

STA 302/1001

Reminder Test is Tues Oct. 25

Regression

Origin of term: Galton, 19th c.

Predicting son's height from father's height

son's height = $26.46 + 0.61152$ father's height

Not $y = x$

At 68.1", model predicts son's height = father's height

Below that: sons of short fathers tend to be short, but not as short as their fathers (tall),

Explan: "Regression to the Mean"

In any test - re-test, expect to get "the regression effect"

Why?

$$b_1 = r \frac{S_y}{S_x}$$

if $S_y \approx S_x$

since $|r| \leq 1$

$$\Rightarrow |b_1| \leq 1$$

Example STA 302, in a previous year there were 2 tests (T_1, T_2)

For students who wrote both tests, use T_1 to predict T_2

Is $T_2 = T_1$ a reasonable guess?

NO

	<u>Mean</u>	<u>SD</u>
T1	62.7	18.4
T2	62.1	18.1

$$\hat{T}_2 = 20.27 + .667 T_1$$

If get 20% on T1, $\hat{T}_2 = 34$
 40% on T1, $\hat{T}_2 = 47$
 75% on T1, $\hat{T}_2 = 70$
 90% on T1, $\hat{T}_2 = \underline{\hspace{1cm}}$

On average, students who score low on 1st ^(high) test will score higher on second. + test

Why? Students with a high T1, are there because:

- That's their skill level
- Random variation caused them to score above their skill level

The regression effect happens because there are more students with skill levels near the mean than away from it.

The regression fallacy occurs when the regression effect is mistaken for a real effect

Examples - Studying that parents' people with extreme values of some measurement and sees how they respond to a treatment expect extremes to be closer to mean when re-measured.

(Solution: use a control group)

- "While overall no effect, people ~~with~~ ~~with~~ baseline measurement

was high, tended to go down and
" " " " " "

Answerable Example: Plot the data

 - evidence of a linear relationship from the
numbers is not enough

Chapter 3: Diagnositis and Transformations for
Simple Linear Regression

From preface: "It makes sense to base inferences or conclusions only on valid models."

Influential Points / Outliers / Leverage Points

Observations whose inclusion or exclusion result in substantial changes to the fitted model (est of coefficients, predicted values) are said to be INFLUENTIAL (eg. NYC in crime example)

Points can be entangling in any (or all or some) of the value of the explanatory variable,

The dependent variable, or its residual,

Outliers w.r.t. the residuals represent

model failures — line doesn't fit that point adequately