b<sub>2</sub>=1, and drop the constants we can drop, we get

 $f(p_1,\ldots,p_n,lpha,eta|\mathsf{data})$ 

$$\propto \prod_{i=1}^{n} p_i^r (1-p_i)^{x_i} \times \prod_{i=1}^{n} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} p_i^{\alpha-1} (1-p_i)^{\beta-1} \times e^{-\alpha} e^{-\beta}$$

31

Suppose our data consist of the n=5 values

 $x_1 = 4$ ,  $x_2 = 10$ ,  $x_3 = 5$ ,  $x_4 = 7$ ,  $x_5 = 3$ (number of failures before r=3 successes)

<u>Problem</u>: We need a way to sample from the 5+2=7 parameters p<sub>1</sub>, ..., p<sub>5</sub>, α, β!

11/3/2016

### <u>Solution</u>: Markov-Chain Monte Carlo (MCMC)

- MCMC is very useful for multivariate distributions, e.g. f(θ<sub>1</sub>,θ<sub>2</sub>, ...,θ<sub>K</sub>)
- Instead of dreaming up a way to make a draw (simulation) of all K variables at once MCMC takes draws one at a time
- We "pay" for this by not getting independent draws. The draws are the states of a <u>Markov Chain.</u>
- The draws will not be "exactly right" right away; the Markov chain has to "burn in" to a stationary distribution; the draws after the "burn-in" are what we want!

#### Summary

- Practical Bayes
- Mastery Learning Example
- A brief taste of JAGS and RUBE
- Hierarchical Form of Bayesian Models
- Extending the "Levels" and the "Slogan"
- Mastery Learning Distribution of "mastery"
- (to be continued...)

11/3/2016