# Lab 4

*Karen Lopez*

*11:59PM March 9, 2019*

Note: the content of this lab is on the midterm exam (March 5) even though the lab itself is due after the midterm exam.

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:

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

Solve for the least squares line by computing $b_0$ and $b_1$ *without* using the functions `mean`, `cor`, `cov`, `var`, `sd` but instead computing it from the $x$ and $y$ quantities manually using base function such as `sum` and other basic operators. See the class notes.

```
mean_x = sum(x) / n
mean_y = sum(y) / n
b_1 = (sum((x * y)) - n * mean_x * mean_y) / (sum(x^2) - n*mean_x^2)
b_0 = mean_y - b_1 * mean_x
```

Verify your computations are correct using the `lm` function in R:

```
lm_mod = lm(y~x)
b_vec = coef(lm_mod)
pacman::p_load(testthat)
expect_equal(b_0, as.numeric(b_vec[1]), tol = 1e-4)
expect_equal(b_1, as.numeric(b_vec[2]), tol = 1e-4)
```

6. 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`.

```
summary(Galton)
```

```
##      parent          child
##  Min.   :64.00   Min.   :61.70
##  1st Qu.:67.50   1st Qu.:66.20
##  Median :68.50   Median :68.20
##  Mean   :68.31   Mean   :68.09
##  3rd Qu.:69.50   3rd Qu.:70.20
##  Max.   :73.00   Max.   :73.70
```

```
str(Galton)
```

```
## 'data.frame':    928 obs. of  2 variables:
##  $ parent: num  70.5 68.5 65.5 64.5 64 67.5 67.5 67.5 66.5 66.5 ...
##  $ child : num  61.7 61.7 61.7 61.7 61.7 62.2 62.2 62.2 62.2 62.2 ...
```

```
?Galton
```

$n = 928$ observations $p = 2$ variables - represented in two columns one named parent which represents the average height of the father and mother and the second column named chil which is the height of the child

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

```
avg_height = (sum(Galton$parent+Galton$child))/(nrow(Galton)*2)
```

```
mean(c(Galton$parent,Galton$child))
```

```
## [1] 68.19833
```

```
(sum(Galton$parent)+sum(Galton$child))/(928*2)
```

```
## [1] 68.19833
```

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

```
y_hat=rep(avg_height,1856) #null model
y_vec=c(Galton$parent,Galton$child)
SSE=sum((y_vec-y_hat)^2)
MSE=SSE/1854
RMSE=sqrt(MSE)
```

```
RMSE
```

```
## [1] 2.186179
```

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.

```
ols1 = lm(child~parent,Galton)
summary(ols1)
```

```
##
## Call:
## lm(formula = child ~ parent, data = Galton)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -7.8050 -1.3661  0.0487  1.6339  5.9264
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 23.94153    2.81088   8.517   <2e-16 ***
## parent       0.64629    0.04114  15.711   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
## 
## Residual standard error: 2.239 on 926 degrees of freedom
## Multiple R-squared:  0.2105, Adjusted R-squared:  0.2096
## F-statistic: 246.8 on 1 and 926 DF,  p-value: < 2.2e-16
```

```r
length(ols1$residuals)
```

```
## [1] 928
```

```r
sse = sum(ols1$residuals^2) #units= inches squared
mse = sse / (928 - 2) #units = inches squared
rmse = sqrt(mse) #units = inches
rmse #
```

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

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

$b_0 = 23.9415$ inches (intercept) $b_1 = 0.6463$ inches (For a one inch increase in the parents height the child high increases by .6463 inches) $R_2 = 0.2105$ (21 percent) (21% of variance in y explained) rmse $= 2.238547$ inches (the model is plus or minus 2.23 inch off of the actual.. the model plus or minus 4.46 is 95% coffidence set of y)

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

Since the $R_2$ represents 21% of the variance in y. Also based on the RMSE your off by +/- 4 inches. THerefore, it does better than the null model but it is not a great model.

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.

This is reasonable to assume because of the way genetics work so it is expected that children reflect their parents because they share genes.

