# Lab 4

## Jonathan Eng

## 11:59PM Feb 29, 2020

We now move on to simple linear modeling using the ordinary least squares algorithm.

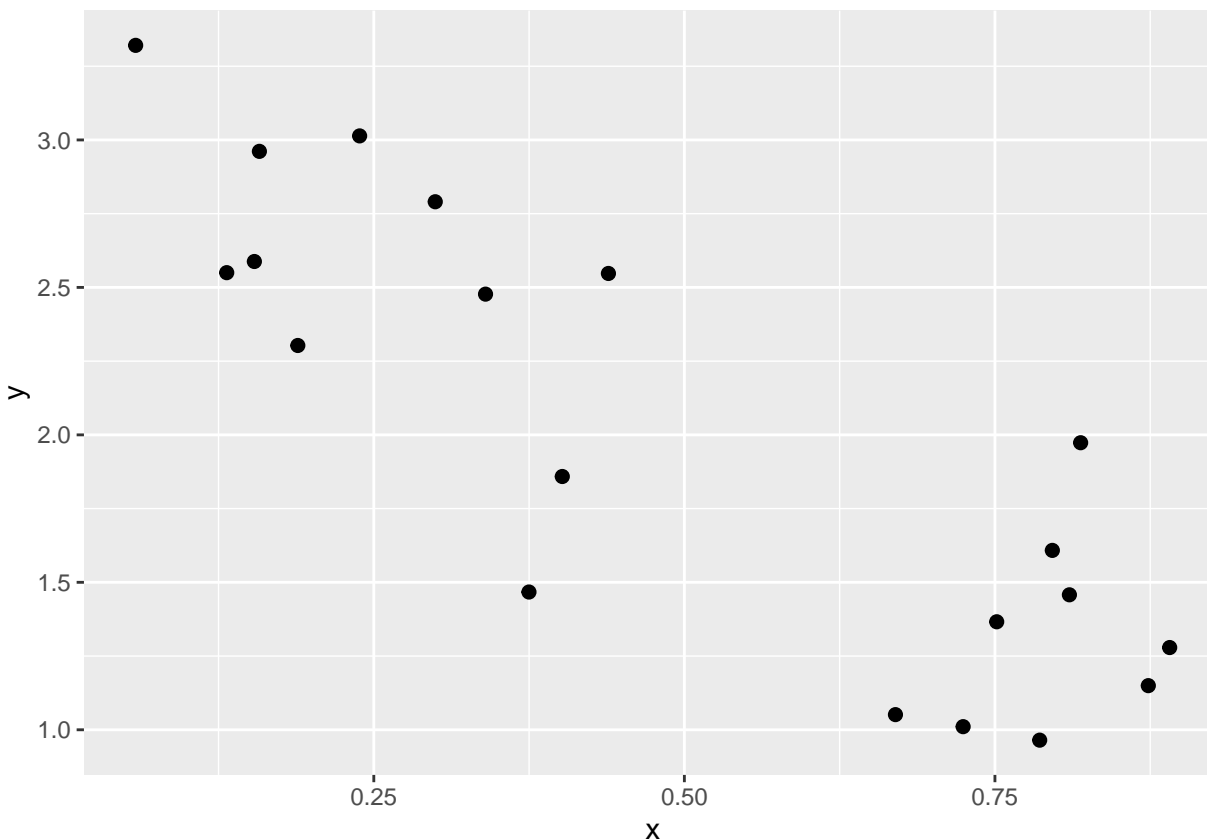Let's quickly recreate the sample data set from practice lecture 7:

```r
n = 20
x = runif(n)
beta_0 = 3
beta_1 = -2
y = beta_0 + beta_1 * x + rnorm(n, mean = 0, sd = 0.33)
```

Rewrite the computation of y so that it is h*(x) + epsilon.

```r
h_star_x = beta_0 + beta_1 * x
epsilon = rnorm(n, mean = 0, sd = 0.33)
y = h_star_x + epsilon
```

Graph the data by running the following chunk:

```r
pacman::p_load(ggplot2)
simple_df = data.frame(x = x, y = y)
simple_viz_obj = ggplot(simple_df, aes(x, y)) +
  geom_point(size = 2)
simple_viz_obj
```

Does this make sense given the values of $beta_0$ and $beta_1$?

Yes.

Write a function `my_simple_ols` that takes in a vector `x` and vector `y` and returns a list that contains the `b_0` (intercept), `b_1` (slope), `yhat` (the predictions), `e` (the residuals), `SSE`, `SST`, `MSE`, `RMSE` and `Rsq` (for the R-squared metric). Internally, you can only use the functions `sum` and `length` and other basic arithmetic operations. You should throw errors if the inputs are non-numeric and not the same length. You should also name the class of the return value 'my_simple_ols_objby using theclass' function as a setter. No need to create ROxygen documentation here.

```r
my_simple_ols = function(x,y){
 if (class(x) != "numeric" | class(y) != "numeric") {stop("argument x or y is not numeric")}
 n = length(x)
 if (n != length(y)){stop("x and y must be same length")}

 y_bar = sum(y)/length(y)
 x_bar = sum(x)/length(x)
 s_x_squared = (1/(n-1) * sum((x - x_bar)^2))
 s_xy = (1/(n-1)) * sum((x - x_bar)*(y - y_bar))
 b1= s_xy/s_x_squared
 b0= y_bar - b1*x_bar
 y_hat = b0 + b1*x

 e = y - y_hat
 SSE = sum(e^2)
 SST = sum((y-y_bar)^2)
```

```r
  Rsq = 1 - SSE/SST # 1 - NA
  MSE = SSE / (n-2)
  RMSE = sqrt(MSE)

  mod = list(b_0 = b0, b_1 = b1, y_hat = y_hat, e = e, SSE = SSE, SST = SST, Rsq = Rsq, MSE = MSE, RMSE
  class(mod) = "my_simple_ols_obj"
  mod
}
```

Verify your computations are correct for the vectors `x` and `y` from the first chunk using the `lm` function in R:

```r
lm_mod = lm(y ~ x)
my_lm_mod = my_simple_ols(x, y)
#run the tests to ensure the function is up to spec
pacman::p_load(testthat)
expect_equal(my_lm_mod$b_0, as.numeric(coef(lm_mod)[1]), tol = 1e-4)
expect_equal(my_lm_mod$b_1, as.numeric(coef(lm_mod)[2]), tol = 1e-4)
expect_equal(my_lm_mod$RMSE, summary(lm_mod)$sigma, tol = 1e-4)
expect_equal(my_lm_mod$Rsq, summary(lm_mod)$r.squared, tol = 1e-4)
```

Verify that the average of the residuals is 0.

```r
expect_equal(mean(my_lm_mod$e), 0, tol = 1e-4)
```

Create the $X$ matrix for this data example.

```r
X = cbind(1, x)
X
```

```
##         x
##  [1,] 1 0.72439818
##  [2,] 1 0.40173039
##  [3,] 1 0.78603297
##  [4,] 1 0.33993088
##  [5,] 1 0.05827157
##  [6,] 1 0.37488552
##  [7,] 1 0.15385093
##  [8,] 1 0.43885504
##  [9,] 1 0.15789404
## [10,] 1 0.89067734
## [11,] 1 0.29945788
## [12,] 1 0.87348045
## [13,] 1 0.66992145
## [14,] 1 0.81900053
## [15,] 1 0.13158586
## [16,] 1 0.80995562
## [17,] 1 0.79622270
## [18,] 1 0.75143082
## [19,] 1 0.18882609
## [20,] 1 0.23857099
```

Use the `model.matrix` function to compute the matrix `X` and verify it is the same as your manual construction.

```
model.matrix(~ x)
```

```
##    (Intercept)          x
## 1            1 0.72439818
## 2            1 0.40173039
## 3            1 0.78603297
## 4            1 0.33993088
## 5            1 0.05827157
## 6            1 0.37488552
## 7            1 0.15385093
## 8            1 0.43885504
## 9            1 0.15789404
## 10           1 0.89067734
## 11           1 0.29945788
## 12           1 0.87348045
## 13           1 0.66992145
## 14           1 0.81900053
## 15           1 0.13158586
## 16           1 0.80995562
## 17           1 0.79622270
## 18           1 0.75143082
## 19           1 0.18882609
## 20           1 0.23857099
## attr(,"assign")
## [1] 0 1
```

