

# Robust Bayesian Regression

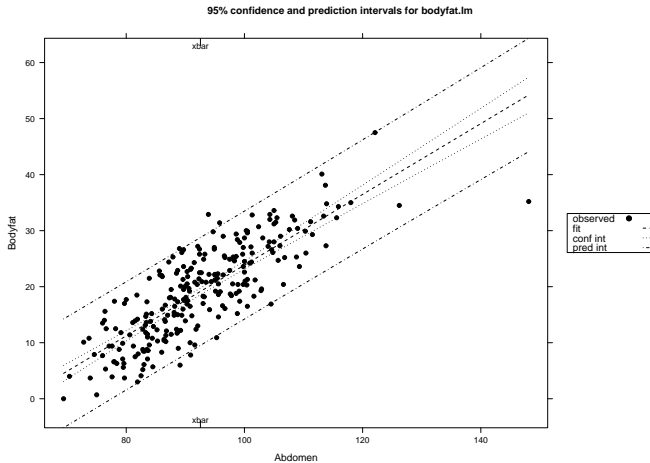
Nov 10, 2015

Readings: Hoff Chapter 9

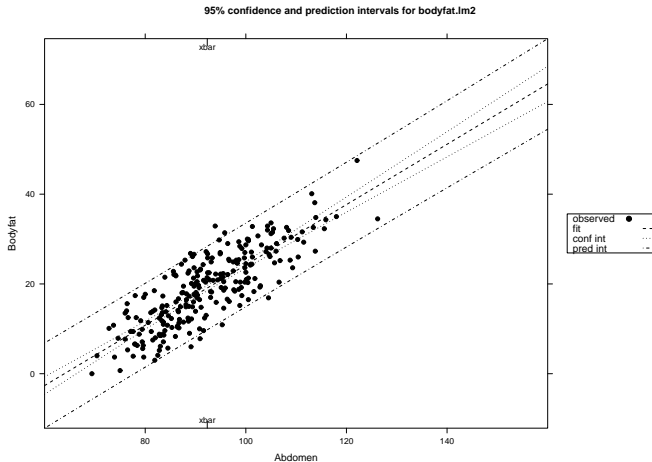
November 10, 2015

# Body Fat Data: Intervals w/ All Data

Response % Body Fat and Predictor Waist Circumference



# Intervals: without case 39



# Interpretations

- ▶ For a given Abdominal circumference, our probability that the mean bodyfat percentage is in the (pointwise) intervals given by the dotted lines is 0.95. For men with 34 inch waist there is a 95% chance that the average bodyfat is between 14.4 to 15.8.
- ▶ For a new man with a given Abdominal circumference, our probability that his bodyfat percentage is in the intervals given by the dashed lines is 0.95. For a man with a 34 inch waist, our probability that his bodyfat is between 5.7 to 24.4 is 0.95.
- ▶ Both have same point estimate
- ▶ Increased uncertainty for prediction of a new observation versus estimating the expected value.

Which analysis do we use? with Case 39 or not – or something different?

# Options for Handling Influential Cases

- ▶ Are there scientific grounds for eliminating the case?
- ▶ Test if the case has a different mean than population
- ▶ Report results with and without the case
- ▶ Model Averaging?
- ▶ Full model  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{I}_n\boldsymbol{\delta} + \boldsymbol{\epsilon}$
- ▶  $2^n$  submodels  $\gamma_i = 0 \Leftrightarrow \delta_i = 0$
- ▶ If  $\gamma_i = 1$  then case  $i$  has a different mean “mean shift” outliers.

# Mean Shift = Variance Inflation

► Model  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{I}_n\delta + \epsilon$

► Prior

$$\delta_i \mid \gamma_i \sim N(0, V\sigma^2\gamma_i)$$

$$\gamma_i \sim \text{Ber}(\pi)$$

Then  $\epsilon_i$  given  $\sigma^2$  is independent of  $\delta_i$  and

$$\epsilon_i^* \equiv \epsilon_i + \delta_i \mid \sigma^2 \begin{cases} N(0, \sigma^2) & \text{wp } (1 - \pi) \\ N(0, \sigma^2(1 + V)) & \text{wp } \pi \end{cases}$$

Model  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \epsilon^*$  “variance inflation”

