

*Homework Notes:* I did not work with anyone else on this homework or refer to resources other than the course notes, textbook, and course Piazza page.

**Problem 1**

**A** Using the normal equation  $\hat{\beta} = (X^T X)^{-1}(X^T Y)$ , we can calculate  $\hat{\beta}$  in  $\mathcal{R}$  with the following program (after loading the data as `lindata`):

Listing 1: R Code for 1A

```
1 X.prime.X = t(X) %*% X
2 X.prime.X.inverse = solve(X.prime.X)
3 X.prime.Y = t(X) %*% Y
4 beta.hat = X.prime.X.inverse %*% X.prime.Y
5 beta.hat ==> 69.942167, 15.119109, 0.321028
6
7 # B
8 linreg = function(X, Y){
```

Letting  $\beta_i$  correspond to  $X_i$  and  $\beta_0$  represent the intercept,  $\hat{\beta}_0 = 69.94$ ,  $\hat{\beta}_1 = 15.12$ ,  $\hat{\beta}_2 = 0.32$ .

**B** We can also estimate  $\beta$  using online stochastic gradient descent using  $\mathcal{R}$ 's `optim` function for minimization:

Listing 2: R Code for 1B

```
9     sq.resid = residuals^2
10     rss = sum(sq.resid)
11     return(rss)
12 }
13 results = optim(rep(0, ncol(X)), RSS,
14     hessian=TRUE, method="BFGS", x=X, y=Y)
15 list(beta=results$par,
16     vcov=solve(results$hessian),
17     converged=results$convergence==0)
18 }
19 model = linreg(X, Y)
20 model$beta ==> 69.942167, 15.119109, 0.321028
21
22 # C
23 ridgereg = function(X, Y, sigma.sq=1, tau.sq=1){
```

Using this method we reach the same results as above:  $\hat{\beta}_0 = 69.94$ ,  $\hat{\beta}_1 = 15.12$ ,  $\hat{\beta}_2 = 0.32$ .

**C** A third method we can use to estimate  $\beta$  is ridge regression, using the following  $\mathcal{R}$  code:

Listing 3: R Code for 1C

```

24  lambda.id = lambda * diag(D)
25  X.prime.X = t(X) %*% X
26  lambda.X.prime.X.inverse = solve(lambda.id + X.prime.X)
27  X.prime.Y = t(X) %*% Y
28  beta.hat = lambda.X.prime.X.inverse %*% X.prime.Y
29  return(beta.hat)
30 }
31 beta.ridge = ridgereg(X, Y)
32 beta.ridge ==> 65.9129285, 14.3521921, 0.3478858
33
34 # D
35 # a

```

From this approach we get the estimates  $\hat{\beta}_1 = 3.00, \hat{\beta}_2 = 0.78$ .

**D** Table 1 presents the residual sum of squares (RSS) for the training data using each of the methods above.

Table 1: Training set RSS for linear regression methods

Method	RSS
Normal equations	98,729.3
Online stochastic gradient descent	98,729.3
Ridge regression	99,005.9

**E** Table 2 presents the residual sum of squares (RSS) for the test data using each of the methods above.

Table 2: Test set RSS for linear regression methods

Method	RSS
Normal equations	11,573.2
Online stochastic gradient descent	11,573.2
Ridge regression	11,523.5

**F** We can now take the predicted blood pressure for a hypothetical female weighing 135 pounds:

Table 3: Predicted values

Method	$\mathbb{E}[Y X = [1, 135]^T, \hat{\beta}]$
Normal equations	128.4
Online stochastic gradient descent	128.4
Ridge regression	127.2

## **G PLOTS HERE**

**H** Of the three methods used above, ridge regression is the least likely to overfit the data. This is because the  $\lambda$  term (using the notation in the MLAPP textbook) helps to regularize the coefficients. This leads to higher RSS but also makes the coefficients less susceptible to outliers in the training data, as evidenced by the lower RSS using the training set data. The other two methods provide no such protection against outliers, and are thus more likely to overfit the training data.

**I** The researchers should use the ridge regression results. One reason for this is its robustness (relative to the other two methods) to overfitting. This is especially pertinent given the small  $n$  of the training data. However, the similarity in the coefficients should help allay any concerns the researchers may have over the differences in the three methods.

**Problem 2**

**A** The following  $\mathcal{R}$  code implements the IRLS algorithm as presented in the MLAPP textbook (after loading the data):

Listing 4: R Code for 2A

```

1 sigm = function(eta){
2   out = exp(eta)/( exp(eta) + 1)
3   return(out)
4 }
5
6 irls = function(X, Y, epsilon=1/1e10){
7   D = ncol(X)
8   N = nrow(X)
9
10  w = matrix(0, nrow=D, ncol=D)
11  w0 = log(mean(Y)/(1-mean(Y)))
12  eta = mu = s = z = matrix(NA, nrow=N, ncol=D)
13
14  converged = FALSE
15  numiters = 0
16  while(!converged){
17    last.w = w
18    for(i in 1:N){
19      xi = matrix(X[i, ], ncol=1)
20      eta[i, ] = w0 + (t(w) %*% xi)
21      mu[i, ] = sigm(eta[i, ])
22      s[i, ] = mu[i, ] * (1-mu[i, ])
23      z[i, ] = eta[i, ] + ((Y[i, ] - mu[i, ])/s[i, ])
24    }
25    S = matrix(0, nrow=N, ncol=N)
26    for(i in 1:N){
27      S[i, i] = s[i]
28    }
29    first.term = solve(t(X) %*% S %*% X)
30    second.term = t(X) %*% S %*% z
31    w = first.term %*% second.term
32    numiters = numiters + 1
33    diff = w - last.w
34    smalldiff = TRUE
35    for(i in 1:D){
36      smalldiff = smalldiff & (diff[i, i] < epsilon)
37    }

```

```
38     if(smalldiff){
39         converged = TRUE
40         output = paste("converged after ", numiters, " iterations", sep="")
41         print(output)
42     }
43     if(numiters > 100){
44         print("reached max iterations")
45         break
46     }
47 }
48 return(w)
49 }
50
51 w = irls(X, Z)
```

**B** The  $\text{RSS}(\hat{\beta})$  for the training data is 141.47.

**C** The  $\text{RSS}(\hat{\beta})$  for the test data is 15.73.

**F** (*sic*)  $\mathbb{E}[Z|X = [1, 135]^T, \hat{\beta}] = 0.088$

**G** **PLOTS HERE**

**H** The correlation between the features  $X_1$  and  $X_2$  is high (0.63). This causes our estimates of  $\beta$  to be biased, since it is difficult to isolate the impact of each feature.