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?
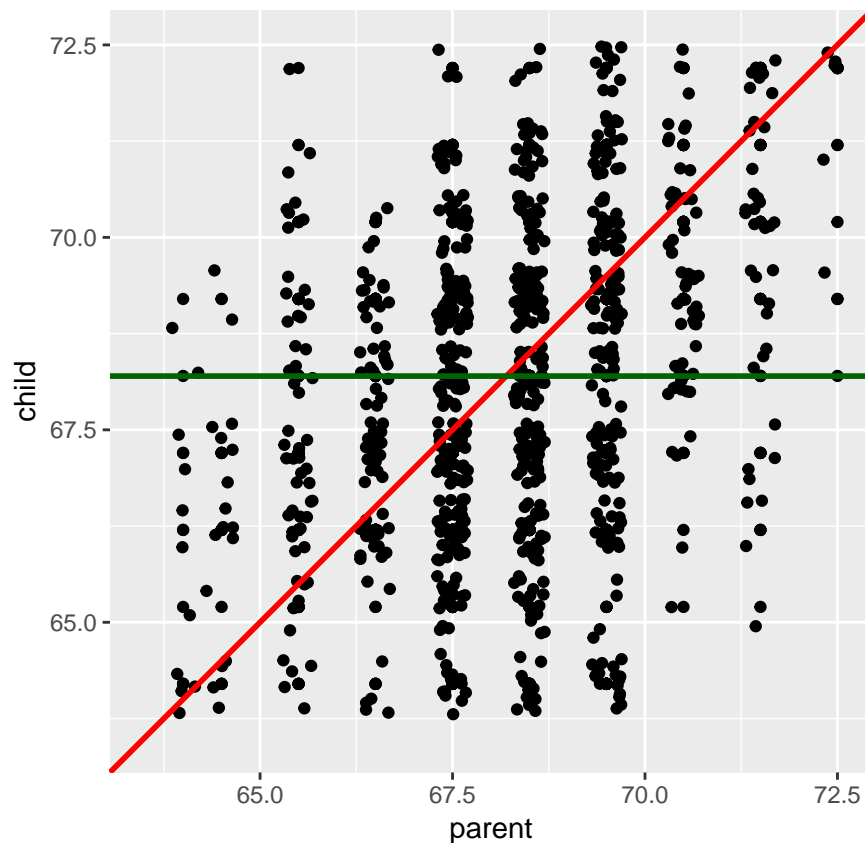
TO-DO

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.

```r
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 88 rows containing missing values (geom_point).
```

Fill in the following sentence:

Children of short parents became taller on average and children of tall parents became shorter on average.

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

Galton called it "Regression towards mediocrity in heredtary stature" becasue essentially the regression line or relationship produced brings us towards the mediocre values essentially average or moderate values of height believed to be a hereditary construction.

Why should this effect be real?

When n is sufficiently large or we have enough observations, heights tend to go towards a central tendency.

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.

TO-DO

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
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.

```
#circle
x = 2^2 + 3^2
```

