

Introduction to Bayesian Statistics

Shin-ichi Mayekawa

mayekawa@hum.titech.ac.jp

www.ms.hum.titech.ac.jp

Outline

1. Introduction to Bayes Theorem
2. Binomial Model and Poisson Model
3. Normal Model
4. Normal Regression Model
5. Bayesian Numerical Calculations
(Markov Chain Monte Carlo)
6. Application of MCMC
7. Model Selection
8. miscellaneous topics

2

Gibbs Sampler: general form

3 : Gibbs Sampler (p blocks)

- 1. Initialize $X_1^{(0)} = x_1^{(0)}, X_2^{(0)} = x_2^{(0)}, \dots, X_p^{(0)} = x_p^{(0)}$
- 2. For $t = 0, 1, \dots$ repeat the following steps:
 - Draw $X_1^{(t+1)} = x_1^{(t+1)}$ according to $f_{X_1}(X_1^{(t+1)} | X_2^{(t)} = x_2^{(t)}, X_3^{(t)} = x_3^{(t)}, \dots, X_p^{(t)} = x_p^{(t)})$
 - Draw $X_2^{(t+1)} = x_2^{(t+1)}$ according to $f_{X_2}(X_2^{(t+1)} | X_1^{(t+1)} = x_1^{(t+1)}, X_3^{(t)} = x_3^{(t)}, \dots, X_p^{(t)} = x_p^{(t)})$
 - Draw $X_3^{(t+1)} = x_3^{(t+1)}$ according to $f_{X_3}(X_3^{(t+1)} | X_1^{(t+1)} = x_1^{(t+1)}, X_2^{(t+1)} = x_2^{(t+1)}, X_4^{(t)} = x_4^{(t)}, \dots, X_p^{(t)} = x_p^{(t)})$
 - \vdots
 - Draw $X_j^{(t+1)} = x_j^{(t+1)}$ according to $f_{X_j}(X_j^{(t+1)} | X_1^{(t+1)} = x_1^{(t+1)}, \dots, X_{j-1}^{(t+1)} = x_{j-1}^{(t+1)}, X_{j+1}^{(t)} = x_{j+1}^{(t)}, \dots, X_p^{(t)} = x_p^{(t)})$
 - \vdots
 - Draw $X_p^{(t+1)} = x_p^{(t+1)}$ according to $f_{X_p}(X_p^{(t+1)} | X_1^{(t+1)} = x_1^{(t+1)}, X_2^{(t+1)} = x_2^{(t+1)}, \dots, X_{p-1}^{(t+1)} = x_{p-1}^{(t+1)})$

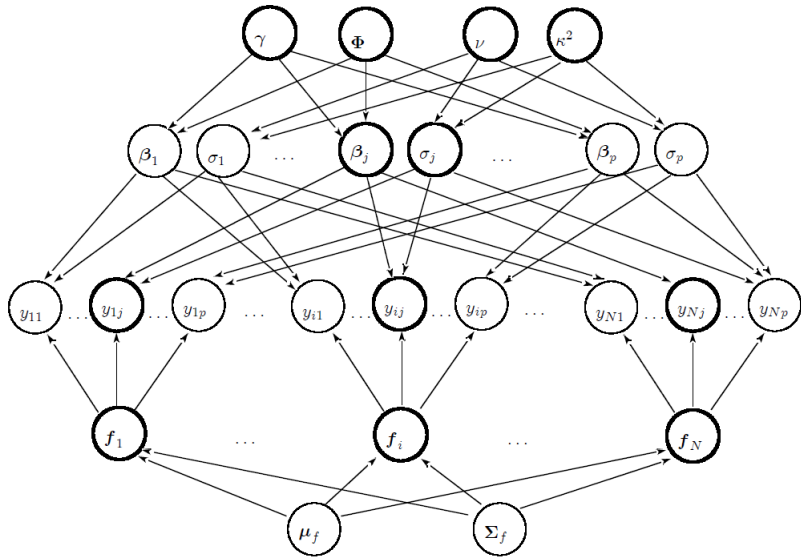
3

How to find the Full Conditional Dist.

