

Lab 4

Janine Lim

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, data = iris )
mod
```

```
##
## Call:
## lm(formula = Petal.Length ~ Species, data = iris)
##
## Coefficients:
##      (Intercept)  Speciesversicolor  Speciesvirginica
##           1.462             2.798             4.090
```

```
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 = c("setosa", "versicolor", "virginica"))) ##returns ybar for each group
```

```
##      1      2      3
## 1.462 4.260 5.552
```

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

```
X = as.matrix(cbind(1,
                    (iris$Species == "versicolor"),
                    (iris$Species == "virginica")
)) ##using setosa as the reference category
```

Find the hat matrix H for this regression.

```
XTX = t(X) %*% X
XTX_inv = solve(XTX)
H = X %*% XTX_inv %*% 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.

```
traceH = sum(diag(H))
```

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_{\perp} .

```
Id = diag(nrow(H))
X_perp = (Id - H) %*% X
t(X_perp) %*% X ##numerically 0, so we see that X_perp and X must be orthogonal
```

```
##           [,1]           [,2]           [,3]
## [1,] -6.835157e-14 -6.952772e-15 -2.600697e-14
## [2,] -1.600109e-14 -1.600109e-14  0.000000e+00
## [3,] -3.494427e-14  0.000000e+00 -3.494427e-14
```

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
yhat = H %*% y ##ybar for each species
e = (diag(nrow(iris)) - H) %*% y
t(e) %*% yhat ##verify orthogonal - should be 0

```

```

##           [,1]
## [1,] -2.2915e-13

```

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

```

ybar=mean(y)
SSE = t(e) %*% e
SST = t(y - ybar)%*% (y - ybar)
SSR = t(yhat - ybar) %*% (yhat - ybar)
Rsqr = 1- SSE/SST
c(SST, SSR, SSE, Rsqr)

```

```

## [1] 464.3254000 437.1028000 27.2226000 0.9413717

```

```
SSR + SSE
```

```

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

```

```
expect_equal(SST, SSR + SSE )
```

Find the angle θ between $y - \bar{y}1$ and $\hat{y} - \bar{y}1$ and then verify that its cosine squared is the same as the R^2 from the previous problem.

```

theta = acos(((t(y - ybar)%*% (yhat - ybar))) / sqrt(SST * SSR))
theta * (180/pi)

```

```

##           [,1]
## [1,] 14.01245

```

```
cos(theta)^2
```

```

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

```

```
expect_equal(cos(theta)^2, Rsqr)
```

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

```

proj1 = (X[,1] %*% t(X[,1]) / as.numeric(t(X[,1]) %*% X[,1])) %*% y
proj2 = (X[,2] %*% t(X[,2]) / as.numeric(t(X[,2]) %*% X[,2])) %*% y
proj3 = (X[,3] %*% t(X[,3]) / as.numeric(t(X[,3]) %*% X[,3])) %*% y

expect_equal(proj1 + proj2 + proj3, yhat) ## expected to fail

```

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

```
X_no_intercept = as.matrix(cbind(
  as.numeric(iris$Species=="setosa"),
  as.numeric(iris$Species=="versicolor"),
  as.numeric(iris$Species=="virginica")),
)
head(X_no_intercept)
```

```
##      [,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 OLS estimates using this design matrix. It should be the sample averages of the petal lengths within species.

```
y_hat = solve(t(X_no_intercept) %*% X_no_intercept) %*% t(X_no_intercept) %*% y
y_hat
```

```
##      [,1]
## [1,] 1.462
## [2,] 4.260
## [3,] 5.552
```

```
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
```

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).

```
H_no_intercept = X_no_intercept %*% solve(t(X_no_intercept) %*% X_no_intercept) %*% t(X_no_intercept)
expect_equal(H_no_intercept, H) ##no error so H and H_no_intercept are equal
```

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

```
proj1 = (X_no_intercept[,1] %*% t(X_no_intercept[,1]) / as.numeric(t(X_no_intercept[,1]) %*% X_no_intercept[,1])) %*% y
proj2 = (X_no_intercept[,2] %*% t(X_no_intercept[,2]) / as.numeric(t(X_no_intercept[,2]) %*% X_no_intercept[,2])) %*% y
proj3 = (X_no_intercept[,3] %*% t(X_no_intercept[,3]) / as.numeric(t(X_no_intercept[,3]) %*% X_no_intercept[,3])) %*% y

sum(proj1, proj2, proj3)
```

```
## [1] 563.7
```

```
sum(yhat) ##are these suppose to be equal?
```

```
## [1] 563.7
```

Convert this design matrix into Q , an orthonormal matrix.

```
qrX = qr(X_no_intercept)
Q = qr.Q(qrX)
```

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

```
proj1_Q = (Q[,1] %*% t(Q[,1]) / as.numeric(t(Q[,1]) %*% Q[,1])) %*% y
proj2_Q = (Q[,2] %*% t(Q[,2]) / as.numeric(t(Q[,2]) %*% Q[,2])) %*% y
proj3_Q = (Q[,3] %*% t(Q[,3]) / as.numeric(t(Q[,3]) %*% Q[,3])) %*% y
sum(proj1_Q, proj2_Q, proj3_Q)
```

```
## [1] 563.7
```

```
sum(yhat)
```

```
## [1] 563.7
```

```
y_hat_q = Q %*% t(Q) %*% y
sum(y_hat_q)
```

```
## [1] 563.7
```

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 ?

