

#1

a) According to p153 of the textbook, $\hat{\beta}_0 = (X^T X)^{-1} X^T Y$

$$\begin{aligned} V[\hat{\beta}_0 | \beta, \sigma^2] &= V[(X^T X)^{-1} X^T Y] \\ &= (X^T X)^{-1} X^T V[Y] X (X^T X)^{-1} \\ &= \frac{1}{(X^T X)^2} \sigma^2 \sum_{i=1}^n \frac{x_i^2}{w_i} \\ &= \frac{\sigma^2}{\left(\sum_{i=1}^n x_i^2\right)^2} \cdot \sum_{i=1}^n \frac{x_i^2}{w_i} \end{aligned}$$

b) $P(y_1, \dots, y_n | \sigma^2, \beta) = \prod_{i=1}^n P(y_i | \sigma^2, \beta)$ since y_i 's are indep.

$$\begin{aligned} &= \prod_{i=1}^n (2\pi \sigma^2 / w_i)^{-\frac{1}{2}} \exp\left\{-\frac{w_i}{2\sigma^2} (y_i - \beta x_i)^2\right\} \\ &= \left\{ \prod_{i=1}^n (2\pi \sigma^2 / w_i)^{-\frac{1}{2}} \right\} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n w_i (y_i - \beta x_i)^2\right\} \end{aligned}$$

$$\begin{aligned} \text{Take log: } L &= \left\{ \sum_{i=1}^n -\frac{1}{2} \cdot 2\pi \sigma^2 / w_i \right\} - \frac{1}{2\sigma^2} \sum_{i=1}^n w_i (y_i - \beta x_i)^2 \\ &= -\pi \sigma^2 \sum_{i=1}^n \frac{1}{w_i} - \frac{1}{2\sigma^2} \sum_{i=1}^n w_i (y_i - \beta x_i)^2 \end{aligned}$$

$$\begin{aligned} \text{Take deriv: } \frac{dL}{d\beta} &= -\frac{1}{2\sigma^2} \sum_{i=1}^n w_i \cdot 2 (y_i - \beta x_i) (-x_i) \\ &= \frac{1}{\sigma^2} \sum_{i=1}^n w_i x_i (y_i - \beta x_i) \end{aligned}$$

$$\begin{aligned} \text{Set to 0: } \frac{dL}{d\beta} = 0 &\Rightarrow \sum_{i=1}^n w_i x_i y_i - \sum_{i=1}^n w_i x_i \beta x_i = 0 \\ \hat{\beta}_n &= \frac{\sum_{i=1}^n w_i x_i y_i}{\sum_{i=1}^n w_i x_i^2} \end{aligned}$$

b) continued

$$\begin{aligned}
 V[\hat{\beta}_m | \beta, \sigma^2] &= \left(\frac{1}{\sum_{i=1}^n w_i x_i^2} \right)^2 V\left[\sum_{i=1}^n w_i x_i y_i \right] \\
 &= \left(\frac{1}{\sum_{i=1}^n w_i x_i^2} \right)^2 \sum_{i=1}^n w_i^2 x_i^2 V[y_i] \\
 &= \frac{1}{\sum_{i=1}^n w_i^2 x_i^4} \sum_{i=1}^n w_i^2 x_i^2 \frac{\sigma^2}{w_i} \\
 &= \frac{\sigma^2}{\sum_{i=1}^n w_i x_i^2}
 \end{aligned}$$

Cauchy-Schwarz inequality: $|x \cdot y| \leq \|x\| \cdot \|y\|$

$$\frac{V[\hat{\beta}_n]}{V[\hat{\beta}_m]} = \frac{\left(\sum_{i=1}^n \frac{x_i^2}{w_i} \right) \left(\sum_{i=1}^n w_i x_i^2 \right)}{\left(\sum_{i=1}^n x_i^2 \right)^2}$$

Let $x = \left[\left(\frac{x_1^2}{w_1} \right)^{\frac{1}{2}}, \dots, \left(\frac{x_n^2}{w_n} \right)^{\frac{1}{2}} \right]^T$, then $\|x\| = \sum_{i=1}^n \frac{x_i^2}{w_i}$

$y = \left[(w_1 x_1^2)^{\frac{1}{2}}, \dots, (w_n x_n^2)^{\frac{1}{2}} \right]^T$, then $\|y\| = \sum_{i=1}^n w_i x_i^2$

$$|x \cdot y| = (x^T y)^2 = [x_1^2, \dots, x_n^2]^2 = \left(\sum_{i=1}^n x_i^2 \right)^2$$

$$\left(\sum_{i=1}^n \frac{x_i^2}{w_i} \right) \left(\sum_{i=1}^n w_i x_i^2 \right) \geq \left(\sum_{i=1}^n x_i^2 \right)^2$$

