



Regression - The Linear Model

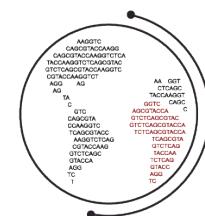
Manuel A. Rivas

Department of Biomedical Data Science

BDS215

Stanford University

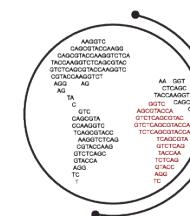
rivaslab.stanford.edu



RIVASLAB

Review concepts learned

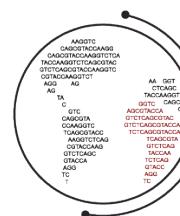
- Multilevel (Hierarchical) Modeling
- Mixture Models



RIVASLAB

Motivating examples for today

- Regression models



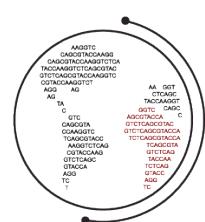
RIVASLAB

Linear model

Linear regression

Linear regression is one of the most widely used statistical tools.

In today's lecture we will focus on Bayesian model building and inference for normal linear models.



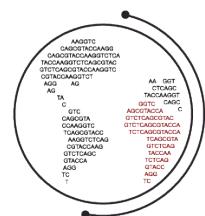
Linear model

Linear regression

Linear regression is one of the most widely used statistical tools.

In today's lecture we will focus on Bayesian model building and inference for normal linear models.

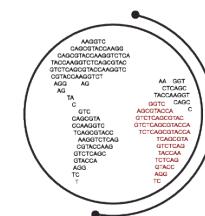
Learning objective: To set up the relevant Bayesian models and draw samples from posterior distributions for parameters θ and future observables \tilde{y} .



RIVASLAB

Motivation

- In many scientific studies concern relations among two or more observable quantities. A common question is: how does a quantity, y , vary as a function of another quantity or vector of quantities, x ?

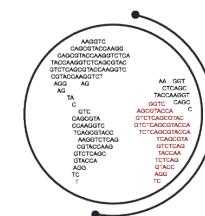


RIVASLAB

Motivation

- In many scientific studies concern relations among two or more observable quantities. A common question is: how does a quantity, y , vary as a function of another quantity or vector of quantities, x ?

We are interested in the conditional distribution of y , given x , denoted as $p(y|\theta, x)$.



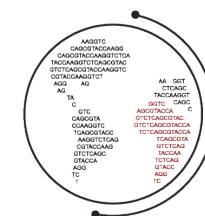
RIVASLAB

Motivation

- In many scientific studies concern relations among two or more observable quantities. A common question is: how does a quantity, y , vary as a function of another quantity or vector of quantities, x ?

We are interested in the conditional distribution of y , given x , denoted as $p(y|\theta, x)$.

- y is called the *response* or *outcome variable*



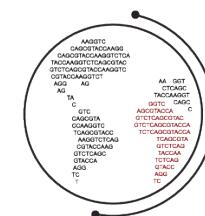
RIVASLAB

Motivation

- In many scientific studies concern relations among two or more observable quantities. A common question is: how does a quantity, y , vary as a function of another quantity or vector of quantities, x ?

We are interested in the conditional distribution of y , given x , denoted as $p(y|\theta, x)$.

- y is called the *response* or *outcome variable*
- $x = (x_1, \dots, x_k)$ are called explanatory variables



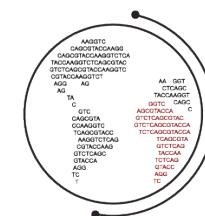
RIVASLAB

Motivation

- In many scientific studies concern relations among two or more observable quantities. A common question is: how does a quantity, y , vary as a function of another quantity or vector of quantities, x ?

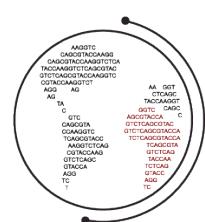
We are interested in the conditional distribution of y , given x , denoted as $p(y|\theta, x)$.

- y is called the *response* or *outcome variable*
- $x = (x_1, \dots, x_k)$ are called explanatory variables
- Sometimes, a single variable x_j may be of primary interest and consider it the *treatment* variable, labeling the other components of x as the *control* variables



The Linear Model

Let's rewrite to better accommodate notation.

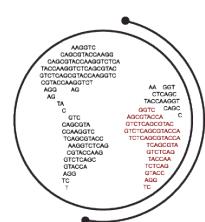


RIVASLAB

The Linear Model

Let's rewrite to better accommodate notation.

$$\mathbf{y} = \mathbb{X}\boldsymbol{\beta} + \epsilon$$



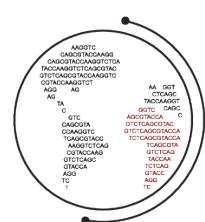
RIVASLAB

The Linear Model

Let's rewrite to better accommodate notation.

$$\mathbf{y} = \mathbb{X}\boldsymbol{\beta} + \epsilon$$

- \mathbf{y} is an $n \times 1$ vector of observations (n is the number of individuals in the study, for example).

