

Lab 5

Kennly Weerasinghe

11:59PM March 18, 2021

Create a 2x2 matrix with the first column 1's and the next column iid normals. Find the absolute value of the angle (in degrees, not radians) between the two columns.

```
norm_vec = function(v){
  sqrt(sum(v^2))
}
X <- matrix(1:1, nrow = 2, ncol=2)
X[,2] = rnorm(2)
cos_theta = t(X[,1]) %*% X[,2]/(norm_vec(X[,1])*norm_vec(X[,2]))
cos_theta
```

```
##           [,1]
## [1,] 0.7644021
abs(90 - acos(cos_theta)*180/pi)
```

```
##           [,1]
## [1,] 49.85383
```

Repeat this exercise Nsim = 1e5 times and report the average absolute angle.

```
Nsim = 1e5
angles = array(NA, Nsim)
for(i in 1:Nsim) {
  X <- matrix(1:1, nrow = 2, ncol=2)
  X[,2] = rnorm(2)
  cos_theta = t(X[,1]) %*% X[,2]/(norm_vec(X[,1])*norm_vec(X[,2]))
  angles[i] = abs(90 - acos(cos_theta)*180/pi)
}

mean(angles)
```

```
## [1] 44.95798
```

Create a nx2 matrix with the first column 1's and the next column iid normals. Find the absolute value of the angle (in degrees, not radians) between the two columns. For n = 10, 50, 100, 200, 500, 1000, report the average absolute angle over Nsim = 1e5 simulations.

```
N_s = c(2, 5, 10, 50, 100, 200, 500, 1000)
Nsim = 1e5
angles = matrix(NA, nrow=Nsim, ncol=length(N_s))
for(j in 1:length(N_s)){
  for(i in 1:Nsim) {
    X <- matrix(1, nrow = N_s[j], ncol=2)
    X[,2] = rnorm(N_s[j])
    cos_theta = t(X[,1]) %*% X[,2]/(norm_vec(X[,1])*norm_vec(X[,2]))
```

```

        angles[i,j] = abs(90 - acos(cos_theta)*180/pi)
    }
}
colMeans(angles)

```

```

## [1] 44.951691 23.211393 15.381952 6.521264 4.590122 3.241769 2.052893
## [8] 1.446414

```

What is this absolute angle converging to? Why does this make sense?

The absolute angle difference from ninety is converging to zero. It makes sense because in a high dimensional space, random directions are orthogonal.

Create a vector y by simulating $n = 100$ standard iid normals. Create a matrix of size 100×2 and populate the first column by all ones (for the intercept) and the second column by 100 standard iid normals. Find the R^2 of an OLS regression of $y \sim X$. Use matrix algebra.

```

n=100
X = cbind(1,rnorm(n))
y = rnorm(n)
H = X %*% solve((t(X) %*% X)) %*% t(X)
y_hat = H %*% y
y_bar = mean(y)

SSR = sum((y_hat - y_bar)^2)
SST = sum((y - y_bar)^2)

Rsqr = (SSR/SST)
Rsqr

```

```
## [1] 0.005046253
```

Write a for loop to each time bind a new column of 100 standard iid normals to the matrix X and find the R^2 each time until the number of columns is 100. Create a vector to save all R^2 's. What happened??

```

Rsqr_s = array(NA, dim=n-2)
for(j in 1: (n-2)){
  X = cbind(X, rnorm(n))
  H = X %*% solve((t(X) %*% X)) %*% t(X)
  y_hat = H %*% y
  y_bar = mean(y)

  SSR = sum((y_hat - y_bar)^2)
  SST = sum((y - y_bar)^2)

  Rsqr_s[j] = (SSR/SST)
}
Rsqr_s

```

```

## [1] 0.005064451 0.017839889 0.021944438 0.026207482 0.088740527 0.151688029
## [7] 0.153065538 0.153296036 0.169859633 0.183578120 0.184395624 0.184435394
## [13] 0.193462424 0.201253453 0.205120194 0.205487164 0.209585357 0.215079434
## [19] 0.221006005 0.222709147 0.226030615 0.226062106 0.226160687 0.241040687
## [25] 0.247058590 0.249064242 0.259603656 0.260326225 0.272323744 0.273384501
## [31] 0.274527784 0.283598919 0.353730454 0.369389194 0.382246809 0.393449138
## [37] 0.419076168 0.437117736 0.437959375 0.449329906 0.462654030 0.462654498