Therefore, $V[\hat{\beta}_n] \geq V[\hat{\beta}_m]$

$$c) \quad p(\beta | y_1, \dots, y_n, \sigma^2) = \frac{p(y_1, \dots, y_n | \beta, \sigma^2) \cdot p(\beta | \sigma^2)}{p(y_1, \dots, y_n)}$$

$$\propto_{\beta} p(y_1, \dots, y_n | \beta, \sigma^2) \cdot p(\beta | \sigma^2)$$

$$\propto_{\beta} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n w_i (y_i - \beta x_i)^2 - \frac{1}{2\tau^2} \beta^2 \right\}$$

$$= \exp \left\{ -\frac{1}{2\sigma^2} \left(\sum_{i=1}^n w_i y_i^2 - 2\beta \sum_{i=1}^n w_i y_i x_i + \beta^2 \sum_{i=1}^n w_i x_i^2 \right) - \frac{1}{2\tau^2} \beta^2 \right\}$$

$$\propto_{\beta} \exp \left\{ -\frac{1}{2} \left[\beta^2 \left(\frac{\sum_{i=1}^n w_i x_i^2}{\sigma^2} + \frac{1}{\tau^2} \right) - \beta \sum_{i=1}^n w_i y_i x_i / \sigma^2 \right] \right\}$$

$$\sim N\left(\frac{b}{a}, \frac{1}{a}\right) \quad \text{where}$$

$$b = \sum_{i=1}^n w_i y_i x_i / \sigma^2 \quad \text{and} \quad a = \sum_{i=1}^n w_i x_i^2 / \sigma^2 + \frac{1}{\tau^2}$$

$$\text{Therefore, } E[\beta | y_1, \dots, y_n, \sigma^2] = \frac{b}{a}$$

$$= \frac{\sum_{i=1}^n w_i y_i x_i / \sigma^2}{\sum_{i=1}^n w_i x_i^2 / \sigma^2 + \frac{1}{\tau^2}}$$

$$\lim_{\tau \rightarrow \infty} = \frac{\sum_{i=1}^n w_i y_i x_i}{\sum_{i=1}^n w_i x_i^2}$$

$$= \hat{\beta}_m$$

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

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

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

$$\propto \exp \left\{ -\frac{1}{2} (y - X\beta)^T (\sigma^2 I)^{-1} (y - X\beta) \right\} \exp \left\{ -\frac{1}{2} \beta^T \left(\frac{\sigma^2}{\lambda} \right)^{-1} \beta \right\}$$

$$\propto \exp \left\{ -\frac{1}{2\sigma^2} (\beta^T X^T X \beta - 2\beta^T X^T y) - \frac{\lambda}{2\sigma^2} \beta^T \beta \right\}$$

$$= \exp \left\{ -\frac{1}{2\sigma^2} \beta^T (X^T X + I\lambda) \beta + \frac{1}{\sigma^2} \beta^T X^T y \right\}$$

$$\sim N_p \left(\frac{b}{a}, \frac{1}{a} \right) \text{ where}$$

$$b = X^T y / \sigma^2 \text{ and}$$

$$a = (X^T X + I\lambda) / \sigma^2$$

$$\text{Therefore, } E[\beta | y, \sigma^2] = \frac{b}{a}$$

$$= \frac{X^T y / \sigma^2}{(X^T X + I\lambda) / \sigma^2}$$

$$= (X^T X + I\lambda)^{-1} X^T y = \hat{\beta}_\lambda$$

$$\lim_{\lambda \rightarrow 0} = (X^T X)^{-1} X^T y = \hat{\beta}_0$$

$$b) \quad X^T X = \begin{bmatrix} x_1^T x_1 & 0 & \dots & 0 \\ 0 & & & \\ & & & \\ 0 & & & x_p^T x_p \end{bmatrix}, \quad \hat{\beta}_0 = (X^T X)^{-1} X^T y, \quad \hat{\beta}_{0i} = (x_i^T x_i)^{-1} x_i^T y$$

$$X^T X + I\lambda = \begin{bmatrix} x_1^T x_1 + \lambda & 0 & \dots & 0 \\ 0 & & & \\ & & & \\ 0 & & & x_p^T x_p + \lambda \end{bmatrix}, \quad \hat{\beta}_\lambda = (X^T X + I\lambda)^{-1} X^T y, \quad \hat{\beta}_{\lambda i} = (x_i^T x_i + \lambda)^{-1} x_i^T y$$

$$\frac{\hat{\beta}_{\lambda i}}{\hat{\beta}_{0i}} = \frac{(x_i^T x_i + \lambda)^{-1} x_i^T y}{(x_i^T x_i)^{-1} x_i^T y} = \frac{x_i^T x_i + \lambda}{x_i^T x_i}$$

$$\hat{\beta}_{\lambda i} = \frac{x_i^T x_i + \lambda}{x_i^T x_i} \hat{\beta}_{0i}$$

Effect of λ : As $\lambda \rightarrow 0$, $\hat{\beta}_{\lambda i} = \hat{\beta}_{0i}$, thus the regressors are similar.

As $\lambda \rightarrow \infty$, $\hat{\beta}_{\lambda i} = \frac{x_i^T y}{x_i^T x_i + \infty} = 0$, the regressor gives 0 as coefficients for x_i .

