

Introduction to Bayesian Psychometric Modeling

PSQF 7375; Spring 2019 Lecture 1 18 January, 2019

PSQF 7375.6: Introduction to Bayesian Psychometric Modeling

Today's Class



An introduction to the course

- > Bayesian Overview
- > MCMC Overview
- > Psychometric Model Overview



AN INTRODUCTION TO BAYESIAN STATISTICS

Bayesian Statistics: The Basics



- Bayesian statistical analysis refers to the use of models where some or all of the parameters are treated as random components
 - > Each parameter comes from some type of distribution
- The likelihood function of the data is then augmented with an additional term that represents the likelihood of the prior distribution for each parameter
 - Think of this as saying each parameter has a certain likelihood the height of the prior distribution
- The final estimates are then considered summaries of the posterior distribution of the parameter, conditional on the data
 - We use these estimates to make inferences, just like we do when using the non-Bayesian approaches (e.g., maximum likelihood/least squares)

Bayesian Statistics: Why It Is Used

- Bayesian methods get used because the <u>relative</u> accessibility of one method of estimation (MCMC – to be discussed shortly)
- There are at least four main reasons why people use MCMC:
- 1. Missing data
 - > Multiple imputation: MCMC is used to estimate model parameters then "impute" data
 - More complicated models for certain types of missing data
- 2. Lack of software capable of handling large sized analyses
 - Have a zero-inflated negative binomial with 21 multivariate outcomes per 18 time points?
- 3. New models/generalizations of models not available in software
 - > Have a new model?
 - Need a certain link function not in software?
- 4. Membership in the "cult" of Bayes
 - Cult members believe philosophical differences exist between numbers from Bayesian analysis and other types of estimators

Bayesian Statistics: Perceptions and Issues

- The use of Bayesian statistics has been controversial
 - The use of certain prior distributions can produce results that are biased or reflect subjective judgment rather than objective science
- Most MCMC estimation methods are computationally intensive with long compute times relative to other estimators
 - Until very recently, very few methods available for those who aren't into programming in Fortran, C, or C++
- Understanding of what Bayesian methods are and how they work is limited outside the field of mathematical statistics

Especially the case in the educational and social sciences

 Over the past 20 years, Bayesian methods have become widespread – making new models estimable and becoming standard in some social science fields (quantitative psychology and educational measurement)

As an Example...

THE UNIVERSITY

Statistical Modelling

These new articles for Statistical Modelling are available online

View online

OnlineFirst Alert

Article

Bayesian residual analysis for spatially correlated data

Viviana GR Lobo, Thaís CO Fonseca Statistical Modelling Jan 13, 2019 | OnlineFirst

A Bayesian approach of analysing semi-continuous longitudinal data with monotone missingness

Jayabrata Biswas, Kiranmoy Das Statistical Modelling Jan 10, 2019 | OnlineFirst

Bayesian semiparametric latent variable model with DP prior for joint analysis: Implementation with nimble

Zhihua Ma, Guanghui Chen Statistical Modelling Jan 8, 2019 | OnlineFirst

Exploring and modelling team performances of the Kaggle European Soccer database

Maurizio Carpita, Enrico Ciavolino, Paola Pasca Statistical Modelling Jan 10, 2019 | OnlineFirst



A PRIMER ON HOW BAYESIAN METHODS WORK

How Bayesian Statistics Work

- The term Bayesian refers to Thomas Bayes (1701-1761)
 Formulated Bayes' Theorem
- Bayesian methods rely on Bayes' Theorem: $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$
 - > P(A) is the prior distribution (pdf) of A \rightarrow WHY THINGS ARE BAYESIAN
 - > P(B) is the marginal distribution (pdf) of B
 - > P(B|A) is the conditional distribution (pdf) of B, given A
 - P(A|B) is the posterior distribution (pdf) of A, given B
- Bayes' Theorem Example...

Imagine a patient takes a test for a rare disease (present 1% of the population) that has a 95% accuracy rate...what is the probability the patient actually has the disease?

Bayes' Theorem Example

THE UNIVERSITY OF IOWA

Imagine a patient takes a test for a rare disease (present 1% of the population) that has a 95% accuracy rate...what is the probability the patient actually has the disease?

- D = the case where the person actually has the disease
- ND = the case where the person does not have the disease
- + = the test for the disease is positive

The question is asking for: P(D|+) From Bayes' Theorem:

$$P(D|+) = \frac{P(+|D)P(D)}{P(+)}$$

What we know:

$$P(D) = .01$$

 $P(+|D) = .95$

Back to Distributions



- We don't know P(+) directly from the problem, but we can figure it out if we recall how distributions work:
- P(+) is a marginal distribution
- P(+|D) is a conditional distribution
- We can get to the marginal by summing across the conditional: P(+) = P(+|D)P(D) + P(+|ND)P(ND) = .95 * .01 + .05 * .99 = .059
- So, to figure out the answer, if a person tests positive for the disease, the posterior probability they actually have the disease is:

$$P(D|+) = \frac{P(+|D)P(D)}{P(+)} = \frac{.01 * .99}{.059} = .17$$

A (Perhaps) More Relevant Example

- The old-fashioned Bayes' Theorem example I've found to be difficult to generalize to your actual data, so...
- Imagine you administer an IQ test to a sample of 50 people
 y_p = person p's IQ test score
- To put this into a linear-models context, the empty model for Y:

$$y_p = \beta_0 + e_p \label{eq:yp}$$
 where $e_p \sim N(0, \sigma_e^2)$

- From this empty model, we know that:
 - > β_0 is the mean of the Y (the mean IQ)
 - > σ_e^2 is the sample variance of Y
 - The conditional distribution of Y is:

$$f(y_p | \beta_0, \sigma_e^2) \sim N(\beta_0, \sigma_e^2)$$

Non-Bayesian Analysis

- Often, least squares (which is equivalent to REML) is used to estimate this model
 - > We could also use ML...
- For ML, we maximized the joint likelihood of the sample with respect to the two unknown parameters β_0 and σ_e^2

$$L(\beta_0, \sigma_e^2) = \prod_{p=1}^N f(y_p | \beta_0, \sigma_e^2) = \prod_{p=1}^N \frac{1}{\sqrt{2\pi\sigma_e^2}} \exp\left(-\frac{(y_p - \beta_0)^2}{2\sigma_e^2}\right)$$

• Here, using gls(), I found:

$$\beta_0 = 102.769$$

 $\sigma_e^2 = 239.490$

• Also, I found:

$$LogL = -207.91$$

Setting up a Bayesian Model



- The (fully) Bayesian model would treat each parameter as a random instance from some prior distribution
- Let's say you know that this version of the IQ test is supposed to have a mean of 100 and a standard deviation of 15
 So β₀ should be 100 and σ_e² should be 225
- Going a step further, let's say you have seen results for administrations of this test that led you to believe that the mean came from a normal distribution with a SD of 2.13
 - > This indicates the prior distribution for the **mean**...or

 $f(\beta_0) \sim N(100, 2.13^2)$

 Let's also say that you don't really have an idea as for the distribution of the variance, but you have seen it range from 200 to 400, so we can come up with a prior distribution for the variance of:

 $f(\sigma_e^2) \sim U(200,\!400)$

 Here the prior is a uniform distribution meaning all values from 200 to 400 are equally likely

More on the Bayesian Approach

- The Bayesian approach is now to seek to find the **posterior** distribution of the parameters given the data:
- $f(\beta_0, \sigma_e^2 | \mathbf{y}_p)$ $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$ We can again use Bayes' Theorem (but for continuous parameters): $f(\alpha_1 | \alpha_2 \alpha_1^2) f(\alpha_2 \alpha_2^2) = f(\alpha_1 | \alpha_2 \alpha_1^2) f(\alpha_1 | \alpha_2 \alpha_1^2) f(\alpha_1 | \alpha_2 \alpha_1^2) f(\alpha_1 | \alpha_1 \alpha_1^2) f(\alpha_1 | \alpha_1$

$$f(\beta_0, \sigma_e^2 | \mathbf{y}_p) = \frac{f(\mathbf{y}_p | \beta_0, \sigma_e^2) f(\beta_0, \sigma_e^2)}{f(\mathbf{y}_p)} = \frac{f(\mathbf{y}_p | \beta_0, \sigma_e^2) f(\beta_0) f(\sigma_e^2)}{f(\mathbf{y}_p)}$$

- Because $f(\mathbf{y}_p)$ essentially is a constant (which involves integrating across β_0 and σ_e^2 to find its value), this term is often referred to as: $f(\beta_0, \sigma_e^2 | \mathbf{y}_p) \propto f(\mathbf{y}_p | \beta_0, \sigma_e^2) f(\beta_0) f(\sigma_e^2)$
- The symbol ∝ is read as "is proportional to" meaning it is the same as when multiplied by a constant
 - $\succ\,$ So it is the same for all values of β_0 and σ_e^2

Unpacking the Posterior Distribution

• $f(\mathbf{y}_p | \beta_0, \sigma_e^2)$ is the **conditional distribution** of the data given the parameters – which comes from our linear model distribution

$$f(\mathbf{y}_p|\beta_0, \sigma_e^2) = \prod_{p=1}^N f(y_p|\beta_0, \sigma_e^2) = \prod_{p=1}^N \frac{1}{\sqrt{2\pi\sigma_e^2}} \exp\left(-\frac{(y_p - \beta_0)^2}{2\sigma_e^2}\right)$$

• $f(\beta_0)$ is the **prior distribution** of β_0 , which we decided would be $N(100,2.13^2)$, giving the height of any β_0 :

$$f(\beta_0) = \frac{1}{\sqrt{2\pi\sigma_{\beta_0}^2}} \exp\left(-\frac{\left(\beta_0 - \mu_{\beta_0}\right)^2}{2\sigma_{\beta_0}^2}\right)$$
$$= \frac{1}{\sqrt{2\pi * 2.13^2}} \exp\left(-\frac{(\beta_0 - 100)^2}{2 * 2.13^2}\right)$$

Unpacking the Posterior Distribution

• $f(\sigma_e^2)$ is the **prior distribution** of σ_e^2 , which we <u>decided</u> would be U(200,400), giving the density function result of any value of σ_e^2 as:

$$f(\sigma_e^2) = \frac{1}{b_{\sigma_e^2} - a_{\sigma_e^2}} = \frac{1}{400 - 200} = \frac{1}{200} = .005$$

