

Lab 4

Frank Palma Gomez

11:59PM March 11, 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)
# response ~ covariate
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] NaN
predict(mod, data.frame(Species = c("setosa")))

##      1
## 1.462
predict(mod, data.frame(Species = c("versicolor")))

##      1
## 4.26
predict(mod, data.frame(Species = c("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.

```
# sum of diag is equal to the rank of a projection
sum(diag(H))
```

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

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} .

```
#TO-DO
```

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

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)
SSR = t(y_hat - y_bar) %>% (y_hat - y_bar)
expect_equal(SSR+SSE, SST)
```

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 - y_bar) %>% (y_hat - y_bar) / sqrt(SST * SSR))
theta * 180 / pi
```

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

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} .

```
proj_1 = (X[, 1] %*% t(X[, 1]) / as.numeric((t(X[, 1])%*% X[, 1])))%*% y
proj_2 = (X[, 2] %*% t(X[, 2]) / as.numeric((t(X[, 1])%*% X[, 1])))%*% y
proj_3 = (X[, 3] %*% t(X[, 3]) / as.numeric((t(X[, 1])%*% X[, 1])))%*% y

expect_equal(proj_1 + proj_2 + proj_3, y_hat, tol=1e4)
```

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

```
X_no_inter = cbind(as.numeric(iris$Species == "setosa"), as.numeric(iris$Species == "versicolor"), as.n
```

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

```
b = solve(t(X_no_inter) %*% X_no_inter) %*% t(X_no_inter) %*% y

expect_equal(b[1], mean(y[iris$Species == "setosa"]))
expect_equal(b[2], mean(y[iris$Species == "versicolor"]))
expect_equal(b[3], mean(y[iris$Species == "virginica"]))
```

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_inter = X_no_inter %*% solve(t(X_no_inter) %*% X_no_inter) %*% t(X_no_inter)

expect_equal(H_no_inter, H)
```

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} .

```
proj_1 = (X_no_inter[, 1] %*% t(X_no_inter[, 1]) / as.numeric((t(X_no_inter[, 1])%*% X_no_inter[, 1])))%*% y
proj_2 = (X_no_inter[, 2] %*% t(X_no_inter[, 2]) / as.numeric((t(X_no_inter[, 2])%*% X_no_inter[, 2])))%*% y
proj_3 = (X_no_inter[, 3] %*% t(X_no_inter[, 3]) / as.numeric((t(X_no_inter[, 3])%*% X_no_inter[, 3])))%*% y

yhat = H_no_inter %*% y

expect_equal(proj_1 + proj_2 + proj_3, y_hat)
```

Convert this design matrix into Q , an orthonormal matrix.

```
qrX_no_inter = qr(X_no_inter)
Q = qr.Q(qrX_no_inter)
R = qr.R(qrX_no_inter)
```

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

```
Qt = t(Q)
yhat_via_Q = Q %*% Qt %*% y
expect_equal(c(yhat), c(yhat_via_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 ?

NO is not the same because each dimension is only one piece of the SSR and only one piece of R^2 as well.

```
modQ = lm(y ~ Q)
modX_no_inter = lm(y ~ X_no_inter)
coef(modQ)
```

```
## (Intercept)      Q1      Q2      Q3
##      5.55200    28.92067    9.13582    NA
```

```
coef(modX_no_inter)
```

```
## (Intercept) X_no_inter1 X_no_inter2 X_no_inter3
##      5.552      -4.090      -1.292      NA
```

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.

```
expect_equal(predict(modX_no_inter, data.frame(X_no_inter)), predict(modQ, data.frame(X_no_inter)))
```

```
## Warning in predict.lm(modX_no_inter, data.frame(X_no_inter)): prediction from a
## rank-deficient fit may be misleading
```

```
## Warning in predict.lm(modQ, data.frame(X_no_inter)): prediction from a rank-
## deficient fit may be misleading
```

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.

```
y = MASS::Boston$medv
X = as.matrix(MASS::Boston[, 1:13])

n = nrow(X)
p_plus_one = ncol(X)

NA_matrix = matrix(data=NA, nrow = p_plus_one, ncol = p_plus_one)
colnames(NA_matrix) = colnames(X)
```

```
X_df = data.frame(X)
for (i in 1:p_plus_one) {
  cols = colnames(X_df)[1:i]
  mod = lm(y ~ ., X_df[cols])
  coefs = coef(mod)
  for (j in 1:length(coefs)) {
    NA_matrix[i, j-1] = coefs[j]
  }
}
```

```
NA_matrix
```

```
##      crim      zn      indus      chas      nox      rm
## [1,] -0.4151903      NA      NA      NA      NA      NA
## [2,] -0.3520783 0.11610909      NA      NA      NA      NA
## [3,] -0.2486283 0.05850082 -0.41557782      NA      NA      NA
## [4,] -0.2287981 0.05928665 -0.44032511 6.894059      NA      NA
## [5,] -0.2185190 0.05511047 -0.38348055 7.026223 -5.424659      NA
## [6,] -0.1769135 0.02128135 -0.14365267 4.784684 -7.184892 7.341586
## [7,] -0.1727607 0.01421402 -0.13089918 4.840730 -4.357411 7.386357
## [8,] -0.1977868 0.06099257 -0.22573089 4.577598 -14.451531 6.752352
## [9,] -0.1780398 0.06095248 -0.21004328 4.536648 -13.342666 6.791184
```

```
## [10,] -0.1795543 0.07145574 -0.10437742 4.110667 -12.591596 6.664084
## [11,] -0.1840321 0.03909990 -0.04232450 3.487528 -22.182110 6.075744
## [12,] -0.1599391 0.03887365 -0.02792186 3.216569 -20.484560 6.123072
## [13,] -0.1080114 0.04642046 0.02055863 2.686734 -17.766611 3.809865
##           age      dis      rad      tax      ptratio      black
## [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,] -0.0236248493      NA      NA      NA      NA      NA
## [8,] -0.0556354540 -1.760312      NA      NA      NA      NA
## [9,] -0.0562612189 -1.748296 -0.04529059      NA      NA      NA
## [10,] -0.0546675064 -1.727933 0.15926305 -0.01434060      NA      NA
## [11,] -0.0451880522 -1.583852 0.25472196 -0.01221262 -0.9962062      NA
## [12,] -0.0459320518 -1.554912 0.28157503 -0.01173838 -1.0142228 0.013620833
## [13,] 0.0006922246 -1.475567 0.30604948 -0.01233459 -0.9527472 0.009311683
##           lstat
## [1,]      NA
## [2,]      NA
## [3,]      NA
## [4,]      NA
## [5,]      NA
## [6,]      NA
## [7,]      NA
## [8,]      NA
## [9,]      NA
## [10,]      NA
## [11,]      NA
## [12,]      NA
## [13,] -0.5247584
```

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

Because every time you're adding a feature and therefore you're increasing your SSR

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

```
#T0-D0
```

Is R^2 monotonically increasing? Why?

Every dimension you add, you explain a little more variance with the model.