Meaning $V[\beta]$ is very small and ≈ 0

STA360 Homework 8 (Ken Ye)

```
library(latex2exp)
library(ggplot2)
library(MASS)
library(ggrepel)
library(mvtnorm)
set.seed(0)
```

Question 3

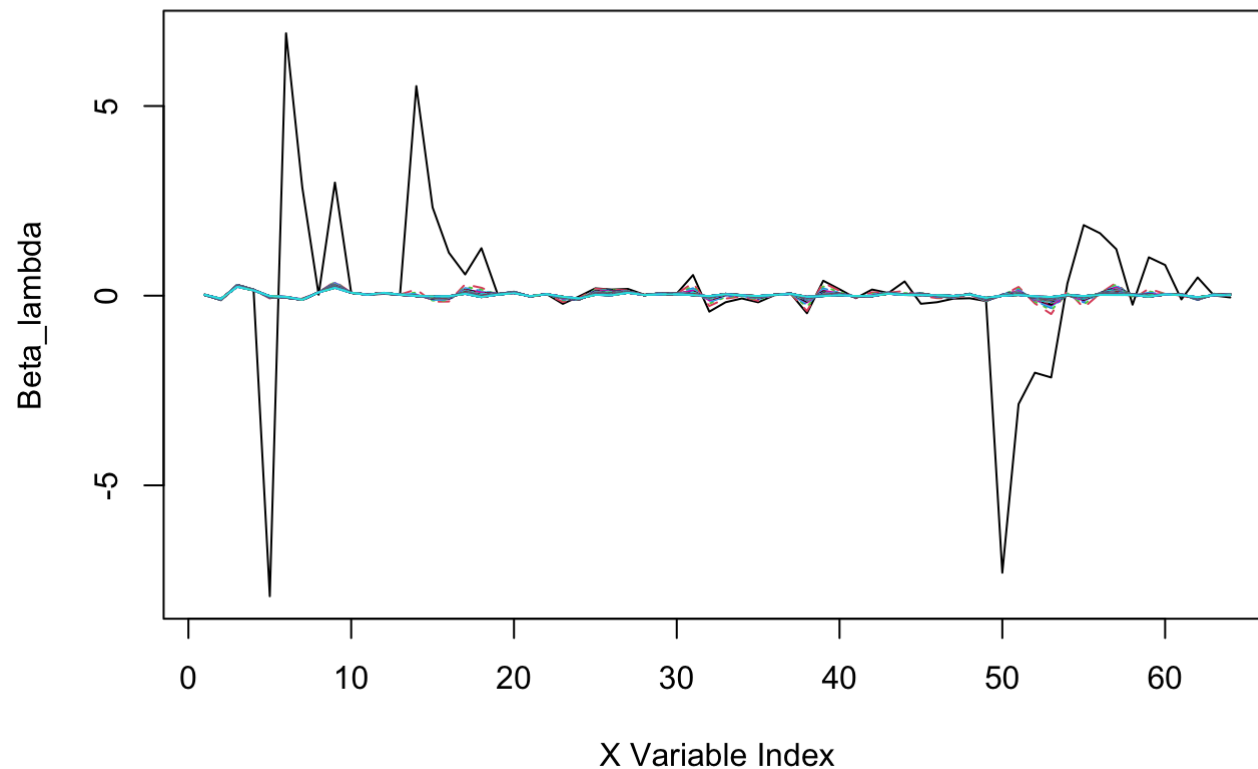
```
# load diabetes data
yX <- dget(url("https://www2.stat.duke.edu/~pdh10/FCBS/Inline/yX.diabetes.train"))
y <- yX[,1]
X <- yX[,-1]
```

Part a

```
# compute beta_lambda hat for each lambda in {0, ..., 100}
lambdas <- seq(0, 100, by = 1)
beta_lambdas <- matrix(0, nrow = 64, ncol = 101)
for (lam in lambdas){
  beta_lambda <- solve(t(X) %*% X + lam * diag(rep(1, 64))) %*% t(X) %*% y
  beta_lambdas[,lam+1] <- beta_lambda
}
```

```
# plot with matplot
matplot(beta_lambdas,
  type = 'l',
  main = 'Beta_lambda vs X Under Different Lambdas',
  ylab = 'Beta_lambda',
  xlab = 'X Variable Index')
```

Beta_lambda vs X Under Different Lambdas



There are 101 lines representing 101 beta_lambda vectors (each 64 by 1) and their respective value for each X variable. It's hard to discern a single beta_lambda vector in the graph because there are so many of them, but this graph shows the beta_lambda estimate for each lambda in $\{0, 1, \dots, 99, 100\}$, and for each beta_lambda, how its beta_lambda_i varies for each X_i, for i from 1 to 64 (there are 64 x variables in total).