- $f(X_j | X_1, X_2, \dots, X_{j-1}, X_{j+1}, X_{j+2}, \dots, X_p)$
- Regard the joint pdf as the pdf of X_j given the rest.
- Use **DAG** : **D**irectional **A**cyclic **G**raph

4

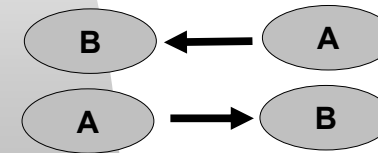
DAG: Directional Acyclic Graph



DAG

- How to draw a DAG.
 - circle \rightarrow random variable
 - arrow \rightarrow dependence
 - Example:
 - B depends on A.
- No arrow indicates (conditional) independence.**

$$\Pr(A,B) = \Pr(B|A)\Pr(A)$$

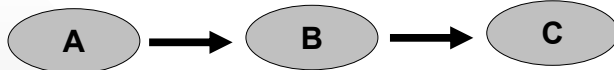


DAG indicates the decomposition of joint probability.

DAG

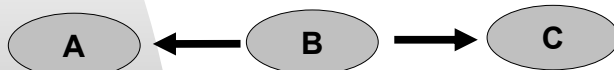
- B depends on A and C depends on B.

$$\Pr(A,B,C) = \Pr(C|B) \Pr(B|A) \Pr(A)$$



- A depends on B and C depends on B.

$$\Pr(A,B,C) = \Pr(C|B) \Pr(A|B) \Pr(B)$$

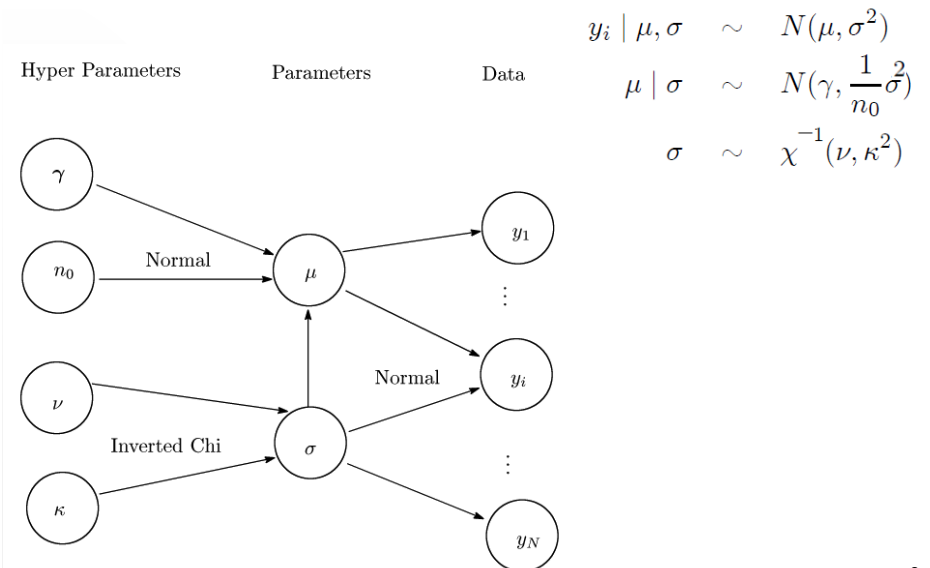


- B depends on both A and C.

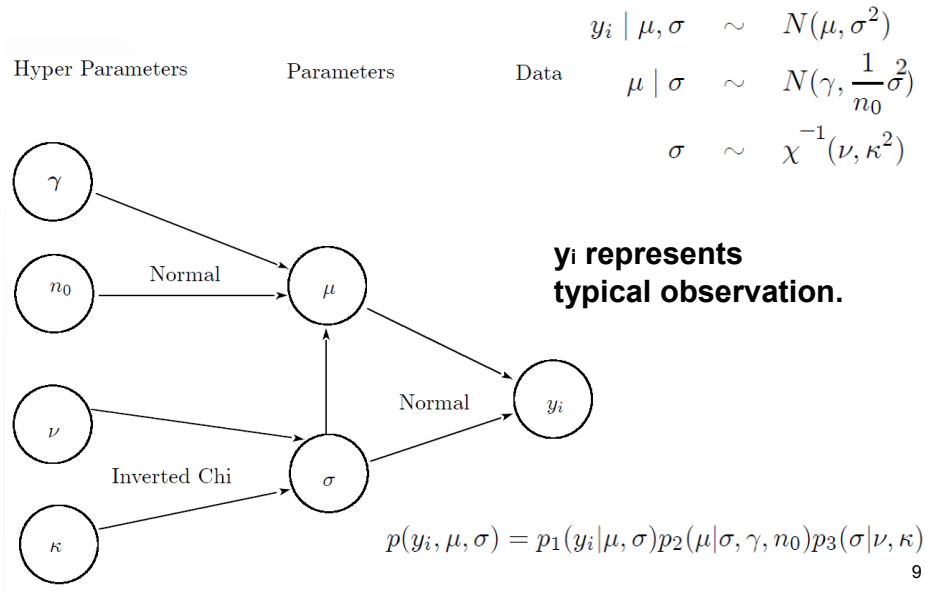
$$\Pr(A,B,C) = \Pr(B|A,C) \Pr(A) \Pr(C)$$