Using matrix algebra, verify the OLS estimate is the same as you computed from the `my_simple_ols` function.

```
XtXinvX = solve(t(X) %*% X) %*% t(X)
b = XtXinvX %*% y
b
```

```
##          [,1]
##     3.080373
## x -2.207916
```

Find the hat matrix $H$.

```
H = X %*% XtXinvX
H
```

```
##                [,1]       [,2]         [,3]       [,4]         [,5]       [,6]
## [1,]   0.082228601 0.03684712  0.090897203 0.02815535 -0.011458514 0.03307153
## [2,]   0.036847117 0.05536785  0.033309356 0.05891507  0.075081965 0.05690871
## [3,]   0.090897203 0.03330936  0.101897418 0.02227974 -0.027989152 0.02851824
## [4,]   0.028155350 0.05891507  0.022279743 0.06480637  0.091656780 0.06147417
## [5,]  -0.011458514 0.07508196 -0.027989152 0.09165678  0.167198660 0.08228184
## [6,]   0.033071531 0.05690871  0.028518241 0.06147417  0.082281840 0.05889189
## [7,]   0.001984208 0.06959583 -0.010930711 0.08254526  0.141563986 0.07522089
## [8,]   0.042068499 0.05323694  0.039935145 0.05537600  0.065125020 0.05416612
## [9,]   0.002552849 0.06936376 -0.010209120 0.08215983  0.140479612 0.07492220
```

```
## [10,]  0.105614869 0.02730290  0.120573728 0.01230406 -0.056055088 0.02078761
## [11,]  0.022463039 0.06123817  0.015056359 0.06866464  0.102511763 0.06446413
## [12,]  0.103196219 0.02828998  0.117504529 0.01394343 -0.051442831 0.02205803
## [13,]  0.074566744 0.03997401  0.081174517 0.03334858  0.003152309 0.03709602
## [14,]  0.095533912 0.03141706  0.107781273 0.01913697 -0.036831150 0.02608275
## [15,] -0.001147254 0.07087381 -0.014904449 0.08466777  0.147535546 0.07686572
## [16,]  0.094261795 0.03193622  0.106166992 0.01999921 -0.034405280 0.02675095
## [17,]  0.092330334 0.03272448  0.103716021 0.02130836 -0.030722069 0.02776547
## [18,]  0.086030597 0.03529548  0.095721830 0.02557834 -0.018708751 0.03107449
## [19,]  0.006903276 0.06758830 -0.004688550 0.07921110  0.132183542 0.07263708
## [20,]  0.013899627 0.06473300  0.004189626 0.07446895  0.118841810 0.06896217
##               [,7]       [,8]        [,9]       [,10]      [,11]       [,12]
##  [1,]  0.001984208 0.04206850  0.002552849  0.10561487 0.022463039  0.103196219
##  [2,]  0.069595827 0.05323694  0.069363758  0.02730290 0.061238168  0.028289978
##  [3,] -0.010930711 0.03993515 -0.010209120  0.12057373 0.015056359  0.117504529
##  [4,]  0.082545260 0.05537600  0.082159833  0.01230406 0.068664641  0.013943430
##  [5,]  0.141563986 0.06512502  0.140479612 -0.05605509 0.102511763 -0.051442831
##  [6,]  0.075220885 0.05416612  0.074922200  0.02078761 0.064464127  0.022058032
##  [7,]  0.121536343 0.06181675  0.120689153 -0.03285783 0.091025950 -0.029254404
##  [8,]  0.061816749 0.05195195  0.061676806  0.03631311 0.056776882  0.036908341
##  [9,]  0.120689153 0.06167681  0.119851995 -0.03187656 0.090540088 -0.028315810
## [10,] -0.032857830 0.03631311 -0.031876561  0.14597108 0.002481200  0.141797366
## [11,]  0.091025950 0.05677688  0.090540088  0.00248120 0.073528300  0.004547758
## [12,] -0.029254404 0.03690834 -0.028315810  0.14179737 0.004547758  0.137805166
## [13,]  0.013399229 0.04395409  0.013832685  0.09239328 0.029009531  0.090549630
## [14,] -0.017838713 0.03879405 -0.017035313  0.12857501 0.011094634  0.125157836
## [15,]  0.126201753 0.06258741  0.125299311 -0.03826160 0.093701553 -0.034423165
## [16,] -0.015943449 0.03910711 -0.015162494  0.12637979 0.012181564  0.123058092
## [17,] -0.013065861 0.03958245 -0.012318985  0.12304679 0.013831855  0.119870041
## [18,] -0.003680197 0.04113282 -0.003044473  0.11217573 0.019214515  0.109471757
## [19,]  0.114207669 0.06060616  0.113447270 -0.02436930 0.086822969 -0.021135038
## [20,]  0.103784152 0.05888435  0.103147197 -0.01229614 0.080845104 -0.009586928
##              [,13]       [,14]       [,15]        [,16]       [,17]
##  [1,] 0.074566744  0.095533912 -0.001147254  0.0942617946  0.092330334
##  [2,] 0.039974014  0.031417058  0.070873815  0.0319362248  0.032724478
##  [3,] 0.081174517  0.107781273 -0.014904449  0.1061669919  0.103716021
##  [4,] 0.033348582  0.019136967  0.084667775  0.0199992121  0.021308363
##  [5,] 0.003152309 -0.036831150  0.147535546 -0.0344052798 -0.030722069
##  [6,] 0.037096016  0.026082753  0.076865725  0.0267509478  0.027765472
##  [7,] 0.013399229 -0.017838713  0.126201753 -0.0159434491 -0.013065861
##  [8,] 0.043954092  0.038794045  0.062587406  0.0391071147  0.039582450
##  [9,] 0.013832685 -0.017035313  0.125299311 -0.0151624940 -0.012318985
## [10,] 0.092393283  0.128575005 -0.038261597  0.1263797929  0.123046792
## [11,] 0.029009531  0.011094634  0.093701553  0.0121815640  0.013831855
## [12,] 0.090549630  0.125157836 -0.034423165  0.1230580919  0.119870041
## [13,] 0.068726376  0.084708920  0.011012224  0.0837392290  0.082266943
## [14,] 0.084708920  0.114332209 -0.022262973  0.1125349089  0.109806060
## [15,] 0.011012224 -0.022262973  0.131171429 -0.0202441053 -0.017178849
## [16,] 0.083739229  0.112534909 -0.020244105  0.1107878216  0.108135211
## [17,] 0.082266943  0.109806060 -0.017178849  0.1081352109  0.105598353
## [18,] 0.077464873  0.100905530 -0.007181076  0.0994833417  0.097324027
## [19,] 0.017148864 -0.010888849  0.118395122 -0.0091877489 -0.006604961
## [20,] 0.022481939 -0.001004112  0.107291812  0.0004208302  0.002584327
##              [,18]       [,19]       [,20]
```

```
## [1,]   0.086030597  0.006903276  0.0138996271
## [2,]   0.035295476  0.067588296  0.0647330000
## [3,]   0.095721830 -0.004688550  0.0041896263
## [4,]   0.025578344  0.079211099  0.0744689498
## [5,]  -0.018708751  0.132183542  0.1188418105
## [6,]   0.031074486  0.072637085  0.0689621650
## [7,]  -0.003680197  0.114207669  0.1037841523
## [8,]   0.041132823  0.060606160  0.0588843488
## [9,]  -0.003044473  0.113447270  0.1031471973
## [10,]   0.112175733 -0.024369304 -0.0122961412
## [11,]   0.019214515  0.086822969  0.0808451037
## [12,]   0.109471757 -0.021135038 -0.0095869279
## [13,]   0.077464873  0.017148864  0.0224819393
## [14,]   0.100905530 -0.010888849 -0.0010041120
## [15,]  -0.007181076  0.118395122  0.1072918116
## [16,]   0.099483342 -0.009187749  0.0004208302
## [17,]   0.097324027 -0.006604961  0.0025843267
## [18,]   0.090281113  0.001819172  0.0096408798
## [19,]   0.001819172  0.107629794  0.0982741340
## [20,]   0.009640880  0.098274134  0.0904372785
```