```

```
## [43] 0.470955705 0.478498715 0.479625361 0.498803830 0.501595654 0.502291645
## [49] 0.502593370 0.504014185 0.504934413 0.524860380 0.524866239 0.526077389
## [55] 0.532833831 0.537551054 0.537584663 0.541250396 0.557602838 0.558007617
## [61] 0.563525529 0.563527129 0.563608159 0.573580147 0.589697657 0.590086611
## [67] 0.594535975 0.598913672 0.629273832 0.649999203 0.651192953 0.652588666
## [73] 0.671170834 0.673508267 0.680818284 0.781334353 0.789022914 0.830368100
## [79] 0.833242393 0.836109205 0.848339594 0.854605918 0.880927983 0.885600092
## [85] 0.894172260 0.898049624 0.898113362 0.899369616 0.900268632 0.951007075
## [91] 0.952571080 0.955410847 0.974665321 0.974745326 0.984240143 0.993687892
## [97] 0.994597951 1.000000000
```

Test that the projection matrix onto this X is the same as I_n . You may have to vectorize the matrices in the `expect_equal` function for the test to work.

```
pacman::p_load(testthat)
dim(X)
```

```
## [1] 100 100
H = X %*% solve(t(X) %*% X) %*% t(X)
I = diag(n)
expect_equal(H, I)
```

Add one final column to X to bring the number of columns to 101. Then try to compute R^2 . What happens?

```
{r}
X = cbind(X, rnorm(n))
dim(X)
H = X %*% solve((t(X) %*% X)) %*% t(X)
y_hat = H %*% y
y_bar = mean(y)
SSR = sum((y_hat - y_bar)^2)
SST = sum((y - y_bar)^2)
Rsqr = (SSR/SST)
Rsqr
```

Why does this make sense? The above chunk failed and this makes sense because $X^T X$ is rank deficient since we added another column, which is linearly dependent, to bring the total columns to 101.

Write a function spec'd as follows:

```
##' Orthogonal Projection
##'
##' Projects vector a onto v.
##'
##' @param a the vector to project
##' @param v the vector projected onto
##'
##' @returns a list of two vectors, the orthogonal projection parallel to v named a_parallel,
##' and the orthogonal error orthogonal to v called a_perpendicular
orthogonal_projection = function(a, v){
  H = v %*% t(v) / norm_vec(v)^2
  a_parallel = H %*% a
  a_perpendicular = a - a_parallel

  list(a_parallel = a_parallel, a_perpendicular = a_perpendicular)
```

```
}
```

Provide predictions for each of these computations and then run them to make sure you're correct.

```
orthogonal_projection(c(1,2,3,4), c(1,2,3,4))
```

```
## $a_parallel
##      [,1]
## [1,]    1
## [2,]    2
## [3,]    3
## [4,]    4
##
## $a_perpendicular
##      [,1]
## [1,]    0
## [2,]    0
## [3,]    0
## [4,]    0
```

```
#prediction: parallels will be itself or c(1,2,3,4) and perpendicular will be 0
orthogonal_projection(c(1, 2, 3, 4), c(0, 2, 0, -1))
```

```
## $a_parallel
##      [,1]
## [1,]    0
## [2,]    0
## [3,]    0
## [4,]    0
##
## $a_perpendicular
##      [,1]
## [1,]    1
## [2,]    2
## [3,]    3
## [4,]    4
```

```
#prediction: parallel will be 0 due to orthogonality and perpendicular will be c(1,2,3,4)
result = orthogonal_projection(c(2, 6, 7, 3), c(1, 3, 5, 7) *37)
t(result$a_parallel) %% result$a_perpendicular
```

```
##      [,1]
## [1,] -3.552714e-15
```

```
#prediction: will result in 0 because they are orthogonal
result$a_parallel + result$a_perpendicular
```

```
##      [,1]
## [1,]    2
## [2,]    6
## [3,]    7
## [4,]    3
```

```
#prediction: gives you back the original vector
result$a_parallel / (c(1, 3, 5, 7)*37)
```

```
##      [,1]
## [1,] 0.02445302
```

```
## [2,] 0.02445302
## [3,] 0.02445302
## [4,] 0.02445302
```

