

STATS 790 Assignment 2

Yiran Zhang
400119421

11 February, 2023

Question 1

(a) In linear regression, the assumed model is $y = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p + \epsilon$, it can be written in matrix form as $\vec{Y} = X\vec{\beta} + \vec{\epsilon}$, where $\vec{\beta}$ is the vector of coefficients (dimension is $p \times 1$, or $(p+1) \times 1$ if consider β_0), Y is the $n \times 1$ vector of responses, X is an $n \times p$ matrix of predictors (or $n \times (p+1)$ if add one column of 1's for β_0), it is also assumed to be full rank. The estimator of \vec{Y} is $\hat{\vec{Y}} = X\hat{\vec{\beta}}$.

Naive Linear Algebra Use the least square estimation:

$$\begin{aligned} RSS(\vec{\beta}) &= \sum_{i=1}^n \epsilon_i^2 \\ &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ &= (\vec{Y} - \hat{\vec{Y}})^T (\vec{Y} - \hat{\vec{Y}}) \\ &= (\vec{Y} - X\vec{\beta})^T (\vec{Y} - X\vec{\beta}) \end{aligned}$$

We want to minimize the residual sum of square value, therefore we take the first derivative of it and set it to 0.

$$\begin{aligned} \frac{\partial RSS(\vec{\beta})}{\partial \vec{\beta}} &= \frac{\partial (\vec{Y} - X\vec{\beta})^T (\vec{Y} - X\vec{\beta})}{\partial \vec{\beta}} \\ &= -2X^T \vec{Y} + 2X^T X \vec{\beta} \\ &= 0 \end{aligned}$$

Therefore, $X^T \vec{Y} = X^T X \vec{\beta}$, solve for $\vec{\beta}$ and get that the coefficients are $\hat{\vec{\beta}} = (X^T X)^{-1} X^T \vec{Y}$. We can also compute the second partial derivative to check if it is minimum: $\frac{\partial^2 RSS(\vec{\beta})}{\partial \vec{\beta}^2} = 2X^T X > 0$. Hence the least square is minimized when $\hat{\vec{\beta}} = (X^T X)^{-1} X^T \vec{Y}$.

QR Decomposition Based on the definition of QR Decomposition, since X is defined as full rank, it can be written as the form $X=QR$, Q is an $n \times (p+1)$ orthogonal matrix (i.e., $Q^T = Q^{-1}$), and R is an $(p+1) \times (p+1)$ invertible upper triangular matrix (Anton & Rorres, 2013). Then we can substitute X in the above equation with Q and R .

$$\begin{aligned} X^T X \vec{\beta} &= X^T \vec{Y} \\ (QR)^T QR \vec{\beta} &= (QR)^T \vec{Y} \\ R^T Q^T QR \vec{\beta} &= R^T Q^T \vec{Y} \\ R^T R \vec{\beta} &= R^T Q^T \vec{Y} \\ R \vec{\beta} &= Q^T \vec{Y} \end{aligned}$$

To obtain $\vec{\beta}$ estimate, use backward substitution since R is upper triangular, and $Q^T \vec{Y}$ is a $(p+1) \times 1$ vector.

SVD X is full rank, therefore we can apply singular value decomposition to X by writing it in the form $X = U\Sigma V^T$, where U is $n \times n$ orthogonal matrix ($U^T U = I$), Σ is $n \times (p+1)$ matrix, V is $(p+1) \times (p+1)$ orthogonal matrix ($V^T V = I$) (Anton & Rorres, 2013).

Then, solve for $X \vec{\beta} = \vec{Y}$

$$\begin{aligned} U\Sigma V^T \vec{\beta} &= \vec{Y} \\ U^T U\Sigma V^T \vec{\beta} &= U^T \vec{Y} \\ \Sigma V^T \vec{\beta} &= U^T \vec{Y} \\ V^T \vec{\beta} &= \Sigma^{-1} U^T \vec{Y} \\ VV^T \vec{\beta} &= V\Sigma^{-1} U^T \vec{Y} \\ \vec{\beta} &= V\Sigma^{-1} U^T \vec{Y} \end{aligned}$$

The coefficients are $\vec{\beta} = V\Sigma^{-1} U^T \vec{Y}$, note that Σ^{-1} is a pseudo inverse of Σ .