Verify that this specific hat matrix is symmetric.

```
expect_equal(H, t(H))
```

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

```
diag(H)
```

```
##  [1] 0.08222860 0.05536785 0.10189742 0.06480637 0.16719866 0.05889189
##  [7] 0.12153634 0.05195195 0.11985200 0.14597108 0.07352830 0.13780517
## [13] 0.06872638 0.11433221 0.13117143 0.11078782 0.10559835 0.09028111
## [19] 0.10762979 0.09043728
```

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

```
## [1] 2
```

Create a prediction method `g` that takes in a vector `x_future` and `my_simple_ols_obj`, an object of type `my_simple_ols_obj` and predicts y values for each entry in `x_future`.

```
g = function(x_future, my_simple_ols_obj){
  my_simple_ols_obj$b_0 + my_simple_ols_obj$b_1 * x_future
}
```

Use this function to verify that when predicting for the average x, you get the average y.

```
expect_equal(g(mean(x), my_lm_mod), mean(y))
```

Create a prediction method `g` that takes in a vector `x_future` and the dataset $\mathbb{D}$ i.e. X where the first column is the one vector and `y` and returns the OLS predictions.

```
g = function(x_future, X, y){
b = solve(t(X) %*% X) %*% t(X) %*% y
b[1] +b[2]*x_future
}
```

In class we spoke about error due to ignorance, misspecification error and estimation error. Show that as n grows, estimation error shrinks. Let us define an error metric that is the difference between $b_0$ and $b_1$ and $\beta_0$ and $\beta_1$. How about $h = ||b - \beta||^2$ where the quantities are now the vectors of size two. Show as n increases, this shrinks.

```
ns = 10^(1:7)
errors = array(dim=length(ns))
beta = c(beta_0, beta_1)
for (i in 1:length(ns)) {
  n = ns[i]
  x = runif(n)
  h_star_x = beta_0 + beta_1 * x
  epsilon = rnorm(n, mean = 0, sd = 0.33)
  y = h_star_x + epsilon

  mod = lm(y ~ x)
  b = coef(mod)
  errors[i] = sum((beta - b)^2)
}
errors
```

```
## [1] 2.821024e-01 1.239159e-02 7.283193e-04 5.419810e-05 2.794452e-05
## [6] 6.252637e-07 3.086186e-07
```

We are now going to repeat one of the first linear model building exercises in history — that of Sir Francis Galton in 1886. First load up package `HistData`.

```
pacman::p_load(HistData)
```

