

## Regression Models

- Bayesian analysis of the “ordinary linear regression” (Chapter 8)
- Example: analysis of radon measurements
- Hierarchical Linear Regression Models (Chapter 13)
- Simple random effect model
- Mixed effect model
- Bayesian Variable Selection

1

2. setting up a prior for  $\theta$  that accurately reflects substantive knowledge

3

## Ordinary Linear Regression

- Question: how does one quantity,  $y$ , vary as a function of another quantity  $x$ ?
- $y = (y_1, \dots, y_n)$  is the continuous outcome
- $x_i = (x_{i1}, \dots, x_{ik})$  is the vector of explanatory variables, discrete or continuous
- $X$  is the  $n \times k$  matrix of the explanatory variables

$$E[y_i | \beta, X] = \beta_1 x_{i1} + \dots + \beta_k x_{ik}$$

$$\text{Var}[y_i | \beta, X] = \sigma^2, \text{ for } i = 1, \dots, n$$

$$\theta = (\beta_1, \dots, \beta_k, \sigma^2)$$

1. define  $x$  and  $y$  (possibly using transformations) so that  $E[y_i | \beta, X]$  is approximately linear, with normal errors

2

## Ordinary Linear Regression, Basic Model

$$y | \beta, \sigma^2, X \sim N(X\beta, \sigma^2 I)$$

$$p(\beta, \sigma^2 | X) \sim \sigma^{-2}$$

conditional distribution of  $\beta$  given  $\sigma$

$$\beta | \sigma^2, y \sim N(\hat{\beta}, V_\beta \sigma^2)$$

$$\hat{\beta} = (X^t X)^{-1} X^t y$$

$$V_\beta = (X^t X)^{-1}$$

conditional distribution of  $\sigma$

$$\sigma^2 | y \sim \chi^{-2}(n - k, s^2)$$

$$s^2 = \frac{1}{n-k} (y - X\hat{\beta})^t (y - X\hat{\beta})$$

$\hat{\beta}$  and  $s^2$  are the MLE.

4

## When the posterior is proper?

For any analysis based on a improper prior distribution, it is important to check that the posterior distribution is proper.

$p(\beta, \sigma^2 | y)$  is a proper as long as

1.  $n > k$
2.  $\text{rank}(X) = k$  (e.g. columns of  $X$  must be linearly independent)

## Sampling from the posterior distribution

- draw  $\sigma^2$  from  $\chi^{-2}(n - k, s^2)$
- draw  $\beta$  conditionally to  $\sigma^2$  from  $N(\hat{\beta}, V_{\beta}\sigma^2)$

5

## Posterior Predictive distribution

- $\tilde{X}$  new data
- $\tilde{y}$  future outcomes
- posterior predictive simulations

To draw a random sample  $\tilde{y}$  from  $p(\tilde{y} | y)$ , we

1. draw  $\sigma^{2(j)}, \beta^{(j)}$  from  $p(\sigma^2, \beta | y)$
2. draw  $\tilde{y} \sim N(\tilde{X}\beta^{(j)}, \sigma^{2(j)}I)$  for  $j = 1, \dots, N$

6

## Analytic form of the posterior predictive distribution

$$\tilde{y} | y \sim t_{n-k}(\tilde{X}\hat{\beta}, s^2(I + \tilde{X}V_{\beta}\tilde{X}^t))$$

$$\begin{aligned} E[\tilde{y} | y] &= E[E(\tilde{y} | \beta, \sigma^2, y) | \sigma^2, y] \\ &= E[\tilde{X}\hat{\beta} | \sigma^2, y] \\ &= \tilde{X}\hat{\beta} \end{aligned}$$

$$\begin{aligned} \text{Var}[\tilde{y} | y, \sigma^2] &= E[\text{Var}(\tilde{y} | \beta, \sigma^2, y) | \sigma^2, y] \\ &+ \text{Var}[E(\tilde{y} | \beta, \sigma^2, y) | \sigma^2, y] = \\ E[\sigma^2 I | \sigma^2, y] &+ \text{Var}[\tilde{X}\hat{\beta} | \sigma^2, y] = \\ \sigma^2 I &+ \sigma^2 \tilde{X}V_{\beta}\tilde{X}^t \end{aligned}$$