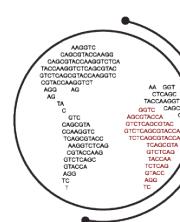


The Linear Model

Let's rewrite to better accommodate notation.

$$\mathbf{y} = \mathbb{X}\boldsymbol{\beta} + \epsilon$$

- \mathbf{y} is an $n \times 1$ vector of observations (n is the number of individuals in the study, for example).
- \mathbb{X} is an $n \times k$ matrix of known coefficients

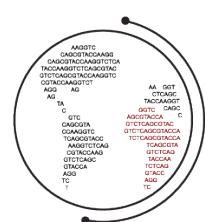


The Linear Model

Let's rewrite to better accommodate notation.

$$\mathbf{y} = \mathbb{X}\boldsymbol{\beta} + \epsilon$$

- \mathbf{y} is an $n \times 1$ vector of observations (n is the number of individuals in the study, for example).
- \mathbb{X} is an $n \times k$ matrix of known coefficients
- $\boldsymbol{\beta}$ is a $k \times 1$ vector of parameters

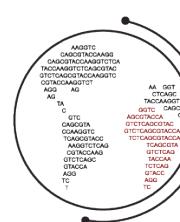


The Linear Model

Let's rewrite to better accommodate notation.

$$\mathbf{y} = \mathbb{X}\beta + \epsilon$$

- \mathbf{y} is an $n \times 1$ vector of observations (n is the number of individuals in the study, for example).
- \mathbb{X} is an $n \times k$ matrix of known coefficients
- β is a $k \times 1$ vector of parameters
- ϵ is an $n \times 1$ vector of random errors



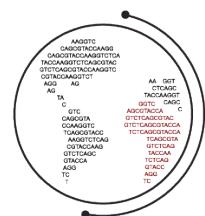
The Linear Model

Let's rewrite to better accommodate notation.

$$\mathbf{y} = \mathbb{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

- \mathbf{y} is an $n \times 1$ vector of observations (n is the number of individuals in the study, for example).
- \mathbb{X} is an $n \times k$ matrix of known coefficients
- $\boldsymbol{\beta}$ is a $k \times 1$ vector of parameters
- $\boldsymbol{\epsilon}$ is an $n \times 1$ vector of random errors

The simplest and most widely used version of the linear model is the *normal linear model*



RIVASLAB

The Linear Model

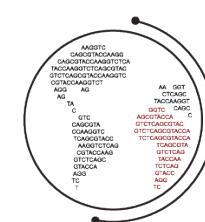
Let's rewrite to better accommodate notation.

$$\mathbf{y} = \mathbb{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

- \mathbf{y} is an $n \times 1$ vector of observations (n is the number of individuals in the study, for example).
- \mathbb{X} is an $n \times k$ matrix of known coefficients
- $\boldsymbol{\beta}$ is a $k \times 1$ vector of parameters
- $\boldsymbol{\epsilon}$ is an $n \times 1$ vector of random errors

The simplest and most widely used version of the linear model is the *normal linear model*

- The elements are assumed to have zero mean (critical and often overlooked in practice)



RIVASLAB

The Linear Model

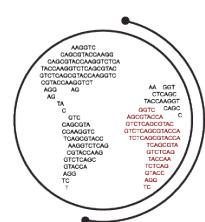
Let's rewrite to better accommodate notation.

$$\mathbf{y} = \mathbb{X}\beta + \epsilon$$

- \mathbf{y} is an $n \times 1$ vector of observations (n is the number of individuals in the study, for example).
- \mathbb{X} is an $n \times k$ matrix of known coefficients
- β is a $k \times 1$ vector of parameters
- ϵ is an $n \times 1$ vector of random errors

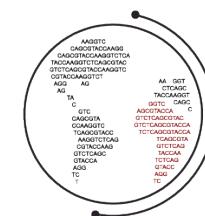
The simplest and most widely used version of the linear model is the *normal linear model*

- The elements are assumed to have zero mean (critical and often overlooked in practice)
- Assumed to be uncorrelated and to have common variance σ^2 , which becomes an additional parameter



The Linear Model

Question: What are the parameters?

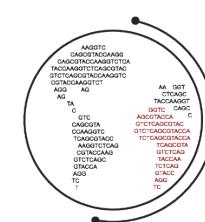


RIVASLAB

The Linear Model

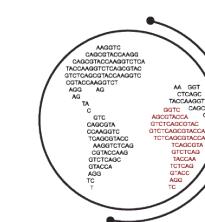
Question: What are the parameters?

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



RIVASLAB

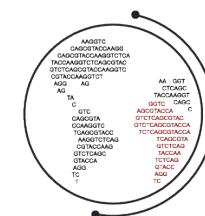
Key statistical modeling issues in normal linear model framework



RIVASLAB

Key statistical modeling issues in normal linear model framework

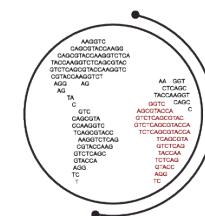
1. Defining the variables x and y (possibly using transformations)



RIVASLAB

Key statistical modeling issues in normal linear model framework

1. Defining the variables x and y (possibly using transformations)
2. Setting up a prior distribution on the model parameters that accurately reflect substantive knowledge

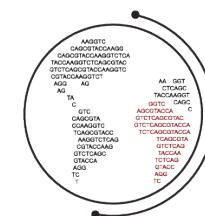


RIVASLAB

Key statistical modeling issues in normal linear model framework

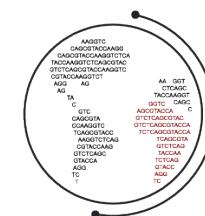
1. Defining the variables x and y (possibly using transformations)
2. Setting up a prior distribution on the model parameters that accurately reflect substantive knowledge

Statistical inference problem is to estimate the parameters θ , conditional on \mathbb{X} and \mathbb{y} .