Cholesky Decomposition Cholesky decomposition states that a positive definite matrix can be written as the product of a lower triangular matrix and its transpose. X is full rank, hence $X^T X$ is positive definite, then it can be written as $X^T X = LL^T$ where L is lower triangular.

$$\begin{aligned}
X\vec{\beta} &= \vec{Y} \\
X^T X \vec{\beta} &= X^T \vec{Y} \\
LL^T \vec{\beta} &= X^T \vec{Y} \\
L(L^T \vec{\beta}) &= X^T \vec{Y}
\end{aligned}$$

To solve for $\vec{\beta}$, we can solve for $L^T \vec{\beta}$ first from the equation $L(L^T \vec{\beta}) = X^T \vec{Y}$ using forward substitution, then solve for $\vec{\beta}$ from $L^T \vec{\beta}$ using backward substitution.

(b) Assume that the predictors X and responses Y are in separate matrices.

Naive Linear Algebra:

```
beta_naive <- function(X,Y){
  # X is the matrix of predictors.
  # Y is the vector of responses
  X_new <- cbind(1, X) # add a column of 1 for beta_0
  solve(t(X_new) %*% X_new) %*% t(X_new) %*% Y # beta estimator
}
```

QR Decomposition:

```
beta_qr <- function(X,Y){
  # X is the matrix of predictors.
  # Y is the vector of responses
  X_new <- cbind(1, X) # add a column of 1 for beta_0
  qr_decomp <- qr(X_new) # Obtain the QR decomposition of X
  Q <- qr.Q(qr_decomp) # Q matrix
  R <- qr.R(qr_decomp) # R matrix
  backsolve(R, t(Q) %*% Y)
}
```

Cholesky Decomposition:

```

beta_chol <- function(X,Y){
  # X is the matrix of predictors.
  # Y is the vector of responses
  X_new <- cbind(1, X) # add a column of 1 for beta_0
  chol_target <- t(X_new)%*%X_new # prepare X^T*X for Cholesky Decomposition
  L <- t(chol(chol_target)) # lower triangular matrix L
  # solve for L^T*beta, L(L^T*beta) = X^T*Y
  L_trans_beta <- forwardsolve(L, t(X_new) %*% Y, upper.tri = FALSE)
  backsolve(t(L), L_trans_beta, upper.tri = TRUE) # solve for beta
}

```

Benchmark:

1. Use $p = 10$ and a range of n values.

```

# Import benchmark library
library(microbenchmark)

p <- 10
n_vec <- round(100*seq(3, 7, by = 0.25)) # a list of n values

# Store the time for each n value
output_table <- matrix(NA, ncol = 4, nrow = length(n_vec))
colnames(output_table) <- c("n", "naive_time", "qr_time", "cholesky_time")
index <- 1

# Run benchmark for different n values
for (n in n_vec){
  X <- matrix(rnorm(n*p), nrow=n, ncol=p) # Simulate X matrix
  Y <- matrix(rnorm(n), ncol=1) # Simulate Y vector
  # Benchmark the three algorithms
  benchmark_result <- microbenchmark(beta_naive(X, Y),
                                     beta_qr(X, Y),
                                     beta_chol(X, Y))

  result <- summary(benchmark_result)
  # Store the mean computational time in the matrix
  output_table[index, ] = c(n, result$mean[1],