$p(\tilde{y} | y)$ , has two components of uncertainty

1. variability of the model not accounted by  $X\beta$  ( $\sigma^2$ )
2. posterior uncertainty in  $\beta$  and  $\sigma^2$  due to the finite sample size of  $y$

7

Table 1: Short-term measurements of radon concentration in a sample of houses in three counties in Minnesota. All measurements were recorded on the basement level of the houses, except for those indicated with \*, which were recorded on the first floor

| County     | Radon Measurements  |
|------------|---|
| Blue Earth | 5.0, 13.0, 7.2, 6.8, 12.8, 5.8*, 9.5, 6.0, 3.8, 14.3*, 1.8, 6.9, 4.7, 9.5     |
| Clay       | 0.9*, 12.9, 2.6, 3.5*, 26.6, 1.5, 13.0, 8.8, 19.5, 2.5*, 9.0, 13.1, 3.6, 6.9* |
| Goodhue    | 14.3, 6.9*, 7.6, 9.8*, 2.6, 43.5, 4.9, 3.5, 4.8, 5.6, 3.5, 3.9, 6.7           |

## Analysis of Radon measurements

- Exercise 8.1: fit a linear regression to the logarithms of the radon measurements, with indicator variables for the three counties and for whether a measurement was recorded on the first floor.

- Basic model

$$\begin{aligned} \log(y_i) &= \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \beta_4 x_{i4} + \epsilon_i \\ \epsilon_i &\sim N(0, \sigma^2) \end{aligned}$$

where

- $x_{i1}$  is the basement indicator

8

- $x_{i2}$ ,  $x_{i3}$  and  $x_{i4}$  are counties indicators
- $\exp(\beta_2)$  denotes the geometric mean of the radon measurements (no in the basement) in the county Blue Earth
- $\exp(\beta_2 + \beta_1)$  denotes the geometric mean of the radon measurements (in the basement) in the county Blue Earth
- Prior

$$\beta_1, \beta_2, \beta_3, \beta_4, \log(\sigma) \sim \text{constant}$$

```
#bayesian approach to ordinary linear regression, radon data-set pag.189 Gelman book,
#ex 8.1(a) exercise 8.1(b) homework
#explanatory variables:basement indicator and three counties indicators
make_indicators_function(x){
  ux_unique(x)
  mat1_matrix(x,nrow=length(x),ncol=length(ux))
  mat2_matrix(ux,nrow=length(x),ncol=length(ux),byrow=T)
  (mat1=mat2)*1}

bayesian_regression_function(nsim=1000){
  y_1_c(5.0, 13.0, 7.2,6.8,12.8,5.8,9.5,6.0,3.8,14.3,1.8,6.9,4.7,9.5)
  y_2_c(0.9,12.9,2.6,3.5,26.6,1.5,13.0,8.8,19.5,2.5,9.0,13.1,3.6,6.9)
  y_3_c(14.3,6.9,7.6,9.8,2.6,43.5,4.9,3.5,4.8,5.6,3.5,3.9,6.7)
  basement_1_c(1,1,1,1,0,1,1,1,0,1,1,1)
  basement_2_c(0,1,1,0,1,1,1,1,0,1,1,1,0)
  basement_3_c(1,0,1,0,1,1,1,1,1,1,1,1)
  counties_rep(1:3,c(length(y_1),length(y_2),length(y_3)))
  y_c(y_1,y_2,y_3)
  x_cbind(c(basement_1,basement_2,basement_3),make_indicators(counties))
  ls_out_lsfit(x,log(y),intercept=F)
  lsd_ls_diag(ls_out)
  n_nrow(x)
  k_ncol(x)
  sigmasqr_rep(NA,nsim)
  beta_matrix(NA,nsim,k)
  for ( i in 1:nsim){
    sigmasqr[i]_lsd$std.dev*(n-k)/rchisq(1,n-k)
    PP_ t(x)%*%x
    VV_solve(PP)
    VV_ .5*(VV+t(VV))
    beta[i,]_simulate.multnorm(as.numeric(ls_out$coef),sigmasqr[i]*VV)
  }
  output_exp(cbind(beta[,2],beta[,1]+beta[,2],beta[,3],beta[,1]+beta[,3],
    beta[,4],beta[,1]+beta[,4],sigmasqr))
  for(i in 1:ncol(output)) print(round(quantile(output[,i],c(.25,.5,.75)),1))
  return(beta)
}
```

