

Exploration librairie swirl : Multiple Testing

> swirl()

Welcome to swirl! Please sign in. If you've been here before, use the same name as you did then. If you are new, call yourself something unique.

*What shall I call you? **jlbellier***

Please choose a course, or type 0 to exit swirl.

- 1: Regression Models*
- 2: Statistical Inference*
- 3: Take me to the swirl course repository!*

*Selection: **1***

Please choose a lesson, or type 0 to return to course menu.

- 1: Introduction*
- 2: Residuals*
- 3: Least Squares Estimation*
- 4: Residual Variation*
- 5: Introduction to Multivariable Regression*
- 6: MultiVar Examples*
- 7: MultiVar Examples2*
- 8: MultiVar Examples3*
- 9: Residuals Diagnostics and Variation*
- 10: Variance Inflation Factors*
- 11: Overfitting and Underfitting*
- 12: Binary Outcomes*
- 13: Count Outcomes*

*Selection: **1***

0%

Introduction to Regression Models. (Slides for this and other Data Science courses may be found at github <https://github.com/DataScienceSpecialization/courses> if you want to use them. They must be downloaded as a zip file and viewed locally. This lesson corresponds to Regression_Models/01_01_introduction.)

===

5%

This is the first lesson on Regression Models. We'll begin with the concept of "regression toward the mean" and illustrate it with some pioneering work of the father of forensic science, Sir Francis Galton.

=====

10%

Sir Francis studied the relationship between heights of parents and their children. His work showed that parents who were taller than average had children who were also tall but closer to the average height. Similarly, parents who were shorter than average had children who were also shorter than average but less so than mom and dad. That is, they were closer to the average height. From one generation to the next the heights moved closer to the average or regressed toward the mean.

=====

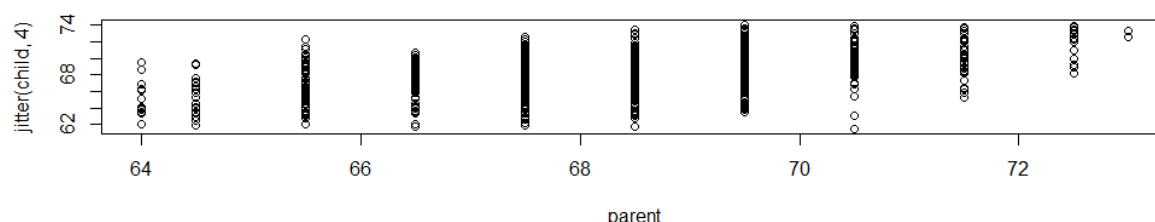
14%

For this lesson we'll use Sir Francis's parent/child height data which we've taken the liberty to load for you as the variable, galton. (Data is from John Verzani's website, <http://wiener.math.csi.cuny.edu/UsingR/>.) So let's get started!

=====

19%

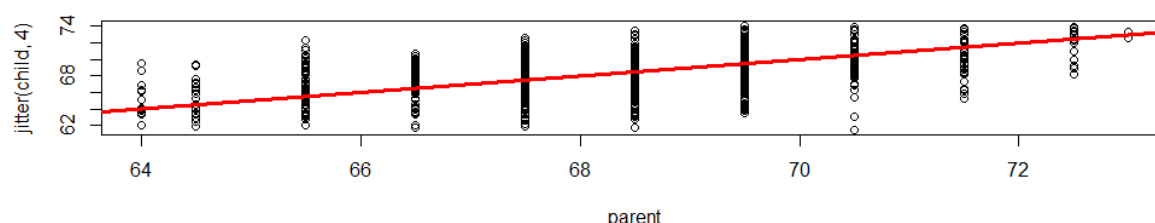
Here is a plot of Galton's data, a set of 928 parent/child height pairs. Moms' and dads' heights were averaged together (after moms' heights were adjusted by a factor of 1.08). In our plot we used the R function "jitter" on the children's heights to highlight heights that occurred most frequently. The dark spots in each column rise from left to right suggesting that children's heights do depend on their parents'. Tall parents have tall children and short parents have short children.



=====

24%

Here we add a red (45 degree) line of slope 1 and intercept 0 to the plot. If children tended to be the same height as their parents, we would expect the data to vary evenly about this line. We see this isn't the case. On the left half of the plot we see a concentration of heights above the line, and on the right half we see the concentration below the line.