```

```

        result$mean[2],
        result$mean[3])

    index <- index + 1
}

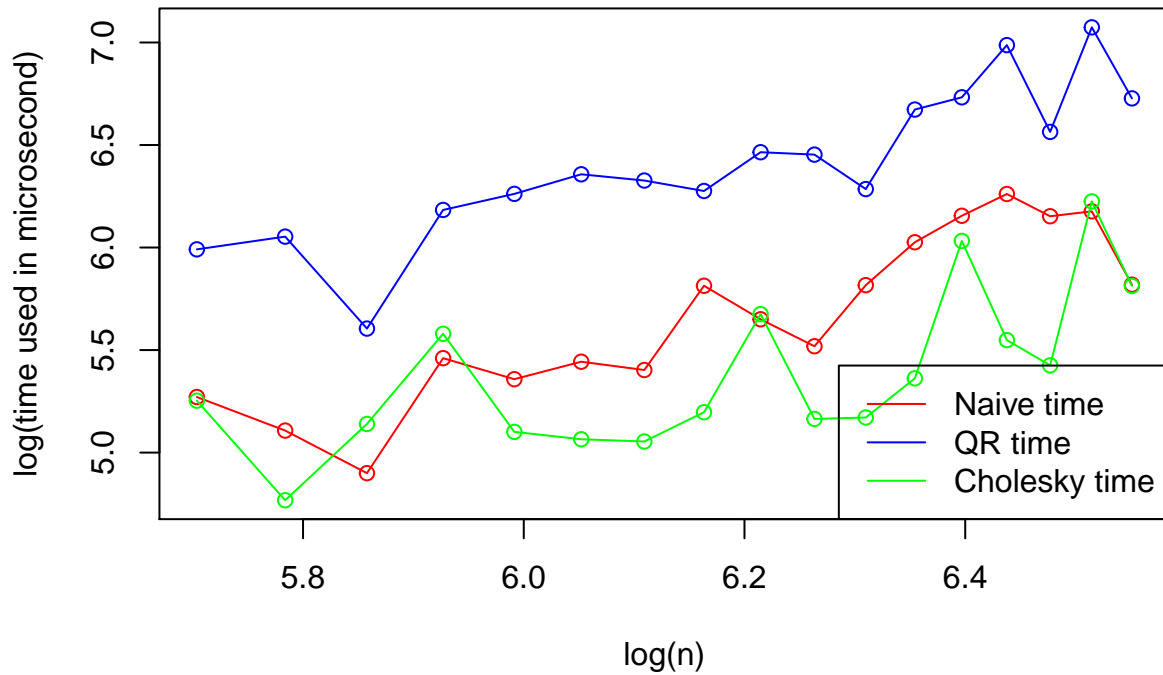
# To plot in log-log scale, take the log of entire table
target <- log(output_table)

# Find limit of y axis for better visualization
ylim_min <- min(target[,2], target[,3], target[,4])
ylim_max <- max(target[,2], target[,3], target[,4])

# Create line graph
plot(target[, 1], target[, 2],
      type = "o", ylim=c(ylim_min, ylim_max), col='red',
      xlab = 'log(n)', ylab = 'log(time used in microsecond)',
      main = 'n VS time in log-log scale (p = 10)')
lines(target[, 1], target[, 3], type = "o", col = 'blue')
lines(target[, 1], target[, 4], type = "o", col = 'green')
legend(x = 'bottomright',
      legend=c("Naive time", "QR time", "Cholesky time"),
      col=c("red", "blue", "green"), lty = 1)

```

n VS time in log-log scale (p = 10)



