Matt Dickenson `mcd31`
STA561/CS571 — Fall 2013      **Homework 2**      Due: 23 September, 2013

*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 $\hat{\beta} = (X^TX)^{-1}(X^TY)$, 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 X = cbind(1, X)
 3 Y = as.matrix(lindata[1:1000  , 3])
 4
 5 normaleqn = function(x, y){
 6   x.prime.x.inverse = solve(t(x) %*% x)
 7   x.prime.y = t(x) %*% y
 8   beta = x.prime.x.inverse %*% x.prime.y
 9   return(beta)
10 }
11
12 beta.hat = normaleqn(X, Y)
13 beta.hat #=> 69.942167, 15.119109, 0.321028
```

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 with $\mathcal{R}$'s `optim` function for minimization:

Listing 2: R Code for 1B

```
14 linreg = function(X, Y){
15   RSS = function(b, x, y){
16     residuals = y - (x %*% b )
17     sq.resid = residuals^2
18     rss = sum(sq.resid)
19     return(rss)
20   }
21   results = optim(rep(0, ncol(X)), RSS,
22     hessian=TRUE, method="BFGS", x=X, y=Y)
23   list(beta=results$par,
24     vcov=solve(results$hessian),
25     converged=results$convergence==0)
26 }
```

```
27 model = linreg(X, Y)
28 model$beta #=> 69.942167, 15.119109, 0.321028
```

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

