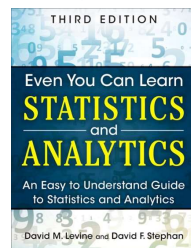


IMGD 2905

Simple Linear Regression

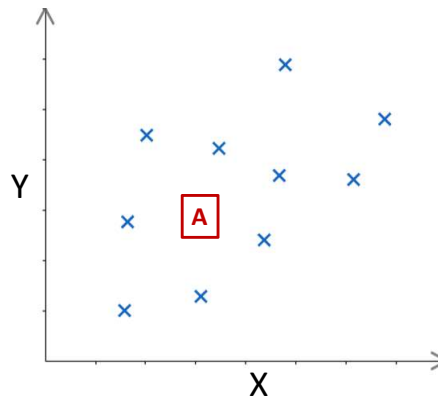
Chapter 10



1

Motivation

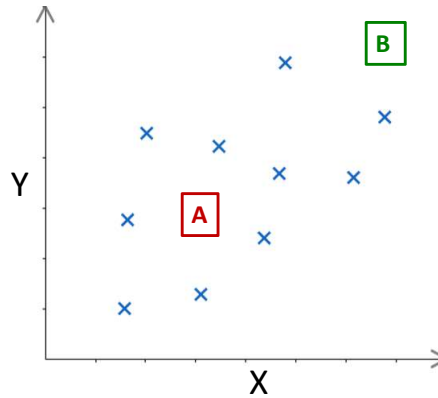
- Have data (sample, x 's)
- Want to know likely value of next observation
 - E.g., *playtime* versus *skins owned*
- A –



2

Motivation

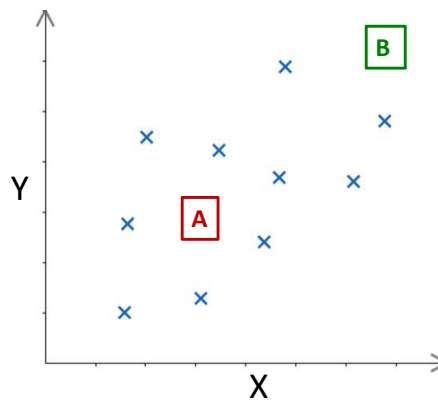
- Have data (sample, x 's)
- Want to know likely value of next observation
 - E.g., *playtime* versus *skins owned*
- **A** – reasonable to compute mean (with confidence interval)
- **B** –



3

Motivation

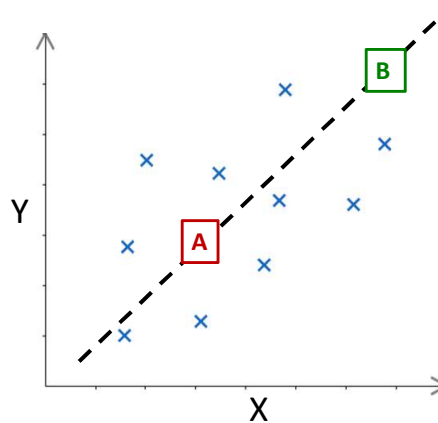
- Have data (sample, x 's)
- Want to know likely value of next observation
 - E.g., *playtime* versus *skins owned*
- **A** – reasonable to compute mean (with confidence interval)
- **B** – could do same, but there appears to be relationship between X and Y!



4

Motivation

- Have data (sample, x 's)
 - Want to know likely value of next observation
 - E.g., *playtime* versus *skins owned*
 - **A** – reasonable to compute mean (with confidence interval)
 - **B** – could do same, but there appears to be relationship between X and Y !
- **Predict B**
e.g., “trendline” (regression)



5

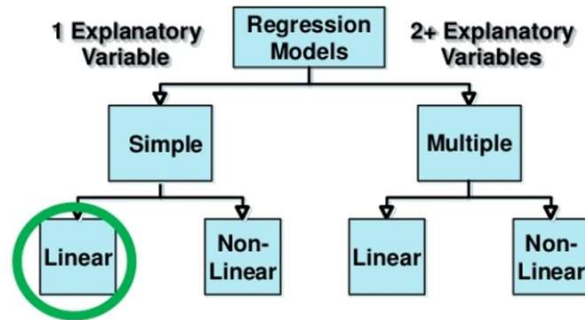
Overview

- Broadly, two types of **prediction techniques**:
 1. **Regression** – mathematical equation to model, use model for predictions
 - We'll discuss **simple linear regression**
 2. **Machine learning** – branch of AI, use computer algorithms to determine relationships (predictions)
 - **CS 453X Machine Learning**



6

Types of Regression Models



- Explanatory variable *explains* dependent variable
 - Variable X (e.g., skill level) explains Y (e.g., KDA)
 - Can have 1 or 2+
- Linear if coefficients added, else Non-linear

7

Outline

- Introduction (done)
- Simple Linear Regression (next)
 - Linear relationship
 - Residual analysis
 - Fitting parameters
- Measures of Variation
- Misc

8

Simple Linear Regression

- Goal – find a **linear** relationship between two values
 - E.g., kills and skill, time and car speed
- First, make sure relationship is linear! How?

9

Simple Linear Regression

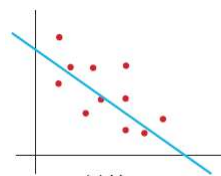
- Goal – find a **linear** relationship between two values
 - E.g., kills and skill, time and car speed
- First, make sure relationship is linear! How?