|  | Posterior quantiles |     |     |
|--|---------------------|-----|-----|
|  | 25%                 | 50% | 75% |
| geometric mean for Blue earth (no basement) $\exp(\beta_2)$        | 3.9                 | 5   | 6.6 |
| geometric mean for Blue earth (basement) $\exp(\beta_1 + \beta_2)$ | 6                   | 7   | 8.2 |
| geometric mean for Clay (no basement) $\exp(\beta_3)$              | 3.7                 | 4.8 | 6.1 |
| geometric mean for Clay (basement) $\exp(\beta_1 + \beta_3)$       | 5.5                 | 6.5 | 7.9 |
| geometric mean for Goodhue (no basement) $\exp(\beta_4)$           | 3.8                 | 4.9 | 6.2 |
| geometric mean for Goodhue (basement) $\exp(\beta_1 + \beta_4)$    | 5.7                 | 6.8 | 7.8 |
| geometric std of predictions, $\exp(\sigma)$                       | 2                   | 2.3 | 2.6 |

## Hierarchical Regression Models (Chapter 13)



Hierarchical Regression Models are useful tools when:

- we have covariate information at different levels of variation
- data are obtained by stratified or cluster sampling

## Posterior Distributions of the coefficients

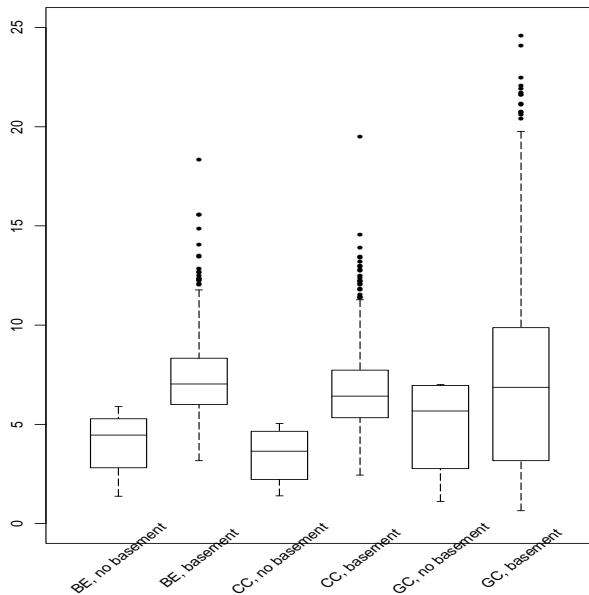


Figure 1: Boxplots of sample of the posterior distributions of the coefficients of “basement” and “no basement” for three counties

For example, in studying scholastic achievement we may have information about

1. individual students (family background)
2. class-level informations (characteristics of the teacher)
3. information about the school (educational policy, type of neighborhood)

rat tumor data set in Gelfand et al. (1990) is an example of hierarchical regression model

$$Y_{ij} \sim N(\alpha_i + \beta_i x_{ij}, \sigma^2), \quad i = 1, \dots, k = 30$$

$$\begin{pmatrix} \alpha_i \\ \beta_i \end{pmatrix} \sim N\left(\begin{bmatrix} \alpha_0 \\ \beta_0 \end{bmatrix}, \Sigma\right), \quad i = 1, \dots, k$$

$$\sigma^2 \sim IG(a, b)$$

$$\begin{pmatrix} \alpha_0 \\ \beta_0 \end{pmatrix} \sim N\left(\begin{bmatrix} \eta_0 \\ \eta_1 \end{bmatrix}, C\right)$$

$$\Sigma^{-1} \sim W((\rho R)^{-1}, \rho), \quad E(\Sigma^{-1}) = R^{-1}, \quad \text{var}(\Sigma) \propto \rho^{-1}$$