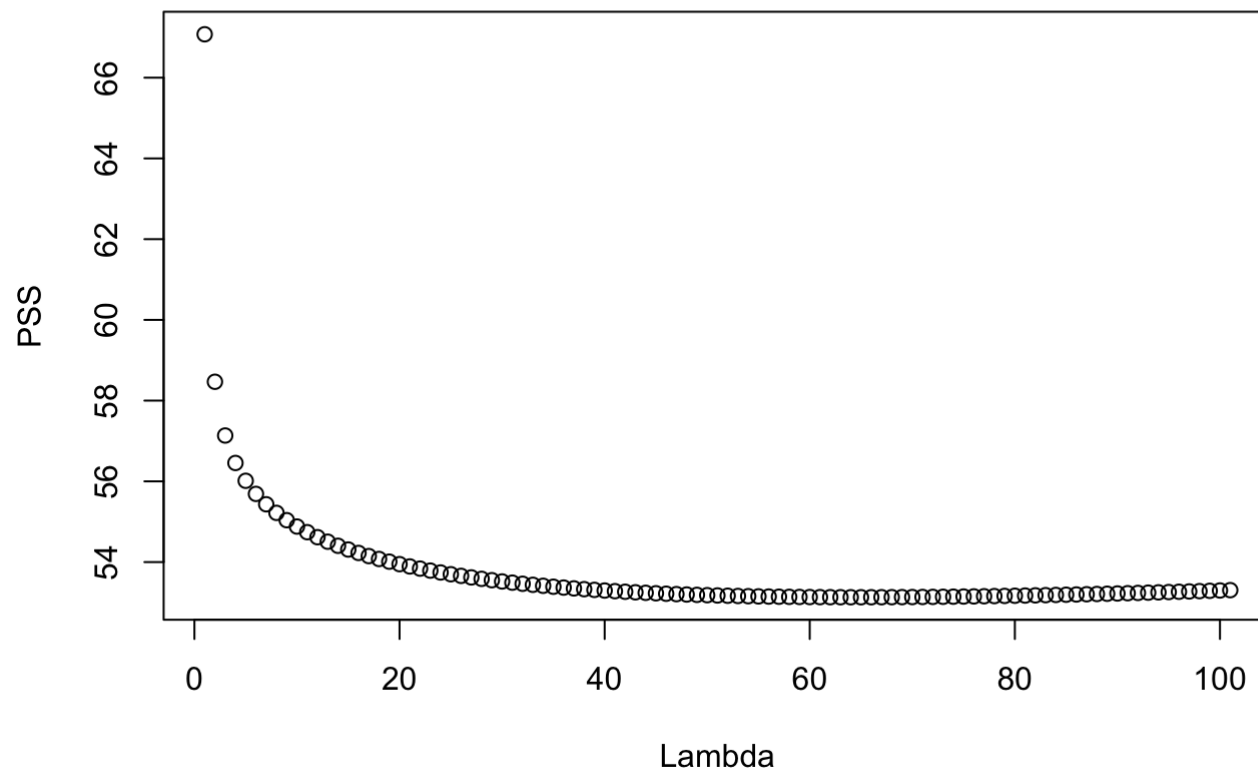
Part b

```
# load diabetes data, test set
yX.test <- dget(url("https://www2.stat.duke.edu/~pdh10/FCBS/Inline/yX.diabetes.test"))
y.test <- yX.test[,1]
X.test <- yX.test[,-1]
```

```
# calculate predictive error sum of squares
PSS <- rep(0,101)
for (i in 1:101){
  PSS[i] <- sum((y.test - (X.test %*% beta_lambdas[,i]))^2)
}
```

```
# plot
plot(PSS,
     main = 'Predictive Error Sum of Squares For Each Lambda',
     ylab = 'PSS',
     xlab = 'Lambda')
```


Predictive Error Sum of Squares For Each Lambda



```
# OLS estimate predictive error sum of squares  
PSS[1]
```

```
## [1] 67.07489
```

We know that $\beta_{OLS} = \beta_{\lambda}$ when $\lambda = 0$. The unbiased OLS estimate for prediction has a predictive error sum of squares of 67.07, which is the highest among all beta estimates (which can be told from the graph). Beta_OLS doesn't perform well.

Part c

```
# identify the value of lambda that has the best predictive performance  
index <- which.min(PSS)  
# this is the index, the value of lambda = index - 1 since lambda starts from 0  
index
```

```
## [1] 65
```

The value of lambda that has the best predictive performance is lambda = 64.

```
# find x-variables that have the largest effects  
beta_lambda.best <- array(beta_lambdas[,65])  
names(beta_lambda.best) <- colnames(X.test)  
sort(beta_lambda.best, decreasing = TRUE)
```