→ Scatterplot

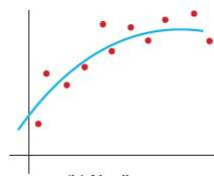
(c) no clear relationship

(b) not a linear relationship

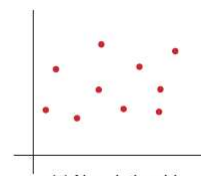
(a) linear relationship – proceed with linear regression



(a) Linear



(b) Nonlinear

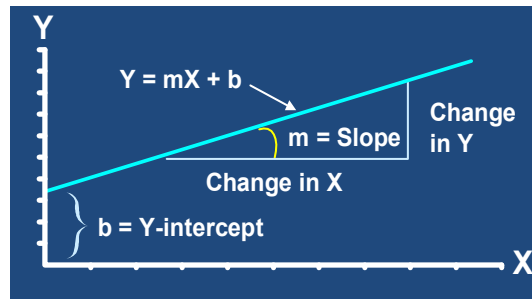


(c) No relationship

10

Linear Relationship

- From algebra: line in form $Y = mX + b$
 - m is slope, b is y-intercept
- Slope (m) is amount Y increases when X increases by 1 unit
- Intercept (b) is where line crosses y-axis, or where y-value when $x = 0$



11

Simple Linear Regression Example

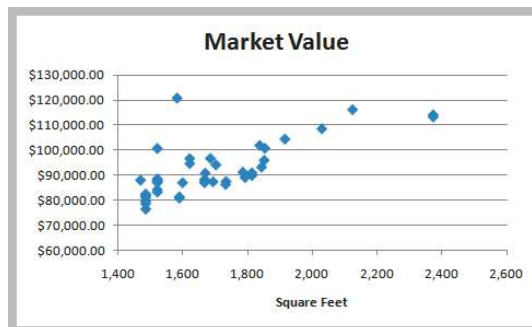
- Size of house related to its market value.

X = square footage
 Y = market value (\$)

- Scatter plot (42 homes) indicates linear trend



	A	B	C
1	Home Market Value		
2			
3	House Age	Square Feet	Market Value
4	33	1,812	\$90,000.00
5	32	1,914	\$104,400.00
6	32	1,842	\$93,300.00
7	33	1,812	\$91,000.00
8	32	1,836	\$101,900.00
9	33	2,028	\$108,500.00
10	32	1,732	\$87,600.00

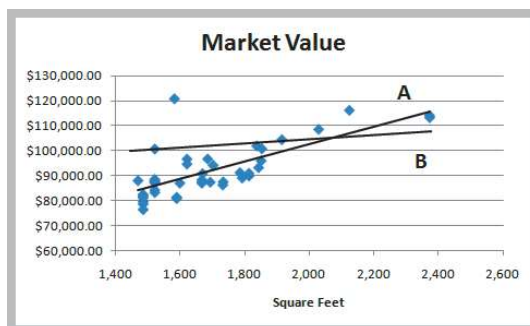


12

Simple Linear Regression Example

- Two possible lines shown below (A and B)
- Want to determine best regression line
- Line A looks a better fit to data
 - But how to know?

$$Y = mX + b$$



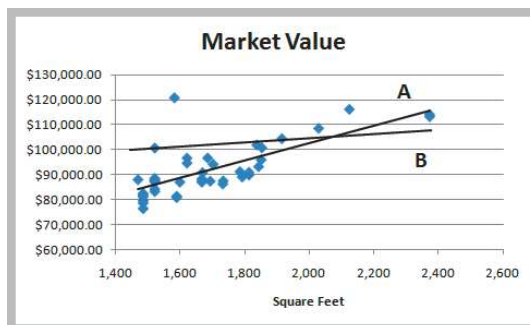
13

Simple Linear Regression Example

- Two possible lines shown below (A and B)
- Want to determine best regression line
- Line A looks a better fit to data
 - But how to know?

$$Y = mX + b$$

Line that gives best fit to data is one that minimizes prediction error
 → Least squares line
 (more later)

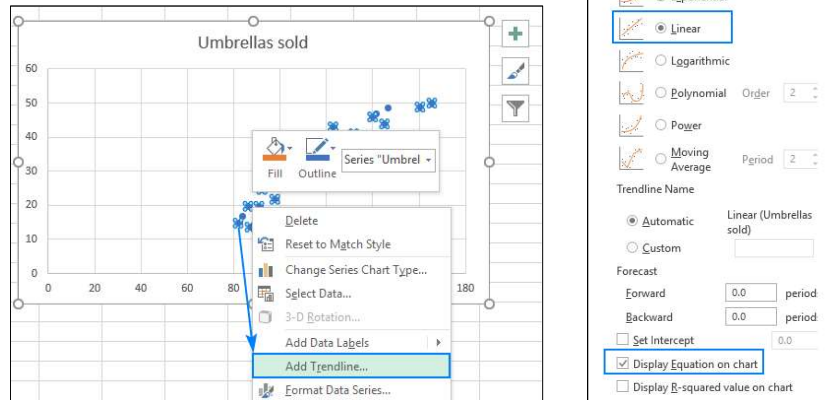


14

Simple Linear Regression Example

Chart

- Scatterplot
- Right click → Add Trendline



15

Simple Linear Regression Example

Formulas

=SLOPE(C4:C45, B4:B45)

- Slope = 35.036

=INTERCEPT(C4:C45, B4:B45)

- Intercept = 32,673

- Estimate Y when X = 1800 square feet

$$Y = 32,673 + 35.036 \times (1800) = \$95,737.80$$

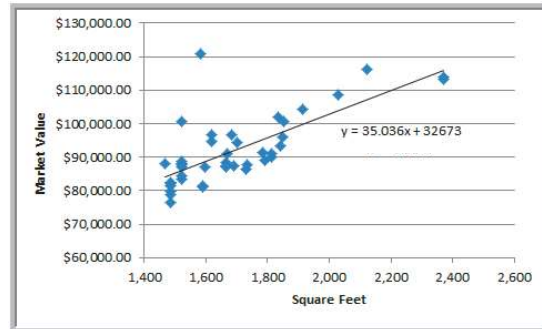
	A	B	C
1	Home Market Value		
2			
3	House Age	Square Feet	Market Value
4	33	1,812	\$90,000.00
5	32	1,914	\$104,400.00
6	32	1,842	\$93,300.00
7	33	1,812	\$91,000.00
8	32	1,836	\$101,900.00
9	33	2,028	\$108,500.00
10	32	1,732	\$87,600.00



