

# SDS 385: Exercises 1 - Preliminaries

August 23, 2016

*Professor James Scott*

Spencer Woody

## Problem 1

(A)

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \sum_{i=1}^N \frac{w_i}{2} (y_i - x_i^T \beta)^2 \quad (1)$$

$$= \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{2} (Y - X\beta)^T W (Y - X\beta) \quad (2)$$

$$\frac{1}{2} (Y - X\beta)^T W (Y - X\beta) = \frac{1}{2} (Y^T - \beta^T X^T) W (Y - X\beta) \quad (3)$$

$$= \frac{1}{2} (Y^T W - \beta^T X^T W) (Y - X\beta) \quad (4)$$

$$= \frac{1}{2} (Y^T W Y - \beta^T X^T W Y - Y^T W X \beta + \beta^T X^T W X \beta) \quad (5)$$

$$= \frac{1}{2} (Y^T W Y - 2(X\beta)^T W Y + \beta^T X^T W X \beta) \quad (6)$$

$$= \frac{1}{2} Y^T W Y - (X\beta)^T W Y + \frac{1}{2} \beta^T X^T W X \beta, \quad (7)$$

because

$$\beta^T X^T W Y = (X\beta)^T W Y, \quad (8)$$

and

$$Y^T W X \beta = (Y^T W X \beta)^T \because Y^T W X \beta \in \mathbb{R}^1 \quad (9)$$

$$(Y^T W X \beta)^T = (W X \beta)^T Y = (X\beta)^T W^T Y = (X\beta)^T W Y. \quad (10)$$

We want to minimize the objective function from Eqn. (7), so we take the gradient with respect to  $\beta$  and set it equal to zero. For each of the three terms, their are respective gradients with respect to  $\beta$  are

(i)

$$\frac{\partial}{\partial \beta} \frac{1}{2} Y^T W Y = 0 \quad (11)$$

(ii)

$$\frac{\partial}{\partial \beta} - (X\beta)^T W Y = -X^T W Y \quad (12)$$

(iii)

$$\frac{\partial}{\partial \beta} \frac{1}{2} \beta^T X^T W X \beta = \frac{1}{2} \beta^T (X^T W X + (X^T W X)^T) \quad (13)$$

$$= X^T W X \beta. \quad (14)$$

Summing these terms and equaling them to zero yields

$$X^T W X \beta - X^T W Y = 0 \therefore \quad (15)$$

$$(X^T W X) \hat{\beta} = X^T W Y \quad (16)$$

(B) The brute force method of solving Eqn. (16) is the *inversion method*, i.e.

$$\hat{\beta} = (X^T W X)^{-1} X^T W y. \quad (17)$$

However, this method is computationally expensive. Therefore I propose an alternative methods to solving this matrix equation using the Cholesky decomposition. **Cholesky Decomposition**  
Let

$$C = X^T W X, \quad D = X^T W y \quad (18)$$

so

$$C \hat{\beta} = D. \quad (19)$$

We decompose matrix  $C$  into a product of a lower-triangular matrix and an upper-triangular matrix, such that  $U = L^T$  so

$$C = LU = LL^T \therefore \quad (20)$$

$$LL^T \hat{\beta} = D. \quad (21)$$

Furthermore we define matrix  $A = L^T \hat{\beta}$ . Thus we are left with two matrix equations to solve.

$$LA = D \quad (22)$$

$$L^T \hat{\beta} = A \quad (23)$$

This method will be much less computationally intensive than the inversion method because of the fact that the two left-matrices  $L$  and  $U = L^T$  are triangular. We still must invert  $L$  and  $L^T$  but this is simpler than taking an inverse of a more complicated matrix  $X^T W X$ . This is similar to an LU decomposition, with the exception that we necessarily have two triangular matrices that are transposes of one another. Therefore, this method gains a computational advantage over LU decomposition from symmetric exploitation.

(C) Code for implementing this method is shown in the appendix to this paper.

(D)

## Problem 2

(A) We have  $y_i \sim \text{Binomial}(m_i, w_i)$ , where

$$w_i = \frac{1}{1 + \exp(-x_i^T \beta)}, \quad 1 - w_i = \frac{\exp(-x_i^T \beta)}{1 + \exp(-x_i^T \beta)}, \quad (24)$$

so the negative log likelihood is

$$\ell(\beta) = -\log \left\{ \prod_{i=1}^N p(y_i | \beta) \right\} \quad (25)$$

$$= -\log \left\{ \prod_{i=1}^N \binom{m_i}{y_i} (w_i)^{y_i} (1 - w_i)^{m_i - y_i} \right\} \quad (26)$$

$$= -\left\{ \sum_{i=1}^N \left( \log \binom{m_i}{y_i} + y_i \log(w_i) + (m_i - y_i) \log(1 - w_i) \right) \right\} \quad (27)$$