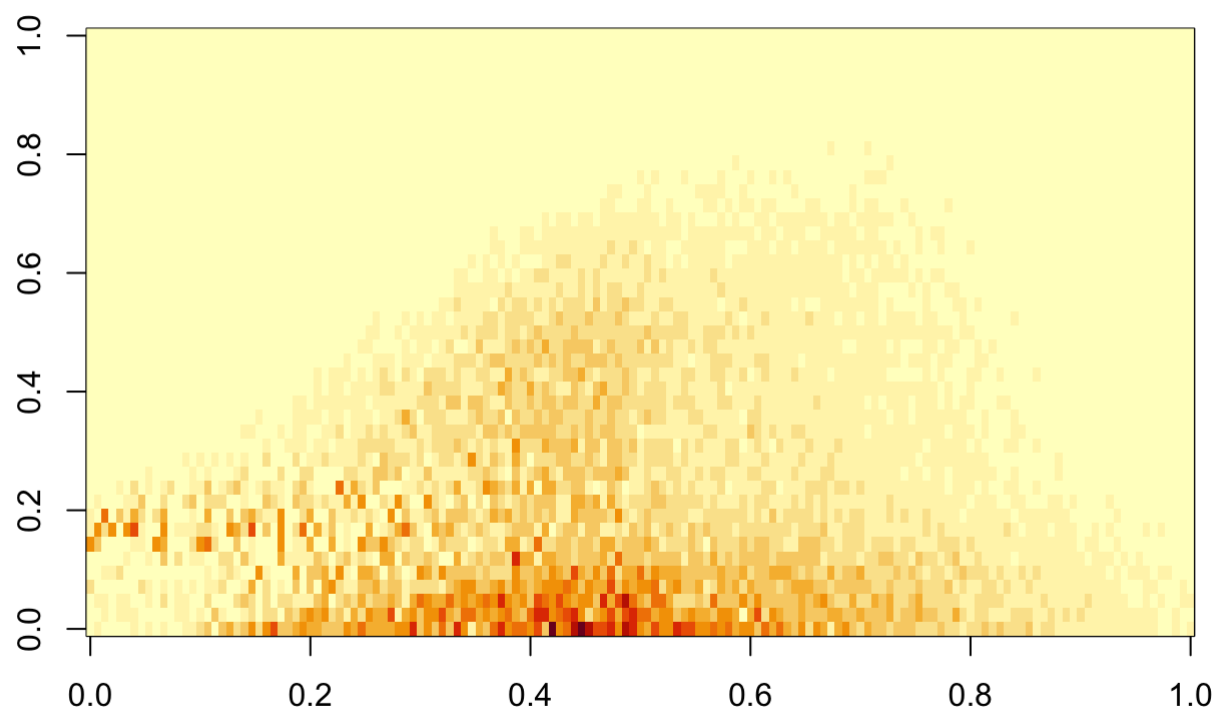
```
##          bmi          ltg          map      age:ltg          tch          glu
## 0.252909731 0.228397916 0.147254764 0.087941087 0.087934685 0.077338474
##      age:sex      bmi^2      tch^2      bmi:glu      bmi:map      sex:hdl
## 0.076324787 0.071152071 0.058732597 0.058276976 0.049170486 0.043500542
##      sex:bmi      sex:map      age:map      ldl:ltg      map:ltg      hdl:ltg
## 0.043490113 0.041124402 0.040800031 0.038342222 0.035184329 0.034370233
##      age:hdl      glu^2      age^2      ldl:tch      map^2      tc:hdl
## 0.033915390 0.031990446 0.028071754 0.026620010 0.024993749 0.024568090
##      map:ldl      sex:glu      ldl:glu          age      map:tc      hdl:glu
## 0.023869197 0.022266772 0.022153006 0.021963242 0.021280595 0.019970286
##      bmi:hdl      sex:tc      ltg:glu      age:tch      tc:glu      age:glu
## 0.018778877 0.018338931 0.017996024 0.017610083 0.017397253 0.016284681
##      tch:glu      map:hdl      sex:tch      tc:ldl      bmi:ldl      tc^2
## 0.014990058 0.009804745 0.009214922 -0.001376226 -0.003008950 -0.007634208
##      map:tch      bmi:ltg      bmi:tch      sex:ltg      age:bmi      ldl:hdl
## -0.011186191 -0.011266451 -0.014018574 -0.014901732 -0.015406394 -0.016465109
##      tc      sex:ldl      hdl^2      ldl^2      hdl:tch      tc:tch
## -0.017435685 -0.017546233 -0.018174450 -0.024323567 -0.025176356 -0.025552851
##      age:tc      ltg^2      bmi:tc      ldl      tc:ltg      tch:ltg
## -0.031782045 -0.034849965 -0.035749440 -0.040005158 -0.047686902 -0.052317906
##      map:glu      age:ldl      sex      hdl
## -0.073321473 -0.080450910 -0.091657677 -0.109873439
```

The x-variables that have the largest effects (top 5) are bmi, ltg, map, age:ltg, and tch. The rest are printed above in decreasing effect order.