16

Simple Linear Regression Example

- Market value = $32673 + 35.036 \times (\text{square feet})$
- Predicts market value better than just average



But before use, examine **residuals**



17

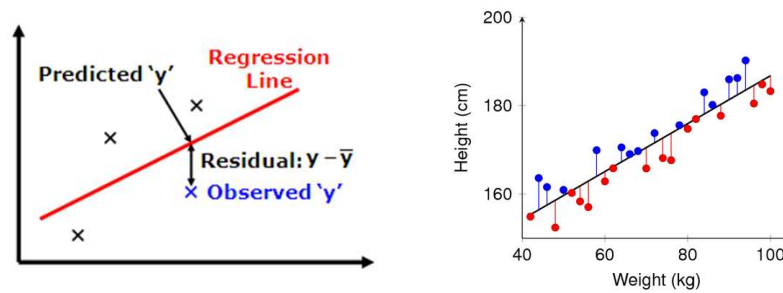
Outline

- Introduction (done)
- Simple Linear Regression
 - Linear relationship (done)
 - Residual analysis (next)
 - Fitting parameters
- Measures of Variation
- Misc

18

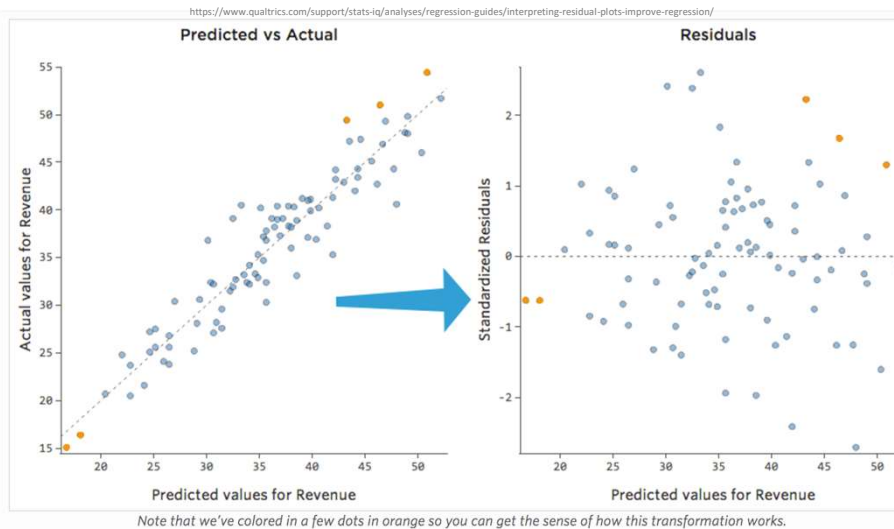
Residual Analysis

- Before predicting, confirm that linear regression assumptions hold
 - Variation around line is normally distributed
 - Variation equal for all X
 - Variation independent for all X
- How? Compute **residuals** (error in prediction) → Chart