#prediction: percentage of the orthogonal projection

Let's use the Boston Housing Data for the following exercises

```
y = MASS::Boston$medv
X = model.matrix(medv ~ ., MASS::Boston)
p_plus_one = ncol(X)
n = nrow(X)
head(X)
```

```
## (Intercept)    crim zn indus chas   nox    rm  age    dis rad tax ptratio
## 1          1 0.00632 18  2.31    0 0.538 6.575 65.2 4.0900   1 296    15.3
## 2          1 0.02731  0  7.07    0 0.469 6.421 78.9 4.9671   2 242    17.8
## 3          1 0.02729  0  7.07    0 0.469 7.185 61.1 4.9671   2 242    17.8
## 4          1 0.03237  0  2.18    0 0.458 6.998 45.8 6.0622   3 222    18.7
## 5          1 0.06905  0  2.18    0 0.458 7.147 54.2 6.0622   3 222    18.7
## 6          1 0.02985  0  2.18    0 0.458 6.430 58.7 6.0622   3 222    18.7
##      black lstat
## 1 396.90  4.98
## 2 396.90  9.14
## 3 392.83  4.03
## 4 394.63  2.94
## 5 396.90  5.33
## 6 394.12  5.21
```

Using your function `orthogonal_projection` orthogonally project onto the column space of `X` by projecting `y` on each vector of `X` individually and adding up the projections and call the sum `yhat_naive`.

```
yhat_naive = rep(0,n)
for(j in 1:p_plus_one){
  yhat_naive = yhat_naive + orthogonal_projection(y,X[,j])$a_parallel
}
```

How much double counting occurred? Measure the magnitude relative to the true LS orthogonal projection.

```
yhat = X %*% solve((t(X) %*% X)) %*% t(X) %*% y
sqrt(sum(yhat_naive^2)) / sqrt(sum(yhat^2))
```

```
## [1] 8.997118
```

Is this ratio expected? Why or why not?

It is expect to be different from 1 because there is a bunch of double counting. Thus `y_hat naive` is not `y_hat`.

Convert `X` into `V` where `V` has the same column space as `X` but has orthogonal columns. You can use the function `orthogonal_projection`. This is the Gram-Schmidt orthogonalization algorithm.

```
V = matrix(NA, nrow = n, ncol = p_plus_one)
V[, 1] = X[, 1]
for (j in 2:p_plus_one) {
  V[,j]= X[,j]
  for(k in 1:(j-1)) {
    V[,j] = V[,j] - orthogonal_projection(X[,j], V[,k])$a_parallel
  }
}
```

```
}
V[,7] %*% V[,9]
```

```
##           [,1]
## [1,] -2.140346e-11
```

Convert V into Q whose columns are the same except normalized

```
Q = matrix(NA, nrow = n, ncol = p_plus_one)
for (j in 1:p_plus_one) {
  Q[,j] = V[,j] /norm_vec(V[,j])
}
```

Verify $Q^T Q$ is $I_{\{p+1\}}$ i.e. Q is an orthonormal matrix.

```
expect_equal(t(Q) %*% Q, diag(p_plus_one))
```

Is your Q the same as what results from R's built-in QR-decomposition function?

```
{r}
Q_from_Rs_builtin = qr.Q(qr(X))
expect_equal(Q, Q_from_Rs_builtin)
```

Is this expected? Why did this happen? Yes this is expected because there are an infinite number of orthonormal basis of any column space and the likelihood of them being equal is highly unlikely.

There are many different orthonormal basis of any column space.

