

# Lab 5

*Tziporah Horowitz*

*11:59PM March 7, 2020*

Load the Boston housing data frame and create the vector  $y$  (the median value) and matrix  $X$  (all other features) from the data frame. Name the columns the same as Boston except for the first name it “(Intercept)”.

```
boston <- MASS::Boston
head(boston)
```

```
##      crim zn indus chas   nox    rm  age    dis rad tax ptratio  black
## 1 0.00632 18  2.31    0 0.538 6.575 65.2 4.0900   1 296    15.3 396.90
## 2 0.02731  0  7.07    0 0.469 6.421 78.9 4.9671   2 242    17.8 396.90
## 3 0.02729  0  7.07    0 0.469 7.185 61.1 4.9671   2 242    17.8 392.83
## 4 0.03237  0  2.18    0 0.458 6.998 45.8 6.0622   3 222    18.7 394.63
## 5 0.06905  0  2.18    0 0.458 7.147 54.2 6.0622   3 222    18.7 396.90
## 6 0.02985  0  2.18    0 0.458 6.430 58.7 6.0622   3 222    18.7 394.12
##   lstat medv
## 1   4.98 24.0
## 2   9.14 21.6
## 3   4.03 34.7
## 4   2.94 33.4
## 5   5.33 36.2
## 6   5.21 28.7
```

```
y <- boston$medv
X <- boston[, 1:13]
```

Run the OLS linear model to get  $b$ , the vector of coefficients. Do not use `lm`.

```
X <- cbind(1, as.matrix(X))
head(X)
```

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

```
b <- solve(t(X) %*% X) %*% t(X) %*% y
```

Find the hat matrix for this regression H. Verify its dimension is correct and verify its rank is correct.

```
H <- X %*% solve(t(X) %*% X) %*% t(X)
dim(H)
```

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

```
pacman::p_load(Matrix)
rankMatrix(X)
```

```
## [1] 14
## attr("method")
## [1] "tolNorm2"
## attr("useGrad")
## [1] FALSE
## attr("tol")
## [1] 1.123546e-13
```

Verify this is a projection matrix by verifying the two sufficient conditions. Use the `testthat` library's `expect_equal(matrix1, matrix2, tolerance = 1e-2)`.

```
pacman::p_load(testthat)
expect_equal(H %*% H, H)
expect(H, t(H))
```

Find the matrix that projects onto the space of residuals `Hcomp` and find its rank. Is this rank expected?

```
Hcomp <- diag(nrow(H)) - H
rankMatrix(H)
```

```
## [1] 14
## attr("method")
## [1] "tolNorm2"
## attr("useGrad")
## [1] FALSE
## attr("tol")
## [1] 1.123546e-13
```

```
rankMatrix(Hcomp, tol = 1e-2)
```

```
## [1] 492
## attr("method")
## [1] "tolNorm2"
## attr("useGrad")
## [1] FALSE
## attr("tol")
## [1] 0.01
```

Verify this is a projection matrix by verifying the two sufficient conditions. Use the `testthat` library.

```
expect_equal(Hcomp %*% Hcomp, Hcomp)
expect(Hcomp, t(Hcomp))
```

Use `diag` to find the trace of both `H` and `Hcomp`.

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

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

```
sum(diag(Hcomp))
```

```
## [1] 492
```

Do you have a conjecture about the trace of an orthogonal projection matrix?

The trace is equal to the rank.

Find the eigendecomposition of both `H` and `Hcomp` as `eigenvals_H`, `eigenvecs_H`, `eigenvals_Hcomp`, `eigenvecs_Hcomp`. Verify these results are the correct dimensions.

```
eigen_H <- eigen(H)
eigen_Hcomp <- eigen(Hcomp)

eigenvals_H <- eigen_H$values
eigenvecs_H <- eigen_H$vectors
eigenvals_Hcomp <- eigen_Hcomp$values
eigenvecs_Hcomp <- eigen_Hcomp$vectors

length(eigenvals_H)
```

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

```
dim(eigenvecs_H)
```

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

```
length(eigenvals_Hcomp)
```

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

```
dim(eigenvecs_Hcomp)
```

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

The eigendecomposition suffers from numerical error which is making them become imaginary. We can coerce imaginary numbers back to real by using the `abs` function. There is also lots of numerical error. Use the `as.numeric` function to coerce to real and the `round` function to round all four objects to the nearest 10 digits.

```
eigenvals_H <- round(Re(eigenvals_H), 10)
eigenvecs_H <- round(Re(eigenvecs_H), 10)
eigenvals_Hcomp <- round(Re(eigenvals_Hcomp), 10)
eigenvecs_Hcomp <- round(Re(eigenvecs_Hcomp), 10)
```

Print out the eigenvalues of both  $\mathbf{H}$  and  $\mathbf{H}_{\text{comp}}$ . Is this expected?

eigenvals\_H

[illegible]

eigenvals\_Hcomp

[illegible]