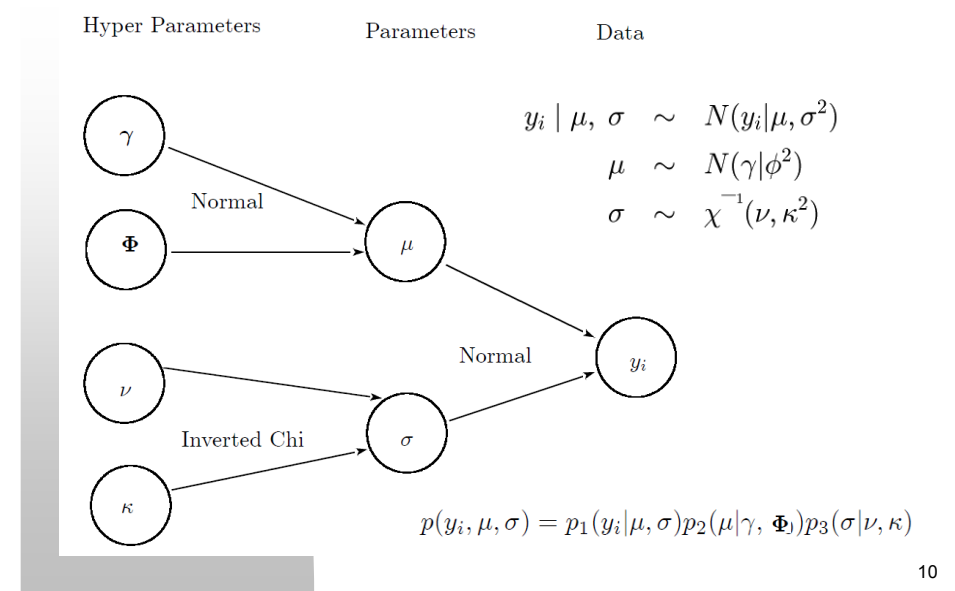
DAG: Normal with Nat. Conj. Prior



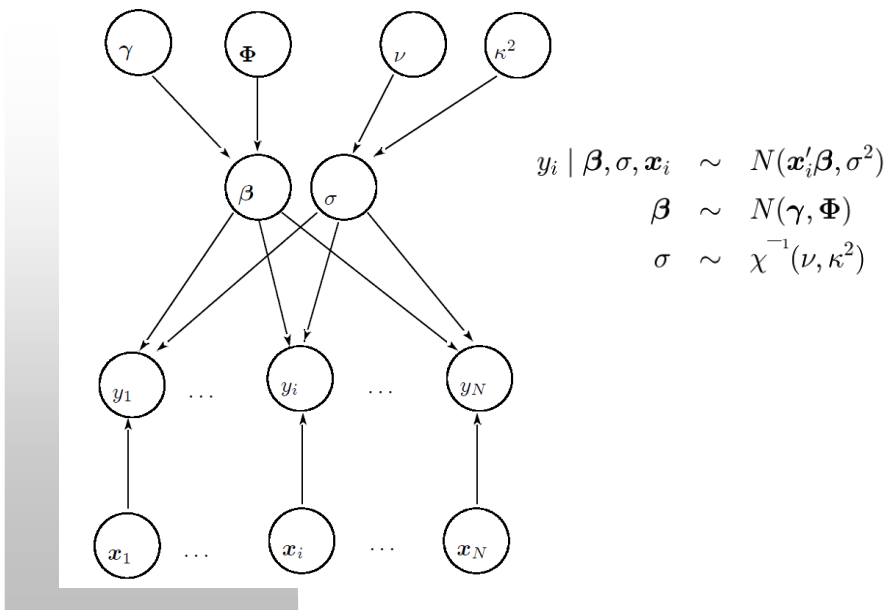
DAG: Normal with Nat. Conj. Prior



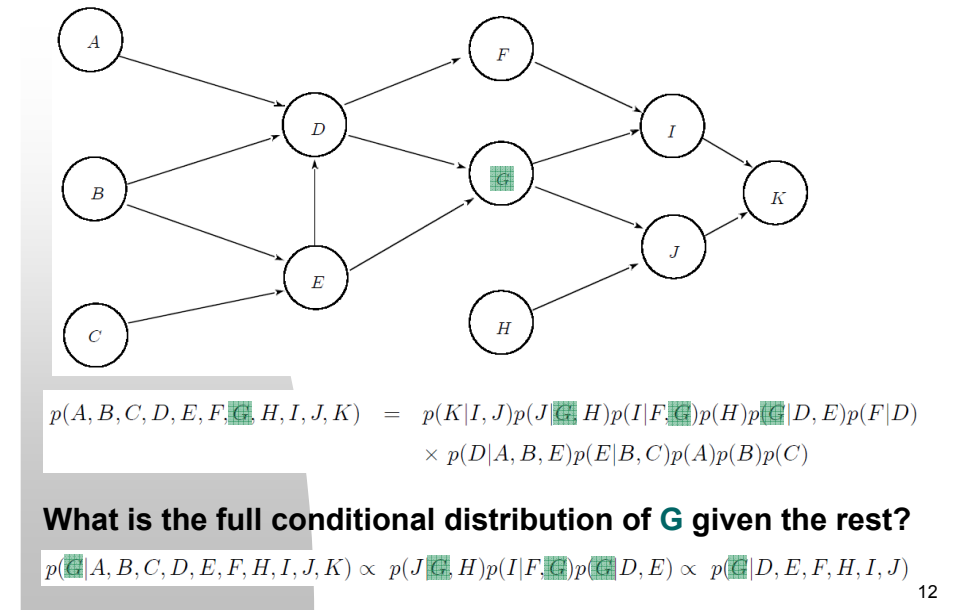
DAG: Normal with Semi Nat. Conj. Prior



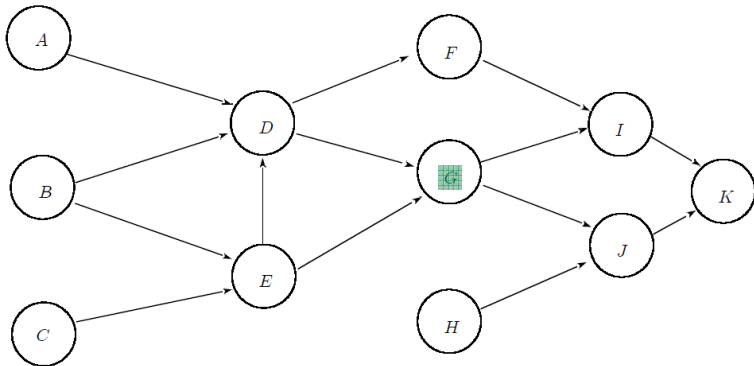
DAG: Regression Model



DAG can be used to find the full conditional distribution of a variable.



DAG can be used to find the full conditional distribution of a variable.



What is the full conditional distribution of **G** given the rest?

$$p(\mathbf{G} | A, B, C, D, E, F, H, I, J, K) = p(\mathbf{G} | D, E, F, H, I, J)$$

Depends on: Parents(D,E), Children(I,J), and Spouses(F,H)
Spouse = Other Parents of the Children

Markov Blanket

13

Markov Blanket

- Given the **Markov Blanket** of a variable, say, **G**, (D,E,F,H,I,J), **G** and the remaining variables, (A,B,C,K), are independent.

$$p(A, B, C, K, \mathbf{G} | D, E, F, H, I, J) =$$

$$p(\mathbf{G} | D, E, F, H, I, J) p(A, B, C, K | D, E, F, H, I, J)$$

- Therefore,

$$p(\mathbf{G} | A, B, C, D, E, F, H, I, J, K)$$

$$= \frac{p(A, B, C, K, \mathbf{G} | D, E, F, H, I, J)}{p(A, B, C, K | D, E, F, H, I, J)}$$