$V + 1 = K = 7$  in the paper by Hoeting et al. package BMA

# Simultaneous Outlier and Variable Selection

```
MC3.REG(all.y = bodyfat$Bodyfat, all.x = as.matrix(bodyfat$Abdom  
num.its = 10000, outliers = TRUE)
```

Model parameters:  $PI=0.02$   $K=7$   $nu=2.58$   $lambda=0.28$   $phi=2.85$

15 models were selected

Best 5 models (cumulative posterior probability = 0.9939):

	prob	model 1	model 2	model 3	model 4	model 5
variables						
all.x	1	x	x	x	x	x
outliers						
39	0.94932	x	x	.	x	.
204	0.04117	.	.	.	x	.
207	0.10427	.	x	.	.	x
post prob		0.815	0.095	0.044	0.035	0.004

## Change Error Assumptions

$$Y_i \stackrel{\text{ind}}{\sim} t(\nu, \alpha + \beta x_i, 1/\phi)$$

$$L(\alpha, \beta, \phi) \propto \prod_{i=1}^n \phi^{1/2} \left( 1 + \frac{\phi(y_i - \alpha - \beta x_i)^2}{\nu} \right)^{-\frac{(\nu+1)}{2}}$$

Use Prior  $p(\alpha, \beta, \phi) \propto 1/\phi$

Posterior distribution

$$p(\alpha, \beta, \phi \mid Y) \propto \phi^{n/2-1} \prod_{i=1}^n \left( 1 + \frac{\phi(y_i - \alpha - \beta x_i)^2}{\nu} \right)^{-\frac{(\nu+1)}{2}}$$



## Bounded Influence - West 1984 (and references within)

Treat  $\sigma^2$  as given, then *influence* of individual observations on the posterior distribution of  $\beta$  in the model where  $E[\mathbf{Y}_i] = \mathbf{x}_i^T \beta$  is investigated through the score function:

$$\frac{d}{d\beta} \log p(\beta \mid \mathbf{Y}) = \frac{d}{d\beta} \log p(\beta) + \sum_{i=1}^n \mathbf{x}_i g(y_i - \mathbf{x}_i^T \beta)$$

where

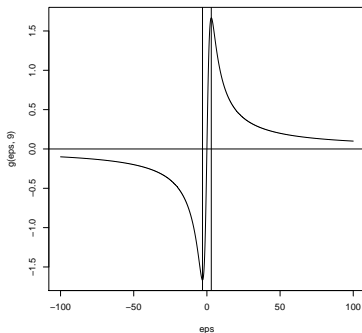
$$g(\epsilon) = -\frac{d}{d\epsilon} \log p(\epsilon)$$

is the influence function of the error distribution (unimodal, continuous, differentiable, symmetric)

An outlying observation  $y_j$  is accommodated if the posterior distribution for  $p(\beta \mid \mathbf{Y}_{(j)})$  converges to  $p(\beta \mid \mathbf{Y})$  for all  $\beta$  as  $|\mathbf{Y}_j| \rightarrow \infty$ . Requires error models with influence functions that go to zero such as the Student  $t$  (O'Hagan, 1979)

## Choice of df

- ▶ Score function for  $t$  with  $\alpha$  degrees of freedom has turning points at  $\pm\sqrt{\alpha}$



- ▶  $g'(\epsilon)$  is negative when  $\epsilon^2 > \alpha$  (standardized errors)
- ▶ Contribution of observation to information matrix is negative and the observation is doubtful
- ▶ Suggest taking  $\alpha = 8$  or  $\alpha = 9$  to reject errors larger than  $\sqrt{8}$  or 3 sd.

# Scale-Mixtures of Normal Representation

$$Z_i \stackrel{\text{iid}}{\sim} t(\nu, 0, \sigma^2) \Leftrightarrow$$

$$Z_i \mid \lambda_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2 / \lambda_i)$$

$$\lambda_i \stackrel{\text{iid}}{\sim} G(\nu/2, \nu/2)$$

Integrate out “latent”  $\lambda$ 's to obtain marginal distribution.