Find the length of all eigenvectors of `H` in one line. Is this expected? What is the convention for eigenvectors in R's `eigen` function?

```
apply(eigenvecs_H, MARGIN = 2, FUN = function(v){
  sqrt(sum(v^2))
})
```

```
##      [1] 1.0000000 1.0000000 1.0000000 1.0000000 1.0000000 1.0000000 1.0000000
```

```

## [8] 1.0000000 1.0000000 1.0000000 1.0000000 1.0000000 1.0000000 1.0000000
## [15] 0.7198274 0.7198274 0.7111798 0.7111798 0.6940586 0.6940586 1.0000000
## [22] 0.7122684 0.7122684 0.7192493 0.7192493 0.7847285 0.7847285 0.8376704
## [29] 0.8376704 0.8282487 0.8282487 0.7586616 0.7586616 0.9528857 0.9528857
## [36] 0.6589249 0.6589249 0.7628320 0.7628320 0.6864124 0.6864124 0.6944800
## [43] 0.6944800 0.8044639 0.8044639 0.9385394 0.9385394 0.7220811 0.7220811
## [50] 0.7623481 0.7623481 0.7133975 0.7133975 0.7339666 0.7339666 0.7071697
## [57] 0.7071697 0.6972923 0.6972923 0.6866897 0.6866897 0.7156528 0.7156528
## [64] 0.7069016 0.7069016 0.7091819 0.7091819 0.7043276 0.7043276 0.7375631
## [71] 0.7375631 0.7200799 0.7200799 0.7067433 0.7067433 1.0000000 0.7291502
## [78] 0.7291502 0.9861821 0.9861821 0.7181425 0.7181425 0.7289067 0.7289067
## [85] 0.7069350 0.7069350 0.7147271 0.7147271 0.7278385 0.7278385 0.7016417
## [92] 0.7016417 0.6985773 0.6985773 0.6978982 0.6978982 0.6983338 0.6983338
## [99] 0.7001943 0.7001943 0.7051763 0.7051763 0.7555366 0.7555366 1.0000000
## [106] 0.6876725 0.6876725 0.7150475 0.7150475 0.7486783 0.7486783 0.7261142
## [113] 0.7261142 0.6985145 0.6985145 0.6836561 0.6836561 1.0000000 0.7187066
## [120] 0.7187066 0.6905980 0.6905980 0.7179283 0.7179283 0.7222461 0.7222461
## [127] 0.7207965 0.7207965 0.7004181 0.7004181 0.7014897 0.7014897 0.6990003
## [134] 0.6990003 0.7115622 0.7115622 0.7237343 0.7237343 0.7703283 0.7703283
## [141] 0.7899225 0.7899225 0.7259246 0.7259246 0.7076892 0.7076892 0.7727441
## [148] 0.7727441 0.7179819 0.7179819 0.7213128 0.7213128 0.6696315 0.6696315
## [155] 0.7021261 0.7021261 0.7045330 0.7045330 0.7363510 0.7363510 1.0000000
## [162] 0.6941660 0.6941660 0.7753994 0.7753994 0.7540773 0.7540773 0.7070886
## [169] 0.7070886 0.6903144 0.6903144 0.6988027 0.6988027 0.7120509 0.7120509
## [176] 0.7158204 0.7158204 0.7312673 0.7312673 1.0000000 0.6920332 0.6920332
## [183] 0.6866946 0.6866946 0.7373332 0.7373332 0.7270449 0.7270449 0.7073476
## [190] 0.7073476 0.7085796 0.7085796 0.7045089 0.7045089 0.7563864 0.7563864
## [197] 0.7618808 0.7618808 0.7044355 0.7044355 0.8477404 0.8477404 0.7109985
## [204] 0.7109985 0.6987366 0.6987366 0.7708438 0.7708438 0.6986430 0.6986430
## [211] 0.7336177 0.7336177 1.0000000 0.6825439 0.6825439 0.6940425 0.6940425
## [218] 0.7314948 0.7314948 0.7171839 0.7171839 0.7281819 0.7281819 1.0000000
## [225] 0.7031402 0.7031402 0.6951803 0.6951803 0.7152899 0.7152899 0.7281453
## [232] 0.7281453 0.7160614 0.7160614 0.7208316 0.7208316 0.7089619 0.7089619
## [239] 0.7120048 0.7120048 0.7095858 0.7095858 0.8786539 0.8786539 0.6955125
## [246] 0.6955125 0.7157227 0.7157227 1.0000000 0.7616379 0.7616379 0.7383818
## [253] 0.7383818 0.7198086 0.7198086 0.7054253 0.7054253 0.7185741 0.7185741
## [260] 0.7113626 0.7113626 0.7210185 0.7210185 0.7061893 0.7061893 0.9134175
## [267] 0.9134175 0.7129986 0.7129986 0.7232002 0.7232002 0.7393014 0.7393014
## [274] 0.7209286 0.7209286 0.7185819 0.7185819 0.7586169 0.7586169 0.7560925
## [281] 0.7560925 0.7144534 0.7144534 0.7047713 0.7047713 0.7334610 0.7334610
## [288] 0.7207679 0.7207679 0.6943252 0.6943252 1.0000000 0.7289015 0.7289015
## [295] 0.7186150 0.7186150 0.7316530 0.7316530 0.8856238 0.8856238 0.7071873
## [302] 0.7071873 0.7283910 0.7283910 0.7209085 0.7209085 0.7375309 0.7375309
## [309] 0.7114813 0.7114813 0.7944529 0.7944529 0.7303951 0.7303951 0.7186130
## [316] 0.7186130 0.8186602 0.8186602 0.6977128 0.6977128 0.7075066 0.7075066
## [323] 0.7688476 0.7688476 0.6874561 0.6874561 0.7118374 0.7118374 1.0000000
## [330] 0.6779750 0.6779750 0.7146818 0.7146818 0.7260265 0.7260265 0.7126152
## [337] 0.7126152 1.0000000 0.7150297 0.7150297 0.7228253 0.7228253 0.7251544
## [344] 0.7251544 0.7112836 0.7112836 0.7131162 0.7131162 0.7254483 0.7254483
## [351] 0.7059673 0.7059673 0.7087931 0.7087931 0.7081581 0.7081581 0.7231071
## [358] 0.7231071 0.7004617 0.7004617 0.8706722 0.8706722 0.7317309 0.7317309
## [365] 0.7945923 0.7945923 1.0000000 0.7322846 0.7322846 0.7326250 0.7326250
## [372] 0.7174917 0.7174917 0.7517336 0.7517336 0.7375921 0.7375921 0.6711866
## [379] 0.6711866 0.6895559 0.6895559 0.6880716 0.6880716 0.7022134 0.7022134