Question 4

```
# load water data  
yX <- readRDS("yXSS.rds")  
y <- yX[,1]  
X <- yX[,-1]
```

```
# view image  
y <- yX[,1]  
image(matrix(y,151,43))
```



Part a

```
# obtain posterior distribution of beta and sigma2 given y with MCMC

n <- dim(X)[1]
p <- dim(X)[2]

# priors
nu.0 <- 1
beta.0 <- rep(1/9, p)
sigma.0 <- matrix(0, p, p)
diag(sigma.0) <- 1 # sigma.0 = I_9
sigma2.0 <- diag(p)
sample.size <- 10000
BETA <- matrix(nrow = sample.size, ncol = p) # posterior storage
SIGMA2 <- rep(NA, sample.size) # posterior storage

# common quantities
X <- as.matrix(X)
y <- as.matrix(y)
isigma.0 <- solve(sigma.0)
XtX <- t(X) %*% X
Xty <- t(X) %*% y

# start values
sigma2 <- var(residuals(lm(y ~ 0 + X)))

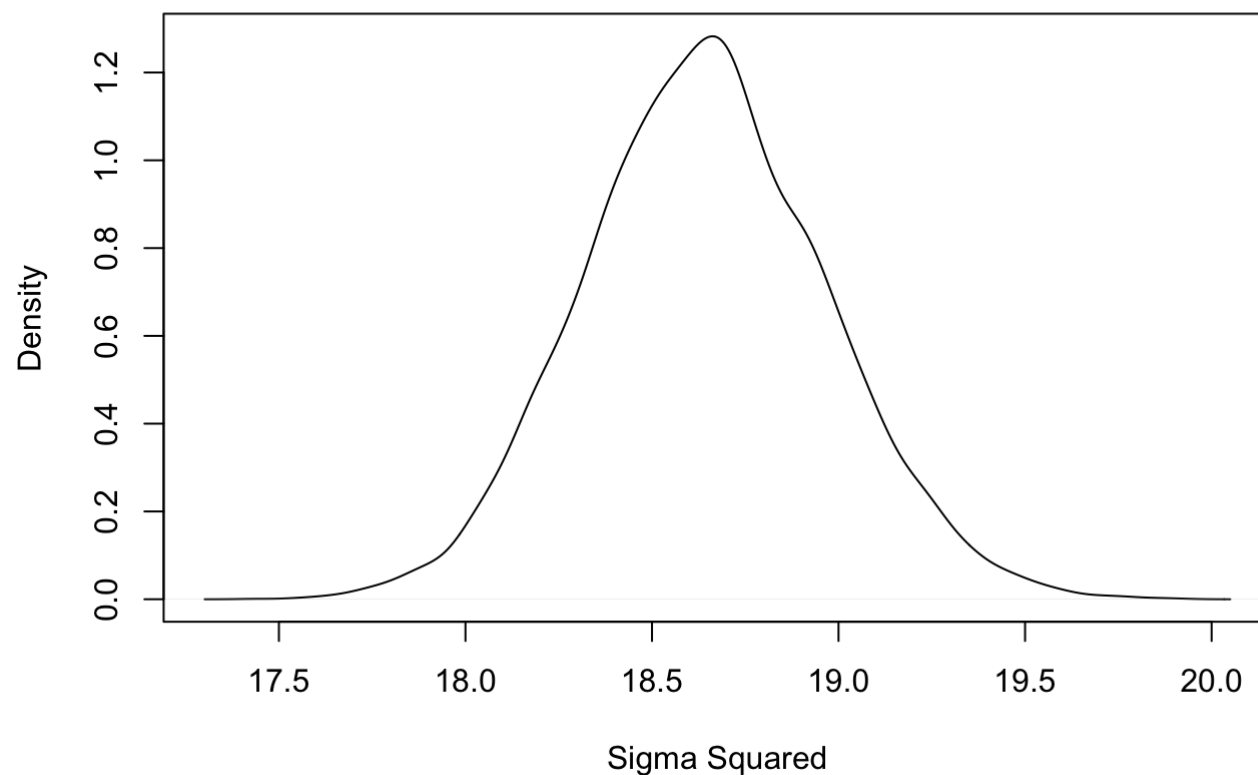
# Gibbs sampling
for (s in 1 : sample.size) {
  # update beta
  beta.V <- solve(isigma.0 + XtX / sigma2)
  beta.E <- beta.V %*% (isigma.0 %*% beta.0 + Xty / sigma2)
  beta <- mvrnorm(1, beta.E, beta.V)

  # update sigma2
  nu.n <- nu.0 + n
  ss.n <- nu.0 * sigma2.0 + sum((y - X %*% beta)^2)
  sigma2 <- 1/rgamma(1, nu.n / 2, ss.n / 2)
```

```
# store sample  
BETA[s,] <- beta  
SIGMA2[s] <- sigma2  
}
```

```
# posterior distribution of sigma2  
plot(density(SIGMA2),  
     main = 'Posterior Distribution of Sigma Squared',  
     xlab = 'Sigma Squared')
```