- Some useful terminology:
 - > The parameters of the model (for the data) get prior distributions
 - The prior distributions each have parameters these parameters are called hyper-parameters
 - The hyper-parameters are not estimated in our example, but could be giving us a case where we would call our priors empirical priors
 - AKA random intercept variance

Up Next: Estimation



- Although MCMC is commonly thought of as the only method for Bayesian estimation, there are several other forms
- The form analogous to ML (where the value of the parameters that maximize the likelihood or log-likelihood) is called Maximum (or Modal) a Posteriori estimation (MAP)
 - The term modal comes from the maximum point coming at the peak (the mode) of the posterior distribution
- In practice, this functions similar to ML, instead of maximizing the joint likelihood of the data, we now have to worry about the prior:

 $f(\beta_0, \sigma_e^2 | \mathbf{y}_p) = \frac{f(\mathbf{y}_p | \beta_0, \sigma_e^2) f(\beta_0) f(\sigma_e^2)}{f(\mathbf{y}_p)} \propto f(\mathbf{y}_p | \beta_0, \sigma_e^2) f(\beta_0) f(\sigma_e^2)$

• Because it is often more easy to work with, the log is often used: $\log \left(f(\beta_0, \sigma_e^2 | \boldsymbol{y}_p) \right) \propto \log f(\boldsymbol{y}_p | \beta_0, \sigma_e^2) + \log f(\beta_0) + \log f(\sigma_e^2)$

Grid Searching for the MAP Estimate of β_0 The University

- To demonstrate, let's imagine we know $\sigma_e^2 = 239.490$
 - > Later we won't know this...when we use MCMC
- We will use Excel to search over a grid of possible values for β_0
- In each, we will use $\log f(\mathbf{y}_p | \beta_0) + \log f(\beta_0)$
- As a comparison, we will also search over the ML log likelihood function $\log f(\pmb{y}_p | \pmb{\beta}_0)$

ML v. Prior for β_0 of N(100, 2.13²)

- Maximum for ML: 102.8
- Maximum for Bayes: 101.4 (estimate is closer to mean of prior)



ML vs. Prior for β_0 of N(100, 1000²)

- Maximum for ML: 102.8
- Maximum for Bayes: 102.8



ML vs. Prior for β_0 of N(100, 0.15²)

- Maximum for ML: 102.8
- Maximum for Bayes: 100



THE UNIVERSITY

OF IOWA

ML vs. Prior for β_0 of U(-1000,1000)

- Maximum for ML: 102.8
- Maximum for Bayes: 102.8



THE UNIVERSITY

OF OWA

Summarizing Bayesian So Far

- THE UNIVERSITY OF IOWA
- Bayesian \rightarrow parameters have prior distributions
- Estimation in Bayesian → MAP estimation is much like estimation in ML, only instead of likelihood of data, now have to add in likelihood for prior of all parameters
 - But...MAP estimation may be difficult as figuring out derivatives for gradient function (for Newton-Raphson) are not always easy
 - > Where they are easy: **Conjugate** priors \rightarrow prior distributions that are the same as the posterior distribution (think multilevel with normal outcomes)
- Priors can be informative (highly peaked) or uninformative (not peaked)
 - Some uninformative priors will give MAP estimates that are equal to ML
- Up next: estimation by brute force: Markov Chain Monte Carlo



MARKOV CHAIN MONTE CARLO ESTIMATION: THE BASICS

How Estimation Works (More or Less)

- Most estimation routines do one of three things:
- Minimize Something: Typically found with names that have "least" in the title. Forms of least squares include "Generalized", "Ordinary", "Weighted", "Diagonally Weighted", "WLSMV", and "Iteratively Reweighted." Typically the estimator of last resort...
- Maximize Something: Typically found with names that have "maximum" in the title. Forms include "Maximum likelihood", "ML", "Residual Maximum Likelihood" (REML), "Robust ML". Historically the gold standard of estimators.
- 3. Use Simulation to Sample from Something: more recent advances in simulation use resampling techniques. Names include "Bayesian Markov Chain Monte Carlo", "Gibbs Sampling", "Metropolis Hastings", "Metropolis Algorithm", and "Monte Carlo". Used for complex models where ML is not available or for methods where prior values are needed.

How MCMC Estimation Works

 MCMC estimation works by taking samples from the posterior distribution of the data given the parameters:

$$f(\beta_0, \sigma_e^2 | \mathbf{y}_p) = \frac{f(\mathbf{y}_p | \beta_0, \sigma_e^2) f(\beta_0) f(\sigma_e^2)}{f(\mathbf{y}_p)}$$

- > How is that possible? We don't know $f(y_p)$...but...we'll see...
- After enough values are drawn, a rough shape of the distribution can be formed
 - > From that shape we can take summaries and make them our parameters (i.e., mean)
- How the sampling mechanism happens comes from several different algorithms that you will hear about, the most popular being:
 - > **Gibbs Sampling**: used when $f(\beta_0, \sigma_e^2 | \mathbf{y}_p)$ is known
 - Parameter values are drawn and kept throughout the chain
 - > **Metropolis-Hastings (within Gibbs):** used when $f(\beta_0, \sigma_e^2 | \mathbf{y}_p)$ is unknown
 - Parameter values are proposed, then either kept or rejected
 - SAS PROC MCMC uses the latter
 - TRIVIA NOTE: The Metropolis algorithm comes from Chemistry (in 1950)
 - > Hybrid MC: Newer versions (1980s; implemented in Stan)
- In some fields (Physics in particular), MCMC estimation is referred to as Monte Carlo estimation

