

Lab 4

Shoara Chowdhury

11:59PM March 10, 2021

Load up the famous iris dataset. We are going to do a different prediction problem. Imagine the only input x is Species and you are trying to predict y which is Petal.Length. A reasonable prediction is the average petal length within each Species. Prove that this is the OLS model by fitting an appropriate `lm` and then using the `predict` function to verify.

```
data(iris)
mod = lm(Petal.Length ~ Species, iris)
mean(iris$Petal.Length[iris$Species=="setosa"])
```

```
## [1] 1.462
```

```
mean(iris$Petal.Length[iris$Species=="versicolor"])
```

```
## [1] 4.26
```

```
mean(iris$Petal.Length[iris$Species=="virginica"])
```

```
## [1] 5.552
```

```
predict(mod,data.frame(Species =("setosa")))
```

```
##      1
## 1.462
```

```
predict(mod,data.frame(Species =("versicolor")))
```

```
##      1
## 4.26
```

```
predict(mod,data.frame(Species =("virginica")))
```

```
##      1
## 5.552
```

Construct the design matrix with an intercept, X without using `model.matrix`.

```
x <- cbind(1,iris$Species == "versicolor", iris$Species == "virginica" )
head(x)
```

```
##      [,1] [,2] [,3]
## [1,]    1    0    0
## [2,]    1    0    0
## [3,]    1    0    0
## [4,]    1    0    0
## [5,]    1    0    0
## [6,]    1    0    0
```

Find the hat matrix H for this regression.

```
H = x %*% solve(t(x) %*% x) %*% t(x)
Matrix::rankMatrix(H)
```

```
## [1] 3
## attr("method")
## [1] "tolNorm2"
## attr("useGrad")
## [1] FALSE
## attr("tol")
## [1] 3.330669e-14
```

Verify this hat matrix is symmetric using the `expect_equal` function in the package `testthat`.

```
pacman::p_load(testthat)
expect_equal(H, t(H))
```

Verify this hat matrix is idempotent using the `expect_equal` function in the package `testthat`.

```
expect_equal(H, H%*%H)
```

Using the `diag` function, find the trace of the hat matrix.

```
diag(H)

##      [1] 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02
##     [16] 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02
##     [31] 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02
##     [46] 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02
##     [61] 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02
##     [76] 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02
##     [91] 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02
##    [106] 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02
##    [121] 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02
##    [136] 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02 0.02
```

It turns out the trace of a hat matrix is the same as its rank! But we don't have time to prove these interesting and useful facts..

For masters students: create a matrix X-perpendicular.

```
sum(diag(H))
```

```
## [1] 3
```

Using the hat matrix, compute the \hat{y} vector and using the projection onto the residual space, compute the e vector and verify they are orthogonal to each other.

```
y = iris$Petal.Length
y_hat = H%% y
e = (diag(nrow(iris))-H) %% y
e
```

```
##           [,1]
## [1,] -0.062
## [2,] -0.062
## [3,] -0.162
## [4,]  0.038
## [5,] -0.062
## [6,]  0.238
## [7,] -0.062
## [8,]  0.038
## [9,] -0.062
## [10,]  0.038
## [11,]  0.038
## [12,]  0.138
## [13,] -0.062
## [14,] -0.362
## [15,] -0.262
## [16,]  0.038
## [17,] -0.162
## [18,] -0.062
## [19,]  0.238
## [20,]  0.038
## [21,]  0.238
## [22,]  0.038
## [23,] -0.462
## [24,]  0.238
## [25,]  0.438
## [26,]  0.138
## [27,]  0.138
## [28,]  0.038
## [29,] -0.062
## [30,]  0.138
## [31,]  0.138
## [32,]  0.038
## [33,]  0.038
## [34,] -0.062
## [35,]  0.038
```

```
## [36,] -0.262
## [37,] -0.162
## [38,] -0.062
## [39,] -0.162
## [40,]  0.038
## [41,] -0.162
## [42,] -0.162
## [43,] -0.162
## [44,]  0.138
## [45,]  0.438
## [46,] -0.062
## [47,]  0.138
## [48,] -0.062
## [49,]  0.038
## [50,] -0.062
## [51,]  0.440
## [52,]  0.240
## [53,]  0.640
## [54,] -0.260
## [55,]  0.340
## [56,]  0.240
## [57,]  0.440
## [58,] -0.960
## [59,]  0.340
## [60,] -0.360
## [61,] -0.760
## [62,] -0.060
## [63,] -0.260
## [64,]  0.440
## [65,] -0.660
## [66,]  0.140
## [67,]  0.240
## [68,] -0.160
## [69,]  0.240
## [70,] -0.360
## [71,]  0.540
## [72,] -0.260
## [73,]  0.640
## [74,]  0.440
## [75,]  0.040
## [76,]  0.140
## [77,]  0.540
## [78,]  0.740
## [79,]  0.240
## [80,] -0.760
## [81,] -0.460
## [82,] -0.560
## [83,] -0.360
## [84,]  0.840
## [85,]  0.240
## [86,]  0.240
## [87,]  0.440
## [88,]  0.140
## [89,] -0.160
```

```
## [90,] -0.260
## [91,]  0.140
## [92,]  0.340
## [93,] -0.260
## [94,] -0.960
## [95,] -0.060
## [96,] -0.060
## [97,] -0.060
## [98,]  0.040
## [99,] -1.260
## [100,] -0.160
## [101,]  0.448
## [102,] -0.452
## [103,]  0.348
## [104,]  0.048
## [105,]  0.248
## [106,]  1.048
## [107,] -1.052
## [108,]  0.748
## [109,]  0.248
## [110,]  0.548
## [111,] -0.452
## [112,] -0.252
## [113,] -0.052
## [114,] -0.552
## [115,] -0.452
## [116,] -0.252
## [117,] -0.052
## [118,]  1.148
## [119,]  1.348
## [120,] -0.552
## [121,]  0.148
## [122,] -0.652
## [123,]  1.148
## [124,] -0.652
## [125,]  0.148
## [126,]  0.448
## [127,] -0.752
## [128,] -0.652
## [129,]  0.048
## [130,]  0.248
## [131,]  0.548
## [132,]  0.848
## [133,]  0.048
## [134,] -0.452
## [135,]  0.048
## [136,]  0.548
## [137,]  0.048
## [138,] -0.052
## [139,] -0.752
## [140,] -0.152
## [141,]  0.048
## [142,] -0.452
## [143,] -0.452
```

```
## [144,] 0.348
## [145,] 0.148
## [146,] -0.352
## [147,] -0.552
## [148,] -0.352
## [149,] -0.152
## [150,] -0.452
```