Posterior Distribution of Sigma Squared



```
# 95% confidence intervals for each element of beta
CI <- NULL
for (i in 1: dim(BETA)[2]){
  print(paste('The 95% CI for Beta', i, "is: "))
  print(quantile(BETA[,i], c(0.0025, 0.975)))
  CI <- c(CI, quantile(BETA[,i], c(0.0025, 0.975)))
}
```

```
## [1] "The 95% CI for Beta 1 is: "
##      0.25%      97.5%
## 0.3794342 0.6672999
## [1] "The 95% CI for Beta 2 is: "
##      0.25%      97.5%
## -0.07597283 0.02150432
## [1] "The 95% CI for Beta 3 is: "
##      0.25%      97.5%
## -0.02184565 0.01643633
## [1] "The 95% CI for Beta 4 is: "
##      0.25%      97.5%
## -0.01622588 0.09902869
## [1] "The 95% CI for Beta 5 is: "
##      0.25%      97.5%
## -0.0315009542 0.0009156621
## [1] "The 95% CI for Beta 6 is: "
##      0.25%      97.5%
## 0.1707828 0.3449928
## [1] "The 95% CI for Beta 7 is: "
##      0.25%      97.5%
## 0.09072786 0.20750827
## [1] "The 95% CI for Beta 8 is: "
##      0.25%      97.5%
## 0.002263687 0.006809795
## [1] "The 95% CI for Beta 9 is: "
##      0.25%      97.5%
## -0.30283697 0.06189633
```

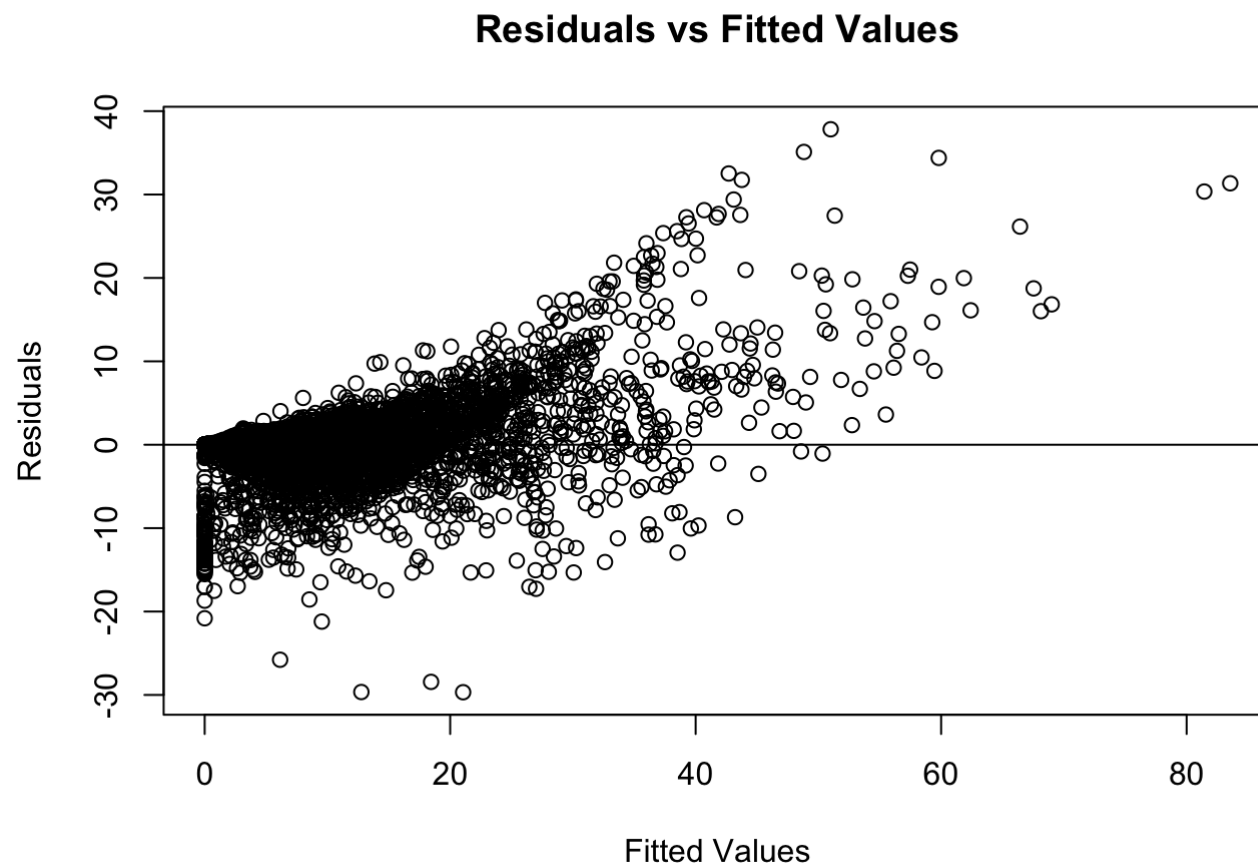
Effluent, soil, street (weak, very small coefficient), and swine are the main sources of the water sample as their 95% CI do not contain 0.