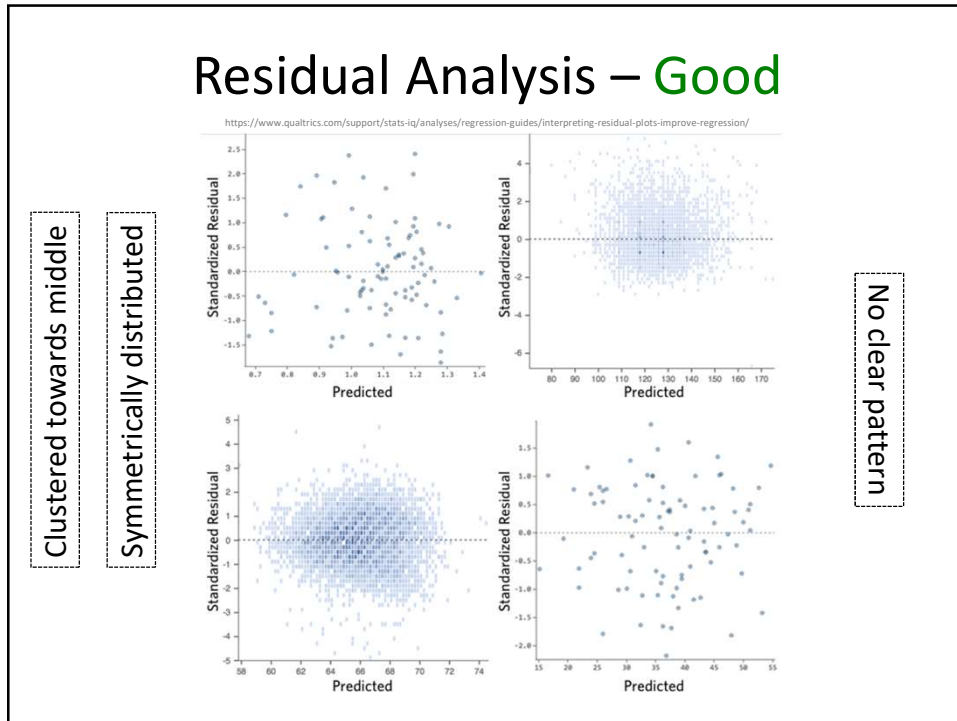


19

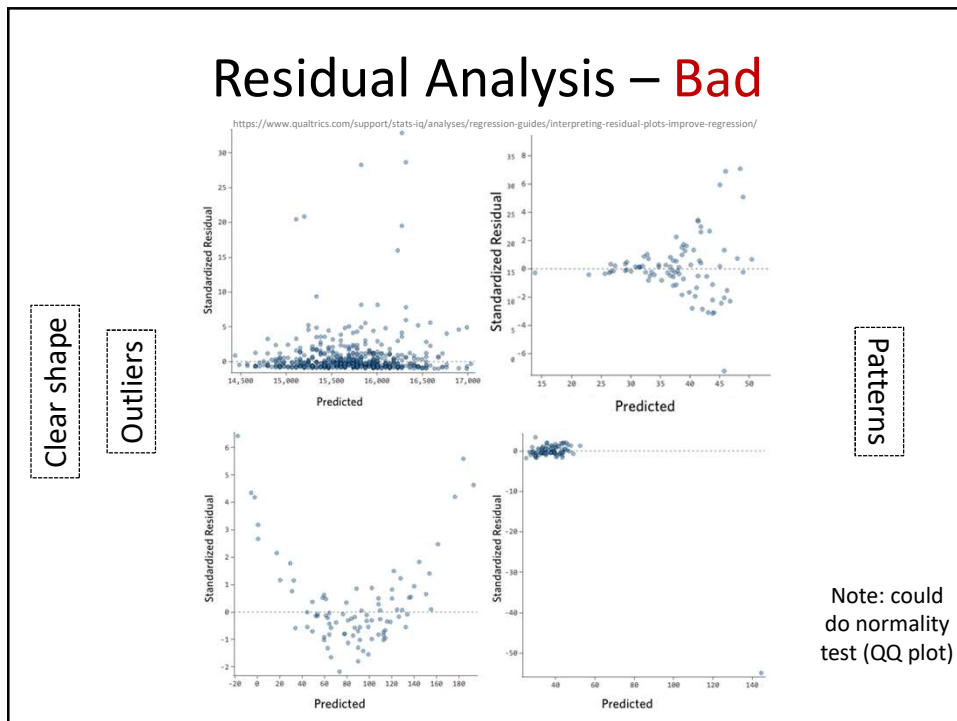
Residual Analysis



20



21

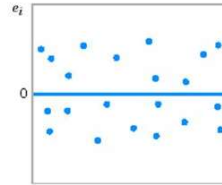


22

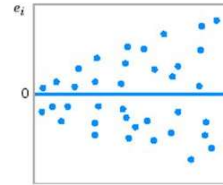
Residual Analysis – Summary

- Regression assumptions:

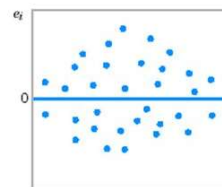
- Normality of variation around regression
- Equal variation for all y values
- Independence of variation



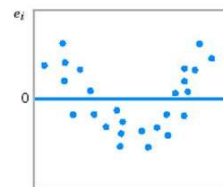
(a)



(b)



(c)



(d)

-
- (a) ok
 - (b) funnel
 - (c) double bow
 - (d) nonlinear

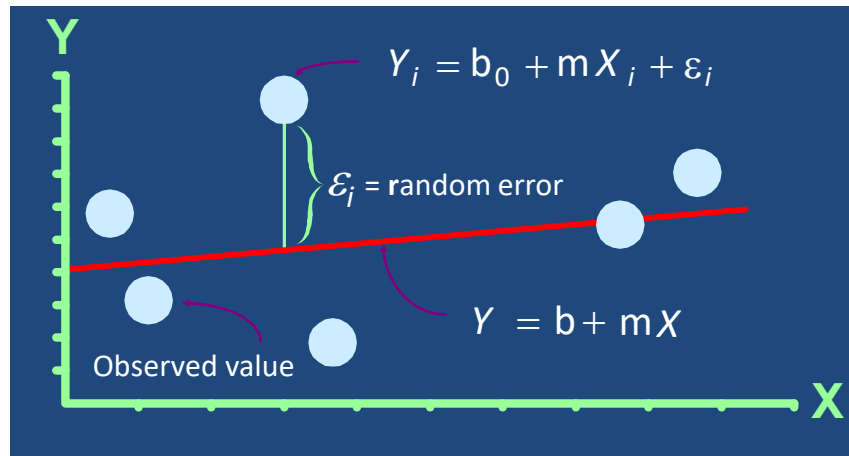
23

Outline

- Introduction (done)
- Simple Linear Regression
 - Linear relationship (done)
 - Residual analysis (done)
 - Fitting parameters (next)
- Measures of Variation
- Misc

24

Linear Regression Model

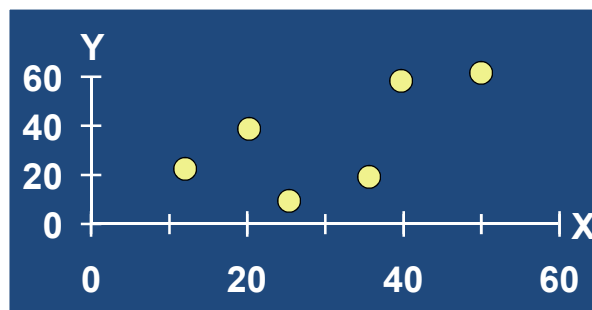


Random error associated with each observation

25

Fitting the Best Line

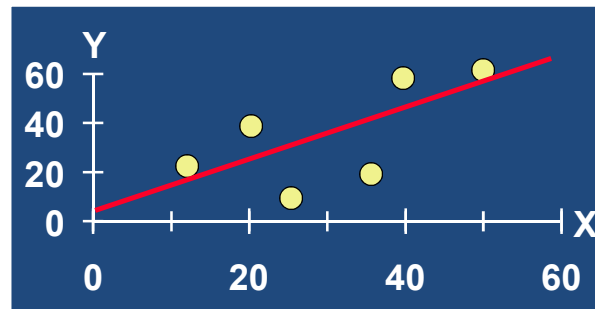
- Plot all (X_i, Y_i) Pairs



26

Fitting the Best Line

- Plot all (X_i, Y_i) Pairs
- Draw a line. But how do we know it is best?