$$= p(\mathbf{G} | D, E, F, H, I, J)$$

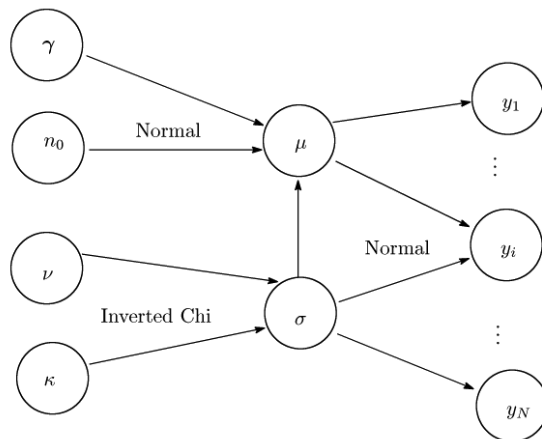
14

DAG: Normal with Nat. Conj. Prior

Hyper Parameters

Parameters

Data



$$y_i | \mu, \sigma \sim N(\mu, \sigma^2)$$

$$\mu | \sigma \sim N\left(\gamma, \frac{1}{n_0} \sigma^2\right)$$

$$\sigma \sim \chi^{-1}(\nu, \kappa^2)$$

$$p(y_i, \mu, \sigma) = p_1(y_i | \mu, \sigma) p_2(\mu | \sigma, \gamma, n_0) p_3(\sigma | \nu, \kappa)$$

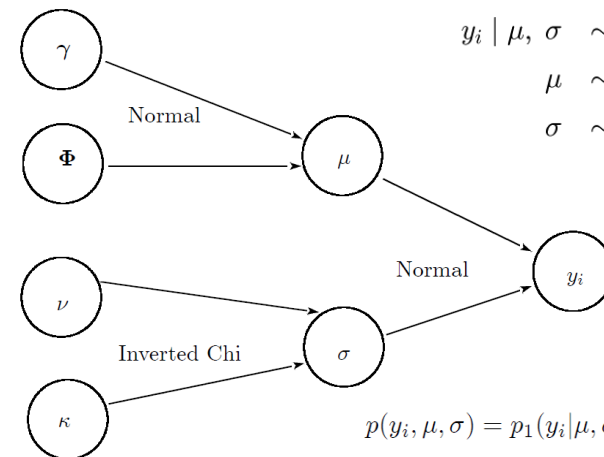
15

DAG: Normal with Semi Nat. Conj. Prior

Hyper Parameters

Parameters

Data



$$y_i | \mu, \sigma \sim N(y_i | \mu, \sigma^2)$$

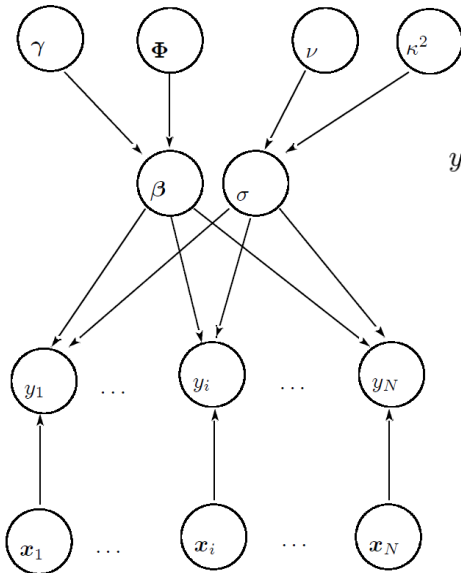
$$\mu \sim N(\gamma | \phi^2)$$

$$\sigma \sim \chi^{-1}(\nu, \kappa^2)$$

$$p(y_i, \mu, \sigma) = p_1(y_i | \mu, \sigma) p_2(\mu | \gamma, \Phi) p_3(\sigma | \nu, \kappa)$$

16

DAG: Regression Model



$$y_i | \beta, \sigma, x_i \sim N(x_i' \beta, \sigma^2)$$

$$\beta \sim N(\gamma, \Phi)$$

$$\sigma \sim \chi^{-1}(\nu, \kappa^2)$$

17

Random Coefficient Model

- Suppose that
 - data (y_k, X_k) , $k=1, 2, \dots, m$, are collected in m classes of the same school.
 - regress y_k on X_k for each class as

$$y_k = X_k \beta_k + E_k, E_k \sim N(0, \sigma_k^2 I_{N_k}), k = 1, 2, \dots, m$$

- Assume that the regression coefficients, β_k , are similar and express this similarity by assuming the same (prior) distributions for the regression coefficients.

$$\beta_k \sim N(\gamma, \Phi), k = 1, 2, \dots, m$$

18

Random Coefficient Model

- m group regression
 - borrowing strength from the neighbors

Hierarchical Linear Model : HLM

Multilevel Model

Slope as Outcome Model

1st stage $y_k = X_k \beta_k + E_k, E_k \sim N(0, \sigma_k^2 I_{N_k}), k = 1, 2, \dots, m$

$$y_k \sim N(X_k \beta_k, \sigma_k^2 I_{N_k}), k = 1, 2, \dots, m$$

2nd stage

$$\beta_k \sim N(\gamma, \Phi), k = 1, 2, \dots, m$$

$$\sigma_k \sim \chi^{-1}(\nu, \kappa^2), k = 1, 2, \dots, m$$

$$h(\beta, \sigma | \mathbf{y}) \propto f(\mathbf{y} | \beta, \sigma) h(\beta | \gamma, \Phi) h(\sigma | \nu, \kappa)$$

$$= \prod_{k=1}^m f(y_k | \beta_k, \sigma_k) h(\beta_k | \gamma, \Phi) h(\sigma_k | \nu, \kappa)$$