Project y onto $\text{colsp}[Q]$ and verify it is the same as the OLS fit. You may have to use the function `unnname` to compare the vectors since they the entries will likely have different names.

```
y_hat = lm(y ~ X)$fitted.values
expect_equal(c(unnname(Q %*% t(Q) %*% y)), unnname(y_hat))
```

Project y onto $\text{colsp}[Q]$ one by one and verify it sums to be the projection onto the whole space.

```
yhat_naive = rep(0,n)

for (j in 1:p_plus_one){
  yhat_naive = yhat_naive + orthogonal_projection(y, Q[, j])$a_parallel
}

H = Q %*% solve(t(Q) %*% Q) %*% t(Q)
expect_equal(H %*% y, yhat_naive)
```

Split the Boston Housing Data into a training set and a test set where the training set is 80% of the observations. Do so at random.

```
K = 5
n_test = round(n * 1 / K)
n_train = n - n_test

test_indices = sample(1:n, 1/K * n)
train_indices = setdiff(1:n, test_indices)

X_train = X[train_indices, ]
y_train = y[train_indices]
X_test = X[test_indices,]
y_test = y[test_indices]
```

```
dim(X_train)
```

```
## [1] 405 14
```

```
dim(X_test)
```

```
## [1] 101 14
```

```
length(y_train)
```

```
## [1] 405
```

```
length(y_test)
```

```
## [1] 101
```

Fit an OLS model. Find the s_e in sample and out of sample. Which one is greater? Note: we are now using s_e and not RMSE since RMSE has the $n-(p+1)$ in the denominator not $n-1$ which attempts to de-bias the error estimate by inflating the estimate when overfitting in high p . Again, we're just using $sd(e)$, the sample standard deviation of the residuals.

```
mod = lm(y_train ~ ., data.frame(X_train))
summary(mod)$sigma
```

```
## [1] 4.8286
```

```
sd(mod$residuals)
```

```
## [1] 4.750277
```

Do these two exercises $N_{sim} = 1000$ times and find the average difference between s_e and $ooss_e$.

```
K = 5
n_test = round(n * 1 / K)
n_train = n - n_test
oosSSE_array = array(NA, dim = Nsim)
s_e_array = array(NA, dim = Nsim)
Nsim = 10000

for(i in 1:Nsim){

  test_indices = sample(1 : n, n_test)
  train_indices = setdiff(1 : n, test_indices)

  X_train = X[train_indices, ]
  y_train = y[train_indices]
  X_test = X[test_indices, ]
  y_test = y[test_indices]

  mod = lm(y_train ~ . + 0, data.frame(X_train))
  y_hat_test = predict(mod, data.frame(X_test))
  oosSSE_array[i] = sd(y_test - y_hat_test)
  s_e_array[i] = sd(mod$residuals)
}

mean(s_e_array - oosSSE_array)
```

```
## [1] NA
```

We'll now add random junk to the data so that $p_{plus_one} = n_{train}$ and create a new data matrix `X_with_junk`.

```
X_with_junk = cbind(X, matrix(rnorm(n * (n_train - p_plus_one)), nrow = n))
dim(X)
```

```
## [1] 506 14
```

```
dim(X_with_junk)
```

```
## [1] 506 405
```