MCMC Estimation with MHG

- The Metropolis-Hastings algorithm works a bit differently than Gibbs sampling:
- 1. Each parameter (here β_0 and σ_e^2) is given an initial value
- 2. In order, a new value is proposed for each model parameter from some distribution:

$$\beta_0^* \sim Q(\beta_0^*|\beta_0); \sigma_e^{2^*} \sim Q(\sigma_e^{2^*}|\sigma_e^2)$$

3. The proposed value is then accepted as the current value with probability $max(r_{MHG}, 1)$:

 $r_{MHG} = \frac{f(\mathbf{y}_{p} | \beta_{0}^{*}, \sigma_{e}^{2^{*}}) f(\beta_{0}^{*}) f(\sigma_{e}^{2^{*}}) Q(\beta_{0} | \beta_{0}^{*}) Q(\sigma_{e}^{2} | \sigma_{e}^{2^{*}})}{f(\mathbf{y}_{p} | \beta_{0}, \sigma_{e}^{2}) f(\beta_{0}) f(\sigma_{e}^{2}) Q(\beta_{0}^{*} | \beta_{0}) Q(\sigma_{e}^{2^{*}} | \sigma_{e}^{2})}$

 The process continues for a pre-specified number of iterations (more is better)

Notes About MHG

• The constant in the denominator of the posterior distribution: $f(\beta_0, \sigma_e^2 | \mathbf{y}_p) = \frac{f(\mathbf{y}_p | \beta_0, \sigma_e^2) f(\beta_0) f(\sigma_e^2)}{f(\mathbf{y}_p)}$

...cancels when the ratio is formed

- The proposal distributions $Q(\beta_0^*|\beta_0)$ and $Q(\sigma_e^{2^*}|\sigma_e^2)$ can literally be any statistical distribution
 - The trick is picking ones that make the chain "converge" quickly
 - Want to find values that lead to moderate number of accepted parameters
 - > SAS PROC MCMC/JAGS don't make you pick these
- Given a long enough chain, the final values of the chain will come from the posterior distribution
 - From that you can get your parameter estimates

Introducing JAGS...

estimation with Bayesian



```
model01Bayes = function(){
   # likelihood
   for (i in 1:n){
     y[i] \sim dnorm(mu, tau)
   }
   #priors
   mu \sim dnorm(100, 1/2.13^2)
   tau \sim dunif(1/400, 1/200)
   sigma2 = 1/tau
 }
 data = list(y = dataIQ$y, n = nrow(dataIQ))
 jags.param = c("mu", "tau", "sigma2")
 fit <- jags.parallel(data=data,</pre>
                         parameters.to.save=jags.param,
                         n.iter=50000, n.chains=2, n.thin=2, n.burnin=40000,
                         model.file=model01Bayes)
```

Examining the Chain and Posteriors

THE UNIVERSITY OF IOWA



Practical Specifics in MCMC Estimation

- A burn-in period is used where a chain is run for a set number of iterations before the sampled parameter values are used in the posterior distribution
- Because of the rejection/acceptance process, any two iterations are likely to have a high correlation (called autocorrelation) → posterior chains use a thinning interval to take every Xth sample to reduce the autocorrelation
 - A high autocorrelation may indicate the standard error of the posterior distribution will be smaller than it should be
- The chain length (and sometimes number of chains) must also be long enough so the rejection/acceptance process can reasonably approximate the posterior distribution
- How does one what values to pick for these? Output diagnostics
 - > Trial. And. Error.

THE UNIVERSITY OF IOWA

Best Output Diagnostics: the Eye Ball Test









Output Statistics and Diagnostics

```
> fit
Inference for Bugs model at "model01Bayes", fit using jags,
2 chains, each with 50000 iterations (first 40000 discarded), n.thin = 2
n.sims = 10000 iterations saved
        mu.vect sd.vect
                           2.5%
                                    25%
                                            50%
                                                    75%
                                                         97.5% Rhat n.eff
        101.312 1.546 98.280 100.279 101.316 102.347 104.349 1.001 10000
mu
        256.627 40.791 202.918 224.724 248.240 280.247 358.019 1.001 10000
sigma2
tau
          0.004 0.001 0.003
                                  0.004 0.004 0.004 0.005 1.001 10000
deviance 417.317 1.404 415.862 416.284 416.869 417.922 421.059 1.001 10000
For each parameter, n.eff is a crude measure of effective sample size,
and Rhat is the potential scale reduction factor (at convergence, Rhat=1).
DIC info (using the rule, pD = var(deviance)/2)
pD = 1.0 and DIC = 418.3
DIC is an estimate of expected predictive error (lower deviance is better).
```

THE UNIVERSITY

OF IOWA

Changing the Prior

- To demonstrate different priors affect the analysis, we will now try a few prior distributions for our parameters
- Prior: $\beta_0 \sim U(-10000, 10000); \sigma_e^2 \sim U(0, 5000)$