3rd stage

$$\gamma \sim N(\zeta, \Psi)$$

$$\Phi \sim W^{-1}(\eta, \mathcal{T})$$

19

Random Coefficient Model

- marginally:

$$y_k \sim N(X_k \gamma, X_k \Phi X_k' + \sigma_k^2 I)$$

$$\gamma \sim N(\zeta, \Psi)$$

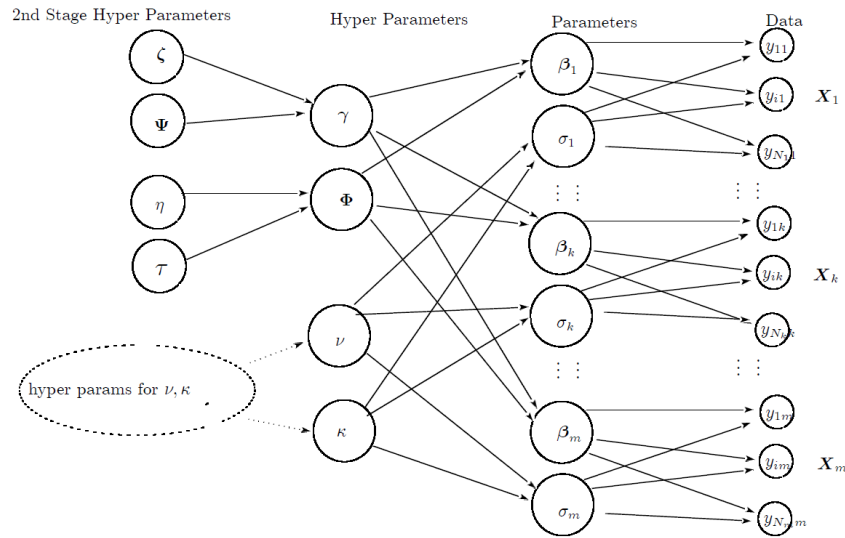
$$\Phi \sim W^{-1}(\eta, \mathcal{T})$$

$$\sigma_k \sim \chi^{-1}(\nu, \kappa^2), k = 1, 2, \dots, m$$

Difficult to generate RV!

20

Random Coefficient Model



Finite Mixture Models

- If we know the membership:

$$d_{iq} = \begin{cases} 1 & \text{if the } i\text{-th obs. belongs class } q. \\ 0 & \text{otherwise.} \end{cases}$$

$$f(y_i | d_{iq} = 1) = f(y_i | \xi_q) \quad \xi_q \sim h(\xi_q | \gamma)$$

- Distribution of the membership

$$h(d_i) \propto \prod_{q=1}^Q \rho_q^{d_{iq}} \quad h(\rho) \propto \prod_{q=1}^Q \rho_q^{\alpha_q - 1}$$

- Since we do not know the membership, marginally:

$$\begin{aligned} f(y_i | \xi) &= \sum_{q=1}^Q f(y_i, d_i | \xi) \\ &= \sum_{q=1}^Q \prod_{i=1}^N (f(y_i | \xi_q) \rho_q)^{d_{iq}} \\ &= \sum_{q=1}^Q f(y_i | \xi_q) \rho_q \end{aligned}$$

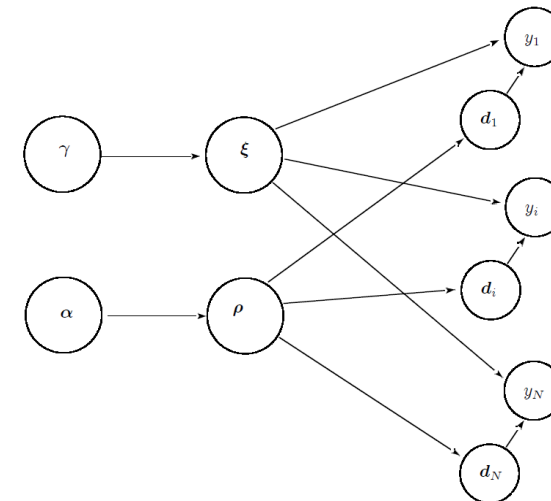
d_{iq} is a latent variable.