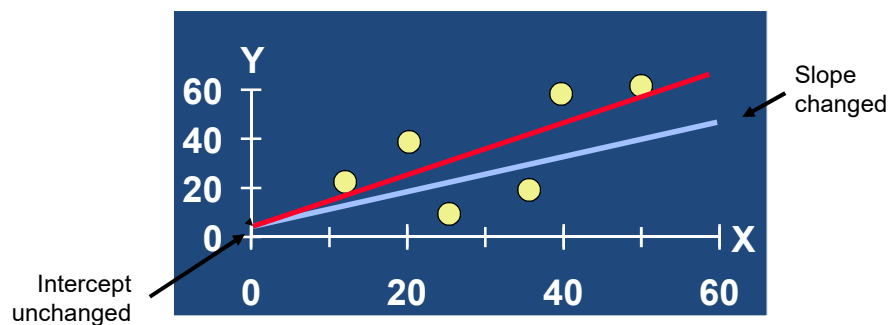


<https://www.scribd.com/presentation/23086725/Fu-Ch11-Linear-Regression>

27

Fitting the Best Line

- Plot all (X_i, Y_i) Pairs
- Draw a line. But how do we know it is best?

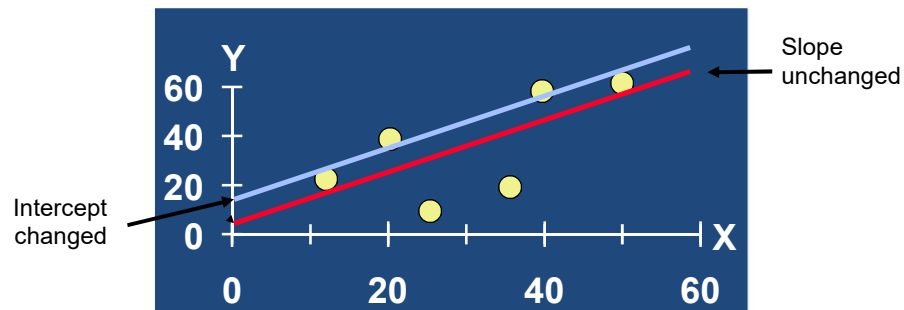


<https://www.scribd.com/presentation/23086725/Fu-Ch11-Linear-Regression>

28

Fitting the Best Line

- Plot all (X_i, Y_i) Pairs
- Draw a line. But how do we know it is best?

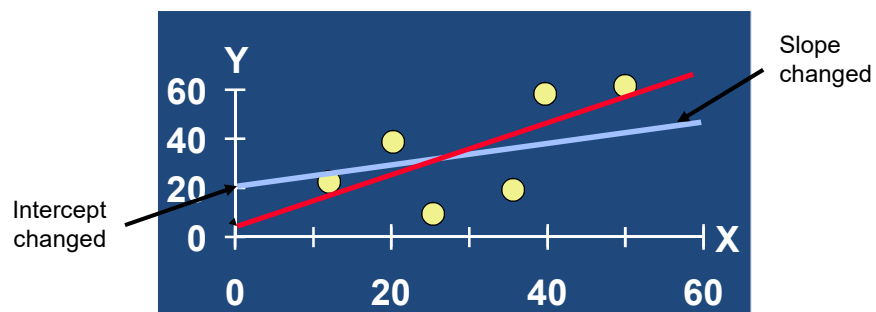


<https://www.scribd.com/presentation/230686725/Fu-Ch11-Linear-Regression>

29

Fitting the Best Line

- Plot all (X_i, Y_i) Pairs
- Draw a line. But how do we know it is best?



<https://www.scribd.com/presentation/230686725/Fu-Ch11-Linear-Regression>

30

Linear Regression Model

- Relationship between variables is linear function

$$Y_i = b_0 + mX_i + \epsilon_i$$

Population
Y-Intercept

Dependent
(response)
Variable
(e.g., kills)

Population
Slope

Independent (explanatory)
Variable
(e.g., skill level)

Random
Prediction
Error

Want error
as small as
possible

31

Least Squares Line

- Want to minimize difference between actual y and predicted \hat{y}
 - Add up ϵ_i for all observed y 's
 - But positive differences offset negative ones
 - (remember when this happened for variance?)
 - Square the errors! Then, minimize (using Calculus)

Simple Linear Regression:

SUM $(y - \hat{y})^2 \rightarrow \min$

Minimize: $\sum (y_i - b_0 - b_1 X_i)^2$

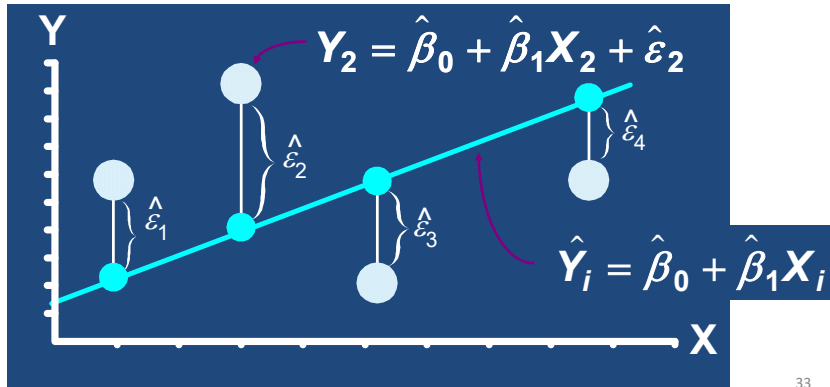
Take derivative

Set to 0 and solve

32

Least Squares Line Graphically

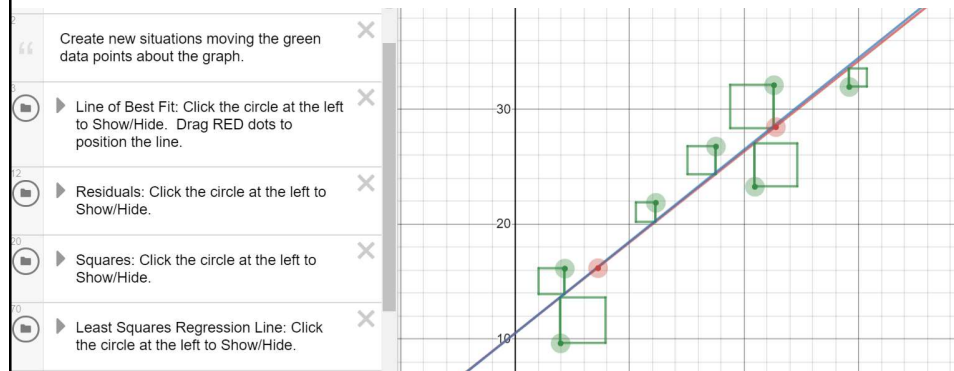
LS minimizes $\sum_{i=1}^n \hat{\varepsilon}_i^2 = \hat{\varepsilon}_1^2 + \hat{\varepsilon}_2^2 + \hat{\varepsilon}_3^2 + \hat{\varepsilon}_4^2$



33

33

Least Squares Line Graphically



<https://www.desmos.com/calculator/zvrc4lg3cr>

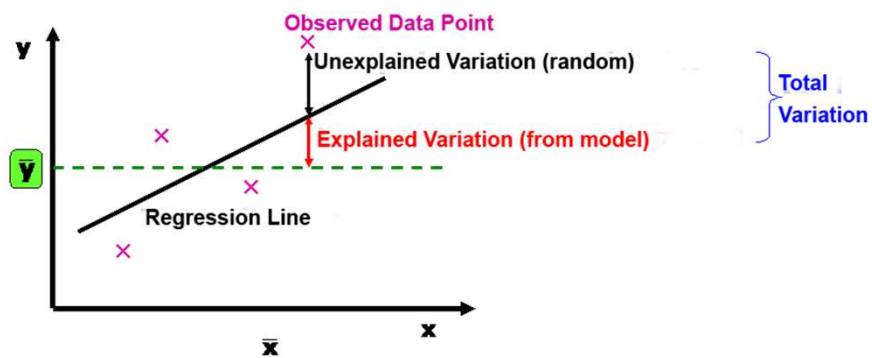
34

Outline

- Introduction (done)
- Simple Linear Regression (done)
- Measures of Variation (next)
 - Coefficient of Determination
 - Correlation
- Misc

35

Measures of Variation

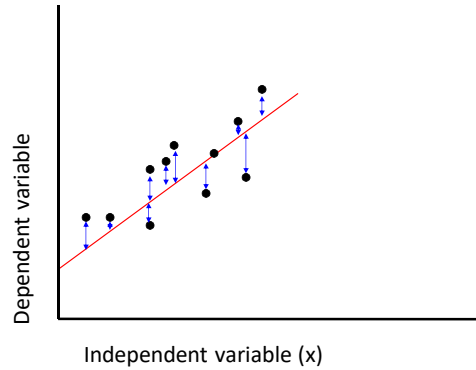


- Several sources of variation in y
 - Error in prediction (unexplained)
 - Variation from model (explained)

Break this
down (next)

36

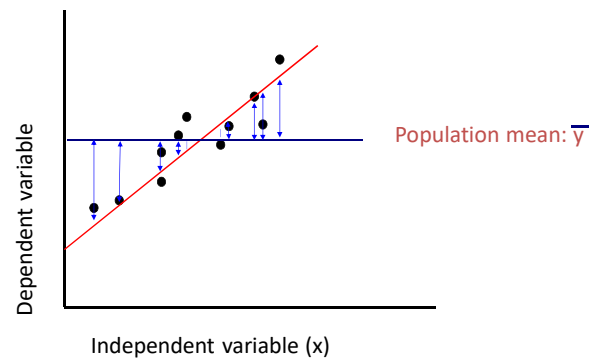
Sum of Squares of Error



- Least squares regression selects line with lowest total sum of squared prediction errors
- Sum of Squares of Error, or **SSE**
- Measure of **unexplained variation**

37

Sum of Squares Regression

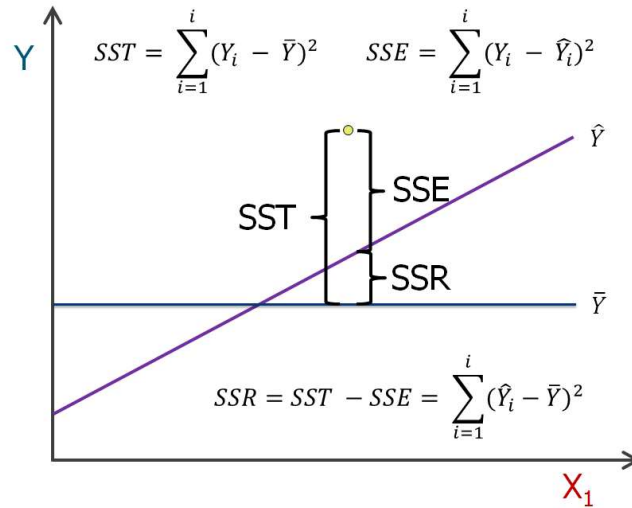


- Differences between prediction and population mean
 - Gets at variation due to X & Y
- Sum of Squares Regression, or **SSR**
- Measure of **explained variation**