2. Use $p = 20$ and a range of n values.

```
p <- 20
n_vec <- round(100*seq(3, 7, by = 0.25)) # a list of n values

# Store the time for each n value
output_table <- matrix(NA, ncol = 4, nrow = length(n_vec))
colnames(output_table) <- c("n", "naive_time", "qr_time", "cholesky_time")
index <- 1

# Run benchmark for different n values
for (n in n_vec){
  X <- matrix(rnorm(n*p), nrow=n, ncol=p) # Simulate X matrix
  Y <- matrix(rnorm(n), ncol=1) # Simulate Y vector
  # Benchmark the three algorithms
  benchmark_result <- microbenchmark(beta_naive(X, Y),
                                     beta_qr(X, Y),
```

```

                                beta_chol(X, Y))

result <- summary(benchmark_result)
# Store the mean computational time in the matrix
output_table[index, ] = c(n, result$mean[1],
                           result$mean[2],
                           result$mean[3])

index <- index + 1
}

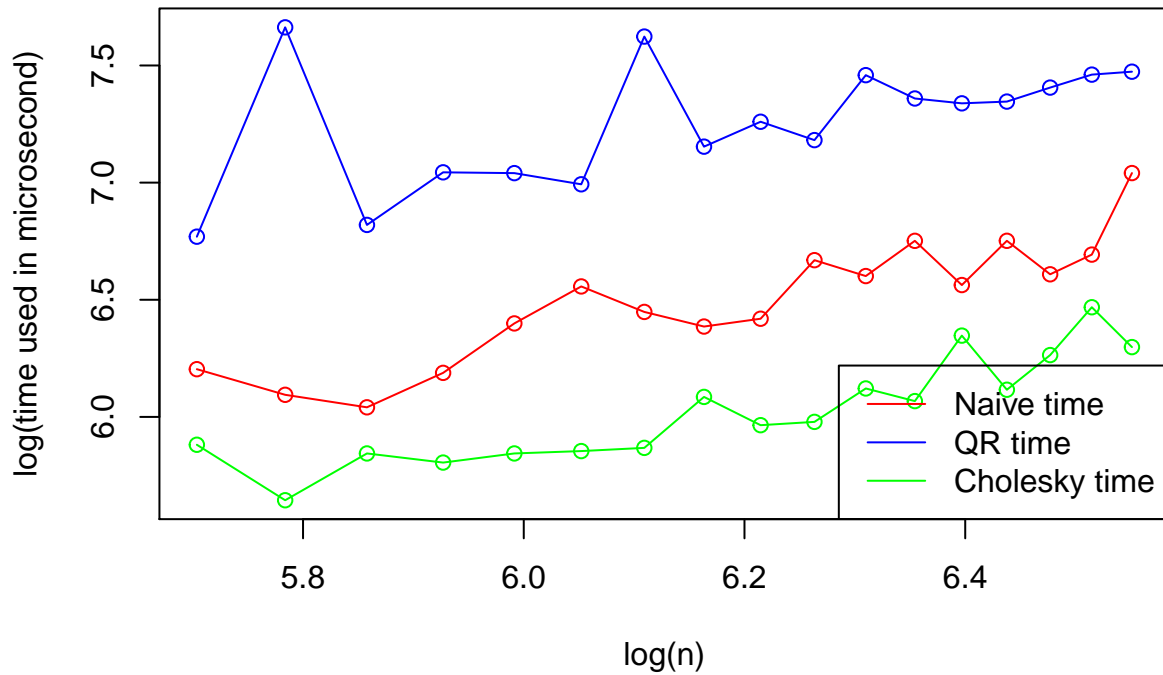
# To plot in log-log scale, take the log of entire table
target <- log(output_table)

# Find limit of y axis for better visualization
ylim_min <- min(target[,2], target[,3], target[,4])
ylim_max <- max(target[,2], target[,3], target[,4])

# Create line graph
plot(target[, 1], target[, 2],
      type = "o", ylim=c(ylim_min, ylim_max), col='red',
      xlab = 'log(n)', ylab = 'log(time used in microsecond)',
      main = 'n VS time in log-log scale (p = 20)')
lines(target[, 1], target[, 3], type = "o", col = 'blue')
lines(target[, 1], target[, 4], type = "o", col = 'green')
legend(x = 'bottomright',
      legend=c("Naive time", "QR time", "Cholesky time"),
      col=c("red", "blue", "green"), lty = 1)

```

n VS time in log-log scale (p = 20)