> †1t2

```
Inference for Bugs model at "model02Bayes", fit using jags,
 2 chains, each with 50000 iterations (first 40000 discarded), n.thin = 2
n.sims = 10000 iterations saved
                                   25%
                                          50% 75% 97.5% Rhat n.eff
        mu.vect sd.vect 2.5%
     102.750 2.229 98.385 101.239 102.758 104.251 107.100 1.001 10000
mu
sigma2 244.899 51.326 164.789 208.259 238.353 273.458 362.843 1.001 10000
    0.004 0.001 0.003 0.004 0.004 0.005 0.006 1.001 10000
tau
deviance 417.856 2.028 415.869 416.415 417.241 418.624 423.285 1.001 3200
For each parameter, n.eff is a crude measure of effective sample size,
and Rhat is the potential scale reduction factor (at convergence, Rhat=1).
DIC info (using the rule, pD = var(deviance)/2)
pD = 2.1 and DIC = 419.9
DIC is an estimate of expected predictive error (lower deviance is better).
```

Chain Plots




Changing the Prior

THE UNIVERSITY

- Prior: $\beta_0 \sim N(0, 100, 000)$;
- $\sigma_e^{-2} \sim gamma(r = .01, \lambda = .01)$

> fit3

Inference for Bugs model at "model03Bayes", fit using jags, 2 chains, each with 50000 iterations (first 40000 discarded), n.thin = 2 n.sims = 10000 iterations saved mu.vect sd.vect 2.5% 25% 50% 75% 97.5% Rhat n.eff 102.784 2.224 98.426 101.274 102.758 104.254 107.175 1.001 7300 mu sigma2 253.996 52.970 171.522 216.053 247.279 284.606 375.192 1.001 9500 0.006 1.001 0.004 0.001 0.003 0.004 0.004 0.005 9500 tau deviance 417.812 1.994 415.872 416.409 417.203 418.572 423.105 1.003 1700

For each parameter, n.eff is a crude measure of effective sample size, and Rhat is the potential scale reduction factor (at convergence, Rhat=1).

```
DIC info (using the rule, pD = var(deviance)/2)
pD = 2.0 and DIC = 419.8
DIC is an estimate of expected predictive error (lower deviance is better).
```

Chain Plots





What About an Informative Prior?

• Prior: $\beta_0 \sim U(102,103); \sigma_e^2 \sim U(238,242)$

> fit4 Inference for Bugs model at "model04Bayes", fit using jags, 2 chains, each with 50000 iterations (first 40000 discarded), n.thin = 2 n.sims = 10000 iterations saved mu.vect sd.vect 2.5% 25% 50% 75% 97.5% Rhat n.eff 102.500 0.289 102.026 102.250 102.502 102.752 102.975 1.001 8200 mu sigma2 239.992 1.155 238.104 238.979 240.011 240.993 241.890 1.001 10000 0.004 0.000 0.004 0.004 0.004 0.004 0.004 1.001 10000 tau deviance 415.853 0.036 415.820 415.823 415.835 415.876 415.935 1.000 1 For each parameter, n.eff is a crude measure of effective sample size, and Rhat is the potential scale reduction factor (at convergence, Rhat=1). DIC info (using the rule, pD = var(deviance)/2) pD = 0.0 and DIC = 415.9DIC is an estimate of expected predictive error (lower deviance is better).

THE UNIVERSITY

OF OWA

Chain Plots





MCMC in R



- R itself does not have an MCMC engine native to the language – but there are many free versions available outside of R
- For instance, if you wanted to estimate a path model with MCMC you can:
 - Install the blavaan package (Bayesian lavaan)
 - > Run the path analysis with MCMC



PSYCHOMETRIC MODELS OVERVIEW

Measurement Models

- Measurement models can be divided into two families of models based on response format alone:
 - Continuous responses Confirmatory Factor Models
 - > Categorical responses Item Response Models
- Both of these families fall under a larger framework:
 Generalized Linear Latent and Mixed Models
 - Provide measurement models for other types of responses
- Other relevant families:
 - Structural Equation Models provides estimates of correlations amongst latent variables in measurement models
 - Path Analysis simultaneous regression amongst multiple observed variables

Confirmatory Factor Analysis Models

- Main idea of CFA: Build a <u>measurement model</u> for response variables that measure the same trait
 - CFA = Linear regression model predicting each continuous observed outcome variable (item, subscale) from a latent trait predictor variable(s)

$$Y_{si} = \mu_i + \lambda_{i1} F_{s1} + e_{si}$$

- i item; s subject; μ_i is the item intercept; λ_{i1} is the item slope (factor loading for factor 1); e_{is} is the error for the item and subject; Y_{is} is the item response (*assumed continuous*)
- Differs from exploratory factor analysis:
 - Number and content of factors is decided a priori
 - > Alternative models are comparable and testable
- Uses of confirmatory factor analysis models:
 - > Analyze relationships among subscales that have normal, continuous distributions
 - > Provide comparability across persons, items, and occasions

CFA Model WITH Factor Means and Item Intercepts





Factor Analysis (Y Observed; F latent)

• The prediction of Y is done using a linear regression:

$$Y_{si} = \mu_i + \lambda_{i1}F_{s1} + e_{si}$$



Item Response Models

 Main idea of IRT: Build a <u>measurement model</u> for response variables that measure the same trait...but the responses are discrete