RIVASLAB

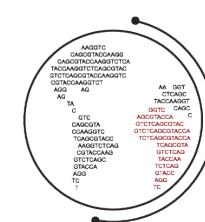
Notation for a basic normal linear model



RIVASLAB

Notation for a basic normal linear model

$$\mathbf{y} | \boldsymbol{\beta}, \sigma^2, \mathbb{X} \sim \mathcal{N}(\mathbb{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$$

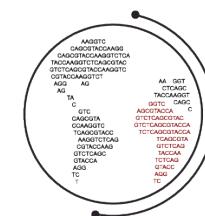


RIVASLAB

Notation for a basic normal linear model

$$\mathbf{y} | \boldsymbol{\beta}, \sigma^2, \mathbb{X} \sim \mathcal{N}(\mathbb{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$$

Multivariate normal distribution where as usual \mathbf{I} represents the identity matrix.



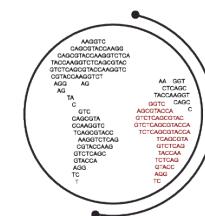
Notation for a basic normal linear model

$$\mathbf{y} | \boldsymbol{\beta}, \sigma^2, \mathbb{X} \sim \mathcal{N}(\mathbb{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$$

Multivariate normal distribution where as usual \mathbf{I} represents the identity matrix.

The likelihood becomes

$$f(\mathbf{y} | \boldsymbol{\beta}, \sigma^2) = (2\pi\sigma^2)^{-n/2} \exp [(\mathbf{y} - \mathbb{X}\boldsymbol{\beta})' (\mathbf{y} - \mathbb{X}\boldsymbol{\beta}) / (2\sigma^2)]$$



RIVASLAB

Notation for a basic normal linear model

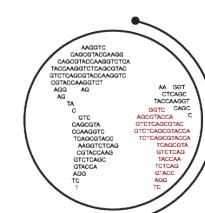
$$\mathbf{y} | \boldsymbol{\beta}, \sigma^2, \mathbb{X} \sim \mathcal{N}(\mathbb{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$$

Multivariate normal distribution where as usual \mathbf{I} represents the identity matrix.

The likelihood becomes

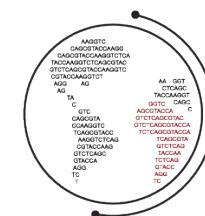
$$f(\mathbf{y} | \boldsymbol{\beta}, \sigma^2) = (2\pi\sigma^2)^{-n/2} \exp [(\mathbf{y} - \mathbb{X}\boldsymbol{\beta})' (\mathbf{y} - \mathbb{X}\boldsymbol{\beta}) / (2\sigma^2)]$$

$(\mathbf{y} - \mathbb{X}\boldsymbol{\beta})' (\mathbf{y} - \mathbb{X}\boldsymbol{\beta})$ is referred to as the *quadratic form*



Least squares estimator of β

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



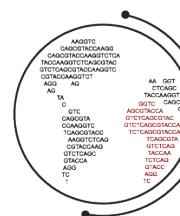
RIVASLAB

Least squares estimator of β

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

Residual sum of squares

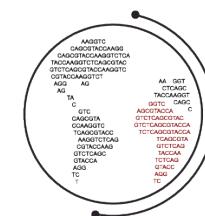
$$S = (\mathbb{y} - \mathbb{X}\hat{\beta})^T (\mathbb{y} - \mathbb{X}\hat{\beta})$$



RIVASLAB

The posterior distribution

We focused on examples of normal distributions with unknown mean and variances in previous lectures.

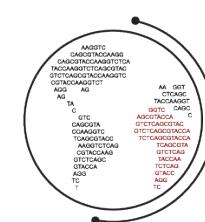


RIVASLAB

The posterior distribution

We focused on examples of normal distributions with unknown mean and variances in previous lectures.

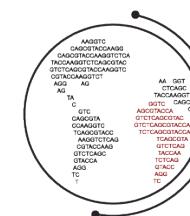
Question: We are interested in the posterior distribution of...



RIVASLAB

The posterior distribution

β, σ^2



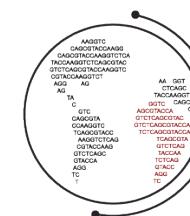
RIVASLAB

The posterior distribution

$$\beta, \sigma^2$$

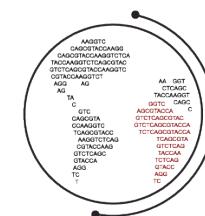
We factor the joint posterior distribution for β and σ^2 as

$$p(\beta, \sigma^2 | y) = p(\beta | \sigma^2, y) p(\sigma^2 | y)$$



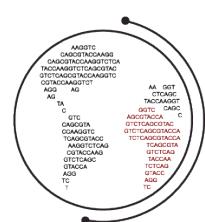
The posterior distribution

Here, we are interested in determining first, the posterior distribution for β , conditioning on σ^2 , and then the marginal posterior distribution for σ^2 .