Finite Mixture Models

- Univariate Normal Mixtures

$$Y_i | d_{iq} \sim N(\mu_q, \sigma_q^2)$$

Finite Mixture Models



t-distributed errors in regression

Errors are distributed as t distribution.

$$y_i = \sum_{j=1}^q x_{ij}\beta_j + e_i$$

$$= \mathbf{x}'_i\boldsymbol{\beta} + e_i, \quad e_i \sim t(\nu, 0, \kappa^2), i = 1, 2, \dots, N, \text{ iid} \quad \text{var}(e_i) = \frac{\nu}{\nu-2}\kappa^2$$

$$y_i \sim t(\nu, \mathbf{x}'_i\boldsymbol{\beta}, \kappa^2), i = 1, 2, \dots, N, \text{ iid}$$

$$f(\mathbf{y}|\boldsymbol{\beta}, \nu, \kappa) \propto \prod_{i=1}^N \frac{\Gamma(\frac{1}{2}(\nu+1))}{\Gamma(\frac{\nu}{2})} \nu^{-\frac{1}{2}\kappa^{-1}} \left(1 + \frac{1}{\nu\kappa^2}(y_i - \mathbf{x}'_i\boldsymbol{\beta})^2\right)^{-\frac{1}{2}(\nu+1)}$$

$$= \left(\frac{\Gamma(\frac{1}{2}(\nu+1))}{\Gamma(\frac{\nu}{2})} \nu^{-\frac{1}{2}\kappa^{-1}}\right)^N \prod_{i=1}^N \left(1 + \frac{1}{\nu\kappa^2}(y_i - \mathbf{x}'_i\boldsymbol{\beta})^2\right)^{-\frac{1}{2}(\nu+1)}$$

No conjugate priors!

t-distributed errors in regression

Alternative formulation of t-errors

If $y_i|\sigma_i \sim N(\mathbf{x}'_i\boldsymbol{\beta}, \sigma_i^2)$

or $\sigma_i \sim \chi^{-1}(\nu, \nu\kappa^2), \text{ iid}$

or $y_i|\sigma_i, \kappa^2 \sim N(\mathbf{x}'_i\boldsymbol{\beta}, \sigma_i^2\kappa^2)$

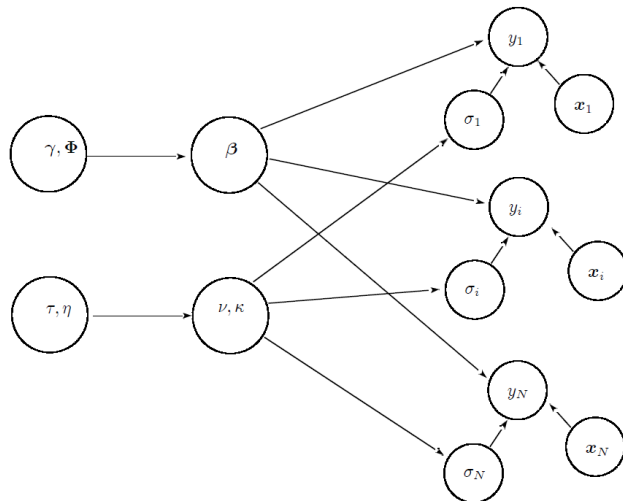
or $\sigma_i \sim \chi^{-1}(\nu, \nu), \text{ iid}$

then, marginally:

$$y_i \sim t(\nu, \mathbf{x}'_i\boldsymbol{\beta}, \kappa^2)$$

σ_i is a **latent variable**.

t-distributed error in regression



Probit Model

- categorical criterion variable: y

$$f(y|\mathbf{p}(\mathbf{x})) = p(\mathbf{x})^y(1 - p(\mathbf{x}))^{1-y}$$

$$\Pr(Y = 1|\mathbf{x}) = p(\mathbf{x})$$

- model for the probability:

$$z = \mathbf{x}'\boldsymbol{\beta} + e, \quad e \sim N(0, \sigma^2) \quad z \sim N(\mathbf{x}'\boldsymbol{\beta}, \sigma^2)$$

$$p(\mathbf{x}) = \Pr(Z > 0) = \Pr\left(\frac{Z - \mathbf{x}'\boldsymbol{\beta}}{\sigma} > -\frac{\mathbf{x}'\boldsymbol{\beta}}{\sigma}\right)$$

$$= 1 - \Phi\left(-\frac{\mathbf{x}'\boldsymbol{\beta}}{\sigma}\right) = \Phi\left(\frac{\mathbf{x}'\boldsymbol{\beta}}{\sigma}\right)$$