THE UNIVERSITY

Binary Item IRT = Logistic regression model predicting each discrete observed outcome variable (item, subscale) from a latent trait predictor variable(s)

$$logit(Y_{si} = 1) = \mu_i + \lambda_{i1}F_{s1}$$
$$= a_i(\theta_s - b_i)$$

- i item; s subject; μ_i is the item intercept; λ_{i1} is the item slope (factor loading for factor 1); e_{is} is the error for the item and subject; Y_{is} is the item response (*assumed continuous*)
- Differs from confirmatory factor analysis:
 - Items are discrete (although CFA has item factor analysis)
 - Traits are normal

Latent Class Models

- THE UNIVERSITY OF IOWA
- Main idea of LCA models: Build a <u>measurement model</u> for response variables that measure the same trait...but the responses are discrete and people belong to a finite set of groups that were not observed (latent classes)
 - Binary Item LCA = Logistic ANOVA model predicting each discrete observed outcome variable (item, subscale) from a latent trait predictor variable(s)

$$logit(Y_{si} = 1) = \mu_{ic} = \mu_i + \lambda_{ic}d_c$$

- i item; s subject; μ_i is the item intercept; λ_{i1} is the item slope (factor loading for factor 1); e_{is} is the error for the item and subject; Y_{is} is the item response (*assumed continuous*)
- Differs from confirmatory factor analysis and IRT:
 - > Traits are discrete--representing a "nominal" level of measurement
 - Model identification is very different (enables connections to AI)
 - > Often used in exploratory analyses
- Note: many parameterizations exist (very flexible models)

DCMs are Confirmatory LCA Models

- DCMs are confirmatory LCA models
 - > Most defined for a set of A dichotomous attributes (α)
 - Attributes are either possessed ($\alpha = 1$) or not ($\alpha = 0$)
 - DCM attributes can have more than two levels
 - > DCMs are LCA models with 2^A latent classes
 - Each possible combination of attribute possession
 - i.e., a test measuring 3 dichotomous attributes has 8 latent classes
- LCA measurement model parameters
 - Items measure only some attributes (so-called Q-matrix indicator)
 - > Equated for classes with equivalent status of measured attributes

The University

Building the LCDM

• To demonstrate the LCDM, consider the item 2+3-1=?

> Measures addition (attribute 1: α_{r1}) and subtraction (attribute 2: α_{r2})

The University

- Only attributes defined by the Q-matrix are modeled for an item
- The LCDM provides the logit of a correct response as a function of the latent attributes mastered by a respondent:

 $Logit(Y_{ri} = 1 | \alpha_r) = \lambda_{i,0} + \lambda_{i,1,(1)} \alpha_{r1} + \lambda_{i,1,(2)} \alpha_{r2} + \lambda_{i,2,(1,2)} \alpha_{r1} \alpha_{r2}$

Bayesian Inference Networks (BayesNets)

- Introduce a new vocabulary for psychometric things
- Nodes: categorical latent variables
 - > Analogs to latent factors in factor analysis or item response theory
- Nodes can be Parents or Children
 - > Parents: Not predicted by anything (we would call this an Exogenous variable)
 - Children: Predicted by parents (we would call this an Endogenous variable)
- Edges: conditional dependencies between:
 - > Nodes
 - Nodes and items
- Often represented with a Directed Acyclic Graph or DAG

Woefully Short Primer on Bayesian Networks

 BINs describe multivariate data using conditional probabilities



Т

т

0.99

0.01

- In the image, three variables observed:
 - > Did it rain?
 - > Were the sprinklers on?
 - > Was the grass wet?
- The BIN includes the set of parameters leading to the probabilities in the tables

The University

Woefully Short Primer

on Bayesian Networks



Joint distribution of Rain, Sprinkler, and Grass Wet given by:

= P(Grass, Sprinkler, Rain)

P(GrassWet = T | Rain, Sprinkler) P(Sprinkler | Rain) P(Rain)

- Conditional/Marginal distribution of each variable: Bernoulli
- This example has all observed variables, but latent variables can also be defined
 - > Hidden/unobserved nodes

The University

Worlds Colliding.... Psychometric Models are BINs

Here are some BINs that may be more familiar in the social sciences...



Conditional/Marginal distribution of each variable: Normal Nodes: Observed variables (or more specifically, X, Y, and M)



Conditional/Marginal distribution of each variable: Normal Nodes: 5 Observed variables (X1 – X5) 1 unobserved variable (G)

THE UNIVERSITY OF LOWA

More BIN Terminology



Network Learning/Training = Estimation of model parameters

 Often done with Bayesian/MCMC where priors are placed on nearly all parameters

- Estimation typically done using cross-validation
 - Estimation on one/several samples of data
 - Prediction done with left-out samples of data
- From Psychometrics: Model fit...not evaluated in same way
 - BIN model fit based on:
 - Prediction of left-out samples
 - Posterior predictive checks
 - Entropy (for categorical hidden nodes)
 - This is like saying your CFA model fits because your Omega reliability coefficient is high



COMMONALITIES ACROSS MODELS

Where (Modern) Test Scores Come From

- Factor scores (by other names) are used in many domains
 > Item response theory (CFA w/categorical items): GRE scores are factor scores
- A factor score is the estimate of a subject's unobserved latent trait
- Because this latent variable is not measured directly, it acts like it is missing data: you really cannot know with certainty its true value
- It is difficult to pin down what the missing data value (factor score value) should be precisely
 - > Each factor score has a posterior distribution of possible values
 - > Often, the mean of the posterior distribution is the "factor score"
 - In CFA, the mean is the most likely value
 - > Depending on the test, there may be a lot of error (variability) in the distribution
- Therefore, the use of factor scores must reflect that the score is not known and is represented by a distribution

Draw Templin, Draw!