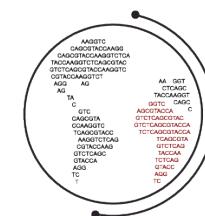
The conditional posterior distribution of β , given σ^2

$$\beta | \sigma^2, \mathbf{y} \sim \mathcal{N}(\hat{\beta}, \mathbf{V}_\beta \sigma^2),$$



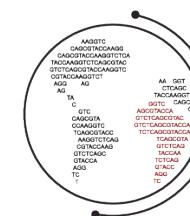
Least squares estimator of β

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



\mathbf{V}_β

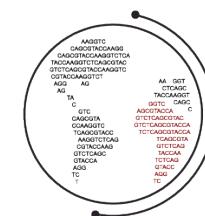
$$\mathbf{V}_\beta = (\mathbb{X}^{\textcolor{blue}{T}} \mathbb{X})^{-1}$$



RIVASLAB

Marginal posterior distribution of σ^2

$$\sigma^2 | \mathbf{y} \sim IG\left(\frac{n-k}{2}, s^2\right),$$



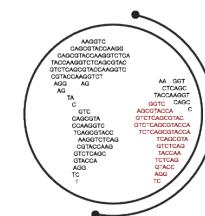
RIVASLAB

Marginal posterior distribution of σ^2

$$\sigma^2 | \mathbf{y} \sim IG\left(\frac{n-k}{2}, s^2\right),$$

What you usually do in practice is to draw inferences by simulation

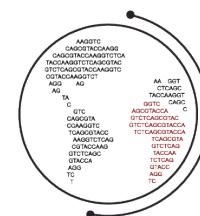
1. Draw simulations of σ^2
2. Then, draw simulations of $\beta | \sigma^2$



non-Bayesian estimates of β and σ^2

- $\hat{\beta}$
- s^2

We obtain the classical standard error estimate for β by setting σ^2 to s^2

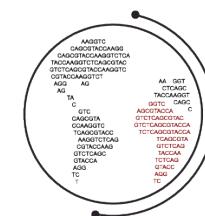


Case study : Trees

Plotting Outcome (Volume) as a function of Covariates (Girth, Height)

Source: http://www.columbia.edu/~cjd11/charles_dimaggio/DIRE/styled-4/styled-11/code-8/

```
data(trees)
attach(trees)
dim(trees)
trees
```

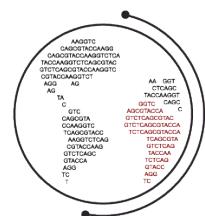


RIVASLAB

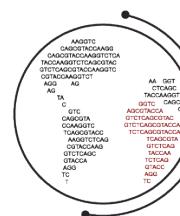
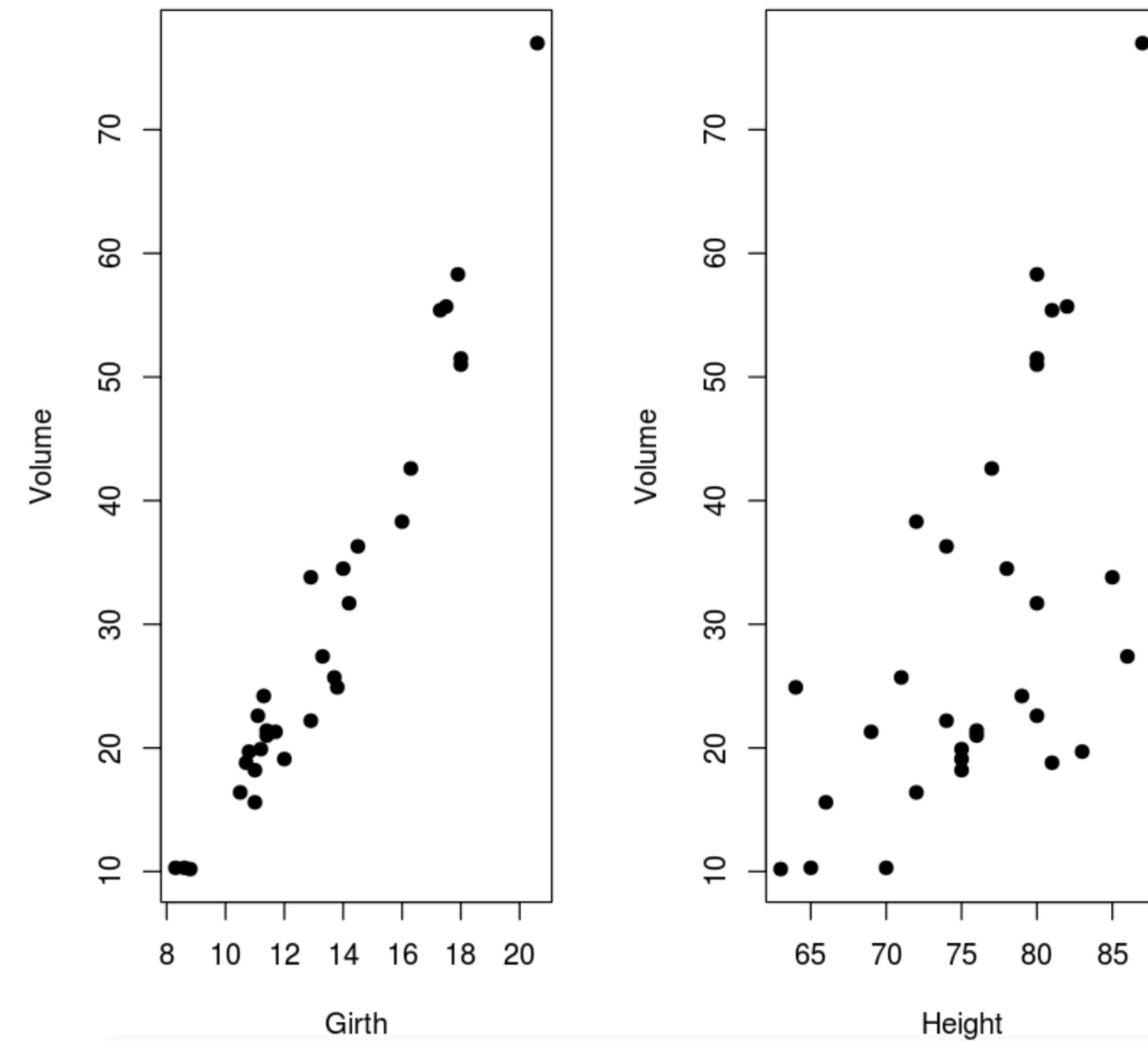
Case study : Trees

