

```

1  Course: Regression_Models
2  Lesson: Residual_Variation
3
4
5  - Class: text
6  Output: "Residual Variation. (Slides for this and other Data Science courses may be
found at github https://github.com/DataScienceSpecialization/courses. If you care to
use them, they must be downloaded as a zip file and viewed locally. This lesson
corresponds to Regression Models/01_06_residualVariation. Galton data is from John
Verzani's website, http://wiener.math.csi.cuny.edu/UsingR/)"
7
8  - Class: text
9  Output: As shown in the slides, residuals are useful for indicating how well data
points fit a statistical model. They "can be thought of as the outcome (Y) with the
linear association of the predictor (X) removed. One differentiates residual
variation (variation after removing the predictor) from systematic variation
(variation explained by the regression model)."
10
11 - Class: text
12 Output: It can also be shown that, given a model, the maximum likelihood estimate of
the variance of the random error is the average squared residual. However, since our
linear model with one predictor requires two parameters we have only (n-2) degrees of
freedom. Therefore, to calculate an "average" squared residual to estimate the
variance we use the formula  $1/(n-2) * (\text{the sum of the squared residuals})$ . If we
divided the sum of the squared residuals by n, instead of n-2, the result would give
a biased estimate.
13
14 - Class: cmd_question
15 Output: To see this we'll use our favorite Galton height data. First regenerate the
regression line and call it fit. Use the R function lm and recall that by default its
first argument is a formula such as "child ~ parent" and its second is the dataset,
in this case galton.
CorrectAnswer: fit <- lm(child ~ parent, galton)
AnswerTests: omnitest(correctExpr='fit <- lm(child ~ parent, galton)')
Hint: Type "fit <- lm(child ~ parent, galton)" at the R prompt.
16
17
18
19
20 - Class: text
21 Output: First, we'll use the residuals (fit$residuals) of our model to estimate the
standard deviation (sigma) of the error. We've already defined n for you as the
number of points in Galton's dataset (928).
22
23 - Class: cmd_question
24 Output: Calculate the sum of the squared residuals divided by the quantity (n-2).
Then take the square root.
CorrectAnswer: sqrt(sum(fit$residuals^2) / (n - 2))
AnswerTests: omnitest(correctExpr='sqrt(sum(fit$residuals^2) / (n - 2))')
Hint: Type "sqrt(sum(fit$residuals^2) / (n - 2))" at the R prompt.
25
26
27
28
29
30 - Class: cmd_question
31 Output: Now look at the "sigma" portion of the summary of fit, "summary(fit)$sigma".
CorrectAnswer: summary(fit)$sigma
AnswerTests: omnitest(correctExpr='summary(fit)$sigma')
Hint: Type "summary(fit)$sigma" at the R prompt.
32
33
34
35
36 - Class: text
37 Output: Pretty cool, huh?
38
39 - Class: cmd_question
40 Output: Another cool thing - take the sqrt of "deviance(fit)/(n-2)" at the R prompt.
CorrectAnswer: sqrt(deviance(fit)/(n-2))
AnswerTests: omnitest(correctExpr='sqrt(deviance(fit)/(n-2))')
Hint: Type "sqrt(deviance(fit)/(n-2))" at the R prompt.
41
42
43
44
45 - Class: text
46 Output: Another useful fact shown in the slides was
47
48 - Class: text
49 Output: Total Variation = Residual Variation + Regression Variation

```

```

50
51 - Class: mult_question
52 Output: Recall the beauty of the slide full of algebra which proved this fact. It
had a bunch of Y's, some with hats and some with bars and several summations of
squared values. The Y's with hats were the estimates provided by the model. (They
were on the regression line.) The Y with the bar was the mean or average of the data.
Which sum of squared term represented Total Variation?
53 AnswerChoices: Yi-mean(Yi); Yi-Yi_hat; Yi_hat-mean(Yi)
54 CorrectAnswer: Yi-mean(Yi)
55 AnswerTests: omnitest(correctVal='Yi-mean(Yi)')
56 Hint: Pick the choice which is independent of the estimated or predicted values, the
(hat terms).