Compute SST, SSR and SSE and R^2 and then show that $SST = SSR + SSE$.

```
SSE= t(e) %*% e
y_bar = mean(y)
SST= t(y-y_bar)%*% (y-y_bar)
Rsqr = 1- SSE/SST
Rsqr
```

```
##          [,1]
## [1,] 0.9413717
```

```
SSR= t(y-y_bar)%*% (y-y_bar)
SSR
```

```
##          [,1]
## [1,] 464.3254
```

```
#expect_equal(SSR+SSE, SST)
var(y)
```

```
## [1] 3.116278
```

```
var(e)
```

```
##          [,1]
## [1,] 0.182702
```

Find the angle theta between $y - \bar{y}$ and $\hat{y} - \bar{y}$ and then verify that its cosine squared is the same as the R^2 from the previous problem.

```
theta = acos(t(y-y_bar) %*% (y_hat-y_bar) / sqrt(SST * SSR))
theta = (180 / pi)
```

Project the y vector onto each column of the X matrix and test if the sum of these projections is the same as \hat{y} .

```
prjt1 = ( x[,1] %*% t(x[,1]) / as.numeric( t(x[,1]) %*% x[,1]) ) %*% y
prjt2 = ( x[,2] %*% t(x[,2]) / as.numeric( t(x[,2]) %*% x[,2]) ) %*% y
prjt3 = ( x[,3] %*% t(x[,3]) / as.numeric( t(x[,3]) %*% x[,3]) ) %*% y
```

Construct the design matrix without an intercept, X, without using `model.matrix`.

```
Mi <- cbind(0,iris$Species == "versicolor", iris$Species == "virginica" )
```

Find the OLS estimates using this design matrix. It should be the sample averages of the petal lengths within species.

```
#expect_equal(Mi,Mi%*%Mi)
```

Verify the hat matrix constructed from this design matrix is the same as the hat matrix constructed from the design matrix with the intercept. (Fact: orthogonal projection matrices are unique).

```
expect_equal(Mi, H%*%Mi)
```

Project the y vector onto each column of the X matrix and test if the sum of these projections is the same as yhat.

```
“{r} prjt4 = ( x[,1] %% t(x[,1]) / as.numeric( t(x[,1]) %% x[,1]) ) %% C prjt5 = ( x[,2] %% t(x[,2]) / as.numeric( t(x[,2]) %% x[,2]) ) %% C prjt6 = ( x[,3] %% t(x[,3]) / as.numeric( t(x[,3]) %% x[,3]) ) %% C
```

Convert this design matrix into Q, an orthonormal matrix.

```
#Q = (diag(nrow(iris))-y) %% C
#Q
```

Project the y vector onto each column of the Q matrix and test if the sum of these projections is the same as yhat.

```
#prjt11 = ( x[,1] %% t(x[,1]) / as.numeric( t(x[,1]) %% x[,1]) ) %% Q
#prjt22 = ( x[,2] %% t(x[,2]) / as.numeric( t(x[,2]) %% x[,2]) ) %% Q
#prjt33 = ( x[,3] %% t(x[,3]) / as.numeric( t(x[,3]) %% x[,3]) ) %% Q
```

Find the p=3 linear OLS estimates if Q is used as the design matrix using the `lm` method. Is the OLS solution the same as the OLS solution for X?

Use the `predict` function and ensure that the predicted values are the same for both linear models: the one created with X as its design matrix and the one created with Q as its design matrix.

Clear the workspace and load the boston housing data and extract X and y. The dimensions are n = 506 and p = 13. Create a matrix that is (p + 1) x (p + 1) full of NA's. Label the columns the same columns as X. Do not label the rows. For the first row, find the OLS estimate of the y regressed on the first column only and put that in the first entry. For the second row, find the OLS estimates of the y regressed on the first and second columns of X only and put them in the first and second entries. For the third row, find the OLS estimates of the y regressed on the first, second and third columns of X only and put them in the first, second and third entries, etc. For the last row, fill it with the full OLS estimates.

```
#TO-DO
```

Why are the estimates changing from row to row as you add in more predictors?

```
#TO-DO
```

Create a vector of length p+1 and compute the R² values for each of the above models.

#TO-DO

Is R^2 monotonically increasing? Why?

#TO-DO