$$= -\left\{ \sum_{i=1}^N \left( \log \binom{m_i}{y_i} + y_i \log \left( \frac{1}{1 + \exp(-x_i^T \beta)} \right) + (m_i - y_i) \log \left( \frac{\exp(-x_i^T \beta)}{1 + \exp(-x_i^T \beta)} \right) \right) \right\} \quad (28)$$

$$= -\left\{ \sum_{i=1}^N \left( \log \binom{m_i}{y_i} - y_i \log(1 + \exp(-x_i^T \beta)) - (m_i - y_i) x_i^T \beta - m_i \log(1 + \exp(-x_i^T \beta)) + y_i \log(1 + \exp(-x_i^T \beta)) \right) \right\} \quad (29)$$

$$= -\left\{ \sum_{i=1}^N \left( \log \binom{m_i}{y_i} - (m_i - y_i) x_i^T \beta - m_i \log(1 + \exp(-x_i^T \beta)) \right) \right\} \quad (30)$$

$$= \sum_{i=1}^N \left( (m_i - y_i) x_i^T \beta + m_i \log(1 + \exp(-x_i^T \beta)) - \log \binom{m_i}{y_i} \right) \quad (31)$$

$$(32)$$

The gradient for this expression is,

$$\frac{\nabla \ell(\beta)}{d\beta} = \sum_{i=1}^N \left( (m_i - y_i) x_i - m_i \frac{1}{1 + \exp(-x_i^T \beta)} \exp(-x_i^T \beta) x_i \right) \quad (33)$$

$$= \sum_{i=1}^N ((m_i - y_i) x_i - m_i w_i \exp(-x_i^T \beta) x_i) \quad (34)$$

$$= \sum_{i=1}^N (m_i - y_i - m_i w_i \exp(-x_i^T \beta)) x_i \quad (35)$$

$$= \sum_{i=1}^N (m_i w_i - y_i) x_i \quad (36)$$

Text here.

(B)

(C)

(D)

(E)

```
#####
##### Created by Spencer Woody on 24 Aug 2016 #####
#####

5 library(Matrix)
  library(microbenchmark)

  ### No. 1 pt C

10 # Set N, P, X, W, and y

  N <- 2000
  P <- 500

15 X <- matrix(rnorm(N * P), nrow = N)
  y <- matrix(rnorm(N), nrow = N)
  W <- diag(rep(1, N))

  # Inversion method

20 Inv.method <- function(X.Inv, W.Inv, y.Inv) {
  XtWX <- (t(X.Inv)*diag(W.Inv)) %*% X.Inv
  XtWY <- (t(X.Inv)*diag(W.Inv)) %*% y.Inv
  bhat.Inv <- solve(XtWX) %*% XtWY
25   return(bhat.Inv)
}

  Cho.decomp <- function(X.Cho, W.Cho, y.Cho) {
  D.Cho <- (t(X.Cho)*diag(W.Cho)) %*% y.Cho
30   C.Cho <- (t(X.Cho)*diag(W.Cho)) %*% X.Cho

  U.Cho <- chol(C.Cho)
  L.Cho <- t(U.Cho)

35   u <- forwardsolve(L.Cho, D.Cho)
  bhat.Cho <- backsolve(U.Cho, u)

  return(bhat.Cho)
}

40 microbenchmark(
  Inv.method(X, W, y),
  Cho.decomp(X, W, y),
  times=5)

45

  ### No. 1 pt D

  N <- 2000
50 P <- 500

  X <- matrix(rnorm(N * P), nrow = N)
  mask <- matrix(rbinom(N * P, 1, 0.05), nrow = N)
```

```
X <- mask * X

55 Inv.methodSPARSE <- function(X.Inv, W.Inv, y.Inv) {
  X <- Matrix(X, sparse = TRUE)
  XtWX <- (t(X.Inv)*diag(W.Inv)) %*% X.Inv
  XtWY <- (t(X.Inv)*diag(W.Inv)) %*% y.Inv
60  bhat.Inv = Matrix::solve(XtWX, XtWY, sparse = TRUE)
  return(bhat.Inv)
}

Inv.methodSPARSE2 <- function(X.Inv, W.Inv, y.Inv) {
65  XtWX <- (t(X.Inv)*diag(W.Inv)) %*% X.Inv
  XtWY <- (t(X.Inv)*diag(W.Inv)) %*% y.Inv
  bhat.Inv = solve(XtWX, XtWY)
  return(bhat.Inv)
}

70 microbenchmark(
  Inv.methodSPARSE(X, W, y),
  Inv.methodSPARSE2(X, W, y),
  Cho.decomp(X, W, y),
75  times=5)

microbenchmark(
  solve(XtWX, XtWY),
80  Matrix::solve(XtWX, XtWY, sparse = TRUE),
  solve(XtWX) %*% XtWY,
  times = 5
)

85
# END
```