```

```
## [386] 0.7630476 0.7630476 0.7072201 0.7072201 0.7656041 0.7656041 0.6909803
## [393] 0.6909803 0.6955879 0.6955879 0.7054902 0.7054902 0.8115084 0.8115084
## [400] 0.6961595 0.6961595 1.0000000 0.6879843 0.6879843 0.7306674 0.7306674
## [407] 0.8083010 0.8083010 0.8121622 0.8121622 0.6934602 0.6934602 0.7164527
## [414] 0.7164527 0.7074909 0.7074909 0.7179766 0.7179766 0.7076450 0.7076450
## [421] 0.7497639 0.7497639 1.0000000 0.8140732 0.8140732 0.7093069 0.7093069
## [428] 0.9117919 0.9117919 0.9473980 0.9473980 0.7074843 0.7074843 0.7350700
## [435] 0.7350700 0.7417871 0.7417871 0.7547830 0.7547830 0.7366121 0.7366121
## [442] 0.7769113 0.7769113 0.7810845 0.7810845 0.6892008 0.6892008 0.8055246
## [449] 0.8055246 0.6814249 0.6814249 0.7272492 0.7272492 0.7020231 0.7020231
## [456] 0.7611428 0.7611428 0.7367486 0.7367486 0.7048714 0.7048714 0.7400683
## [463] 0.7400683 1.0000000 0.7952697 0.7952697 0.6912670 0.6912670 0.8435088
## [470] 0.8435088 0.7646925 0.7646925 0.7883719 0.7883719 0.8329977 0.8329977
## [477] 0.7629826 0.7629826 1.0000000 0.7070124 0.7070124 0.7649891 0.7649891
## [484] 0.6532954 0.6532954 0.7107735 0.7107735 1.0000000 0.7442381 0.7442381
## [491] 1.0000000 0.7999522 0.7999522 1.0000000 0.7369916 0.7369916 0.7339597
## [498] 0.7339597 0.9876976 0.9876976 0.9150144 0.9150144 0.8738908 0.8738908
## [505] 1.0000000 1.0000000
```

The convention is length 1.

The first  $p+1$  eigenvectors are the columns of  $X$  but they are in arbitrary order. Find the column that represents the one-vector.

```
head(eigenvecs_H[, 3])
```

```
## [1] 0.04561114 0.04803979 0.04756035 0.04872081 0.04788055 0.04985263
```

Why is it not exactly 506 1's?

Numeric error.

Use the first  $p+1$  eigenvectors as a model matrix and run the OLS model of medv on that model matrix.

```
mod1 <- lm(y ~ 0 + X)
mod2 <- lm(y ~ eigenvecs_H[, 1:14])
coef(mod1)
```

```
##           X           Xcrim           Xzn           Xindus           Xchas
## 3.645949e+01 -1.080114e-01 4.642046e-02 2.055863e-02 2.686734e+00
##           Xnox           Xrm           Xage           Xdis           Xrad
## -1.776661e+01 3.809865e+00 6.922246e-04 -1.475567e+00 3.060495e-01
##           Xtax           Xptratio           Xblack           Xlstat
## -1.233459e-02 -9.527472e-01 9.311683e-03 -5.247584e-01
```

```
coef(mod2)
```

```
##           (Intercept) eigenvecs_H[, 1:14]1 eigenvecs_H[, 1:14]2
##           287393655           123831293           -290861139
## eigenvecs_H[, 1:14]3 eigenvecs_H[, 1:14]4 eigenvecs_H[, 1:14]5
##           -6559216960           604392624           96659807
## eigenvecs_H[, 1:14]6 eigenvecs_H[, 1:14]7 eigenvecs_H[, 1:14]8
##           -299961037           271606996           16309421
```

```
## eigenvecs_H[, 1:14]9 eigenvecs_H[, 1:14]10 eigenvecs_H[, 1:14]11
##          -969027434          -8130072          -203027856
## eigenvecs_H[, 1:14]12 eigenvecs_H[, 1:14]13 eigenvecs_H[, 1:14]14
##          155231970          357020877          -31894038
```

Is  $\mathbf{b}$  about the same above (in arbitrary order)?