In it, there is a dataset called `Galton`. Load it up.

```
data(Galton)
```

You now should have a data frame in your workspace called `Galton`. Summarize this data frame and write a few sentences about what you see. Make sure you report $n$, $p$ and a bit about what the columns represent and how the data was measured. See the help file `?Galton`. p is 1 and n is 928 the number of observations

```
Galton
```

```
##      parent child
## 1      70.5  61.7
## 2      68.5  61.7
## 3      65.5  61.7
## 4      64.5  61.7
## 5      64.0  61.7
## 6      67.5  62.2
```

```
## 7      67.5  62.2
## 8      67.5  62.2
## 9      66.5  62.2
## 10     66.5  62.2
## 11     66.5  62.2
## 12     64.5  62.2
## 13     70.5  63.2
## 14     69.5  63.2
## 15     68.5  63.2
## 16     68.5  63.2
## 17     68.5  63.2
## 18     68.5  63.2
## 19     68.5  63.2
## 20     68.5  63.2
## 21     68.5  63.2
## 22     67.5  63.2
## 23     67.5  63.2
## 24     67.5  63.2
## 25     67.5  63.2
## 26     67.5  63.2
## 27     66.5  63.2
## 28     66.5  63.2
## 29     66.5  63.2
## 30     65.5  63.2
## 31     65.5  63.2
## 32     65.5  63.2
## 33     65.5  63.2
## 34     65.5  63.2
## 35     65.5  63.2
## 36     65.5  63.2
## 37     65.5  63.2
## 38     65.5  63.2
## 39     64.5  63.2
## 40     64.5  63.2
## 41     64.5  63.2
## 42     64.5  63.2
## 43     64.0  63.2
## 44     64.0  63.2
## 45     69.5  64.2
## 46     69.5  64.2
## 47     69.5  64.2
## 48     69.5  64.2
## 49     69.5  64.2
## 50     69.5  64.2
## 51     69.5  64.2
## 52     69.5  64.2
## 53     69.5  64.2
## 54     69.5  64.2
## 55     69.5  64.2
## 56     69.5  64.2
## 57     69.5  64.2
## 58     69.5  64.2
## 59     69.5  64.2
## 60     69.5  64.2
```

```
## 61    68.5  64.2
## 62    68.5  64.2
## 63    68.5  64.2
## 64    68.5  64.2
## 65    68.5  64.2
## 66    68.5  64.2
## 67    68.5  64.2
## 68    68.5  64.2
## 69    68.5  64.2
## 70    68.5  64.2
## 71    68.5  64.2
## 72    67.5  64.2
## 73    67.5  64.2
## 74    67.5  64.2
## 75    67.5  64.2
## 76    67.5  64.2
## 77    67.5  64.2
## 78    67.5  64.2
## 79    67.5  64.2
## 80    67.5  64.2
## 81    67.5  64.2
## 82    67.5  64.2
## 83    67.5  64.2
## 84    67.5  64.2
## 85    67.5  64.2
## 86    66.5  64.2
## 87    66.5  64.2
## 88    66.5  64.2
## 89    66.5  64.2
## 90    66.5  64.2
## 91    65.5  64.2
## 92    65.5  64.2
## 93    65.5  64.2
## 94    65.5  64.2
## 95    65.5  64.2
## 96    64.5  64.2
## 97    64.5  64.2
## 98    64.5  64.2
## 99    64.5  64.2
## 100   64.0  64.2
## 101   64.0  64.2
## 102   64.0  64.2
## 103   64.0  64.2
## 104   71.5  65.2
## 105   70.5  65.2
## 106   69.5  65.2
## 107   69.5  65.2
## 108   69.5  65.2
## 109   69.5  65.2
## 110   68.5  65.2
## 111   68.5  65.2
## 112   68.5  65.2
## 113   68.5  65.2
## 114   68.5  65.2
```

```
## 115   68.5   65.2
## 116   68.5   65.2
## 117   68.5   65.2
## 118   68.5   65.2
## 119   68.5   65.2
## 120   68.5   65.2
## 121   68.5   65.2
## 122   68.5   65.2
## 123   68.5   65.2
## 124   68.5   65.2
## 125   68.5   65.2
## 126   67.5   65.2
## 127   67.5   65.2
## 128   67.5   65.2
## 129   67.5   65.2
## 130   67.5   65.2
## 131   67.5   65.2
## 132   67.5   65.2
## 133   67.5   65.2
## 134   67.5   65.2
## 135   67.5   65.2
## 136   67.5   65.2
## 137   67.5   65.2
## 138   67.5   65.2
## 139   67.5   65.2
## 140   67.5   65.2
## 141   66.5   65.2
## 142   66.5   65.2
## 143   65.5   65.2
## 144   65.5   65.2
## 145   65.5   65.2
## 146   65.5   65.2
## 147   65.5   65.2
## 148   65.5   65.2
## 149   65.5   65.2
## 150   64.5   65.2
## 151   64.0   65.2
## 152   71.5   66.2
## 153   71.5   66.2
## 154   71.5   66.2
## 155   70.5   66.2
## 156   69.5   66.2
## 157   69.5   66.2
## 158   69.5   66.2
## 159   69.5   66.2
## 160   69.5   66.2
## 161   69.5   66.2
## 162   69.5   66.2
## 163   69.5   66.2
## 164   69.5   66.2
## 165   69.5   66.2
## 166   69.5   66.2
## 167   69.5   66.2
## 168   69.5   66.2
```

```
## 169    69.5   66.2
## 170    69.5   66.2
## 171    69.5   66.2
## 172    69.5   66.2
## 173    68.5   66.2
## 174    68.5   66.2
## 175    68.5   66.2
## 176    68.5   66.2
## 177    68.5   66.2
## 178    68.5   66.2
## 179    68.5   66.2
## 180    68.5   66.2
## 181    68.5   66.2
## 182    68.5   66.2
## 183    68.5   66.2
## 184    68.5   66.2
## 185    68.5   66.2
## 186    68.5   66.2
## 187    68.5   66.2
## 188    68.5   66.2
## 189    68.5   66.2
## 190    68.5   66.2
## 191    68.5   66.2
## 192    68.5   66.2
## 193    68.5   66.2
## 194    68.5   66.2
## 195    68.5   66.2
## 196    68.5   66.2
## 197    68.5   66.2
## 198    67.5   66.2
## 199    67.5   66.2
## 200    67.5   66.2
## 201    67.5   66.2
## 202    67.5   66.2
## 203    67.5   66.2
## 204    67.5   66.2
## 205    67.5   66.2
## 206    67.5   66.2
## 207    67.5   66.2
## 208    67.5   66.2
## 209    67.5   66.2
## 210    67.5   66.2
## 211    67.5   66.2
## 212    67.5   66.2
## 213    67.5   66.2
## 214    67.5   66.2
## 215    67.5   66.2
## 216    67.5   66.2
## 217    67.5   66.2
## 218    67.5   66.2
## 219    67.5   66.2
## 220    67.5   66.2
## 221    67.5   66.2
## 222    67.5   66.2
```

```
## 223    67.5   66.2
## 224    67.5   66.2
## 225    67.5   66.2
## 226    67.5   66.2
## 227    67.5   66.2
## 228    67.5   66.2
## 229    67.5   66.2
## 230    67.5   66.2
## 231    67.5   66.2
## 232    67.5   66.2
## 233    67.5   66.2
## 234    66.5   66.2
## 235    66.5   66.2
## 236    66.5   66.2
## 237    66.5   66.2
## 238    66.5   66.2
## 239    66.5   66.2
## 240    66.5   66.2
## 241    66.5   66.2
## 242    66.5   66.2
## 243    66.5   66.2
## 244    66.5   66.2
## 245    66.5   66.2
## 246    66.5   66.2
## 247    66.5   66.2
## 248    66.5   66.2
## 249    66.5   66.2
## 250    66.5   66.2
## 251    65.5   66.2
## 252    65.5   66.2
## 253    65.5   66.2
## 254    65.5   66.2
## 255    65.5   66.2
## 256    65.5   66.2
## 257    65.5   66.2
## 258    65.5   66.2
## 259    65.5   66.2
## 260    65.5   66.2
## 261    65.5   66.2
## 262    64.5   66.2
## 263    64.5   66.2
## 264    64.5   66.2
## 265    64.5   66.2
## 266    64.5   66.2
## 267    64.0   66.2
## 268    64.0   66.2
## 269    71.5   67.2
## 270    71.5   67.2
## 271    71.5   67.2
## 272    71.5   67.2
## 273    70.5   67.2
## 274    70.5   67.2
## 275    70.5   67.2
## 276    69.5   67.2
```

```
## 277   69.5   67.2
## 278   69.5   67.2
## 279   69.5   67.2
## 280   69.5   67.2
## 281   69.5   67.2
## 282   69.5   67.2
## 283   69.5   67.2
## 284   69.5   67.2
## 285   69.5   67.2
## 286   69.5   67.2
## 287   69.5   67.2
## 288   69.5   67.2
## 289   69.5   67.2
## 290   69.5   67.2
## 291   69.5   67.2
## 292   69.5   67.2
## 293   69.5   67.2
## 294   69.5   67.2
## 295   69.5   67.2
## 296   69.5   67.2
## 297   69.5   67.2
## 298   69.5   67.2
## 299   69.5   67.2
## 300   69.5   67.2
## 301   69.5   67.2
## 302   69.5   67.2
## 303   68.5   67.2
## 304   68.5   67.2
## 305   68.5   67.2
## 306   68.5   67.2
## 307   68.5   67.2
## 308   68.5   67.2
## 309   68.5   67.2
## 310   68.5   67.2
## 311   68.5   67.2
## 312   68.5   67.2
## 313   68.5   67.2
## 314   68.5   67.2
## 315   68.5   67.2
## 316   68.5   67.2
## 317   68.5   67.2
## 318   68.5   67.2
## 319   68.5   67.2
## 320   68.5   67.2
## 321   68.5   67.2
## 322   68.5   67.2
## 323   68.5   67.2
## 324   68.5   67.2
## 325   68.5   67.2
## 326   68.5   67.2
## 327   68.5   67.2
## 328   68.5   67.2
## 329   68.5   67.2
## 330   68.5   67.2
```

```
## 331    68.5   67.2
## 332    68.5   67.2
## 333    68.5   67.2
## 334    67.5   67.2
## 335    67.5   67.2
## 336    67.5   67.2
## 337    67.5   67.2
## 338    67.5   67.2
## 339    67.5   67.2
## 340    67.5   67.2
## 341    67.5   67.2
## 342    67.5   67.2
## 343    67.5   67.2
## 344    67.5   67.2
## 345    67.5   67.2
## 346    67.5   67.2
## 347    67.5   67.2
## 348    67.5   67.2
## 349    67.5   67.2
## 350    67.5   67.2
## 351    67.5   67.2
## 352    67.5   67.2
## 353    67.5   67.2
## 354    67.5   67.2
## 355    67.5   67.2
## 356    67.5   67.2
## 357    67.5   67.2
## 358    67.5   67.2
## 359    67.5   67.2
## 360    67.5   67.2
## 361    67.5   67.2
## 362    67.5   67.2
## 363    67.5   67.2
## 364    67.5   67.2
## 365    67.5   67.2
## 366    67.5   67.2
## 367    67.5   67.2
## 368    67.5   67.2
## 369    67.5   67.2
## 370    67.5   67.2
## 371    67.5   67.2
## 372    66.5   67.2
## 373    66.5   67.2
## 374    66.5   67.2
## 375    66.5   67.2
## 376    66.5   67.2
## 377    66.5   67.2
## 378    66.5   67.2
## 379    66.5   67.2
## 380    66.5   67.2
## 381    66.5   67.2
## 382    66.5   67.2
## 383    66.5   67.2
## 384    66.5   67.2
```

```
## 385    66.5   67.2
## 386    66.5   67.2
## 387    66.5   67.2
## 388    66.5   67.2
## 389    65.5   67.2
## 390    65.5   67.2
## 391    65.5   67.2
## 392    65.5   67.2
## 393    65.5   67.2
## 394    65.5   67.2
## 395    65.5   67.2
## 396    65.5   67.2
## 397    65.5   67.2
## 398    65.5   67.2
## 399    65.5   67.2
## 400    64.5   67.2
## 401    64.5   67.2
## 402    64.5   67.2
## 403    64.5   67.2
## 404    64.5   67.2
## 405    64.0   67.2
## 406    64.0   67.2
## 407    72.5   68.2
## 408    71.5   68.2
## 409    71.5   68.2
## 410    71.5   68.2
## 411    70.5   68.2
## 412    70.5   68.2
## 413    70.5   68.2
## 414    70.5   68.2
## 415    70.5   68.2
## 416    70.5   68.2
## 417    70.5   68.2
## 418    70.5   68.2
## 419    70.5   68.2
## 420    70.5   68.2
## 421    70.5   68.2
## 422    70.5   68.2
## 423    69.5   68.2
## 424    69.5   68.2
## 425    69.5   68.2
## 426    69.5   68.2
## 427    69.5   68.2
## 428    69.5   68.2
## 429    69.5   68.2
## 430    69.5   68.2
## 431    69.5   68.2
## 432    69.5   68.2
## 433    69.5   68.2
## 434    69.5   68.2
## 435    69.5   68.2
## 436    69.5   68.2
## 437    69.5   68.2
## 438    69.5   68.2
```

```