```
29 ridgereg = function(X, Y, sigma.sq=1, tau.sq=1){
30   lambda = sigma.sq / tau.sq
31   D = ncol(X)
32   lambda.id = lambda * diag(D)
33   X.prime.X = t(X) %*% X
34   lambda.X.prime.X.inverse = solve(lambda.id + X.prime.X)
35   X.prime.Y = t(X) %*% Y
36   beta.hat = lambda.X.prime.X.inverse %*% X.prime.Y
37   return(beta.hat)
38 }
39 beta.ridge = ridgereg(X, Y)
40 beta.ridge #=> 65.9129285, 14.3521921, 0.3478858
```

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.

**F**    We can now take the predicted blood pressure for a hypothetical female weighing 135 pounds for the three different methods (Table 3):

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 |

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**   Bivariate plots of $X_1$ and $X_2$ with $Y$ are given in Figures 1 and 2 below, along with regression lines for the three methods used above. Note that the lines for the normal equation method and the stochastic gradient descent (SGD) method overlap in both cases. (The intercept and $\hat{\beta}_i$ are used in the plot for $X_i$, while $\hat{\beta}_j, j \neq i$ is ignored.)
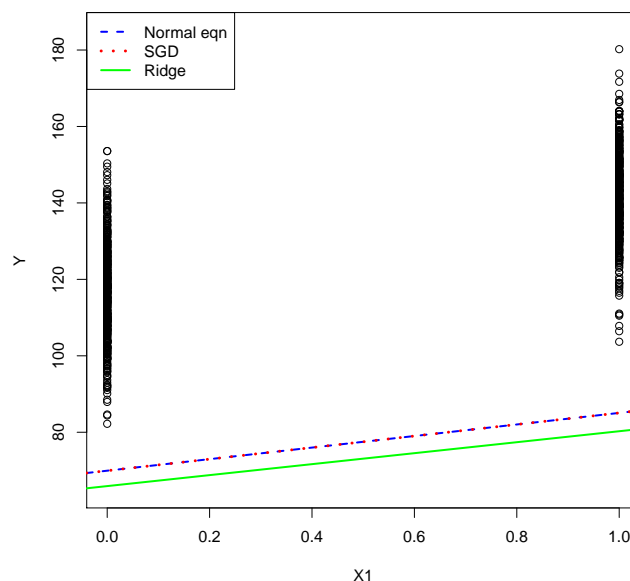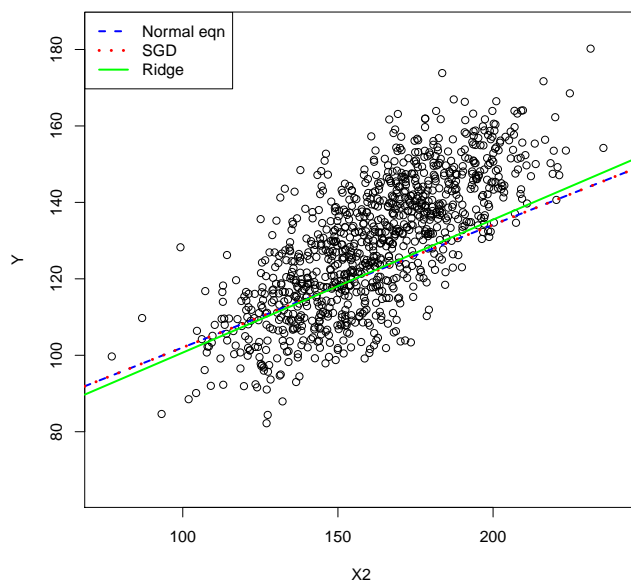


Figure 1: $X_1$ and $Y$ with Regression Lines

Figure 2: $X_2$ and $Y$ with Regression Lines

**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 in the training data but also makes the coefficients less susceptible to outliers, as evidenced by the lower RSS using the test 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 regresion 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 RSS($\hat{\beta}$) for the training data is 141.47.

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

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

**G**    Figures 3 and 4 below plot $X_1$ and $X_2$ with $Z$ and the respective regression line. Because $X_1$ and $Z$ are binary variables, a small random jitter is added to make the points more apparent.

**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.

**I**    A model with only $X_2$ would do almost as well as the model with both features. This is due to the correlation between the two features and the relatively minimal impact of $X_1$ on the predicted $Z$ values, as shown in Figure 3. As Figure 4 shows, risk for $Z$ increases noticeably when $X_2 > 160$, so this should be an important component of the researchers' decision rule.
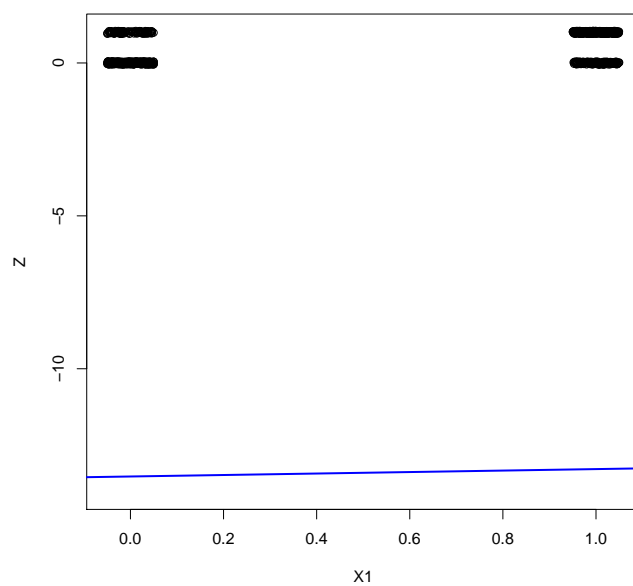
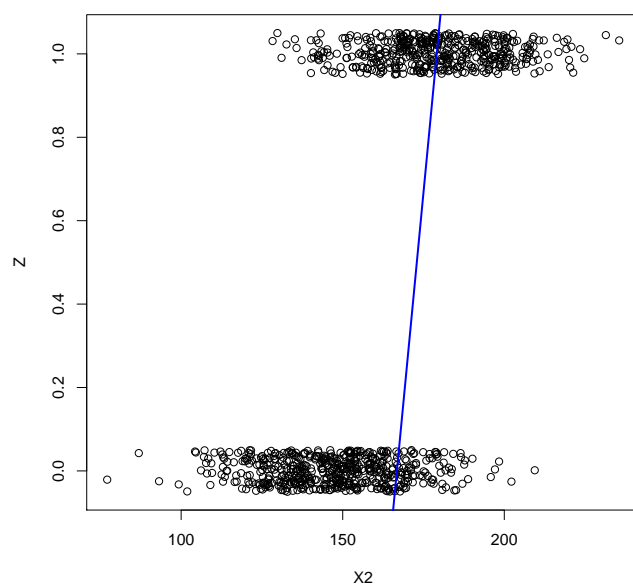Figure 3: $X_1$ and $Z$ with Regression Lines



Figure 4: $X_2$ and $Z$ with Regression Lines