No,  $\mathbf{b}$  is now scaled by the eigendecomposition.

Calculate  $\hat{\mathbf{y}}$  using the hat matrix.

```
yhat <- H %*% y
yhat
```

```
##          [,1]
## 1  30.0038434
## 2  25.0255624
## 3  30.5675967
## 4  28.6070365
## 5  27.9435242
## 6  25.2562845
## 7  23.0018083
## 8  19.5359884
## 9  11.5236369
## 10 18.9202621
## 11 18.9994965
## 12 21.5867957
## 13 20.9065215
## 14 19.5529028
## 15 19.2834821
## 16 19.2974832
## 17 20.5275098
## 18 16.9114013
## 19 16.1780111
## 20 18.4061360
## 21 12.5238575
## 22 17.6710367
## 23 15.8328813
## 24 13.8062853
## 25 15.6783383
## 26 13.3866856
## 27 15.4639765
## 28 14.7084743
## 29 19.5473729
## 30 20.8764282
## 31 11.4551176
## 32 18.0592329
## 33  8.8110574
## 34 14.2827581
## 35 13.7067589
## 36 23.8146353
## 37 22.3419371
## 38 23.1089114
## 39 22.9150261
```

## 40 31.3576257  
## 41 34.2151023  
## 42 28.0205641  
## 43 25.2038663  
## 44 24.6097927  
## 45 22.9414918  
## 46 22.0966982  
## 47 20.4232003  
## 48 18.0365509  
## 49 9.1065538  
## 50 17.2060775  
## 51 21.2815254  
## 52 23.9722228  
## 53 27.6558508  
## 54 24.0490181  
## 55 15.3618477  
## 56 31.1526495  
## 57 24.8568698  
## 58 33.1091981  
## 59 21.7753799  
## 60 21.0849356  
## 61 17.8725804  
## 62 18.5111021  
## 63 23.9874286  
## 64 22.5540887  
## 65 23.3730864  
## 66 30.3614836  
## 67 25.5305651  
## 68 21.1133856  
## 69 17.4215379  
## 70 20.7848363  
## 71 25.2014886  
## 72 21.7426577  
## 73 24.5574496  
## 74 24.0429571  
## 75 25.5049972  
## 76 23.9669302  
## 77 22.9454540  
## 78 23.3569982  
## 79 21.2619827  
## 80 22.4281737  
## 81 28.4057697  
## 82 26.9948609  
## 83 26.0357630  
## 84 25.0587348  
## 85 24.7845667  
## 86 27.7904920  
## 87 22.1685342  
## 88 25.8927642  
## 89 30.6746183  
## 90 30.8311062  
## 91 27.1190194  
## 92 27.4126673  
## 93 28.9412276



## 94 29.0810555  
## 95 27.0397736  
## 96 28.6245995  
## 97 24.7274498  
## 98 35.7815952  
## 99 35.1145459  
## 100 32.2510280  
## 101 24.5802202  
## 102 25.5941347  
## 103 19.7901368  
## 104 20.3116713  
## 105 21.4348259  
## 106 18.5399401  
## 107 17.1875599  
## 108 20.7504903  
## 109 22.6482911  
## 110 19.7720367  
## 111 20.6496586  
## 112 26.5258674  
## 113 20.7732364  
## 114 20.7154831  
## 115 25.1720888  
## 116 20.4302559  
## 117 23.3772463  
## 118 23.6904326  
## 119 20.3357836  
## 120 20.7918087  
## 121 21.9163207  
## 122 22.4710778  
## 123 20.5573856  
## 124 16.3666198  
## 125 20.5609982  
## 126 22.4817845  
## 127 14.6170663  
## 128 15.1787668  
## 129 18.9386859  
## 130 14.0557329  
## 131 20.0352740  
## 132 19.4101340  
## 133 20.0619157  
## 134 15.7580767  
## 135 13.2564524  
## 136 17.2627773  
## 137 15.8784188  
## 138 19.3616395  
## 139 13.8148390  
## 140 16.4488147  
## 141 13.5714193  
## 142 3.9888551  
## 143 14.5949548  
## 144 12.1488148  
## 145 8.7282236  
## 146 12.0358534  
## 147 15.8208206