## 439   69.5   68.2
## 440   69.5   68.2
## 441   69.5   68.2
## 442   69.5   68.2
## 443   68.5   68.2
## 444   68.5   68.2
## 445   68.5   68.2
## 446   68.5   68.2
## 447   68.5   68.2
## 448   68.5   68.2
## 449   68.5   68.2
## 450   68.5   68.2
## 451   68.5   68.2
## 452   68.5   68.2
## 453   68.5   68.2
## 454   68.5   68.2
## 455   68.5   68.2
## 456   68.5   68.2
## 457   68.5   68.2
## 458   68.5   68.2
## 459   68.5   68.2
## 460   68.5   68.2
## 461   68.5   68.2
## 462   68.5   68.2
## 463   68.5   68.2
## 464   68.5   68.2
## 465   68.5   68.2
## 466   68.5   68.2
## 467   68.5   68.2
## 468   68.5   68.2
## 469   68.5   68.2
## 470   68.5   68.2
## 471   68.5   68.2
## 472   68.5   68.2
## 473   68.5   68.2
## 474   68.5   68.2
## 475   68.5   68.2
## 476   68.5   68.2
## 477   67.5   68.2
## 478   67.5   68.2
## 479   67.5   68.2
## 480   67.5   68.2
## 481   67.5   68.2
## 482   67.5   68.2
## 483   67.5   68.2
## 484   67.5   68.2
## 485   67.5   68.2
## 486   67.5   68.2
## 487   67.5   68.2
## 488   67.5   68.2
## 489   67.5   68.2
## 490   67.5   68.2
## 491   67.5   68.2
## 492   67.5   68.2
```

```
## 493    67.5    68.2
## 494    67.5    68.2
## 495    67.5    68.2
## 496    67.5    68.2
## 497    67.5    68.2
## 498    67.5    68.2
## 499    67.5    68.2
## 500    67.5    68.2
## 501    67.5    68.2
## 502    67.5    68.2
## 503    67.5    68.2
## 504    67.5    68.2
## 505    66.5    68.2
## 506    66.5    68.2
## 507    66.5    68.2
## 508    66.5    68.2
## 509    66.5    68.2
## 510    66.5    68.2
## 511    66.5    68.2
## 512    66.5    68.2
## 513    66.5    68.2
## 514    66.5    68.2
## 515    66.5    68.2
## 516    66.5    68.2
## 517    66.5    68.2
## 518    66.5    68.2
## 519    65.5    68.2
## 520    65.5    68.2
## 521    65.5    68.2
## 522    65.5    68.2
## 523    65.5    68.2
## 524    65.5    68.2
## 525    65.5    68.2
## 526    64.0    68.2
## 527    72.5    69.2
## 528    72.5    69.2
## 529    71.5    69.2
## 530    71.5    69.2
## 531    71.5    69.2
## 532    71.5    69.2
## 533    71.5    69.2
## 534    70.5    69.2
## 535    70.5    69.2
## 536    70.5    69.2
## 537    70.5    69.2
## 538    70.5    69.2
## 539    70.5    69.2
## 540    70.5    69.2
## 541    70.5    69.2
## 542    70.5    69.2
## 543    70.5    69.2
## 544    70.5    69.2
## 545    70.5    69.2
## 546    70.5    69.2
```

```
## 547    70.5   69.2
## 548    70.5   69.2
## 549    70.5   69.2
## 550    70.5   69.2
## 551    70.5   69.2
## 552    69.5   69.2
## 553    69.5   69.2
## 554    69.5   69.2
## 555    69.5   69.2
## 556    69.5   69.2
## 557    69.5   69.2
## 558    69.5   69.2
## 559    69.5   69.2
## 560    69.5   69.2
## 561    69.5   69.2
## 562    69.5   69.2
## 563    69.5   69.2
## 564    69.5   69.2
## 565    69.5   69.2
## 566    69.5   69.2
## 567    69.5   69.2
## 568    69.5   69.2
## 569    69.5   69.2
## 570    69.5   69.2
## 571    69.5   69.2
## 572    69.5   69.2
## 573    69.5   69.2
## 574    69.5   69.2
## 575    69.5   69.2
## 576    69.5   69.2
## 577    69.5   69.2
## 578    69.5   69.2
## 579    69.5   69.2
## 580    69.5   69.2
## 581    69.5   69.2
## 582    69.5   69.2
## 583    69.5   69.2
## 584    69.5   69.2
## 585    68.5   69.2
## 586    68.5   69.2
## 587    68.5   69.2
## 588    68.5   69.2
## 589    68.5   69.2
## 590    68.5   69.2
## 591    68.5   69.2
## 592    68.5   69.2
## 593    68.5   69.2
## 594    68.5   69.2
## 595    68.5   69.2
## 596    68.5   69.2
## 597    68.5   69.2
## 598    68.5   69.2
## 599    68.5   69.2
## 600    68.5   69.2
```

```
## 601   68.5   69.2
## 602   68.5   69.2
## 603   68.5   69.2
## 604   68.5   69.2
## 605   68.5   69.2
## 606   68.5   69.2
## 607   68.5   69.2
## 608   68.5   69.2
## 609   68.5   69.2
## 610   68.5   69.2
## 611   68.5   69.2
## 612   68.5   69.2
## 613   68.5   69.2
## 614   68.5   69.2
## 615   68.5   69.2
## 616   68.5   69.2
## 617   68.5   69.2
## 618   68.5   69.2
## 619   68.5   69.2
## 620   68.5   69.2
## 621   68.5   69.2
## 622   68.5   69.2
## 623   68.5   69.2
## 624   68.5   69.2
## 625   68.5   69.2
## 626   68.5   69.2
## 627   68.5   69.2
## 628   68.5   69.2
## 629   68.5   69.2
## 630   68.5   69.2
## 631   68.5   69.2
## 632   68.5   69.2
## 633   67.5   69.2
## 634   67.5   69.2
## 635   67.5   69.2
## 636   67.5   69.2
## 637   67.5   69.2
## 638   67.5   69.2
## 639   67.5   69.2
## 640   67.5   69.2
## 641   67.5   69.2
## 642   67.5   69.2
## 643   67.5   69.2
## 644   67.5   69.2
## 645   67.5   69.2
## 646   67.5   69.2
## 647   67.5   69.2
## 648   67.5   69.2
## 649   67.5   69.2
## 650   67.5   69.2
## 651   67.5   69.2
## 652   67.5   69.2
## 653   67.5   69.2
## 654   67.5   69.2
```

```
## 655    67.5   69.2
## 656    67.5   69.2
## 657    67.5   69.2
## 658    67.5   69.2
## 659    67.5   69.2
## 660    67.5   69.2
## 661    67.5   69.2
## 662    67.5   69.2
## 663    67.5   69.2
## 664    67.5   69.2
## 665    67.5   69.2
## 666    67.5   69.2
## 667    67.5   69.2
## 668    67.5   69.2
## 669    67.5   69.2
## 670    67.5   69.2
## 671    66.5   69.2
## 672    66.5   69.2
## 673    66.5   69.2
## 674    66.5   69.2
## 675    66.5   69.2
## 676    66.5   69.2
## 677    66.5   69.2
## 678    66.5   69.2
## 679    66.5   69.2
## 680    66.5   69.2
## 681    66.5   69.2
## 682    66.5   69.2
## 683    66.5   69.2
## 684    65.5   69.2
## 685    65.5   69.2
## 686    65.5   69.2
## 687    65.5   69.2
## 688    65.5   69.2
## 689    65.5   69.2
## 690    65.5   69.2
## 691    64.5   69.2
## 692    64.5   69.2
## 693    64.0   69.2
## 694    72.5   70.2
## 695    71.5   70.2
## 696    71.5   70.2
## 697    71.5   70.2
## 698    71.5   70.2
## 699    71.5   70.2
## 700    71.5   70.2
## 701    71.5   70.2
## 702    71.5   70.2
## 703    71.5   70.2
## 704    71.5   70.2
## 705    70.5   70.2
## 706    70.5   70.2
## 707    70.5   70.2
## 708    70.5   70.2
```

```
## 709   70.5   70.2
## 710   70.5   70.2
## 711   70.5   70.2
## 712   70.5   70.2
## 713   70.5   70.2
## 714   70.5   70.2
## 715   70.5   70.2
## 716   70.5   70.2
## 717   70.5   70.2
## 718   70.5   70.2
## 719   69.5   70.2
## 720   69.5   70.2
## 721   69.5   70.2
## 722   69.5   70.2
## 723   69.5   70.2
## 724   69.5   70.2
## 725   69.5   70.2
## 726   69.5   70.2
## 727   69.5   70.2
## 728   69.5   70.2
## 729   69.5   70.2
## 730   69.5   70.2
## 731   69.5   70.2
## 732   69.5   70.2
## 733   69.5   70.2
## 734   69.5   70.2
## 735   69.5   70.2
## 736   69.5   70.2
## 737   69.5   70.2
## 738   69.5   70.2
## 739   69.5   70.2
## 740   69.5   70.2
## 741   69.5   70.2
## 742   69.5   70.2
## 743   69.5   70.2
## 744   68.5   70.2
## 745   68.5   70.2
## 746   68.5   70.2
## 747   68.5   70.2
## 748   68.5   70.2
## 749   68.5   70.2
## 750   68.5   70.2
## 751   68.5   70.2
## 752   68.5   70.2
## 753   68.5   70.2
## 754   68.5   70.2
## 755   68.5   70.2
## 756   68.5   70.2
## 757   68.5   70.2
## 758   68.5   70.2
## 759   68.5   70.2
## 760   68.5   70.2
## 761   68.5   70.2
## 762   68.5   70.2
```

```
## 763    68.5   70.2
## 764    68.5   70.2
## 765    67.5   70.2
## 766    67.5   70.2
## 767    67.5   70.2
## 768    67.5   70.2
## 769    67.5   70.2
## 770    67.5   70.2
## 771    67.5   70.2
## 772    67.5   70.2
## 773    67.5   70.2
## 774    67.5   70.2
## 775    67.5   70.2
## 776    67.5   70.2
## 777    67.5   70.2
## 778    67.5   70.2
## 779    67.5   70.2
## 780    67.5   70.2
## 781    67.5   70.2
## 782    67.5   70.2
## 783    67.5   70.2
## 784    66.5   70.2
## 785    66.5   70.2
## 786    66.5   70.2
## 787    66.5   70.2
## 788    65.5   70.2
## 789    65.5   70.2
## 790    65.5   70.2
## 791    65.5   70.2
## 792    65.5   70.2
## 793    72.5   71.2
## 794    72.5   71.2
## 795    71.5   71.2
## 796    71.5   71.2
## 797    71.5   71.2
## 798    71.5   71.2
## 799    70.5   71.2
## 800    70.5   71.2
## 801    70.5   71.2
## 802    70.5   71.2
## 803    70.5   71.2
## 804    70.5   71.2
## 805    70.5   71.2
## 806    69.5   71.2
## 807    69.5   71.2
## 808    69.5   71.2
## 809    69.5   71.2
## 810    69.5   71.2
## 811    69.5   71.2
## 812    69.5   71.2
## 813    69.5   71.2
## 814    69.5   71.2
## 815    69.5   71.2
## 816    69.5   71.2
```

```
## 817   69.5   71.2
## 818   69.5   71.2
## 819   69.5   71.2
## 820   69.5   71.2
## 821   69.5   71.2
## 822   69.5   71.2
## 823   69.5   71.2
## 824   69.5   71.2
## 825   69.5   71.2
## 826   68.5   71.2
## 827   68.5   71.2
## 828   68.5   71.2
## 829   68.5   71.2
## 830   68.5   71.2
## 831   68.5   71.2
## 832   68.5   71.2
## 833   68.5   71.2
## 834   68.5   71.2
## 835   68.5   71.2
## 836   68.5   71.2
## 837   68.5   71.2
## 838   68.5   71.2
## 839   68.5   71.2
## 840   68.5   71.2
## 841   68.5   71.2
## 842   68.5   71.2
## 843   68.5   71.2
## 844   67.5   71.2
## 845   67.5   71.2
## 846   67.5   71.2
## 847   67.5   71.2
## 848   67.5   71.2
## 849   67.5   71.2
## 850   67.5   71.2
## 851   67.5   71.2
## 852   67.5   71.2
## 853   67.5   71.2
## 854   67.5   71.2
## 855   65.5   71.2
## 856   65.5   71.2
## 857   73.0   72.2
## 858   72.5   72.2
## 859   72.5   72.2
## 860   72.5   72.2
## 861   72.5   72.2
## 862   72.5   72.2
## 863   72.5   72.2
## 864   72.5   72.2
## 865   71.5   72.2
## 866   71.5   72.2
## 867   71.5   72.2
## 868   71.5   72.2
## 869   71.5   72.2
## 870   71.5   72.2
```

```
## 871   71.5   72.2
## 872   71.5   72.2
## 873   71.5   72.2
## 874   70.5   72.2
## 875   70.5   72.2
## 876   70.5   72.2
## 877   70.5   72.2
## 878   69.5   72.2
## 879   69.5   72.2
## 880   69.5   72.2
## 881   69.5   72.2
## 882   69.5   72.2
## 883   69.5   72.2
## 884   69.5   72.2
## 885   69.5   72.2
## 886   69.5   72.2
## 887   69.5   72.2
## 888   69.5   72.2
## 889   68.5   72.2
## 890   68.5   72.2
## 891   68.5   72.2
## 892   68.5   72.2
## 893   67.5   72.2
## 894   67.5   72.2
## 895   67.5   72.2
## 896   67.5   72.2
## 897   65.5   72.2
## 898   73.0   73.2
## 899   73.0   73.2
## 900   73.0   73.2
## 901   72.5   73.2
## 902   72.5   73.2
## 903   71.5   73.2
## 904   71.5   73.2
## 905   70.5   73.2
## 906   70.5   73.2
## 907   70.5   73.2
## 908   69.5   73.2
## 909   69.5   73.2
## 910   69.5   73.2
## 911   69.5   73.2
## 912   68.5   73.2
## 913   68.5   73.2
## 914   68.5   73.2
## 915   72.5   73.7
## 916   72.5   73.7
## 917   72.5   73.7
## 918   72.5   73.7
## 919   71.5   73.7
## 920   71.5   73.7
## 921   70.5   73.7
## 922   70.5   73.7
## 923   70.5   73.7
## 924   69.5   73.7
```

```
## 925    69.5  73.7
## 926    69.5  73.7
## 927    69.5  73.7
## 928    69.5  73.7
```