Part b

```
# check constant variance condition

# obtain residual
beta.mean <- apply(BETA, 2, mean)
y.pred <- X %*% beta.mean
residual <- y - y.pred

# residual plot
plot(y, residual,
     main = 'Residuals vs Fitted Values',
     ylab = 'Residuals',
     xlab = 'Fitted Values')
abline(0,0)
```

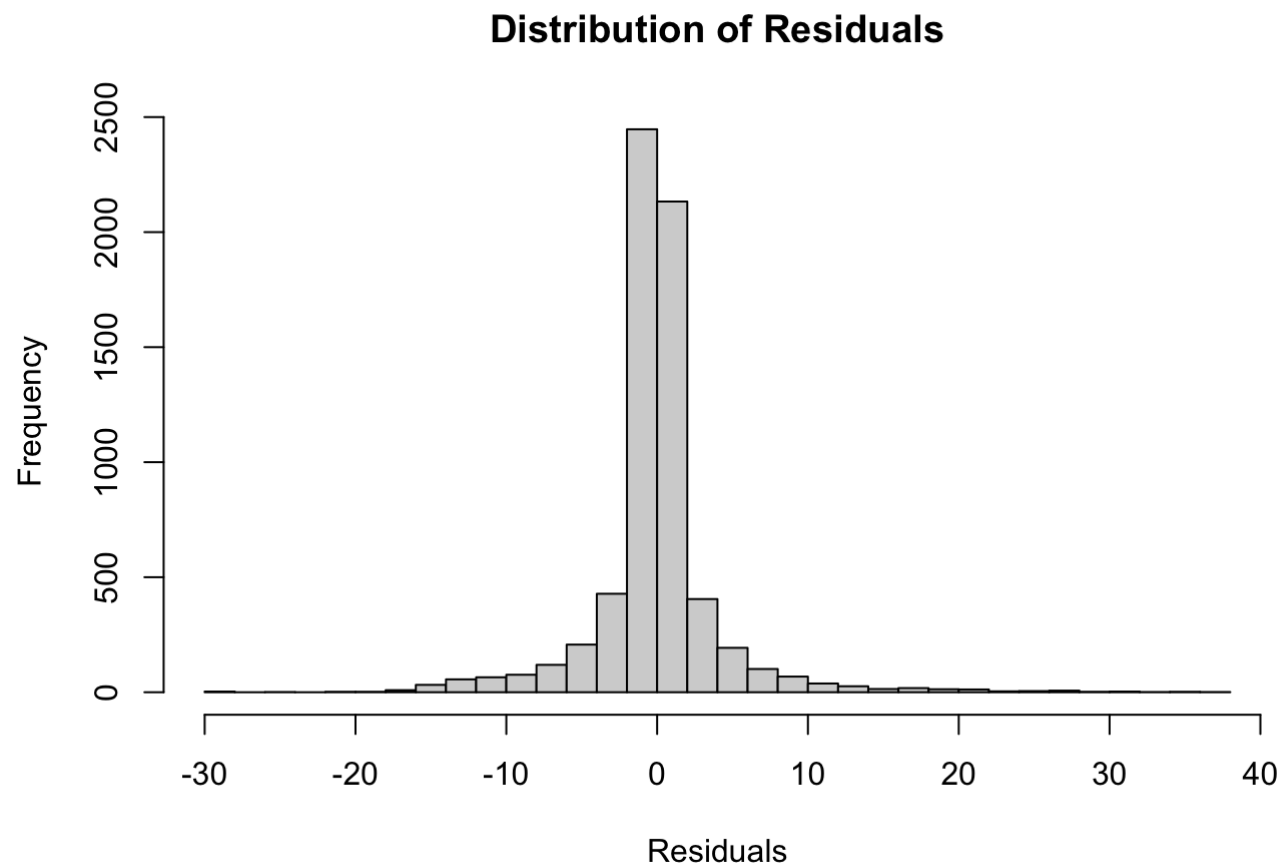


The constant variance condition seems to be violated as the residuals are not randomly spread across all fitted values of y . In fact, as the fitted value increases, the residual increases.

Check Independence Condition

The independence condition seems to be satisfied since y is the vectorization of a spectroscopy image of a water sample taken from the Neuse River in North Carolina. It is reasonable to assume the entries of the vectors are uncorrelated.

```
# check normal condition  
# distribution of y  
hist(residual,  
      main = 'Distribution of Residuals',  
      xlab = 'Residuals',  
      breaks = 30)
```



The distribution of the residuals approximately normal, and the sample size is larger than 30, thus the normal condition is satisfied.

Part c

Since it doesn't make sense for the coefficients of beta to be negative, a modification to the prior distribution for beta could be picking a distribution with a positive support, such as beta or gamma.

Suppose we use beta distribution as the prior distribution for beta. We first need to derive the posterior distribution of beta. For the Gibbs sampler, with starting values for σ^2 and beta, in each iteration we update and get a new beta from the posterior calculation formula, and get a new σ^2 (just as we did in part a). In the end, we would get a storage vector BETA of all positive values, since beta distribution lies in the first quadrant.