31 3

Girth	Height	Volume
8.3	70	10.3
8.6	65	10.3
8.8	63	10.2
10.5	72	16.4
10.7	81	18.8
10.8	83	19.7
11.0	66	15.6
11.0	75	18.2
11.1	80	22.6
11.2	75	19.9
11.3	79	24.2
11.4	76	21.0
11.4	76	21.4
11.7	69	21.3
12.0	75	19.1
12.9	74	22.2

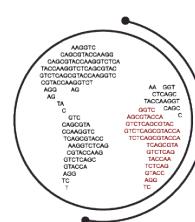


Case study : Trees



Case study : Trees

```
## MLE fit of Regression Model  
model <- lm(Volume~Girth+Height)  
summary(model)
```



Case study : Trees

Call:

```
lm(formula = Volume ~ Girth + Height)
```

Residuals:

Min	1Q	Median	3Q	Max
-6.4065	-2.6493	-0.2876	2.2003	8.4847

Coefficients:

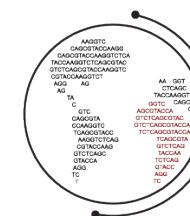
	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-57.9877	8.6382	-6.713	2.75e-07 ***
Girth	4.7082	0.2643	17.816	< 2e-16 ***
Height	0.3393	0.1302	2.607	0.0145 *

Signif. codes:	0	'***'	0.001	'**'
	0.01	'*'	0.05	'. '
	0.1	' '	1	

Residual standard error: 3.882 on 28 degrees of freedom

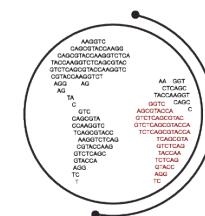
Multiple R-squared: 0.948, Adjusted R-squared: 0.9442

F-statistic: 255 on 2 and 28 DF, p-value: < 2.2e-16



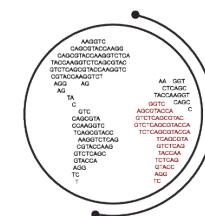
Sampling from posterior distribution of coefficients beta and variance sigsq

```
beta.hat <- model$coef  
n <- length(Volume)  
k <- length(beta.hat)  
s2 <- (n-k)*summary(model)$sigma^2  
V.beta <- summary(model)$cov.unscaled
```

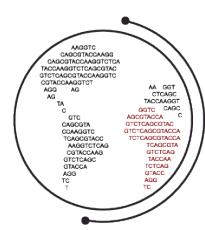
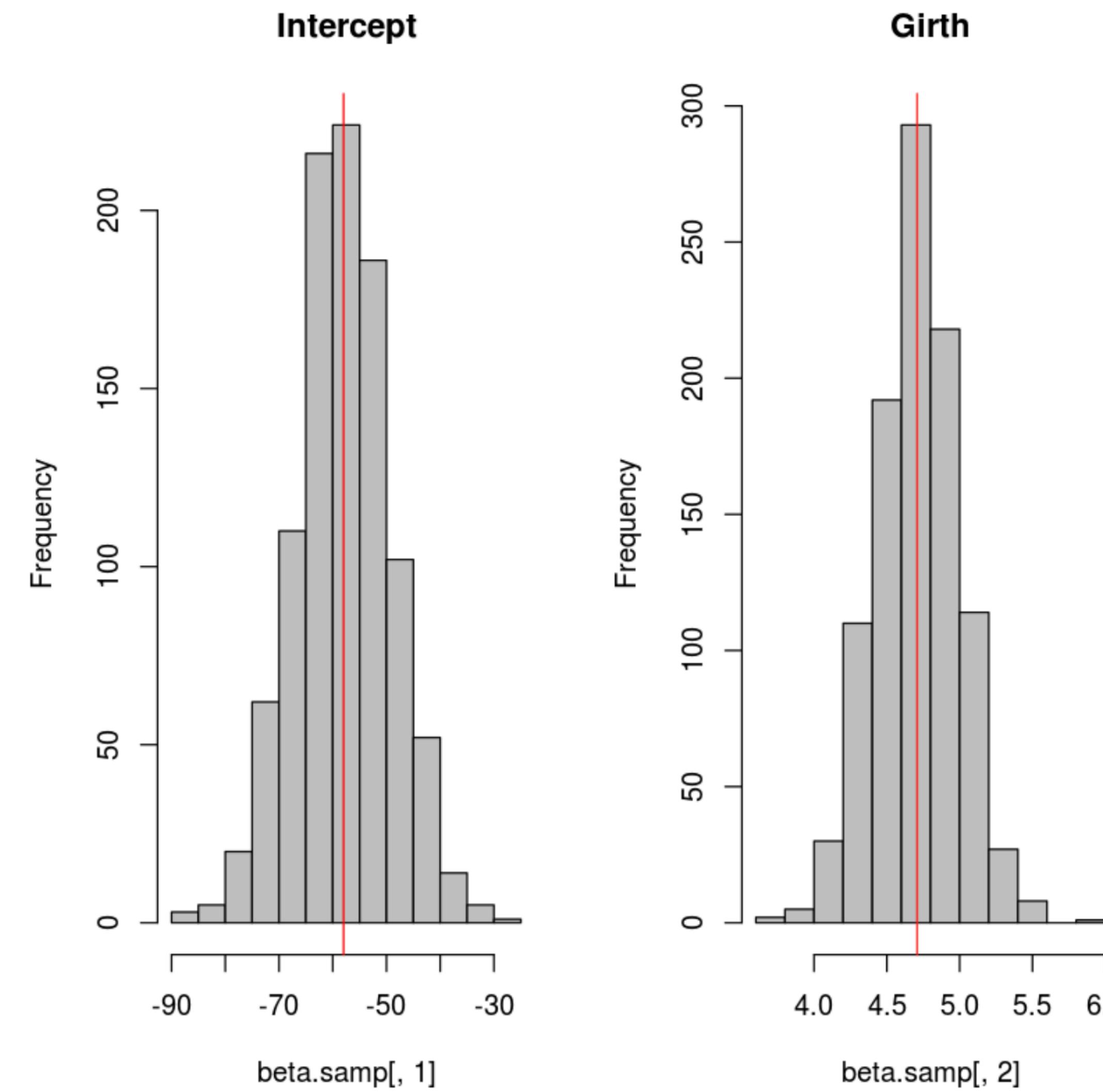