TO-DO

Find the average height (include both parents and children in this computation).

```
avg_height = (mean(Galton$parent) + mean(Galton$child)) / 2
```

If you were to use the null model, what would the RMSE be of this model be?

```
rmse_null = sqrt(mean((Galton$child - avg_height) ^ 2))
```

Note that in Math 241 you learned that the sample average is an estimate of the "mean", the population expected value of height. We will call the average the "mean" going forward since it is probably correct to the nearest tenth of an inch with this amount of data.

Run a linear model attempting to explain the childrens' height using the parents' height. Use `lm` and use the R formula notation. Compute and report $b_0$, $b_1$, RMSE and $R^2$. Use the correct units to report these quantities.

```
mod = lm(child ~ parent, Galton)
b_0 = mod$coefficients[1]
b_1 = mod$coefficients[2]
b_0
```

```
## (Intercept)
##    23.94153
```

```
b_1
```

```
##    parent
## 0.6462906
```

```
summary(mod)$sigma
```

```
## [1] 2.238547
```

```
summary(mod)$r.sq
```

```
## [1] 0.2104629
```

Interpret all four quantities: $b_0$, $b_1$, RMSE and $R^2$.

b_0: b_0 is the intercept. b_1: for every unit increase in parents heigh the average height of the child increases by 0.6462906 (b_1) R^2: Shows the amount of error explained in the model. Since R^2 is only ~21% the model has a lot of variance unaccounted for and is likely a bad model.