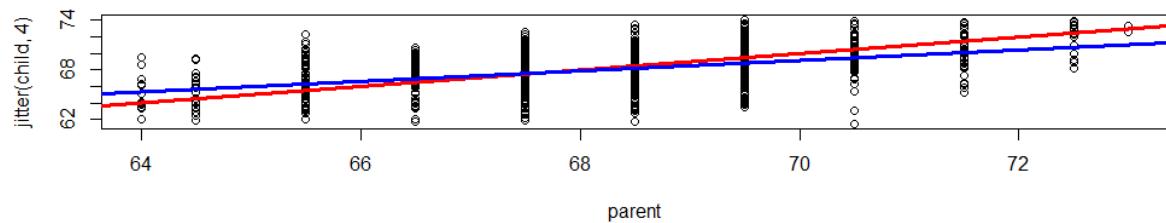


=====

29%

Now we've added a blue regression line to the plot. This is the line which has the minimum variation of the data around it. (For theory see the slides.) Its slope is greater than zero

indicating that parents' heights do affect their children's. The slope is also less than 1 as would have been the case if children tended to be the same height as their parents.



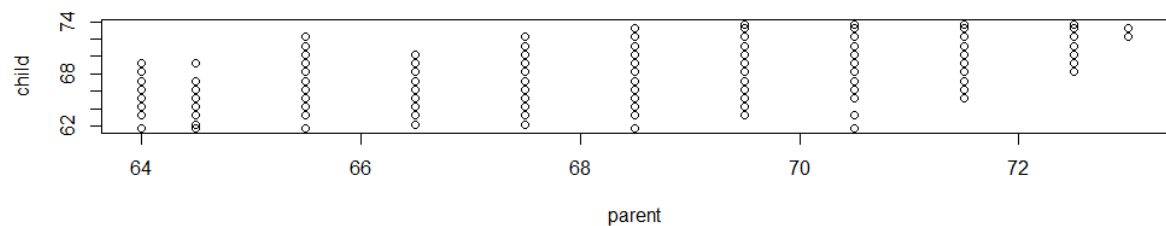
=====

33%

Now's your chance to plot in R. Type "`plot(child ~ parent, galton)`" at the R prompt.

> `plot(child ~ parent, galton)`

Excellent job!



=====

38%

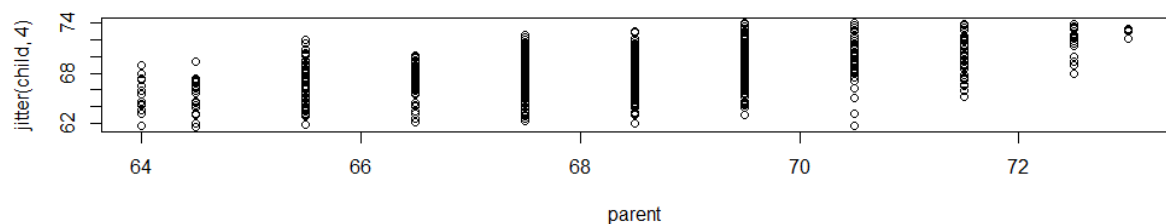
You'll notice that this plot looks a lot different than the original we displayed. Why? Many people are the same height to within measurement error, so points fall on top of one another. You can see that some circles appear darker than others. However, by using R's function "jitter" on the children's heights, we can spread out the data to simulate the measurement errors and make high frequency heights more visible.

=====

43%

Now it's your turn to try. Just type "`plot(jitter(child,4) ~ parent,galton)`" and see the magic.

> `plot(jitter(child,4) ~ parent,galton)`



You got it!

=====

48%

Now for the regression line. This is quite easy in R. The function `lm` (linear model) needs a formula and dataset. You can type `?formula` for more information, but, in simple terms, we just need to specify the dependent variable (children's heights) ~ the independent variable (parents' heights).

=====

52%

So generate the regression line and store it in the variable `regline`. Type `"regline <- lm(child ~ parent, galton)"`

```
> regline <- lm(child ~ parent, galton)
```

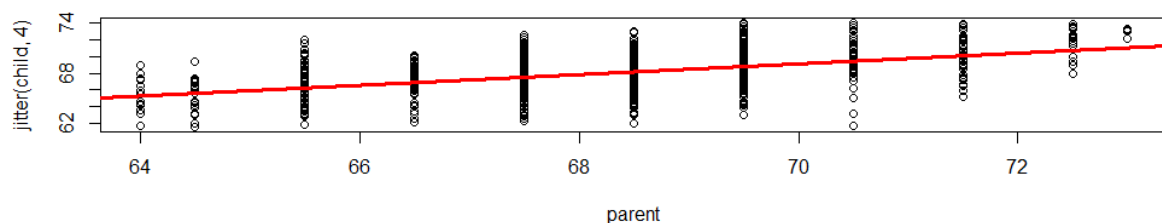
All that hard work is paying off!