38

Sum of Squares Total

- Total Sum of Squares, or SST = SSR + SSE



39

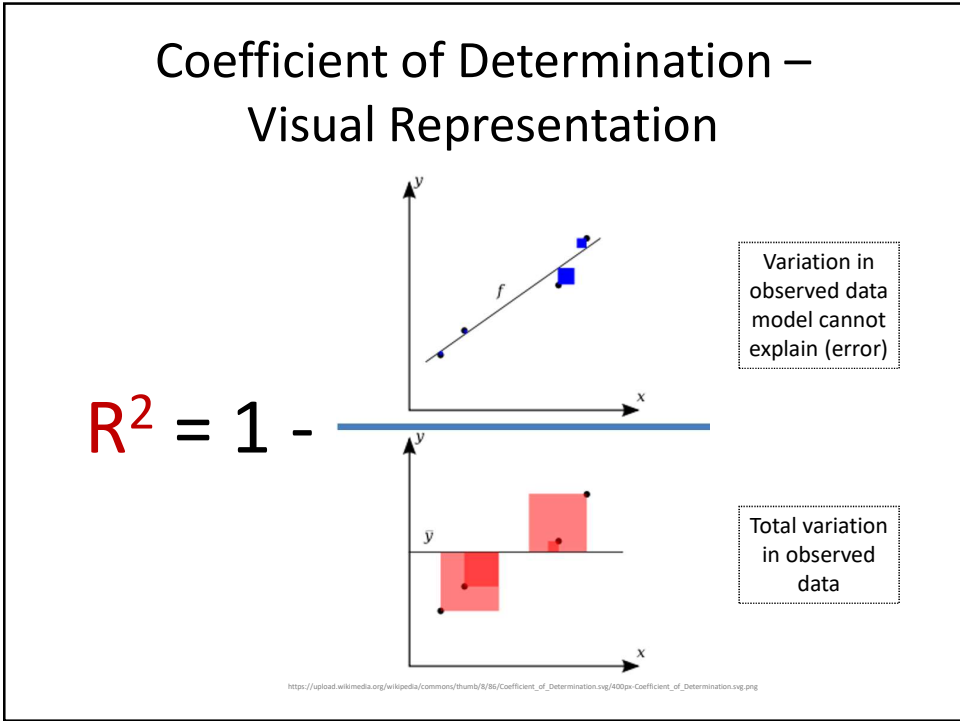
Coefficient of Determination

- Proportion of total variation (SST) explained by the regression (SSR) is known as the **Coefficient of Determination** (R^2)

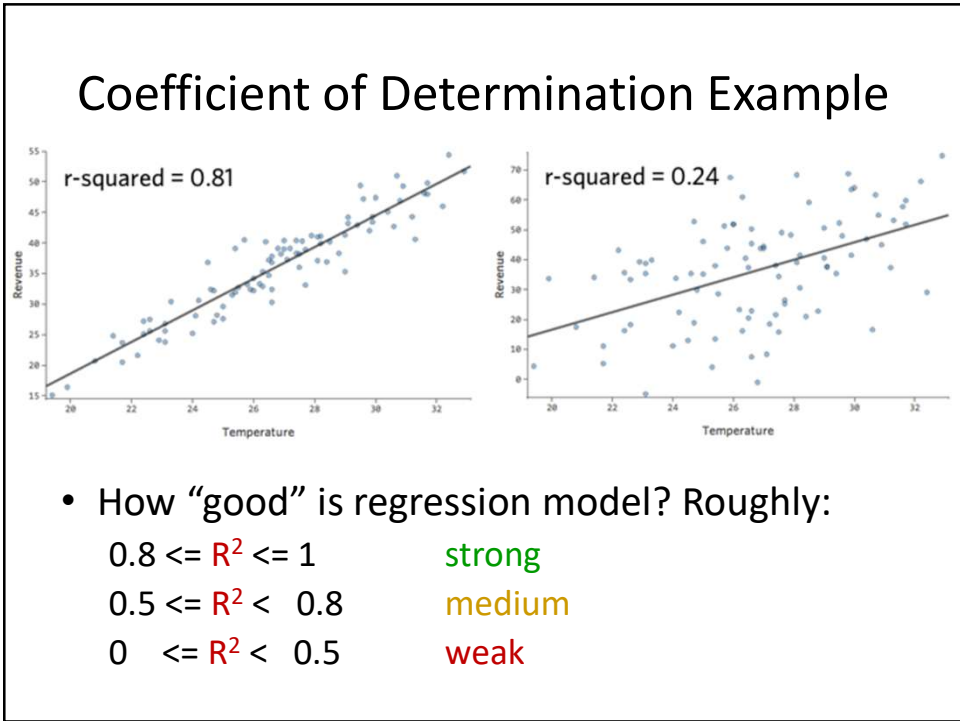
$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

- Ranges from 0 to 1 (often said as a percent)
 - 1 – regression explains all of variation
 - 0 – regression explains none of variation

40

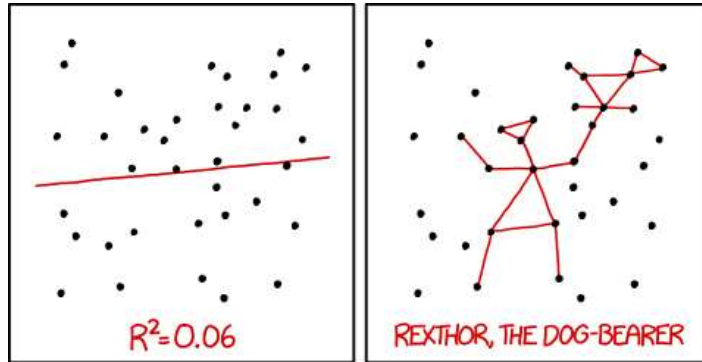


41



42

How "good" is the Regression Model?

