

*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 = as.matrix(lindata[1:1000, 1:2])
2 Y = as.matrix(lindata[1:1000, 3])
3 X.prime.X = t(X) %*% X
4 X.prime.X.inverse = solve(X.prime.X)
5 X.prime.Y = t(X) %*% Y
6 beta.hat = X.prime.X.inverse %*% X.prime.Y
7 beta.hat #=> 3.0117879, 0.7832938
```

Letting  $\beta_i$  correspond to  $X_i$ ,  $\hat{\beta}_1 = 3.01$ ,  $\hat{\beta}_2 = 0.78$ .

**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

```
8 linreg = function(X, Y){
9   RSS = function(b, x, y){
10    residuals = y - (x %*% b)
11    sq.resid = residuals^2
12    rss = sum(sq.resid)
13    return(rss)
14  }
15  results = optim(rep(0, ncol(X)), RSS,
16    hessian=TRUE, method="BFGS", x=X, y=Y)
17  list(beta=results$par,
18    vcov=solve(results$hessian),
19    converged=results$convergence==0)
20 }
21 model = linreg(X, Y)
22 model$beta #=> 3.0117879, 0.7832938
```

Using this method we reach the same results as above:  $\hat{\beta}_1 = 3.01$ ,  $\hat{\beta}_2 = 0.78$ .

**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

```

23 ridgereg = function(X, Y, sigma.sq=1, tau.sq=1){
24   lambda = sigma.sq / tau.sq
25   D = ncol(X)
26   lambda.id = lambda * diag(D)
27   print(lambda.id)
28   X.prime.X = t(X) %*% X
29   lambda.X.prime.X.inverse = solve(lambda.id + X.prime.X)
30   X.prime.Y = t(X) %*% Y
31   beta.hat = lambda.X.prime.X.inverse %*% X.prime.Y
32   return(beta.hat)
33 }
34 beta.ridge = ridgereg(X, Y)
35 beta.ridge  $\Rightarrow$  2.9972982, 0.7833412

```

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	181767.9
Online stochastic gradient descent	181767.9
Ridge regression	181768

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

Table 2: Test set RSS for linear regression methods

Method	RSS
Normal equations	16963.35
Online stochastic gradient descent	16963.35
Ridge regression	16963.59

**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	108.76
Online stochastic gradient descent	108.76
Ridge regression	108.75

## 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. The other two methods provide no such protection against outliers.

**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 relatively 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.

Matt Dickenson mcd31

STA561/CS571 — Fall 2013

## Homework 2

Due: 23 September, 2013

---

### Problem 2