A different version of factor model identification would change the numbers on the Xaxis, but the shapes and order of the distributions would not change

Factor scores provide a weak ordering of people (weak because of error)

How Distributions get Summarized into Scores

THE UNIVERSITY OF IOWA

- There are two ways of providing a score from the factor score posterior distribution:
 - > Expected a posteriori (EAP): the mean of the distribution
 - > Maximum a posteriori (MAP): the most likely score from the distribution
- In CFA factor score distributions are normal (so EAP=MAP)



Additional Information on Factor Scores

For EAP factor scores:

 $\widehat{F}_{p} = E\left(f(F_{p}|\mathbf{Y})\right)$ $F_{p} = \int Var\left(f(F_{p}|\mathbf{Y})\right)$

• For MAP factor scores:

$$F_{p} = \arg \max_{F_{p}} f(F_{p} | \mathbf{Y})$$

$$SE(\hat{F}_{p}) = \left[\frac{\partial^{2}}{\partial F_{p}^{2}} f(F_{p} | \mathbf{Y})\right|_{\hat{F}_{p}}^{-\frac{1}{2}} \text{ (square root of Fisher's information)}$$

- For CFA (Normal Data/Normal Factor) measurement models:
 - \rightarrow MAP = EAP
 - > Variance is identical across all people, regardless of score
- For non-CFA measurement models:
 - > MAP \neq EAP (but does with infinite items)
 - Standard error is a function of the factor score

Scores Are Empirical Bayes Estimates

- For most (if not all) latent variable techniques, the factor scores come from Empirical Bayes estimation—meaning there is a prior distribution present
 - Empirical = some or all of the parameters of the distribution of the latent variable are estimated (i.e., factor mean and variance)
 - Bayes = comes from the use of Bayes' Theorem
- Prior == Assumed factor distribution with mean/variance
- This is true for all CFA, IRT, mixed/multilevel/hierarchical models
 - > And is true for models that don't have a label (e.g., Poisson Factor Analysis?)

Bayes' Theorem

 Bayes' Theorem states the conditional distribution of a variable A (soon to be our factor score) given values of a variable B (soon to be our data) is:

For Categorical A, replace integral with sum

$$f(A|B) = \frac{f(B|A)f(A)}{f(B)} = \frac{f(B|A)f(A)}{\int_{a \in A} f(B|A = a)f(A = a)da}$$

- f(A|B) is the **distribution** of A, conditional on B
 - > We will come to know this as the posterior distribution of the factor score, conditional on the data observed or $f(\mathbf{F}|\mathbf{Y})$
- f(B|A) is the **distribution** of B, conditional on A
 - > We will come to know this as our measurement model or $f(\mathbf{Y}|\mathbf{F})$
- f(A) is the marginal distribution of A
 - > We will come to know this as the prior distribution of the factor or $f(\mathbf{F})$

Putting Together the Pieces of Empirical Bayes Factor Scores



- For $f(\mathbf{Y}|\mathbf{F})$, consider the measurement model (here CFA) for one item: $Y_{pi} = \mu_i + \lambda_i F_p + e_{pi}$ Where: $e_{pi} \sim N(0, \psi_i^2)$
- Using expected values, we can show the distribution for this one item is: $f(Y_{pi}|F_p) \sim N(\mu_i + \lambda_i F_p, \psi_i^2)$
- Therefore, for all *I* items, our conditional distribution is: $f(\mathbf{Y}|F_p) \sim N_I(\boldsymbol{\mu} + \boldsymbol{\Lambda} F_p, \boldsymbol{\Psi})$
- With multiple factors, this becomes:

 $f(\mathbf{Y}|\mathbf{F}) \sim N_I(\boldsymbol{\mu} + \boldsymbol{\Lambda}\mathbf{F}, \boldsymbol{\Psi})$

Putting Together the Pieces of Empirical Bayes Factor Scores



- For $f(\mathbf{F})$, consider the distribution assumed by the factor:
 - For one factor

$$f(F_p) \sim N(\mu_F, \sigma_F^2)$$

For multiple factors K

$$f(\mathbf{F}) \sim N_K(\mathbf{\mu}_F, \mathbf{\Phi})$$

- We must pick an identification method which determines if certain parameters of $\mu_{\it F}$ and Φ are fixed or are estimated
 - > Any method identification works, so we keep μ_F and Φ throughout

Putting Together the Pieces of Empirical Bayes Factor Scores



$$f(A|B) = \frac{f(B|A)f(A)}{f(B)} = f(\mathbf{F}|\mathbf{Y}) = \frac{f(\mathbf{Y}|\mathbf{F})f(\mathbf{F})}{f(\mathbf{Y})}$$

 For f(Y), we return to the model-implied mean vector and covariance matrix:

 $f(\mathbf{Y}) \sim N_I(\mathbf{\mu} + \mathbf{\Lambda}^T \mathbf{\mu}_F, \mathbf{\Lambda} \mathbf{\Phi} \mathbf{\Lambda}^T + \mathbf{\Psi})$