How good is this model? How well does it predict? Discuss.

The model is not good because the R^2 value is only about 21%. Because of this, it is likely the model in unable to predict accuractly due to the amount of variance unaccounted for.

It is reasonable to assume that parents and their children have the same height? Explain why this is reasonable using basic biology and common sense.

It is reasonable to assume that the height of the child will match the height the match of the parents because they parents pass on genetics that code for that specific height. Though the height won't be exactly equal due to other factors influencing height such as diet, it is highly unlikely the child would be a significantly different height compared to the parents due to genetics.

If they were to have the same height and any differences were just random noise with expectation 0, what would the values of $\beta_0$ and $\beta_1$ be?

beta_0 would be 0 beta_1 would be 1

Let's plot (a) the data in $\mathbb{D}$ as black dots, (b) your least squares line defined by $b_0$ and $b_1$ in blue, (c) the theoretical line $\beta_0$ and $\beta_1$ if the parent-child height equality held in red and (d) the mean height in green.

```
pacman::p_load(ggplot2)
ggplot(Galton, aes(x = parent, y = child)) +
  geom_point() +
  geom_jitter() +
  geom_abline(intercept = b_0, slope = b_1, color = "blue", size = 1) +
  geom_abline(intercept = 0, slope = 1, color = "red", size = 1) +
  geom_abline(intercept = avg_height, slope = 0, color = "darkgreen", size = 1) +
  xlim(63.5, 72.5) +
  ylim(63.5, 72.5) +
  coord_equal(ratio = 1)
```

```
## Warning: Removed 76 rows containing missing values (geom_point).
```

```
## Warning: Removed 89 rows containing missing values (geom_point).
```

Fill in the following sentence:

TO-DO: Children of short parents became ... on average and children of tall parents became ... on average.

Why did Galton call it "Regression towards mediocrity in hereditary stature" which was later shortened to "regression to the mean"?

The average height of the children "regressed" to the "mean" of the height of the parents. The children of taller parents ended up being shorter, averaging towards the mean height. The children of shorter parrents ended up being taller, averaging towards the mean about.

Why should this effect be real?

This effect likely occured due to sampling error. Since results are taken from children that are tall and short, the two will average out towards the mean. If taking a similar sample size from two opposing values they will always regress towards the means because their values would average each other out.

You now have unlocked the mystery. Why is it that when modeling with $y$ continuous, everyone calls it "regression"? Write a better, more descriptive and appropriate name for building predictive models with $y$ continuous.

It is called a regression because of Galton's findings that the average heights of short children and tall children regressed towards the average height, coining the term "regression" and the terminology stayed. A better descriptive name for this model would be "Best Fit Line Estimator"

Create a dataset $\mathbb{D}$ which we call `Xy` such that the linear model as $R^2$ about 50% and RMSE approximately 1.

```
x = #TO DO
y = #TO DO
```

27

```
Xy = data.frame(x = x, y = y)
```

Create a dataset $\mathbb{D}$ which we call `Xy` such that the linear model as $R^2$ about 0% but x, y are clearly associated.

```
x = seq(1,1000)
y = x * ((-1)^x)

lm(y~x)
```

```
##
## Call:
## lm(formula = y ~ x)
##
## Coefficients:
## (Intercept)            x
##   -1.003003     0.003003
```

```
summary(lm(y~x))
```

```
##
## Call:
## lm(formula = y ~ x)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -1001.0  -500.0     1.5   500.0   998.0
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -1.003003  36.606216  -0.027    0.978
## x            0.003003   0.063356   0.047    0.962
##
## Residual standard error: 578.4 on 998 degrees of freedom
## Multiple R-squared:  2.251e-06,  Adjusted R-squared:  -0.0009998
## F-statistic: 0.002247 on 1 and 998 DF,  p-value: 0.9622
```

```
Xy = data.frame(x = x, y = y)
```

Extra credit: create a dataset $\mathbb{D}$ and a model (hint: not a linear model) that can give you $R^2$ arbitrarily close to 1 but RMSE arbitrarily high.

```
#TO-DO
```

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. What would a reasonable, naive prediction be under all Species? Hint: it's what we did in class.

```
g0 = function(x){
  prediction = mean( (iris[iris$Species == x, ])$Petal.Length)
  prediction
}

g0("setosa")
```

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

```r
g0("versicolor")
```

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

```r
g0("virginica")
```

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

Prove that this is the OLS model by fitting an appropriate `lm` and then using the predict function to verify you get the same answers as you wrote previously. Show this by doing a linear regression with and without the intercept.

```r
?predict()
```

```
## starting httpd help server ... done
```

```r
x = iris$Species
y = iris$Petal.Length
```

```r
lm(y ~ x)
```

```
##
## Call:
## lm(formula = y ~ x)
##
## Coefficients:
## (Intercept)  xversicolor    xvirginica
##       1.462        2.798         4.090
```

```r
lm(y ~ 0 + x)
```

```
##
## Call:
## lm(formula = y ~ 0 + x)
##
## Coefficients:
##     xsetosa  xversicolor    xvirginica
##       1.462        4.260         5.552
```

```r
predict(lm(y ~ x))
```

```
##     1     2     3     4     5     6     7     8     9    10    11    12    13
## 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462
##    14    15    16    17    18    19    20    21    22    23    24    25    26
## 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462
##    27    28    29    30    31    32    33    34    35    36    37    38    39
## 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462
```

```
##    40    41    42    43    44    45    46    47    48    49    50    51    52
## 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 4.260 4.260
##    53    54    55    56    57    58    59    60    61    62    63    64    65
## 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260
##    66    67    68    69    70    71    72    73    74    75    76    77    78
## 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260
##    79    80    81    82    83    84    85    86    87    88    89    90    91
## 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260
##    92    93    94    95    96    97    98    99   100   101   102   103   104
## 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 5.552 5.552 5.552 5.552
##   105   106   107   108   109   110   111   112   113   114   115   116   117
## 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552
##   118   119   120   121   122   123   124   125   126   127   128   129   130
## 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552
##   131   132   133   134   135   136   137   138   139   140   141   142   143
## 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552
##   144   145   146   147   148   149   150
## 5.552 5.552 5.552 5.552 5.552 5.552 5.552
```

```r
predict(lm(y ~ 0 + x))
```

```
##     1     2     3     4     5     6     7     8     9    10    11    12    13
## 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462
##    14    15    16    17    18    19    20    21    22    23    24    25    26
## 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462
##    27    28    29    30    31    32    33    34    35    36    37    38    39
## 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462
##    40    41    42    43    44    45    46    47    48    49    50    51    52
## 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 1.462 4.260 4.260
##    53    54    55    56    57    58    59    60    61    62    63    64    65
## 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260
##    66    67    68    69    70    71    72    73    74    75    76    77    78
## 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260
##    79    80    81    82    83    84    85    86    87    88    89    90    91
## 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260
##    92    93    94    95    96    97    98    99   100   101   102   103   104
## 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 4.260 5.552 5.552 5.552 5.552
##   105   106   107   108   109   110   111   112   113   114   115   116   117
## 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552
##   118   119   120   121   122   123   124   125   126   127   128   129   130
## 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552
##   131   132   133   134   135   136   137   138   139   140   141   142   143
## 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552 5.552
##   144   145   146   147   148   149   150
## 5.552 5.552 5.552 5.552 5.552 5.552 5.552
```