# Latent Variable Model

$$\begin{aligned} Y_i \mid \alpha, \beta, \phi, \lambda &\stackrel{\text{ind}}{\sim} N\left(\alpha + \beta x_i, \frac{1}{\phi \lambda_i}\right) \\ \lambda_i &\stackrel{\text{iid}}{\sim} G(\nu/2, \nu/2) \\ p(\alpha, \beta, \phi) &\propto 1/\phi \end{aligned}$$

Joint Posterior Distribution:

$$\begin{aligned} p((\alpha, \beta, \phi, \lambda_1, \dots, \lambda_n \mid Y) \propto & \phi^{n/2} \exp \left\{ -\frac{\phi}{2} \sum \lambda_i (y_i - \alpha - \beta x_i)^2 \right\} \times \\ & \phi^{-1} \\ & \prod_{i=1}^n \lambda_i^{\nu/2-1} \exp(-\lambda_i \nu/2) \end{aligned}$$

# Single Component Gibbs Sampler

Start with  $(\alpha^{(0)}, \beta^{(0)}, \phi^{(0)}, \lambda_1^{(0)}, \dots, \lambda_n^{(0)})$

For  $t = 1, \dots, T$ , generate from the following sequence of Full Conditional distributions:

- ▶  $p(\alpha \mid \beta^{(t-1)}, \phi^{(t-1)}, \lambda_1^{(t-1)}, \dots, \lambda_n^{(t-1)}, Y)$
- ▶  $p(\beta \mid \alpha^{(t)}, \phi^{(t-1)}, \lambda_1^{(t-1)}, \dots, \lambda_n^{(t-1)}, Y)$
- ▶  $p(\phi \mid \alpha^{(t)}, \beta^{(t)}, \phi^{(t-1)}, \lambda_1^{(t-1)}, \dots, \lambda_n^{(t-1)}, Y)$
- ▶  $p(\lambda_j \mid \alpha^{(t)}, \beta^{(t)}, \phi^{(t)}, \lambda_{(-j)}^{(t-1)}, Y)$  for  $j = 1, \dots, n$

$\lambda_{(-j)}$  is the vector of  $\lambda$ s excluding the  $j$ th component

Easy to find and sample! (work out)

# Programs

## BUGS: Bayesian inference Using Gibbs Sampling

- ▶ WinBUGS is the Windows implementation
  - ▶ can be called from R with R2WinBUGS package
  - ▶ can be run on any intel-based computer using VMware, wine
- ▶ OpenBUGS open source version of WinBUGS
- ▶ LinBUGS is the Linux implementation of OpenBUGS.
- ▶ JAGS: Just Another Gibbs Sampler is an alternative program that uses the (almost) same model description as BUGS (Linux, MAC OS X, Windows) Can call from R using `library(R2jags)`

Include more than just Gibbs Sampling

# JAGS

- ▶ Model
- ▶ Data
- ▶ Initial values (optional)

May do this through ordinary text files or use the functions in R2jags to specify model, data, and initial values then call jags.

## Model Specification via R2jags

```
rr.model = function() {  
  for (i in 1:n) {  
    mu[i] <- alpha0 + alpha1*(X[i] - Xbar)  
    lambda[i] ~ dgamma(9/2, 9/2)  
    prec[i] <- phi*lambda[i]  
    Y[i] ~ dnorm(mu[i], prec[i])  
  }  
  phi ~ dgamma(1.0E-6, 1.0E-6)  
  alpha0 ~ dnorm(0, 1.0E-6)  
  alpha1 ~ dnorm(0, 1.0E-6)  
}
```



# Notes on Models

- ▶ Distributions of stochastic “nodes” are specified using  $\sim$
- ▶ Assignment of deterministic “nodes” uses  $\leftarrow$  (NOT  $=$ )
- ▶ JAGS allows expressions as arguments in distributions (WinBUGS does not)
- ▶ Normal distributions are parameterized using *precisions*, so `dnorm(0, 1.0E-6)` is a  $N(0, 1.0 \times 10^6)$
- ▶ uses for loop structure as in R for model description but coded in C++ so is fast!