## 148 8.5149902  
## 149 9.7184414  
## 150 14.8045137  
## 151 20.8385815  
## 152 18.3010117  
## 153 20.1228256  
## 154 17.2860189  
## 155 22.3660023  
## 156 20.1037592  
## 157 13.6212589  
## 158 33.2598270  
## 159 29.0301727  
## 160 25.5675277  
## 161 32.7082767  
## 162 36.7746701  
## 163 40.5576584  
## 164 41.8472817  
## 165 24.7886738  
## 166 25.3788924  
## 167 37.2034745  
## 168 23.0874875  
## 169 26.4027396  
## 170 26.6538211  
## 171 22.5551466  
## 172 24.2908281  
## 173 22.9765722  
## 174 29.0719431  
## 175 26.5219434  
## 176 30.7220906  
## 177 25.6166931  
## 178 29.1374098  
## 179 31.4357197  
## 180 32.9223157  
## 181 34.7244046  
## 182 27.7655211  
## 183 33.8878732  
## 184 30.9923804  
## 185 22.7182001  
## 186 24.7664781  
## 187 35.8849723  
## 188 33.4247672  
## 189 32.4119915  
## 190 34.5150995  
## 191 30.7610949  
## 192 30.2893414  
## 193 32.9191871  
## 194 32.1126077  
## 195 31.5587100  
## 196 40.8455572  
## 197 36.1277008  
## 198 32.6692081  
## 199 34.7046912  
## 200 30.0934516  
## 201 30.6439391

## 202 29.2871950  
## 203 37.0714839  
## 204 42.0319312  
## 205 43.1894984  
## 206 22.6903480  
## 207 23.6828471  
## 208 17.8544721  
## 209 23.4942899  
## 210 17.0058772  
## 211 22.3925110  
## 212 17.0604275  
## 213 22.7389292  
## 214 25.2194255  
## 215 11.1191674  
## 216 24.5104915  
## 217 26.6033477  
## 218 28.3551871  
## 219 24.9152546  
## 220 29.6865277  
## 221 33.1841975  
## 222 23.7745666  
## 223 32.1405196  
## 224 29.7458199  
## 225 38.3710245  
## 226 39.8146187  
## 227 37.5860575  
## 228 32.3995325  
## 229 35.4566524  
## 230 31.2341151  
## 231 24.4844923  
## 232 33.2883729  
## 233 38.0481048  
## 234 37.1632863  
## 235 31.7138352  
## 236 25.2670557  
## 237 30.1001074  
## 238 32.7198716  
## 239 28.4271706  
## 240 28.4294068  
## 241 27.2937594  
## 242 23.7426248  
## 243 24.1200789  
## 244 27.4020841  
## 245 16.3285756  
## 246 13.3989126  
## 247 20.0163878  
## 248 19.8618443  
## 249 21.2883131  
## 250 24.0798915  
## 251 24.2063355  
## 252 25.0421582  
## 253 24.9196401  
## 254 29.9456337  
## 255 23.9722832

## 256 21.6958089  
## 257 37.5110924  
## 258 43.3023904  
## 259 36.4836142  
## 260 34.9898859  
## 261 34.8121151  
## 262 37.1663133  
## 263 40.9892850  
## 264 34.4463409  
## 265 35.8339755  
## 266 28.2457430  
## 267 31.2267359  
## 268 40.8395575  
## 269 39.3179239  
## 270 25.7081791  
## 271 22.3029553  
## 272 27.2034097  
## 273 28.5116947  
## 274 35.4767660  
## 275 36.1063916  
## 276 33.7966827  
## 277 35.6108586  
## 278 34.8399338  
## 279 30.3519266  
## 280 35.3098070  
## 281 38.7975697  
## 282 34.3312319  
## 283 40.3396307  
## 284 44.6730834  
## 285 31.5968909  
## 286 27.3565923  
## 287 20.1017415  
## 288 27.0420667  
## 289 27.2136458  
## 290 26.9139584  
## 291 33.4356331  
## 292 34.4034963  
## 293 31.8333982  
## 294 25.8178324  
## 295 24.4298235  
## 296 28.4576434  
## 297 27.3626700  
## 298 19.5392876  
## 299 29.1130984  
## 300 31.9105461  
## 301 30.7715945  
## 302 28.9427587  
## 303 28.8819102  
## 304 32.7988723  
## 305 33.2090546  
## 306 30.7683179  
## 307 35.5622686  
## 308 32.7090512  
## 309 28.6424424

## 310 23.5896583  
## 311 18.5426690  
## 312 26.8788984  
## 313 23.2813398  
## 314 25.5458025  
## 315 25.4812006  
## 316 20.5390990  
## 317 17.6157257  
## 318 18.3758169  
## 319 24.2907028  
## 320 21.3252904  
## 321 24.8868224  
## 322 24.8693728  
## 323 22.8695245  
## 324 19.4512379  
## 325 25.1178340  
## 326 24.6678691  
## 327 23.6807618  
## 328 19.3408962  
## 329 21.1741811  
## 330 24.2524907  
## 331 21.5926089  
## 332 19.9844661  
## 333 23.3388800  
## 334 22.1406069  
## 335 21.5550993  
## 336 20.6187291  
## 337 20.1609718  
## 338 19.2849039  
## 339 22.1667232  
## 340 21.2496577  
## 341 21.4293931  
## 342 30.3278880  
## 343 22.0473498  
## 344 27.7064791  
## 345 28.5479412  
## 346 16.5450112  
## 347 14.7835964  
## 348 25.2738008  
## 349 27.5420512  
## 350 22.1483756  
## 351 20.4594409  
## 352 20.5460542  
## 353 16.8806383  
## 354 25.4025351  
## 355 14.3248663  
## 356 16.5948846  
## 357 19.6370469  
## 358 22.7180661  
## 359 22.2021889  
## 360 19.2054806  
## 361 22.6661611  
## 362 18.9319262  
## 363 18.2284680