3. Use $p = 50$ and a range of n values.

```
p <- 50
n_vec <- round(100*seq(3, 7, by = 0.25)) # a list of n values

# Store the time for each n value
output_table <- matrix(NA, ncol = 4, nrow = length(n_vec))
colnames(output_table) <- c("n", "naive_time", "qr_time", "cholesky_time")
index <- 1

# Run benchmark for different n values
for (n in n_vec){
  X <- matrix(rnorm(n*p), nrow=n, ncol=p) # Simulate X matrix
  Y <- matrix(rnorm(n), ncol=1) # Simulate Y vector
  # Benchmark the three algorithms
  benchmark_result <- microbenchmark(beta_naive(X, Y),
                                     beta_qr(X, Y),
```



```

                                beta_chol(X, Y))

result <- summary(benchmark_result)
# Store the mean computational time in the matrix
output_table[index, ] = c(n, result$mean[1],
                           result$mean[2],
                           result$mean[3])

index <- index + 1
}

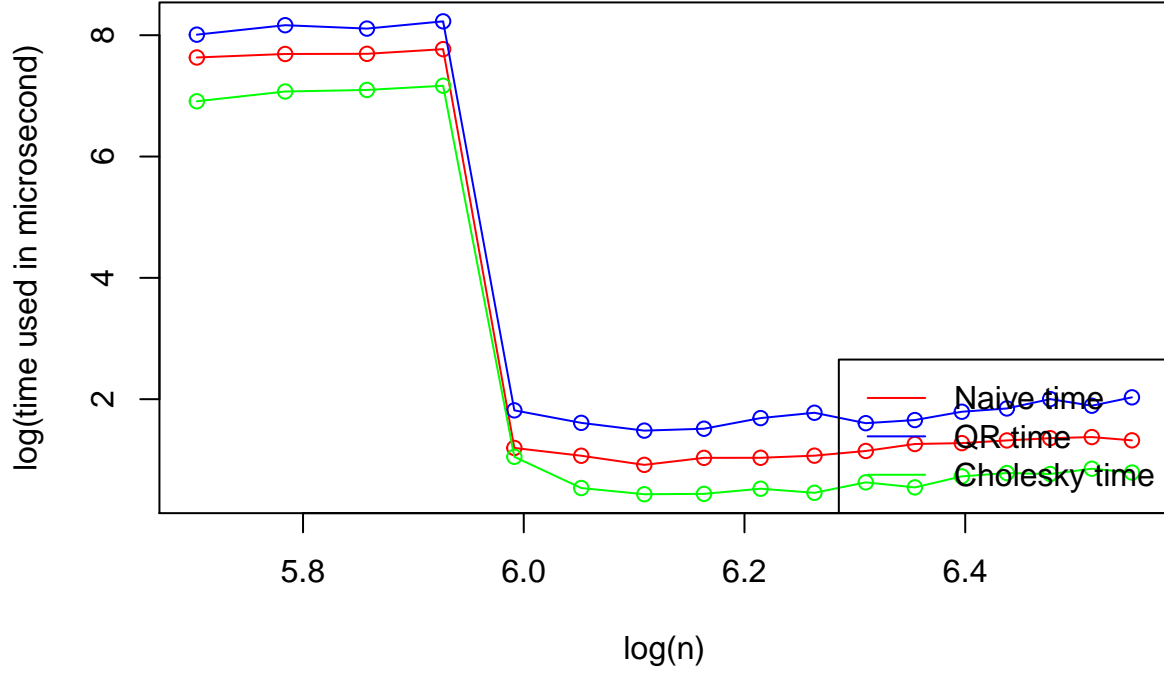
# To plot in log-log scale, take the log of entire table
target <- log(output_table)

# Find limit of y axis for better visualization
ylim_min <- min(target[,2], target[,3], target[,4])
ylim_max <- max(target[,2], target[,3], target[,4])

# Create line graph
plot(target[, 1], target[, 2],
      type = "o", ylim=c(ylim_min, ylim_max), col='red',
      xlab = 'log(n)', ylab = 'log(time used in microsecond)',
      main = 'n VS time in log-log scale (p = 50)')
lines(target[, 1], target[, 3], type = "o", col = 'blue')
lines(target[, 1], target[, 4], type = "o", col = 'green')
legend(x = 'bottomright',
      legend=c("Naive time", "QR time", "Cholesky time"),
      col=c("red", "blue", "green"), lty = 1)