=====

57%

Now add the regression line to the plot with `"abline"`. Make the line wide and red for visibility. Type `"abline(regline, lwd=3, col='red')"`

```
> abline(regline, lwd=3, col='red')
```



Excellent work!

=====

62%

The regression line will have a slope and intercept which are estimated from data. Estimates are not exact. Their accuracy is gauged by theoretical techniques and expressed in terms of standard error." You can use `"summary(regline)"` to examine the Galton regression line. Do this now.

```
> summary(regline)
```

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

Perseverance, that's the answer.

=====

67%

The slope of the line is the estimate of the coefficient, or multiplier, of "parent", the independent variable of our data (in this case, the parents' heights). From the output of "summary" what is the slope of the regression line?

- 1: 23.94153
- 2: .64629
- 3: .04114

Selection: **1**

Not quite, but you're learning! Try again.

Look at the line labelled "parent" and the column "Estimate"

- 1: .64629
- 2: .04114
- 3: 23.94153

Selection: **1**

That's a job well done!

=====

71%

What is the standard error of the slope?

- 1: .64629
- 2: .04114
- 3: 23.94153

Selection: **2**

You got it right!

=====

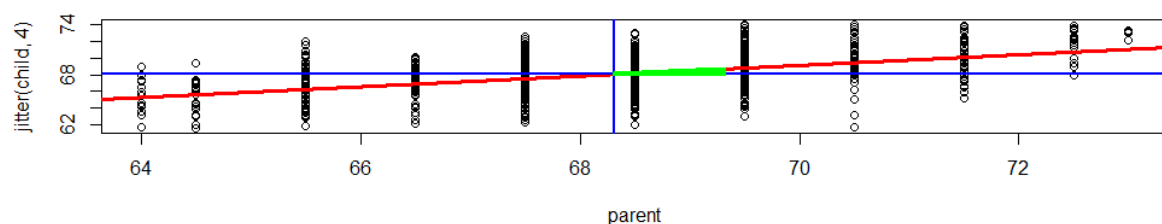
76%

A coefficient will be within 2 standard errors of its estimate about 95% of the time. This means the slope of our regression is significantly different than either 0 or 1 since $(.64629) \pm (2 \cdot .04114)$ is near neither 0 nor 1.

=====

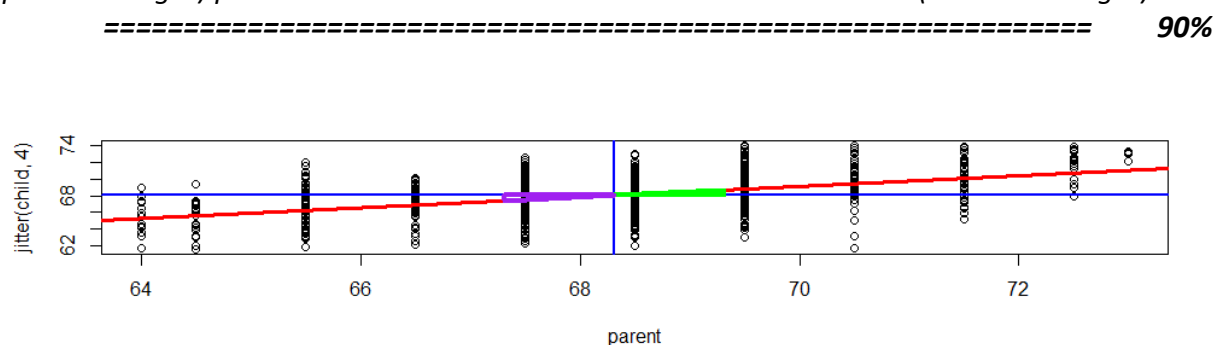
81%

We're now adding two blue lines to indicate the means of the children's heights (horizontal) and the parents' (vertical). Note that these lines and the regression line all intersect in a point. Pretty cool, huh? We'll talk more about this in a later lesson. (Something you can look forward to.)



===== 86%

The slope of a line shows how much of a change in the vertical direction is produced by a change in the horizontal direction. So, parents "1 inch" above the mean in height tend to have children who are only .65 inches above the mean. The green triangle illustrates this point. From the mean, moving a "1 inch distance" horizontally to the right (increasing the parents' height) produces a ".65 inch" increase in the vertical direction (children's height).



===== 95%

Similarly, parents who are 1 inch below average in height have children who are only .65 inches below average height. The purple triangle illustrates this. From the mean, moving a "1 inch distance" horizontally to the left (decreasing the parents' height) produces a ".65 inch" decrease in the vertical direction (children's height).

===== 95%

This concludes our lesson on regression toward the mean. We hope you found it above average!

===== 100%