## 364 20.2315081  
## 365 37.4944739  
## 366 14.2819073  
## 367 15.5428625  
## 368 10.8316232  
## 369 23.8007290  
## 370 32.6440736  
## 371 34.6068404  
## 372 24.9433133  
## 373 25.9998091  
## 374 6.1263250  
## 375 0.7777981  
## 376 25.3071306  
## 377 17.7406106  
## 378 20.2327441  
## 379 15.8333130  
## 380 16.8351259  
## 381 14.3699483  
## 382 18.4768283  
## 383 13.4276828  
## 384 13.0617751  
## 385 3.2791812  
## 386 8.0602217  
## 387 6.1284220  
## 388 5.6186481  
## 389 6.4519857  
## 390 14.2076474  
## 391 17.2122518  
## 392 17.2988727  
## 393 9.8911664  
## 394 20.2212419  
## 395 17.9418118  
## 396 20.3044578  
## 397 19.2955908  
## 398 16.3363278  
## 399 6.5516232  
## 400 10.8901678  
## 401 11.8814587  
## 402 17.8117451  
## 403 18.2612659  
## 404 12.9794878  
## 405 7.3781636  
## 406 8.2111586  
## 407 8.0662619  
## 408 19.9829479  
## 409 13.7075637  
## 410 19.8526845  
## 411 15.2230830  
## 412 16.9607198  
## 413 1.7185181  
## 414 11.8057839  
## 415 -4.2813107  
## 416 9.5837674  
## 417 13.3666081

## 418 6.8956236  
## 419 6.1477985  
## 420 14.6066179  
## 421 19.6000267  
## 422 18.1242748  
## 423 18.5217713  
## 424 13.1752861  
## 425 14.6261762  
## 426 9.9237498  
## 427 16.3459065  
## 428 14.0751943  
## 429 14.2575624  
## 430 13.0423479  
## 431 18.1595569  
## 432 18.6955435  
## 433 21.5272830  
## 434 17.0314186  
## 435 15.9609044  
## 436 13.3614161  
## 437 14.5207938  
## 438 8.8197601  
## 439 4.8675110  
## 440 13.0659131  
## 441 12.7060970  
## 442 17.2955806  
## 443 18.7404850  
## 444 18.0590103  
## 445 11.5147468  
## 446 11.9740036  
## 447 17.6834462  
## 448 18.1269524  
## 449 17.5183465  
## 450 17.2274251  
## 451 16.5227163  
## 452 19.4129110  
## 453 18.5821524  
## 454 22.4894479  
## 455 15.2800013  
## 456 15.8208934  
## 457 12.6872558  
## 458 12.8763379  
## 459 17.1866853  
## 460 18.5124761  
## 461 19.0486053  
## 462 20.1720893  
## 463 19.7740732  
## 464 22.4294077  
## 465 20.3191185  
## 466 17.8861625  
## 467 14.3747852  
## 468 16.9477685  
## 469 16.9840576  
## 470 18.5883840  
## 471 20.1671944

```
## 472 22.9771803
## 473 22.4558073
## 474 25.5782463
## 475 16.3914763
## 476 16.1114628
## 477 20.5348160
## 478 11.5427274
## 479 19.2049630
## 480 21.8627639
## 481 23.4687887
## 482 27.0988732
## 483 28.5699430
## 484 21.0839878
## 485 19.4551620
## 486 22.2222591
## 487 19.6559196
## 488 21.3253610
## 489 11.8558372
## 490 8.2238669
## 491 3.6639967
## 492 13.7590854
## 493 15.9311855
## 494 20.6266205
## 495 20.6124941
## 496 16.8854196
## 497 14.0132079
## 498 19.1085414
## 499 21.2980517
## 500 18.4549884
## 501 20.4687085
## 502 23.5333405
## 503 22.3757189
## 504 27.6274261
## 505 26.1279668
## 506 22.3442123
```

Calculate  $e$  two ways: (1) the difference of  $y$  and  $\hat{y}$  and (2) the projection onto the space of the residuals. Verify the two means of calculating the residuals provide the same results via `expect_equal`.

```
e2 <- Hcomp %*% y
e1 <- y - yhat
expect_equal(e1, e2)
```

Calculate  $R^2$  using the angle relationship between the responses and their predictions.

```
len_vec <- function(v){sqrt(sum(v^2))}

y_avg_adj <- y - mean(y)
h_yhat_adj <- yhat - mean(y)

cos(sum(y * yhat)/(len_vec(y) * len_vec(yhat)))

## [1] 0.5559107
```



Find the cosine-squared of  $y - \bar{y}$  and  $\hat{y} - \bar{y}$  and verify it is the same as  $R^2$ . This empirically demonstrates what I missed in class: that the angle between  $y$  and  $\hat{y}$  is equal to the angle between  $y - \bar{y}$  and  $\hat{y} - \bar{y}$ .

```
summary(mod1)$r.squared
```

```
## [1] 0.9630247
```