Use the `model.matrix` function to compute the matrix `X` for the regression with the intercept and without the intercept. What is different?

```r
x = iris$Species
y = iris$Petal.Length

model.matrix(y ~ x)
```

```
##    (Intercept) xversicolor xvirginica
## 1            1           0          0
## 2            1           0          0
## 3            1           0          0
## 4            1           0          0
## 5            1           0          0
## 6            1           0          0
## 7            1           0          0
## 8            1           0          0
## 9            1           0          0
## 10           1           0          0
## 11           1           0          0
## 12           1           0          0
## 13           1           0          0
## 14           1           0          0
## 15           1           0          0
## 16           1           0          0
## 17           1           0          0
## 18           1           0          0
## 19           1           0          0
## 20           1           0          0
## 21           1           0          0
## 22           1           0          0
## 23           1           0          0
## 24           1           0          0
## 25           1           0          0
## 26           1           0          0
## 27           1           0          0
## 28           1           0          0
## 29           1           0          0
## 30           1           0          0
## 31           1           0          0
## 32           1           0          0
## 33           1           0          0
## 34           1           0          0
## 35           1           0          0
## 36           1           0          0
## 37           1           0          0
## 38           1           0          0
## 39           1           0          0
## 40           1           0          0
## 41           1           0          0
## 42           1           0          0
## 43           1           0          0
## 44           1           0          0
## 45           1           0          0
## 46           1           0          0
## 47           1           0          0
## 48           1           0          0
## 49           1           0          0
## 50           1           0          0
## 51           1           1          0
## 52           1           1          0
## 53           1           1          0
```

```
## 54               1           1           0
## 55               1           1           0
## 56               1           1           0
## 57               1           1           0
## 58               1           1           0
## 59               1           1           0
## 60               1           1           0
## 61               1           1           0
## 62               1           1           0
## 63               1           1           0
## 64               1           1           0
## 65               1           1           0
## 66               1           1           0
## 67               1           1           0
## 68               1           1           0
## 69               1           1           0
## 70               1           1           0
## 71               1           1           0
## 72               1           1           0
## 73               1           1           0
## 74               1           1           0
## 75               1           1           0
## 76               1           1           0
## 77               1           1           0
## 78               1           1           0
## 79               1           1           0
## 80               1           1           0
## 81               1           1           0
## 82               1           1           0
## 83               1           1           0
## 84               1           1           0
## 85               1           1           0
## 86               1           1           0
## 87               1           1           0
## 88               1           1           0
## 89               1           1           0
## 90               1           1           0
## 91               1           1           0
## 92               1           1           0
## 93               1           1           0
## 94               1           1           0
## 95               1           1           0
## 96               1           1           0
## 97               1           1           0
## 98               1           1           0
## 99               1           1           0
## 100              1           1           0
## 101              1           0           1
## 102              1           0           1
## 103              1           0           1
## 104              1           0           1
## 105              1           0           1
## 106              1           0           1
## 107              1           0           1
```

```
## 108              1            0            1
## 109              1            0            1
## 110              1            0            1
## 111              1            0            1
## 112              1            0            1
## 113              1            0            1
## 114              1            0            1
## 115              1            0            1
## 116              1            0            1
## 117              1            0            1
## 118              1            0            1
## 119              1            0            1
## 120              1            0            1
## 121              1            0            1
## 122              1            0            1
## 123              1            0            1
## 124              1            0            1
## 125              1            0            1
## 126              1            0            1
## 127              1            0            1
## 128              1            0            1
## 129              1            0            1
## 130              1            0            1
## 131              1            0            1
## 132              1            0            1
## 133              1            0            1
## 134              1            0            1
## 135              1            0            1
## 136              1            0            1
## 137              1            0            1
## 138              1            0            1
## 139              1            0            1
## 140              1            0            1
## 141              1            0            1
## 142              1            0            1
## 143              1            0            1
## 144              1            0            1
## 145              1            0            1
## 146              1            0            1
## 147              1            0            1
## 148              1            0            1
## 149              1            0            1
## 150              1            0            1
## attr(,"assign")
## [1] 0 1 1
## attr(,"contrasts")
## attr(,"contrasts")$x
## [1] "contr.treatment"
```

```r
model.matrix(y ~ 0 + x)
```

```
##      xsetosa xversicolor xvirginica
## 1          1            0           0
## 2          1            0           0
```

```
## 3            1               0               0
## 4            1               0               0
## 5            1               0               0
## 6            1               0               0
## 7            1               0               0
## 8            1               0               0
## 9            1               0               0
## 10           1               0               0
## 11           1               0               0
## 12           1               0               0
## 13           1               0               0
## 14           1               0               0
## 15           1               0               0
## 16           1               0               0
## 17           1               0               0
## 18           1               0               0
## 19           1               0               0
## 20           1               0               0
## 21           1               0               0
## 22           1               0               0
## 23           1               0               0
## 24           1               0               0
## 25           1               0               0
## 26           1               0               0
## 27           1               0               0
## 28           1               0               0
## 29           1               0               0
## 30           1               0               0
## 31           1               0               0
## 32           1               0               0
## 33           1               0               0
## 34           1               0               0
## 35           1               0               0
## 36           1               0               0
## 37           1               0               0
## 38           1               0               0
## 39           1               0               0
## 40           1               0               0
## 41           1               0               0
## 42           1               0               0
## 43           1               0               0
## 44           1               0               0
## 45           1               0               0
## 46           1               0               0
## 47           1               0               0
## 48           1               0               0
## 49           1               0               0
## 50           1               0               0
## 51           0               1               0
## 52           0               1               0
## 53           0               1               0
## 54           0               1               0
## 55           0               1               0
## 56           0               1               0
```

```
## 57           0              1              0
## 58           0              1              0
## 59           0              1              0
## 60           0              1              0
## 61           0              1              0
## 62           0              1              0
## 63           0              1              0
## 64           0              1              0
## 65           0              1              0
## 66           0              1              0
## 67           0              1              0
## 68           0              1              0
## 69           0              1              0
## 70           0              1              0
## 71           0              1              0
## 72           0              1              0
## 73           0              1              0
## 74           0              1              0
## 75           0              1              0
## 76           0              1              0
## 77           0              1              0
## 78           0              1              0
## 79           0              1              0
## 80           0              1              0
## 81           0              1              0
## 82           0              1              0
## 83           0              1              0
## 84           0              1              0
## 85           0              1              0
## 86           0              1              0
## 87           0              1              0
## 88           0              1              0
## 89           0              1              0
## 90           0              1              0
## 91           0              1              0
## 92           0              1              0
## 93           0              1              0
## 94           0              1              0
## 95           0              1              0
## 96           0              1              0
## 97           0              1              0
## 98           0              1              0
## 99           0              1              0
## 100          0              1              0
## 101          0              0              1
## 102          0              0              1
## 103          0              0              1
## 104          0              0              1
## 105          0              0              1
## 106          0              0              1
## 107          0              0              1
## 108          0              0              1
## 109          0              0              1
## 110          0              0              1
```

```
## 111          0               0               1
## 112          0               0               1
## 113          0               0               1
## 114          0               0               1
## 115          0               0               1
## 116          0               0               1
## 117          0               0               1
## 118          0               0               1
## 119          0               0               1
## 120          0               0               1
## 121          0               0               1
## 122          0               0               1
## 123          0               0               1
## 124          0               0               1
## 125          0               0               1
## 126          0               0               1
## 127          0               0               1
## 128          0               0               1
## 129          0               0               1
## 130          0               0               1
## 131          0               0               1
## 132          0               0               1
## 133          0               0               1
## 134          0               0               1
## 135          0               0               1
## 136          0               0               1
## 137          0               0               1
## 138          0               0               1
## 139          0               0               1
## 140          0               0               1
## 141          0               0               1
## 142          0               0               1
## 143          0               0               1
## 144          0               0               1
## 145          0               0               1
## 146          0               0               1
## 147          0               0               1
## 148          0               0               1
## 149          0               0               1
## 150          0               0               1
## attr(,"assign")
## [1] 1 1 1
## attr(,"contrasts")
## attr(,"contrasts")$x
## [1] "contr.treatment"
```

Including the intercept turned "setosa" into a reference variable, eliminating it as a category since it is now the default species and the rows added up to 1 or 2. Without the intercept there is no default species and the rows added up to 1.