## Hierarchical Regression Model

We assume a normal linear regression model for the likelihood

$$y \mid \beta, X, \Sigma_y \sim N(X\beta, \Sigma_y)$$

$$\dim(X) = n \times J \text{ explanatory variables}$$

$$\beta = (\beta_1, \dots, \beta_J) \text{ regression coefficient:}$$

$$E(y_i \mid \beta, X) = \beta_1 x_{i1} + \dots + \beta_J x_{iJ}$$

$$\text{Var}(y_i \mid \beta, X) = [\Sigma_y]_{ii}$$

$$\text{Cov}(y_i, y_j \mid \beta, X) = [\Sigma_y]_{ij}$$

### Simple random effect model

$$y \mid \beta, X, \Sigma_y \sim N_n(X\beta, \Sigma_y)$$

$$\beta \mid \alpha, \sigma_\beta^2 \sim N_J(\mathbf{1}\alpha, \sigma_\beta^2 I)$$

- $\mathbf{1}$  is a vector of ones
- $\sigma_\beta^2 = 0 \rightarrow$  all  $\beta'_j$ s equal
- $\sigma_\beta^2 = \infty \rightarrow$  unrelated  $\beta'_j$ s
- $p(\log \sigma_\beta) = \text{constant} \rightarrow$  improper posterior (sec.5.4) with all its mass on  $\sigma_\beta = 0$
- we can use  $\sigma_\beta^2 \sim \chi^2$  scaled with  $df \leq 2$
- always check that the posterior inferences are not sensitive to the choice of the prior distribution for  $\sigma^2$

17

Positive intraclass correlation in a linear regression can be subsumed into a random effect model by augmenting the regression with  $J$  indicator variables whose coefficients have the population distribution  $\beta \sim N(\mathbf{1}\alpha, \sigma_\beta^2 I)$

19

### Intraclass correlation

- data  $y_1, \dots, y_n$  fall into  $J$  batches

$$y \mid \alpha, \Sigma_y \sim N(\alpha, \Sigma_y)$$

$$\text{Var}(y_i) = \sigma^2, \text{ for all } i$$

$$\text{Cov}(y_{i_1}, y_{i_2}) = \begin{cases} \rho\sigma^2 & i_1, i_2 \text{ are in the same batch} \\ 0 & \text{otherwise} \end{cases}$$

- if  $\rho \geq 0$ , this is equivalent to the model

$$y \sim N_n(X\beta, \sigma^2 I)$$

$$\beta \sim N(\mathbf{1}\alpha, \sigma_\beta^2 I)$$

where  $X$  is a  $n \times J$  matrix of indicator variables such that  $X_{ij} = 1$  if unit  $i$  is in batch  $j$  and 0 otherwise.

- the equivalence of the two models occurs when  $\rho = \sigma_\beta^2 / (\sigma^2 + \sigma_\beta^2)$ .

18

### Mixed effect model

$$\beta = \begin{bmatrix} \beta_{J_1} \\ \beta_{J_2} \end{bmatrix}$$

$$\beta_{J_1} \sim \text{Unif}(-\infty, \infty)$$

$$\beta_{J_2} \sim N_{J_2}(\mathbf{1}\alpha, \sigma_\beta^2 I), \quad J_2 = J - J_1$$

- $\beta_{J_1}$  are exchangeable with infinite prior variance and they are labeled *fixed effects*
- the random effect model with the school means normally distributed and a uniform prior density assumed for their mean  $\alpha$  is equivalent to a mixed-effect model with a single constant *fixed-effect* and a set of *random effects* of mean 0.