Sampling from posterior distribution of coefficients beta and variance sigsq

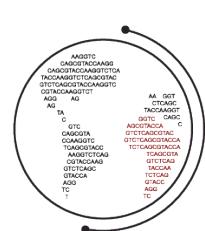
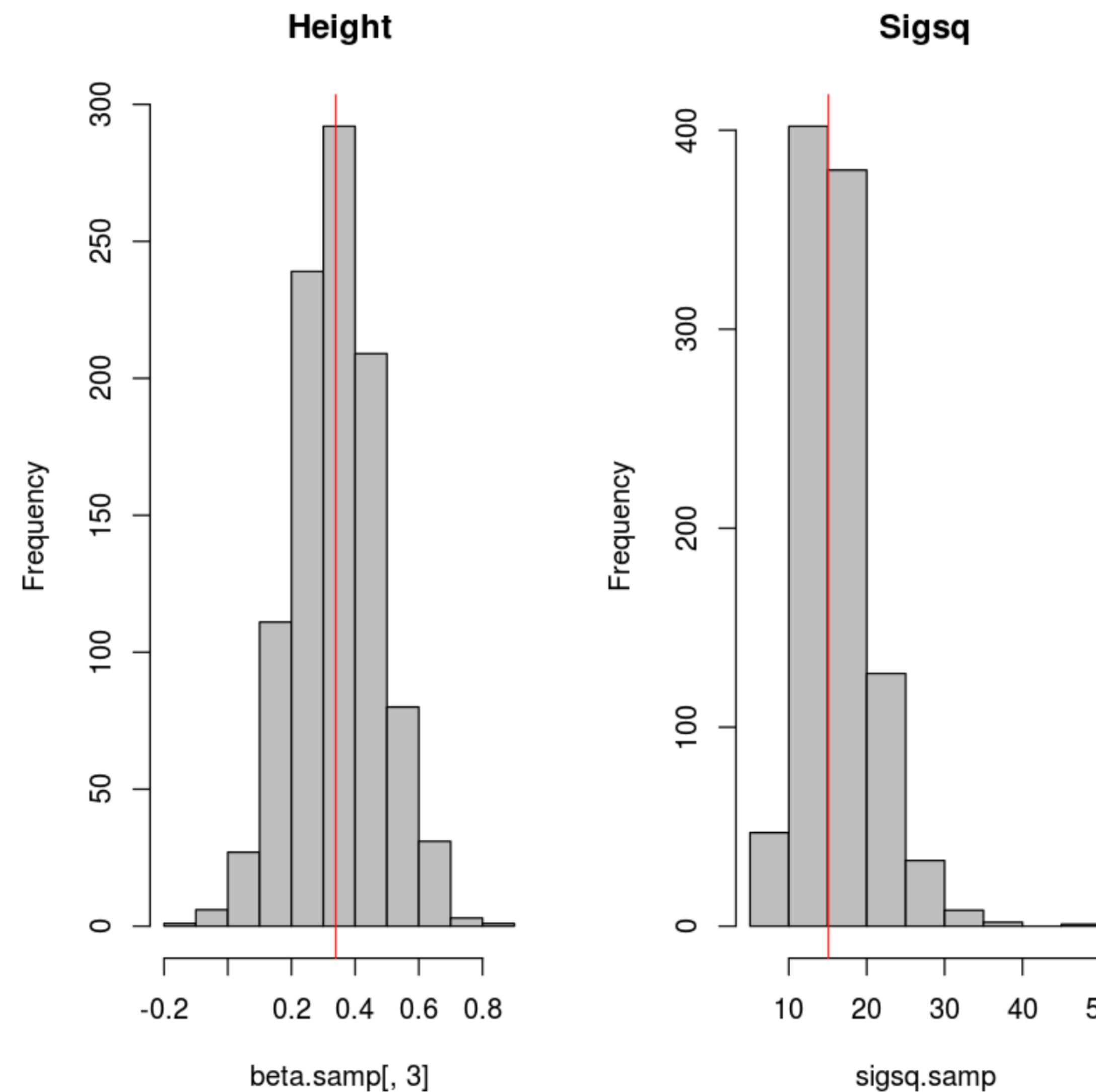
```
numsamp <- 1000
beta.samp <- matrix(NA,nrow=numsamp,ncol=k)
sigsq.samp <- rep(NA,numsamp)
for (i in 1:numsamp){
  temp <- rgamma(1,shape=(n-k)/2,rate=s2/2)
  cursigsq <- 1/temp # sampling from IG
  curvarbeta <- cursigsq*V.beta # sigma^2 *V_{beta}
  curvarbeta.chol <- t(chol(curvarbeta))
  z <- rnorm(k,0,1)
  curbeta <- beta.hat+curvarbeta.chol%*%z # sampling from normal
  sigsq.samp[i] <- cursigsq
  beta.samp[i,] <- curbeta
}
```