Probit Model

likelihood

$$f(\mathbf{y}|\beta, \sigma) = \prod_{i=1}^N p(\mathbf{x}_i)^{y_i} (1 - p(\mathbf{x}_i))^{1-y_i}$$

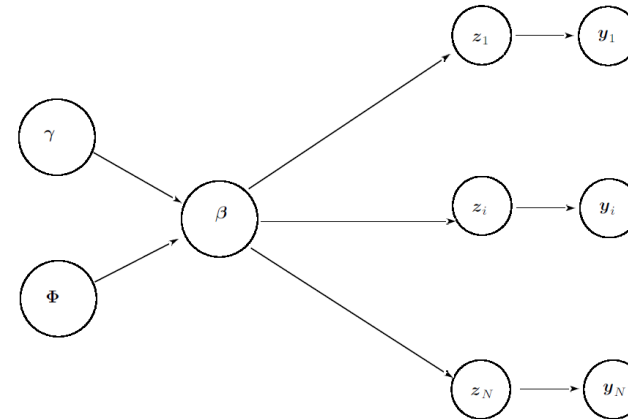
$$= \prod_{i=1}^N \Phi(\mathbf{x}'_i \beta)^{y_i} (1 - \Phi(\mathbf{x}'_i \beta))^{1-y_i}$$

No conjugate priors!

Treat z as the latent variable.

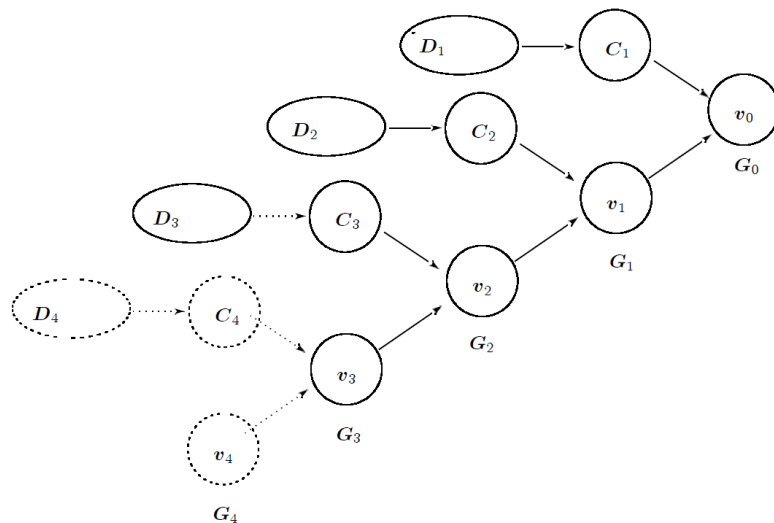
29

Probit Model



30

Hierarchical Model



31

Leftovers

- model comparison
 - ◆ M1: $y = X\beta + e$, $e \sim N(0, \sigma^2)$
 - ◆ M2: $y = X\beta + e$, $e \sim t(0, \nu, c^2)$
 - ◆ posterior odds of models, Bayes Factor, etc
- What is Markov Chain?
 - ◆ Sequence of random numbers, X_1, X_2, \dots , in which $f(X_i | X_1, X_2, \dots, X_{i-1}) = f(X_i | X_{i-1})$.
- What do we do when the full-conditional distribution cannot be identified and therefore we are unable to generate random numbers from it?
 - ◆ Use Metropolis-Hastings method which is another form of MCMC algorithm.

32

Leftovers

- multivariate analysis
 - ◆ multivariate regression, analysis of variance
- latent variable models
 - ◆ factor analysis
 - ◆ item response theory
- comparison with the classical statistics
- how to assess prior distribution
- Bayesian decision theory
- etc etc

33

Exam Topics

- Outline the general procedure for the Bayesian inference using the following keywords.
 - ◆ keywords:
model, parameters, model distribution, subjective probability, prior distribution, likelihood, posterior distribution, posterior mean, posterior standard deviation, natural conjugate prior, semi-natural conjugate prior, Monte Carlo method, Markov Chain Monte Carlo, gibbs sampler, DAG, etc etc...
- What do you think about the Bayesian statistics?
- Due on Aug. 16, in pdf format, by e-mail

34