```

n VS time in log-log scale (p = 50)



We benchmark the three algorithms using microbenchmark command for $p = 10, 20, 50$ respectively and using n from the range (300, 700). The log-log scale plots above shows that the algorithm using QR Decomposition has the highest computational time, linear regression with naive linear algebra has the second highest computational time, and linear regression with Cholesky Decomposition has the lowest computational time. The general trends show that the computational time increases as n increases, and increases as p increases.

(c) Time complexity for naive linear algebra is said to be $O(np^2 + p^3)$, for QR decomposition is $O(np^2)$, for SVD is $O(n^2p)$, and for Choleskey decomposition is $O(np^2)$.

Question 2

The formula for ridge regression is $\hat{\beta}^{ridge} = \operatorname{argmin}[\sum_{i=1}^n (y_i - \beta_0 - \sum_{j=1}^p x_{ij}\beta_j)^2 + \lambda \sum_{j=1}^p \beta_j^2]$. In matrix form $\hat{\beta}^{ridge} = (\vec{Y} - X\vec{\beta})^T(\vec{Y} - X\vec{\beta}) + \lambda \vec{\beta}^T \vec{\beta}$ (Hastie et al., 2009).

Minimize the RSS function by taking the first derivative and set to 0:

$$\begin{aligned}\frac{\partial RSS(\vec{\beta}^{ridge})}{\partial \vec{\beta}^{ridge}} &= \frac{\partial (\vec{Y} - X\vec{\beta})^T (\vec{Y} - X\vec{\beta}) + \lambda \vec{\beta}^T \vec{\beta}}{\partial \vec{\beta}} \\ &= -2X^T \vec{Y} + 2X^T X \vec{\beta} + 2\lambda \vec{\beta} \\ &= 0\end{aligned}$$

$\vec{\beta}^{ridge} = (X^T X + \lambda I_p)^{-1} X^T \vec{Y}$, we will implement this estimation into the R code below.

```
ridge_augmented <- function(X, Y, lambda=0){
  X <- cbind(1, X) # add one column of 1
  p <- ncol(X)
  solve(t(X) %*% X + lambda * diag(p)) %*% t(X) %*% Y # compute the coefficients
}
```

Apply the above algorithm and the glmnet package algorithm to the prostate cancer dataset to compare the output and timing.

```
# Import cancer dataset.
cancer <- read.table('https://hastie.su.domains/ElemStatLearn/datasets/prostate.data')

# Split train and test set.
train_cancer <- cancer[cancer$train == TRUE, ]
test_cancer <- cancer[cancer$train == FALSE, ]

# Split the design matrix and response vector for train and test set.
x_train <- as.matrix(scale(train_cancer[,1:8])) # scale x as required.
y_train <- train_cancer$lpsa

x_test <- as.matrix(scale(test_cancer[,1:8]))
y_test <- test_cancer$lpsa

# Calculate the coefficients for the newly implemented ridge regression
my_ridge <- ridge_augmented(x_train, y_train, lambda = 0.087)

# Augment the design matrix for test set for prediction
x_test_aug <- cbind(1, x_test)
```

```

# Predict on the test set
my_ridge_pred <- x_test_aug %*% my_ridge

# Calculate the mean squared error
mse_my_ridge <- mean((my_ridge_pred - y_test)^2)
mse_my_ridge