Sampling from posterior distribution of coefficients beta and variance sigsq



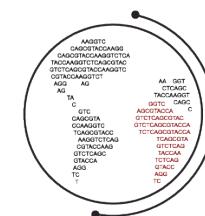
Sampling from posterior distribution of coefficients beta and variance sigsq



Posterior predictive distribution for new data

Suppose we apply the regression model to a new set of data, for which we have observed the matrix $\tilde{\mathbf{X}}$ of explanatory variables, and we wish to predict the outcomes, $\tilde{\mathbf{y}}$.

- If β and σ^2 were known exactly, the vector $\tilde{\mathbf{y}}$ would have a normal distribution with mean $\tilde{\mathbf{X}}\beta$ and variance matrix $\sigma^2\mathbf{I}$.
- Instead, it is summarized by our posterior distribution.



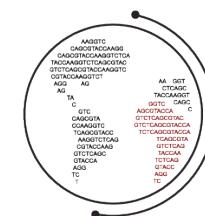
Posterior predictive distribution for new data

Suppose we apply the regression model to a new set of data, for which we have observed the matrix $\tilde{\mathbf{X}}$ of explanatory variables, and we wish to predict the outcomes, $\tilde{\mathbf{y}}$.

- If β and σ^2 were known exactly, the vector $\tilde{\mathbf{y}}$ would have a normal distribution with mean $\tilde{\mathbf{X}}\beta$ and variance matrix $\sigma^2\mathbf{I}$.
- Instead, it is summarized by our posterior distribution.

To draw a random sample $\tilde{\mathbf{y}}$ from its posterior predictive distribution

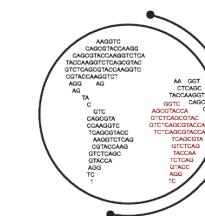
- We first draw (β, σ^2) from the joint posterior distribution
- Then, we draw $\tilde{\mathbf{y}} \sim \mathcal{N}(\tilde{\mathbf{X}}\beta, \sigma^2\mathbf{I})$