## Initial Values

Function to calculate initial values for parameters as a list

```
rr.inits = function() {  
  bf.lm = lm(bf.data$Y ~ bf.data$X)  
  coefs = coef(bf.lm)  
  alpha1=coefs[2]  
  alpha0 = coefs[1] - alpha1*bf.data$Xbar  
  phi = (1/summary(bf.lm)$sigma)^2  
  lambda = rep(1, bf.data$n)  
  return(list(alpha0=alpha0, alpha1 = alpha1,  
              phi=phi, lambda=lambda))  
}
```

# Data

A list or rectangular data structure for all data and summaries of data used in the model

```
bf.data = list(Y = bodyfat$Bodyfat,  
               X=bodyfat$Abdomen)  
bf.data$n = length(bf.data$Y)  
bf.data$Xbar = mean(bf.data$X)
```

## Specifying which Parameters to Save

The parameters to be monitored and returned to R are specified with the variable `parameters`

```
parameters = c("beta0", "beta1", "sigma",  
               "mu34", "y34", "lambda[39]")
```

- ▶ All of the above (except `lambda`) are calculated from the other parameters. (See R-code for definitions of these parameters.)
- ▶ `lambda[39]` saves only the 39th case of  $\lambda$
- ▶ To save a whole vector (for example all `lambdas`, just give the vector name)

## Running jags from R

Write the model out as a text file, then call `jags()`

```
model.file = "rr-model.txt"
write.model(rr.model, model.file)

bf.sim = jags(bf.data, rr.inits, parameters,
              model.file=model.file,
              n.chains=2, n.iter=5000,
              )
```

## Output

	mean	sd	2.5%	50%	97.5%
beta0	-41.70	2.75	-46.91	-41.67	-36.40
beta1	0.66	0.03	0.60	0.66	0.71
sigma	4.48	0.23	4.05	4.46	4.96
mu34	15.10	0.35	14.43	15.10	15.82
y34	14.94	5.15	4.37	15.21	24.65
lambda[39]	0.33	0.16	0.11	0.30	0.72

95% HPD interval for expected bodyfat (14.5, 15.8)

95% HPD interval for bodyfat (5.1, 25.3)

## Comparison

- ▶ 95% Probability Interval for  $\beta$  is (0.60, 0.71) with  $t_9$  errors
- ▶ 95% Confidence Interval for  $\beta$  is (0.58, 0.69) (all data normal model)
- ▶ 95% Confidence Interval for  $\beta$  is (0.61, 0.73) ( normal model without case 39)

Results intermediate without having to remove any observations  
Case 39 down weighted by  $\lambda_{39}$

## Full Conditional for $\lambda_j$

$$\begin{aligned} p(\lambda_j \mid \text{rest}, Y) &\propto p(\alpha, \beta, \phi, \lambda_1, \dots, \lambda_n \mid Y) \\ &\propto \phi^{n/2-1} \prod_{i=1}^n \exp \left\{ -\frac{\phi}{2} \lambda_i (y_i - \alpha - \beta x_i)^2 \right\} \times \\ &\quad \prod_{i=1}^n \lambda_i^{\frac{\nu+1}{2}-1} \exp(-\lambda_i \frac{\nu}{2}) \end{aligned}$$

Ignore all terms except those that involve  $\lambda_j$

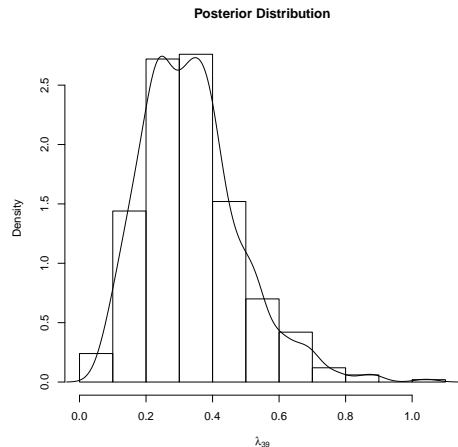
$$\lambda_j \mid \text{rest}, Y \sim G \left( \frac{\nu+1}{2}, \frac{\phi(y_j - \alpha - \beta x_j)^2 + \nu}{2} \right)$$



# Weights

Under prior  $E[\lambda_i] = 1$

Under posterior, large residuals are down-weighted (approximately those bigger than  $\sqrt{\nu}$ )



# Prior Distributions on Parameter

As a general recommendation, the prior distribution should have “heavier” tails than the likelihood

- ▶ with  $t_9$  errors use a  $t_\alpha$  with  $\alpha < 9$
- ▶ also represent via scale mixture of normals
- ▶ Horseshoe, Double Pareto, Cauchy all have heavier tails
- ▶ See Stack-loss code