These are currently not the same estimates, but will give us the same OLS solution as the OLS solution for X .

```
mod2 = lm(y~0+Q)
coef(mod2)
```

```
##           Q1           Q2           Q3
## -10.33790 -30.12275 -39.25857
```

```
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
```

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.

```
colnames(X)<-c("setosa", "versicolor", "virginica")
mod3 = lm(y~0+X_no_intercept)
unique(predict(mod3, data.frame(X_no_intercept)))
```

```
## [1] 1.462 4.260 5.552
```

```
mod4 = lm(y~0+Q)
unique(predict(mod4, data.frame(Q)))
```

```
## [1] 1.462 1.462 4.260 5.552
```

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 $(p + 1) \times (p + 1)$ 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.

```
rm(list=ls())
Boston = MASS::Boston
intercept = rep(1, nrow(Boston))
X = as.matrix(cbind(intercept, Boston[, 1:13]))
y = Boston[,14]

matrix1 = matrix(data = NA, nrow = 14, ncol = 14)
colnames(matrix1) = c(colnames(X))

for (i in 1:ncol(matrix1)){
  b=array(NA, dim = ncol(matrix1))
  X_new = X[, 1:i]
  X_new = as.matrix(X_new)
  XTX_inv = solve(t(X_new) %*% X_new)
  b[1:i] = XTX_inv %*% t(X_new) %*% y
  matrix1[i, ] <- b
}
matrix1
```

```

##      intercept      crim      zn      indus      chas      nox
## [1,] 22.5328063      NA      NA      NA      NA      NA
## [2,] 24.0331062 -0.4151903      NA      NA      NA      NA
## [3,] 22.4856281 -0.3520783 0.11610909      NA      NA      NA
## [4,] 27.3946468 -0.2486283 0.05850082 -0.41557782      NA      NA
## [5,] 27.1128031 -0.2287981 0.05928665 -0.44032511 6.894059      NA
## [6,] 29.4899406 -0.2185190 0.05511047 -0.38348055 7.026223 -5.424659
## [7,] -17.9546350 -0.1769135 0.02128135 -0.14365267 4.784684 -7.184892
## [8,] -18.2649261 -0.1727607 0.01421402 -0.13089918 4.840730 -4.357411
## [9,] 0.8274820 -0.1977868 0.06099257 -0.22573089 4.577598 -14.451531
## [10,] 0.1553915 -0.1780398 0.06095248 -0.21004328 4.536648 -13.342666
## [11,] 2.9907868 -0.1795543 0.07145574 -0.10437742 4.110667 -12.591596
## [12,] 27.1523679 -0.1840321 0.03909990 -0.04232450 3.487528 -22.182110
## [13,] 20.6526280 -0.1599391 0.03887365 -0.02792186 3.216569 -20.484560
## [14,] 36.4594884 -0.1080114 0.04642046 0.02055863 2.686734 -17.766611
##      rm      age      dis      rad      tax      ptratio
## [1,]      NA      NA      NA      NA      NA      NA
## [2,]      NA      NA      NA      NA      NA      NA
## [3,]      NA      NA      NA      NA      NA      NA
## [4,]      NA      NA      NA      NA      NA      NA
## [5,]      NA      NA      NA      NA      NA      NA
## [6,]      NA      NA      NA      NA      NA      NA
## [7,] 7.341586      NA      NA      NA      NA      NA
## [8,] 7.386357 -0.0236248493      NA      NA      NA      NA
## [9,] 6.752352 -0.0556354540 -1.760312      NA      NA      NA
## [10,] 6.791184 -0.0562612189 -1.748296 -0.04529059      NA      NA
## [11,] 6.664084 -0.0546675064 -1.727933 0.15926305 -0.01434060      NA
## [12,] 6.075744 -0.0451880522 -1.583852 0.25472196 -0.01221262 -0.9962062
## [13,] 6.123072 -0.0459320518 -1.554912 0.28157503 -0.01173838 -1.0142228
## [14,] 3.809865 0.0006922246 -1.475567 0.30604948 -0.01233459 -0.9527472
##      black      lstat
## [1,]      NA      NA
## [2,]      NA      NA
## [3,]      NA      NA
## [4,]      NA      NA
## [5,]      NA      NA
## [6,]      NA      NA
## [7,]      NA      NA
## [8,]      NA      NA
## [9,]      NA      NA
## [10,]      NA      NA
## [11,]      NA      NA
## [12,]      NA      NA
## [13,] 0.013620833      NA
## [14,] 0.009311683 -0.5247584

```

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

The estimates are changing from row to row as we add more predictors because the new predictors will change how much each predictor will contribute to the median price of the house. Even if they only change slightly, adding more predictors will change the “weight” of each input feature.

Create a vector of length $p + 1$ and compute the R^2 values for each of the above models.

```

R_sq_vec = array(NA, dim = 14)
ybar = mean(y)
SST = sum((y - ybar)^2)
for (i in 1:nrow(matrix1)){
  b = c(matrix1[i,1:i],rep(0, nrow(matrix1)-i))
  yhat = X %*% b
  SSR = sum((yhat - ybar)^2)
  Rsq = SSR / SST
  R_sq_vec[i] = Rsq
}

R_sq_vec

```

```

## [1] 5.382448e-30 1.507805e-01 2.339884e-01 2.937136e-01 3.295277e-01
## [6] 3.313127e-01 5.873770e-01 5.894902e-01 6.311488e-01 6.319479e-01
## [11] 6.396628e-01 6.703141e-01 6.842043e-01 7.406427e-01

```

Is R^2 monotonically increasing? Why?

Yes, R^2 is monotonically increasing because we see that as we fit more features, the algorithm believes the input data is valid when it is not, and begins overfitting the model with these features and thus, will make R^2 increase.