Verify  $\hat{y}$  and  $e$  are orthogonal.

```
sum(yhat * e1)
```

```
## [1] -4.991219e-08
```

Verify  $\hat{y} - \bar{y}$  and  $e$  are orthogonal.

```
sum((yhat - mean(y)) * e1)
```

```
## [1] 2.832455e-09
```

Verify the sum of squares identity which we learned was due to the Pythagorean Theorem (applies since the projection is specifically orthogonal). You need to compute all three quantities first: SST, SSR and SSE.

```
SST <- len_vec(y_avg_adj)^2
SSR <- len_vec(h_yhat_adj)^2
SSE <- len_vec(e1)^2

expect_equal(SST, SSR + SSE)
```

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.

```
M <- matrix(NA, ncol(X), ncol(X))
colnames(M) <- colnames(X)

for (j in 1:ncol(X)){
  X_j <- X[, 1:j, drop = FALSE]
  b <- solve(t(X_j) %*% X_j) %*% t(X_j) %*% y
  M[j, 1:j] <- b
}

M
```

```
##           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
```

Examine this matrix. Why are the estimates changing from row to row as you add in more predictors?

Because the weight of the missing predictors is included in the present predictors.

Clear the workspace and load the diamonds dataset in the package `ggplot2`.

```
rm(list = ls())

pacman::p_load(ggplot2)
data("diamonds")
head(diamonds)
```

```
## # A tibble: 6 x 10
```

```
##   carat cut      color clarity depth table price      x      y      z
##   <dbl> <ord>      <ord> <ord>   <dbl> <dbl> <int> <dbl> <dbl> <dbl>
## 1 0.23 Ideal      E      SI2     61.5   55   326  3.95  3.98  2.43
## 2 0.21 Premium    E      SI1     59.8   61   326  3.89  3.84  2.31
## 3 0.23 Good       E      VS1     56.9   65   327  4.05  4.07  2.31
## 4 0.290 Premium   I      VS2     62.4   58   334  4.2   4.23  2.63
## 5 0.31 Good       J      SI2     63.3   58   335  4.34  4.35  2.75
## 6 0.24 Very Good J      VVS2     62.8   57   336  3.94  3.96  2.48
```

Extract  $y$ , the price variable and  $col$ , the nominal variable “color” as vectors.

```
y <- diamonds$price
col <- diamonds$color
```

Convert the  $col$  vector to  $X$  which contains an intercept and an appropriate number of dummies. Let the color G be the reference category as it is the modal color. Name the columns of  $X$  appropriately. The first should be “(Intercept)”. Delete  $col$ .

```
level <- levels(col)
X <- rep(1, length(y))
temp <- c()
for (i in 1:length(level)){
  temp <- ifelse(col == level[i], 1, 0)
  X <- cbind(X, temp)
}

colnames(X) <- c("(Intercept)", levels(col))
X <- X[, -5]
```

Repeat the iterative exercise above we did for Boston here.

```
M <- matrix(NA, ncol(X), ncol(X))
colnames(M) <- colnames(X)

for (j in 1:ncol(X)){
  X_j <- X[, 1:j, drop = FALSE]
  b <- solve(t(X_j) %*% X_j) %*% t(X_j) %*% y
  M[j, 1:j] <- b
}

M
```

```
##      (Intercept)      D      E      F      H      I
## [1,]    3932.800      NA      NA      NA      NA      NA
## [2,]    4042.378 -872.4243      NA      NA      NA      NA
## [3,]    4295.543 -1125.5885 -1218.7901      NA      NA      NA
## [4,]    4491.230 -1321.2760 -1414.4776 -766.3437      NA      NA
## [5,]    4493.170 -1323.2160 -1416.4176 -768.2837 -6.50092      NA
## [6,]    4262.945 -1092.9907 -1186.1923 -538.0584 223.72444 828.9302
## [7,]    3999.136 -829.1816 -922.3832 -274.2493 487.53352 1092.7393
##           J
## [1,]      NA
## [2,]      NA
```

```
## [3,]      NA
## [4,]      NA
## [5,]      NA
## [6,]      NA
## [7,] 1324.682
```

Why didn't the estimates change as we added more and more features?

It did change because we added more weight each time.

Model `price` with both `color` and `clarity` with and without an intercept and report the coefficients.

```
mod1 <- lm(price ~ color + clarity, data = diamonds)
mod2 <- lm(price ~ 0 + color + clarity, data = diamonds)

coef(mod1)
```

```
## (Intercept)      color.L      color.Q      color.C      color^4
## 3813.005017 2137.537090  16.275100 -221.703368  77.356602
##      color^5      color^6      clarity.L      clarity.Q      clarity.C
## -259.950208   2.119947 -1718.615415 -594.151103  681.362698
##      clarity^4      clarity^5      clarity^6      clarity^7
## -248.597841   806.884319 -228.711112   191.108568
```

```
coef(mod2)
```

```
##      colorD      colorE      colorF      colorG      colorH      colorI
## 2747.6576 2757.0769 3462.3053 3841.9086 4167.6095 4780.8262
##      colorJ      clarity.L      clarity.Q      clarity.C      clarity^4      clarity^5
## 4933.6510 -1718.6154 -594.1511  681.3627 -248.5978  806.8843
##      clarity^6      clarity^7
## -228.7111   191.1086
```

Which coefficients did not change between the models and why?