4

```
y = 13
Xy = data.frame(x = x, y = y)
```

Load up the famous iris dataset and drop the data for Species "virginica".

```
data("iris")
summary(iris)
```

```
##   Sepal.Length    Sepal.Width    Petal.Length    Petal.Width
## Min.   :4.300   Min.   :2.000   Min.   :1.000   Min.   :0.100
## 1st Qu.:5.100   1st Qu.:2.800   1st Qu.:1.600   1st Qu.:0.300
## Median :5.800   Median :3.000   Median :4.350   Median :1.300
## Mean   :5.843   Mean   :3.057   Mean   :3.758   Mean   :1.199
## 3rd Qu.:6.400   3rd Qu.:3.300   3rd Qu.:5.100   3rd Qu.:1.800
## Max.   :7.900   Max.   :4.400   Max.   :6.900   Max.   :2.500
##        Species
## setosa    :50
## versicolor:50
## virginica :50
##
##
##
```

```
iris=iris[iris$Species != "virginica",]
iris
```

```
##    Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## 1           5.1         3.5          1.4         0.2  setosa
## 2           4.9         3.0          1.4         0.2  setosa
## 3           4.7         3.2          1.3         0.2  setosa
## 4           4.6         3.1          1.5         0.2  setosa
## 5           5.0         3.6          1.4         0.2  setosa
## 6           5.4         3.9          1.7         0.4  setosa
## 7           4.6         3.4          1.4         0.3  setosa
## 8           5.0         3.4          1.5         0.2  setosa
## 9           4.4         2.9          1.4         0.2  setosa
## 10          4.9         3.1          1.5         0.1  setosa
## 11          5.4         3.7          1.5         0.2  setosa
## 12          4.8         3.4          1.6         0.2  setosa
## 13          4.8         3.0          1.4         0.1  setosa
## 14          4.3         3.0          1.1         0.1  setosa
## 15          5.8         4.0          1.2         0.2  setosa
## 16          5.7         4.4          1.5         0.4  setosa
## 17          5.4         3.9          1.3         0.4  setosa
## 18          5.1         3.5          1.4         0.3  setosa
## 19          5.7         3.8          1.7         0.3  setosa
## 20          5.1         3.8          1.5         0.3  setosa
## 21          5.4         3.4          1.7         0.2  setosa
## 22          5.1         3.7          1.5         0.4  setosa
## 23          4.6         3.6          1.0         0.2  setosa
## 24          5.1         3.3          1.7         0.5  setosa
## 25          4.8         3.4          1.9         0.2  setosa
## 26          5.0         3.0          1.6         0.2  setosa
## 27          5.0         3.4          1.6         0.4  setosa
## 28          5.2         3.5          1.5         0.2  setosa
```

```
## 29           5.2         3.4         1.4         0.2      setosa
## 30           4.7         3.2         1.6         0.2      setosa
## 31           4.8         3.1         1.6         0.2      setosa
## 32           5.4         3.4         1.5         0.4      setosa
## 33           5.2         4.1         1.5         0.1      setosa
## 34           5.5         4.2         1.4         0.2      setosa
## 35           4.9         3.1         1.5         0.2      setosa
## 36           5.0         3.2         1.2         0.2      setosa
## 37           5.5         3.5         1.3         0.2      setosa
## 38           4.9         3.6         1.4         0.1      setosa
## 39           4.4         3.0         1.3         0.2      setosa
## 40           5.1         3.4         1.5         0.2      setosa
## 41           5.0         3.5         1.3         0.3      setosa
## 42           4.5         2.3         1.3         0.3      setosa
## 43           4.4         3.2         1.3         0.2      setosa
## 44           5.0         3.5         1.6         0.6      setosa
## 45           5.1         3.8         1.9         0.4      setosa
## 46           4.8         3.0         1.4         0.3      setosa
## 47           5.1         3.8         1.6         0.2      setosa
## 48           4.6         3.2         1.4         0.2      setosa
## 49           5.3         3.7         1.5         0.2      setosa
## 50           5.0         3.3         1.4         0.2      setosa
## 51           7.0         3.2         4.7         1.4 versicolor
## 52           6.4         3.2         4.5         1.5 versicolor
## 53           6.9         3.1         4.9         1.5 versicolor
## 54           5.5         2.3         4.0         1.3 versicolor
## 55           6.5         2.8         4.6         1.5 versicolor
## 56           5.7         2.8         4.5         1.3 versicolor
## 57           6.3         3.3         4.7         1.6 versicolor
## 58           4.9         2.4         3.3         1.0 versicolor
## 59           6.6         2.9         4.6         1.3 versicolor
## 60           5.2         2.7         3.9         1.4 versicolor
## 61           5.0         2.0         3.5         1.0 versicolor
## 62           5.9         3.0         4.2         1.5 versicolor
## 63           6.0         2.2         4.0         1.0 versicolor
## 64           6.1         2.9         4.7         1.4 versicolor
## 65           5.6         2.9         3.6         1.3 versicolor
## 66           6.7         3.1         4.4         1.4 versicolor
## 67           5.6         3.0         4.5         1.5 versicolor
## 68           5.8         2.7         4.1         1.0 versicolor
## 69           6.2         2.2         4.5         1.5 versicolor
## 70           5.6         2.5         3.9         1.1 versicolor
## 71           5.9         3.2         4.8         1.8 versicolor
## 72           6.1         2.8         4.0         1.3 versicolor
## 73           6.3         2.5         4.9         1.5 versicolor
## 74           6.1         2.8         4.7         1.2 versicolor
## 75           6.4         2.9         4.3         1.3 versicolor
## 76           6.6         3.0         4.4         1.4 versicolor
## 77           6.8         2.8         4.8         1.4 versicolor
## 78           6.7         3.0         5.0         1.7 versicolor
## 79           6.0         2.9         4.5         1.5 versicolor
## 80           5.7         2.6         3.5         1.0 versicolor
## 81           5.5         2.4         3.8         1.1 versicolor
## 82           5.5         2.4         3.7         1.0 versicolor
```

```
## 83             5.8          2.7          3.9          1.2 versicolor
## 84             6.0          2.7          5.1          1.6 versicolor
## 85             5.4          3.0          4.5          1.5 versicolor
## 86             6.0          3.4          4.5          1.6 versicolor
## 87             6.7          3.1          4.7          1.5 versicolor
## 88             6.3          2.3          4.4          1.3 versicolor
## 89             5.6          3.0          4.1          1.3 versicolor
## 90             5.5          2.5          4.0          1.3 versicolor
## 91             5.5          2.6          4.4          1.2 versicolor
## 92             6.1          3.0          4.6          1.4 versicolor
## 93             5.8          2.6          4.0          1.2 versicolor
## 94             5.0          2.3          3.3          1.0 versicolor
## 95             5.6          2.7          4.2          1.3 versicolor
## 96             5.7          3.0          4.2          1.2 versicolor
## 97             5.7          2.9          4.2          1.3 versicolor
## 98             6.2          2.9          4.3          1.3 versicolor
## 99             5.1          2.5          3.0          1.1 versicolor
## 100            5.7          2.8          4.1          1.3 versicolor
```

If 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 both Species? Hint: it's what we did in class.

```
mean_versicolor = iris[mean(iris$Species == "versicolor")]
mean_setosa = iris[mean(iris$Species == "setosa")]
y = mean(iris$Petal.Length)


#g = mean_versicolor + (mean_versicolor - mean_setosa)
```

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.

```
#TO-DO
```