57
58 - Class: mult_question
59 Output: Which sum of squared term represents Residual Variation?
60 AnswerChoices: Yi-Yi_hat; Yi-mean(Yi); Yi_hat-mean(Yi)
61 CorrectAnswer: Yi-Yi_hat
62 AnswerTests: omnitest(correctVal='Yi-Yi_hat')
63 Hint: Residuals represent the vertical distance between actual values and estimated
(hat) values.
64
65 - Class: text
66 Output: The term  $R^2$  represents the percent of total variation described by the
model, the regression variation (the term we didn't ask about in the preceding
multiple choice questions). Also, since it is a percent we need a ratio or fraction
of sums of squares. Let's do this now for our Galton data.
67
68 - Class: cmd_question
69 Output: We'll start with easy steps. Calculate the mean of the children's heights and
store it in a variable called mu. Recall that we reference the children's heights
with the expression 'galton$child' and the parents' heights with the expression
'galton$parent'.
70 CorrectAnswer: mu <- mean(galton$child)
71 AnswerTests: omnitest(correctExpr='mu <- mean(galton$child)')
72 Hint: Type "mu <- mean(galton$child)" at the R prompt.
73
74 - Class: cmd_question
75 Output: Recall that centering data means subtracting the mean from each data point.
Now calculate the sum of the squares of the centered children's heights and store
the result in a variable called sTot. This represents the Total Variation of the data.
76 CorrectAnswer: sTot <- sum((galton$child-mu)^2)
77 AnswerTests: ANY_of_exprs('sTot <- sum((galton$child-mu)^2)', 'sTot <-
sum((galton$child-mu)*(galton$child-mu))')
78 Hint: Type "sTot <- sum((galton$child-mu)^2)" at the R prompt.
79
80 - Class: cmd_question
81 Output: Now create the variable sRes. Use the R function deviance to calculate the
sum of the squares of the residuals. These are the distances between the children's
heights and the regression line. This represents the Residual Variation.
82 CorrectAnswer: sRes <- deviance(fit)
83 AnswerTests: omnitest(correctExpr='sRes <- deviance(fit)')
84 Hint: Type "sRes <- deviance(fit)" at the R prompt.
85
86 - Class: cmd_question
87 Output: Finally, the ratio sRes/sTot represents the percent of total variation
contributed by the residuals. To find the percent contributed by the model, i.e., the
regression variation, subtract the fraction sRes/sTot from 1. This is the value  $R^2$ .
88 CorrectAnswer: 1-sRes/sTot
89 AnswerTests: omnitest(correctExpr='1-sRes/sTot')
90 Hint: Type "1-sRes/sTot" at the R prompt.
91
92 - Class: cmd_question
93 Output: For fun you can compare your result to the values shown in
summary(fit)$r.squared to see if it looks familiar. Do this now.
94 CorrectAnswer: summary(fit)$r.squared
95 AnswerTests: omnitest(correctExpr='summary(fit)$r.squared')
96 Hint: Type "summary(fit)$r.squared" at the R prompt.
97
98 - Class: cmd_question

```

```
99      Output: To see some real magic, compute the square of the correlation of the galton
100      data, the children and parents. Use the R function cor.
101      CorrectAnswer: cor(galton$parent,galton$child)^2
102      AnswerTests:
103      ANY_of_exprs('cor(galton$parent,galton$child)^2','cor(galton$child,galton$parent)^2')
104      Hint: Type "cor(galton$parent,galton$child)^2" at the R prompt.
105
106      - Class: text
107      Output: We'll now summarize useful facts about R^2. It is the percentage of variation
108      explained by the regression model. As a percentage it is between 0 and 1. It also
109      equals the sample correlation squared. However, R^2 doesn't tell the whole story.
110
111      - Class: text
112      Output: Congrats! You've finished this lesson on Residual Variation.
113
114      - Class: mult_question
115      Output: "Would you like to receive credit for completing this course on
116      Coursera.org?"
117      CorrectAnswer: NULL
118      AnswerChoices: Yes;No
119      AnswerTests: coursera_on_demand()
120      Hint: ""
```