```

```
## [1] 0.5488212
```

```

# Naive implementation by glmnet package in R
library(glmnet)

# Build ridge regression model
ridge_naive <- glmnet(x_train, y_train, alpha=0)

# Predict on the test set
naive_pred <- predict(ridge_naive, x_test)

# Calculate mean squared error
mse_naive <- mean((naive_pred - y_test)^2)
mse_naive

```

```
## [1] 0.7605016
```

```

# Compare the timing between two methods
a <-microbenchmark(my_ridge, ridge_naive, times = 10000)

```

Based on the above outputs, the new implemented ridge regression has less mean squared error than the build-in ridge function in glmnet. The microbenchmark function is ran multiple times on both functions and sometimes the new ridge regression function is faster, sometimes the build-in function is faster, but generally they have similar computational time.

Question 3

ESL 3.6 Given prior $\vec{\beta} \sim N(0, \tau I)$, then $\pi(\vec{\beta}) \propto \exp(-\frac{1}{2\tau} \vec{\beta}^T \vec{\beta})$. Given the sampling model is also Gaussian, $\vec{Y} | \vec{\beta} \sim N(X\vec{\beta}, \sigma^2 I)$, then $\pi(\vec{Y} | \vec{\beta}) \propto \exp(-\frac{1}{2\sigma^2} (\vec{Y} - X\vec{\beta})^T (\vec{Y} - X\vec{\beta}))$.

From Bayesian inference, we calculate the posterior distribution to be:

$$\begin{aligned}
\pi(\vec{\beta}|\vec{Y}) &\propto \pi(\vec{\beta})\pi(\vec{Y}|\vec{\beta}) \\
&= \exp(-\frac{1}{2\tau}\vec{\beta}^T\vec{\beta} - \frac{1}{2\sigma^2}(\vec{Y} - X\vec{\beta})^T(\vec{Y} - X\vec{\beta})) \\
&= \exp\{-\frac{1}{2}(\frac{\vec{\beta}^T\vec{\beta}}{\tau} + \frac{\vec{Y}^T\vec{Y} - 2\vec{Y}^TX\vec{\beta} + \vec{\beta}^TX^TX\vec{\beta}}{\sigma^2})\} \\
&\propto \exp\{-\frac{1}{2}[(\frac{1}{\tau}I + \frac{X^TX}{\sigma^2})\vec{\beta}^T\vec{\beta} - \frac{2X^T\vec{Y}}{\sigma^2}\vec{\beta}]\} \\
&= \exp\{-\frac{1}{2}(\frac{1}{\tau}I + \frac{X^TX}{\sigma^2})[\vec{\beta}^T\vec{\beta} - 2\frac{\frac{X^T\vec{Y}}{\sigma^2}}{\frac{1}{\tau}I + \frac{X^TX}{\sigma^2}}\vec{\beta}]\} \\
&\propto \exp\{-\frac{1}{2}(\frac{1}{\tau}I + \frac{X^TX}{\sigma^2})(\vec{\beta} - E[\vec{\beta}|\vec{Y}])^2\}
\end{aligned}$$

The posterior distribution is also Gaussian, with mean equal to $\frac{\frac{X^T\vec{Y}}{\sigma^2}}{\frac{1}{\tau}I + \frac{X^TX}{\sigma^2}}$.

Simplify it and we get $(X^TX + \frac{\sigma^2}{\tau}I)^{-1}X^T\vec{Y}$, if taking the regularization parameter $\lambda = \frac{\sigma^2}{\tau}$, then the mean of the posterior is $(X^TX + \lambda I)^{-1}X^T\vec{Y}$, which is exactly the formula for ridge regression estimate. Moreover, since the posterior is Gaussian, the mode is the same as the mean, which is also the ridge regression estimate.

ESL 3.19

ESL 3.28

ESL 3.30

Reference

- Hastie, T., Tibshirani, R., Friedman, J. (2009). The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition. Springer.
- Anton, H., & Rorres, C. (2013). Elementary Linear Algebra: Applications Version. Wiley.