The clarity coefficients did not change between the models because it was not included in the intercept.

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.

```
m <- matrix(c(1, 1, rnorm(2)), 2, 2)
theta_in_rad <- acos((m[, 1] %*% m[, 2]) / (sqrt(sum(m[, 1]^2)) * sqrt(sum(m[, 2]^2))))
theta_in_rad * 180 / pi
```

```
##      [,1]
## [1,] 8.760682
```

Repeat this exercise  $N_{sim} = 1e5$  times and report the average absolute angle.

```
Nsim <- 1e5
abs_angle <- c()

for (i in 1:Nsim){
```

```

m <- matrix(c(1, 1, rnorm(2)), 2, 2)
theta_in_rad <- acos((m[, 1] %*% m[, 2]) / (sqrt(sum(m[, 1]^2)) * sqrt(sum(m[, 2]^2))))
abs_angle <- c(abs_angle, theta_in_rad * 180 / pi)
}

mean(abs_angle)

```

```
## [1] 90.13522
```

Create a  $n \times 2$  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  $n \in 10, 50, 100, 200, 500, 1000$ , report the average absolute angle over  $Nsim = 1e5$  simulations.

```

n <- c(10, 50, 100, 200, 500, 1000)
avg <- c()
for (i in 1:length(n)){
  m <- cbind(rep(1, n[i]), rnorm(n[i]))
  for (j in 1:Nsim){
    theta_in_rad <- acos((m[, 1] %*% m[, 2]) / (sqrt(sum(m[, 1]^2)) * sqrt(sum(m[, 2]^2))))
    abs_angle <- c(abs_angle, theta_in_rad * 180 / pi)
  }
  avg <- c(avg, mean(abs_angle))
}
avg

```

```
## [1] 83.62079 88.14742 90.21100 91.47835 91.30413 90.69217
```

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

The angle is converging to 90 degrees because the columns are orthogonal

Create a vector  $y$  by simulating  $n = 100$  standard iid normals. Create a matrix of size  $100 \times 2$  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
y <- as.matrix(rnorm(n))
X <- cbind(rep(1, n), rnorm(n))

H <- X %*% solve(t(X) %*% X) %*% t(X)
yhat <- H %*% y

len_vec <- function(v){sqrt(sum(v^2))}
Rsqr <- (len_vec(yhat - mean(y))^2)/(len_vec(y - mean(y))^2)
Rsqr

```

```
## [1] 0.02036548
```

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??

```

R2s <- c(Rsq)
SST <- len_vec(y - mean(y))^2
for (i in 3:100){
  X <- cbind(X, rnorm(n))
  H <- X %*% solve(t(X) %*% X) %*% t(X)
  yhat <- H %*% y

  Rsq <- (len_vec(yhat - mean(y))^2)/SST # SSR/SST
  R2s <- c(R2s, Rsq)
}

R2s

```

```

## [1] 0.02036548 0.02395415 0.02397333 0.03445465 0.04586071 0.04952533
## [7] 0.10289810 0.10311230 0.10480411 0.10480701 0.10777721 0.10887180
## [13] 0.11001862 0.11619717 0.11841232 0.13947058 0.15286702 0.17515671
## [19] 0.17781993 0.18240392 0.18245525 0.19053069 0.19260500 0.21163225
## [25] 0.21164144 0.21441605 0.21543327 0.24789709 0.25076376 0.25136198
## [31] 0.25186815 0.25409995 0.25456782 0.27373545 0.30372799 0.30380112
## [37] 0.30703995 0.30704436 0.30855809 0.30887335 0.30902992 0.31184685
## [43] 0.31357303 0.31465793 0.31958338 0.32083198 0.32332016 0.32479675
## [49] 0.32859305 0.33533715 0.43518740 0.46047421 0.49029190 0.49389148
## [55] 0.51066182 0.51211432 0.51448677 0.53357754 0.53976017 0.56337677
## [61] 0.56559039 0.59751332 0.61233312 0.61526058 0.62364621 0.62365269
## [67] 0.62605707 0.63027003 0.64996877 0.65596404 0.67309744 0.68556168
## [73] 0.72298632 0.72417445 0.73050694 0.73097949 0.73748435 0.73868913
## [79] 0.73942723 0.74049968 0.76228999 0.76300417 0.80864915 0.81071283
## [85] 0.82003463 0.83392521 0.84550251 0.85696065 0.86098500 0.86098563
## [91] 0.87731172 0.89776120 0.90994759 0.91170068 0.92111436 0.94456238
## [97] 0.99644908 0.99798478 1.00000000

```

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

```

# X <- cbind(X, rnorm(n))
# H <- X %*% solve(t(X) %*% X) %*% t(X) # X is no longer invertible
# yhat <- H %*% y
# Rsq <- (len_vec(yhat - mean(y))^2)/(len_vec(y - mean(y))^2)
# Rsq

```