20

Several sets of random effects

$$\beta = \begin{bmatrix} \beta_{J_1} \\ \vdots \\ \beta_{J_K} \end{bmatrix}$$

$$\beta_{J_k} \sim N_{J_k}(1\alpha_k, \sigma_{\beta_k}^2 I)$$

- a mixed effect is obtained by setting the variance to  $\infty$  for one of the clusters of random effects

21

General notation and computation for hierarchical linear models

$$y \mid X, \beta, \Sigma \sim N_n(X\beta, \Sigma_y) \quad \textit{likelihood}$$

$$\beta \mid X, \alpha, \Sigma_\beta \sim N_J(X_\beta\alpha, \Sigma_\beta) \quad \textit{population dist.}$$

$$\alpha \mid \alpha_0, \Sigma_\alpha \sim N_K(\alpha_0, \Sigma_\alpha) \quad \textit{hyperprior dist.}$$

- $n$  data points  $y_i$
- $J$  parameters  $\beta_j$
- $K$  parameters  $\alpha_k$

22

Interpretation as a single regression

We consider the hierarchical model as a single normal regression model using a “larger” data set that includes as observations the information added by the population and hyperprior distributions

$$y_\star \mid X_\star, \gamma, \Sigma_\star \sim N(X_\star\gamma, \Sigma_\star)$$

where

- $\gamma = (\beta, \alpha)$ ,  $\dim(\gamma) = J + K$
- $y_\star = (y, 0, \alpha_0)^t$

$$\bullet X_\star = \begin{pmatrix} X & 0 \\ I_j & -X_\beta \\ 0 & I_K \end{pmatrix}, \quad \Sigma_\star^{-1} = \begin{pmatrix} \Sigma_y^{-1} & 0 & 0 \\ 0 & \Sigma_y^{-1} & 0 \\ 0 & 0 & \Sigma_\alpha^{-1} \end{pmatrix}$$

23

If any of the components of  $\beta$  or  $\alpha$  have non informative prior distributions, the corresponding rows in  $y_\star$  and  $X_\star$ , as well as the corresponding rows and columns in  $\Sigma_\star^{-1}$ , can be eliminated, because they correspond to observations with infinite variance. The resulting regression then has  $n + J_\star + K_\star$  “observations” where  $J_\star$  and  $K_\star$  are the number of components of  $\beta$  and  $\alpha$  with informative prior distributions.

24

## Bayesian Variable Selection

(George and McCulloch JASA 1993)

- The problem is to find and fit the “best” model of the form

$$Y = X_1^* \beta_1 + X_2^* \beta_2 + \dots + X_q^* \beta_q + \epsilon$$

- where  $(X_1^*, X_2^*, \dots, X_q^*)$  is a “selected” subset of  $X_1, \dots, X_p$ .

- If we consider the canonical regression setup

$$y \sim N_n(X\beta, \sigma^2 I)$$

where  $\beta = (\beta_1, \dots, \beta_p)$ .

- Selecting a subset of predictor is equivalent to setting to 0 those  $\beta_i$ s corresponding to non selected predictors.

25

- Bayesian variable selection introduce a binary latent variable  $\gamma_i$  and assumes

$$\beta_i \gamma_i \sim (1 - \gamma_i)N(0, \tau_i^2) + \gamma_i N(0, c_i^2 \tau_i^2)$$
$$p(\gamma_i = 1) = 1 - p(\gamma_i = 0) = p_i$$

- $p_i$  is the prior probability that  $\beta_i$  will require a non zero estimate, or equivalently that  $X_i$  will be included in the model.
- promising subset of predictors can be identified as those with higher posterior probability
- In practice, a common sense Bayesian perspective indicates that the key is to use substantive knowledge, either as a formal prior distribution or more informally, in choosing which variable to include.

26