Repeat the exercise above measuring the average `s_e` and `ooss_e` but this time record these metrics by number of features used. That is, do it for the first column of `X_with_junk` (the intercept column), then do it for the first and second columns, then the first three columns, etc until you do it for all columns of `X_with_junk`. Save these in `s_e_by_p` and `ooss_e_by_p`.

```
K = 5 # The test set is one fifth of the entire historical dataset
```

```
n_test = round(n * 1 / K)
```

```
n_train = n - n_test
```

```
ooss_e_by_p = array(NA, dim = ncol(X_with_junk))
```

```
s_e_by_p = array(NA, dim = ncol(X_with_junk))
```

```
Nsim = 100
```

```
for(j in 1:ncol(X_with_junk)){
  oosSSE_array = array(NA, dim = Nsim)
  s_e_array = array(NA, dim = Nsim)
  for(i in 1:Nsim){

    test_indices = sample(1 : n, n_test)
    train_indices = setdiff(1 : n, test_indices)

    X_train = X_with_junk[train_indices, 1:j, drop = FALSE ]
    y_train = y[train_indices]
    X_test = X_with_junk[test_indices, 1:j, drop = FALSE]
    y_test = y[test_indices]

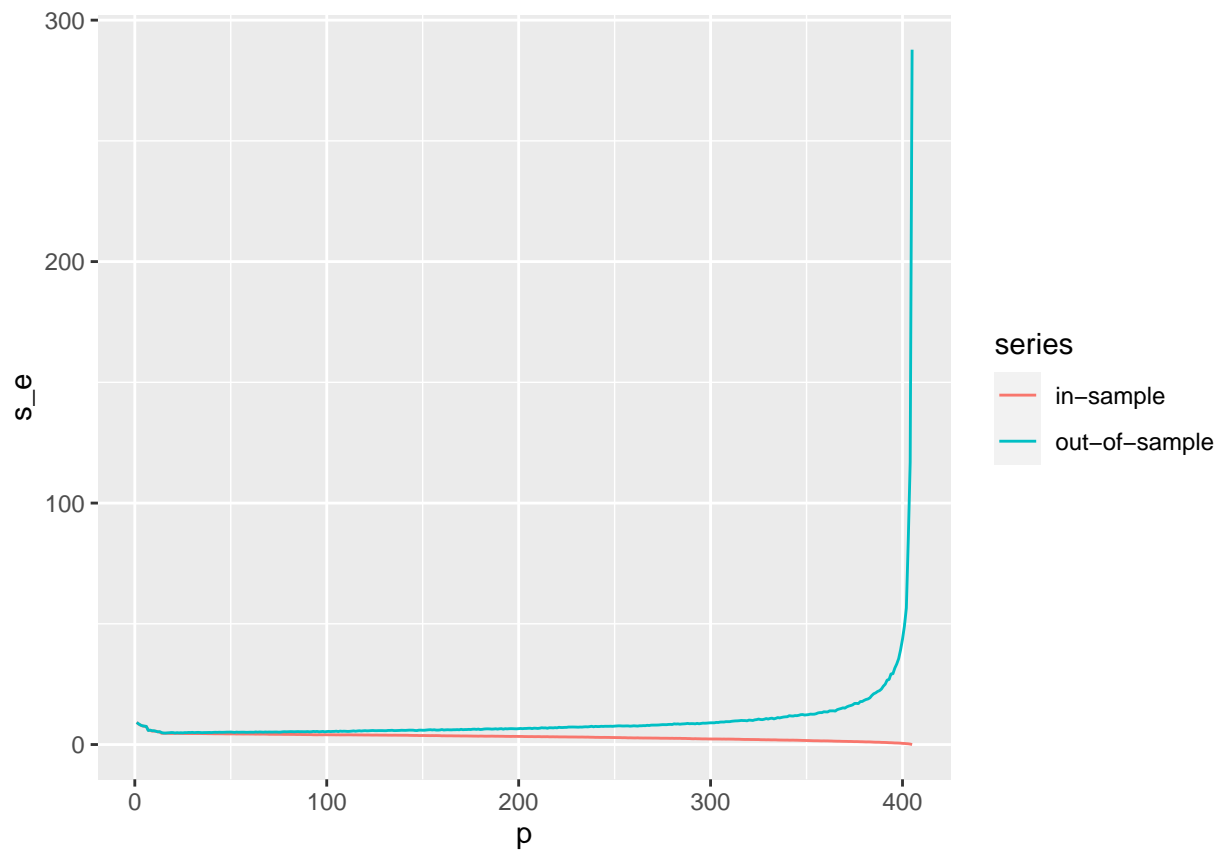
    mod = lm(y_train ~ . + 0, data.frame(X_train))
    y_hat_test = predict(mod, data.frame(X_test))
    oosSSE_array[i] = sd(y_test - y_hat_test)
    s_e_array[i] = sd(mod$residuals) #s_e

  }

  ooss_e_by_p[j] = mean(oosSSE_array)
  s_e_by_p[j] = mean(s_e_array)
}
```

You can graph them here:

```
pacman::p_load(ggplot2)
ggplot(
  rbind(
    data.frame(s_e = s_e_by_p, p = 1 : 1 : n_train, series = "in-sample"),
    data.frame(s_e = ooss_e_by_p, p = 1 : n_train, series = "out-of-sample")
  ) +
  geom_line(aes(x = p, y = s_e, col = series))
```

Is this shape expected? Explain.

Yes this shape is expected because as we add more features the in-sample error will decrease due to the the model fitting the additional features and data. However, the out of sample error will start to get worse due to the over-fitting that is occurring. This will lead to a worse model that produces worse predictions for data that is out of sample.