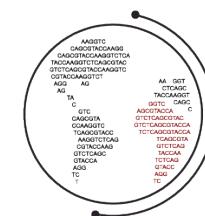
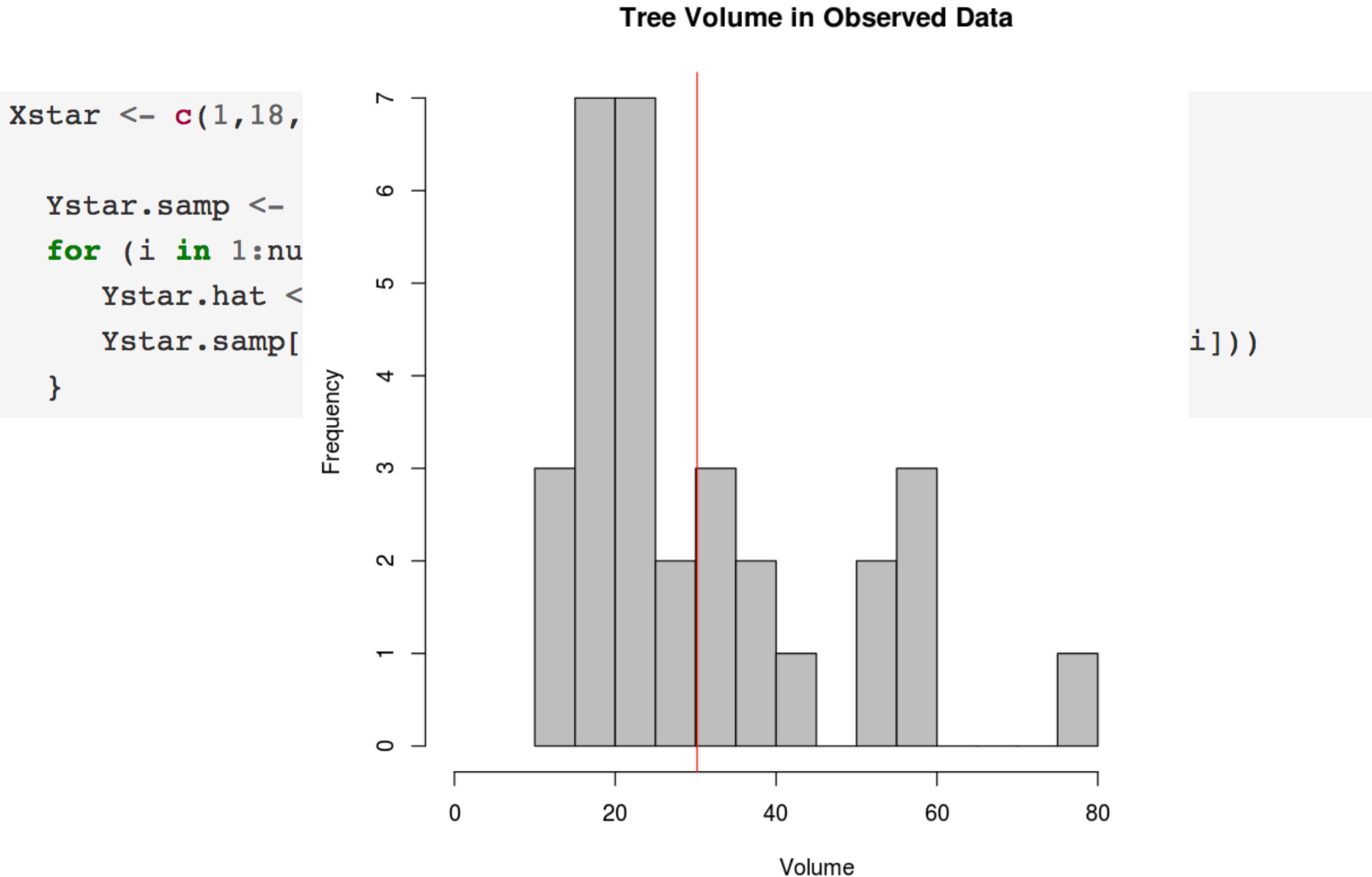
Posterior Predictive Sampling for new tree with girth = 18 and height = 80

```
Xstar <- c(1,18,80) # new tree with girth = 18 and height = 80

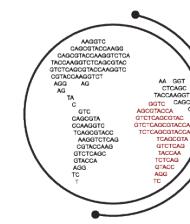
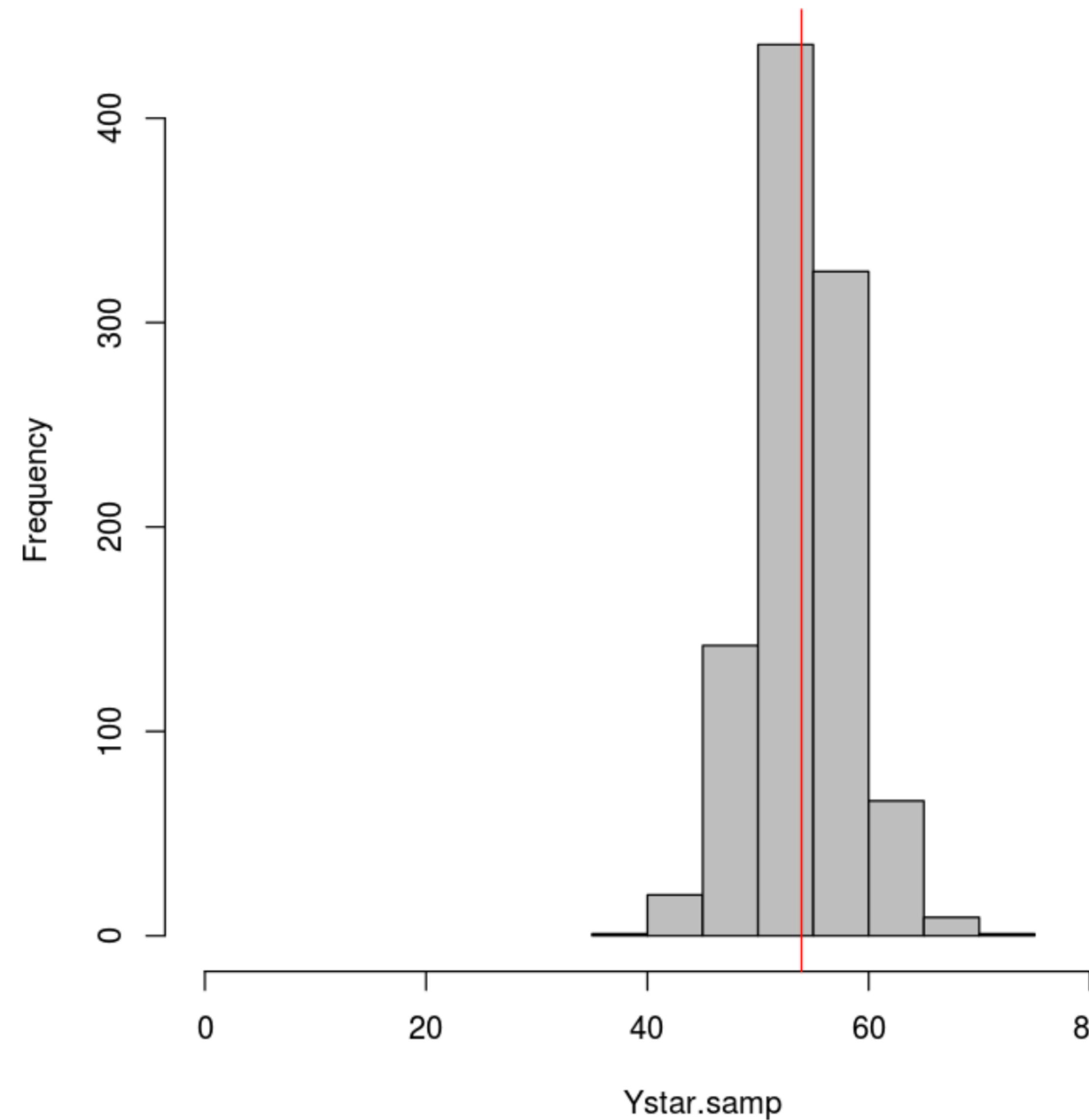
Ystar.samp <- rep(NA,numsamp)
for (i in 1:numsamp){
  Ystar.hat <- sum(beta.samp[i,]*Xstar)
  Ystar.samp[i] <- rnorm(1,mean=Ystar.hat,sd=sqrt(sigsq.samp[i]))
}
```



Posterior Predictive Sampling for new tree with girth = 18 and height = 80



Predicted Volume of Tree with girth = 85 and height = 80



Predicted Volume of Tree with girth = 85 and height = 80