A Quick Reminder About Types of Distributions

- For two random variables x and z, a conditional distribution is written as: f(z|x)
- The conditional distribution is also equal to the joint distribution divided by the marginal distribution of the conditioning random variable

$$f(z|x) = \frac{f(z,x)}{f(x)}$$

- Therefore, the joint distribution can be found by the product of the conditional and marginal distributions: f(z, x) = f(z|x)f(x)
- We can use this result in our analysis:

A Quick Reminder about Multivariate Normal Distributions



- If X is distributed multivariate normally:
 Conditional distributions of X are multivariate normal
- We can show that f(Y, F), the joint distribution of the data and the factors, is multivariate normal
- We can then use the result above (shown on the next slides) to show that our posterior distribution of the factor scores is also multivariate normal
 - This result <u>only</u> applies for measurement models assuming normally distributed data and normally distributed factors: CFA
 - For IRT (and other measurement models), this result will not hold—but this distribution is asymptotically normal as the number of items gets large

Conditional Distributions of MVN Variables are Multivariate Normal

- The conditional distribution of sets of variables from a MVN is also MVN
- If we were interested in the distribution of the first *q* variables, we partition three matrices:

> The data: $[\mathbf{X}_{1:(N \times q)} \mid \mathbf{X}_{2:(N \times p-q)}]$

> The mean vector:
$$\begin{bmatrix} \boldsymbol{\mu}_{1:(q \ x \ 1)} \\ \boldsymbol{\mu}_{2:(p-q \ x \ 1)} \end{bmatrix}$$

> The covariance matrix:
$$\begin{bmatrix} \Sigma_{11:(q \times q)} & \Sigma_{12:(q \times p-q)} \\ \overline{\Sigma}_{21:(p-q \times q)} & \overline{\Sigma}_{22:(p-q \times p-q)} \end{bmatrix}$$

The University

Conditional Distributions of MVN Variables

• The, $f(\mathbf{X}_1|\mathbf{X}_2)$, conditional distribution of \mathbf{X}_1 given the values of $\mathbf{X}_2 = \mathbf{x}_2$ is then: $\mathbf{X}_1|\mathbf{X}_2 \sim N_q(\boldsymbol{\mu}^*, \boldsymbol{\Sigma}^*)$

Where (using our partitioned matrices):

$$\boldsymbol{\mu}^* = \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2^T - \boldsymbol{\mu}_2)$$

And:

$$\boldsymbol{\Sigma}^* = \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21}$$

Derive, Templin, Derive!

- The joint distribution of all *I* items and *K* factor scores is $f(\mathbf{Y}, \mathbf{F}) = f\left(\begin{bmatrix}\mathbf{Y}\\\mathbf{F}\end{bmatrix}\right)$ $= N_{I+K}\left(\begin{bmatrix}\boldsymbol{\mu} + \boldsymbol{\Lambda}^T \boldsymbol{\mu}_F\\ - \overline{\boldsymbol{\mu}}_F\end{bmatrix}, \begin{bmatrix}\boldsymbol{\Lambda} \boldsymbol{\Phi} \boldsymbol{\Lambda}^T + \boldsymbol{\Psi} & \boldsymbol{\Lambda} \boldsymbol{\Phi}\\ - \overline{\boldsymbol{\Phi}} \boldsymbol{\Lambda}^T & - \overline{\boldsymbol{\Phi}} \end{bmatrix}\right)$
- Using the conditional distributions of MVNs result: $f(\mathbf{F}_p | \mathbf{Y}_p)$ is MVN: With mean: $\boldsymbol{\mu}_F + \boldsymbol{\Phi} \boldsymbol{\Lambda}^T (\boldsymbol{\Lambda} \boldsymbol{\Phi} \boldsymbol{\Lambda}^T + \boldsymbol{\Psi})^{-1} (\mathbf{Y}_p^T - \boldsymbol{\mu})$ And Covariance: $\boldsymbol{\Phi} - \boldsymbol{\Phi} \boldsymbol{\Lambda}^T (\boldsymbol{\Lambda} \boldsymbol{\Phi} \boldsymbol{\Lambda}^T + \boldsymbol{\Psi})^{-1} \boldsymbol{\Lambda} \boldsymbol{\Phi}$

#WTFTemplin

- When using measurement models assuming normally distributed data and normally distributed factors (CFA):
 - > The posterior distribution of the factor scores is MVN
 - > Therefore, the most likely factor score (MAP) and the expected factor score (EAP) is given by the mean from the previous slides
 - > The factor score is a function of the model parameter estimates and the data

The University

Factor Scores via Metropolis Hastings

• Place holder for R script...


WRAPPING UP

Wrapping Up



- Today was an introduction to Bayesian statistics
 - Bayes = use of prior distributions on parameters
- We used two methods for estimation:
 - MAP estimation far less common
 - MCMC estimation
 - Commonly, people will say Bayesian and mean MCMC but Bayesian is just the addition of priors. MCMC is one way of estimating Bayesian models!
- MCMC is effective for most Bayesian models:

Model likelihood and prior likelihood are all that are needed

• MCMC is estimation by brute force:

> Can be very slow, computationally intensive, and disk-space intensive