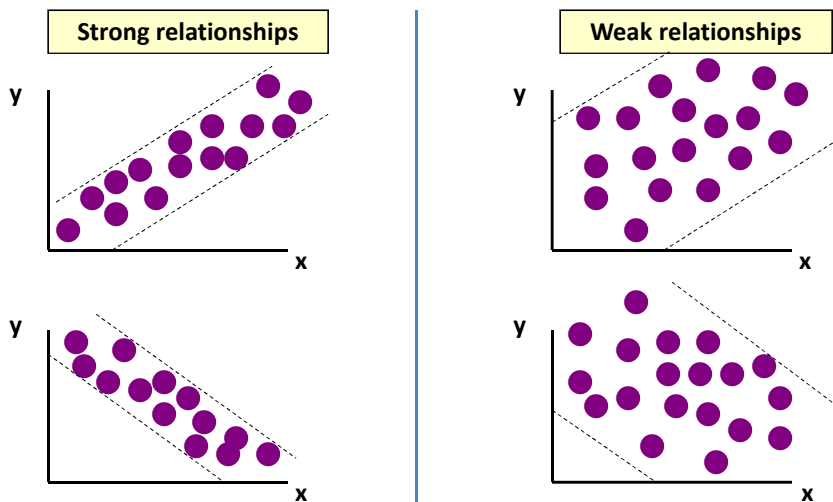


I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER TO GUESS THE DIRECTION OF THE CORRELATION FROM THE SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

<https://xkcd.com/1725/>

43

Relationships Between X & Y



44

Relationship Strength and Direction – Correlation

- **Correlation** measures strength and direction of linear relationship
 - 1 perfect neg. to +1 perfect pos.
 - Sign is same as regression slope
 - Denoted R . Why? $R = \sqrt{R^2}$

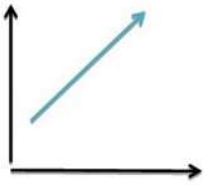
Pearson's Correlation Coefficient

$$r = \frac{\sum(x-\bar{x})(y-\bar{y})}{\sqrt{\sum(x-\bar{x})^2} \sqrt{\sum(y-\bar{y})^2}}$$

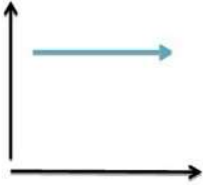
Where, \bar{x} - mean of X variable
 \bar{y} - mean of Y variable

Vary together

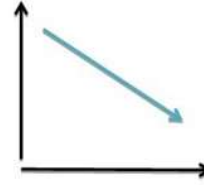
Vary Separately



POSITIVE CORRELATION



ZERO CORRELATION

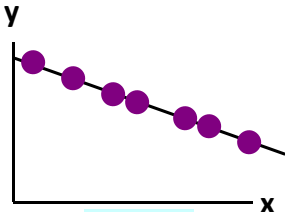


NEGATIVE CORRELATION

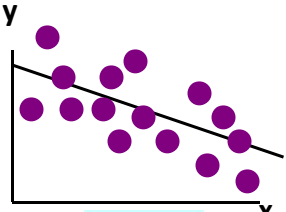
https://www.mbakool.com/2013_images/stories/doc_images/pearson-coeff-0001.jpg

45

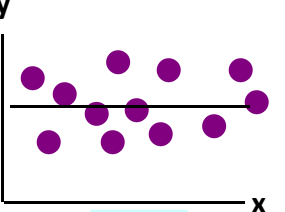
Correlation Examples (1 of 3)



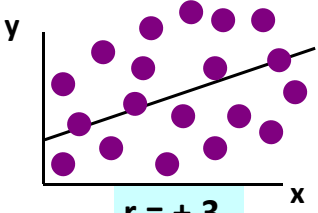
$r = -1$



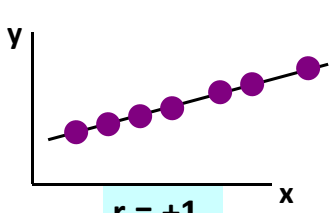
$r = -.6$



$r = 0$

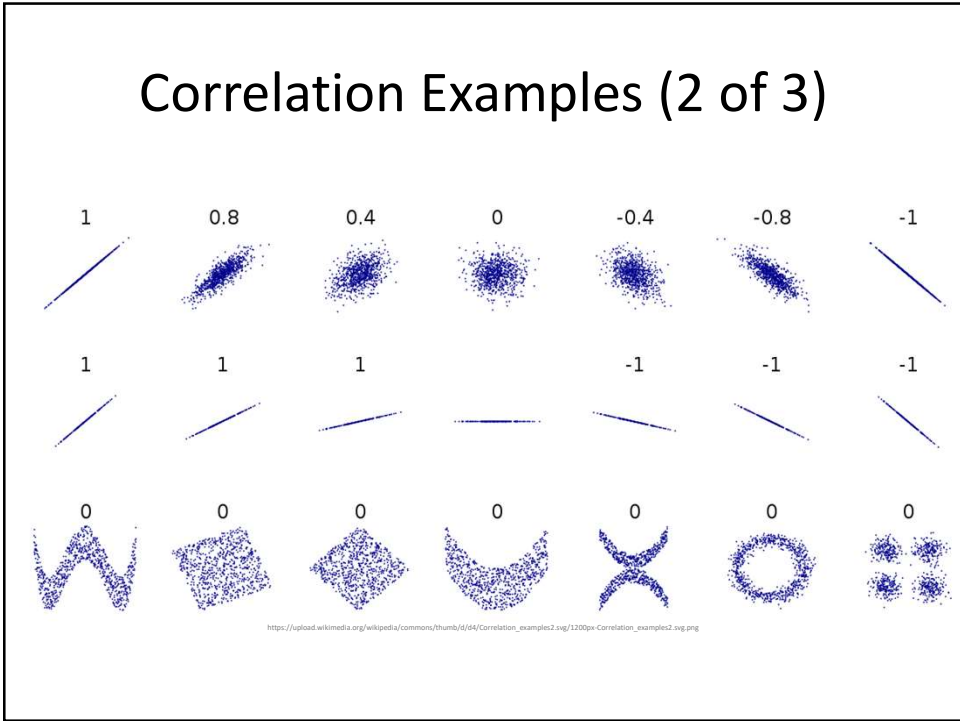


$r = +.3$

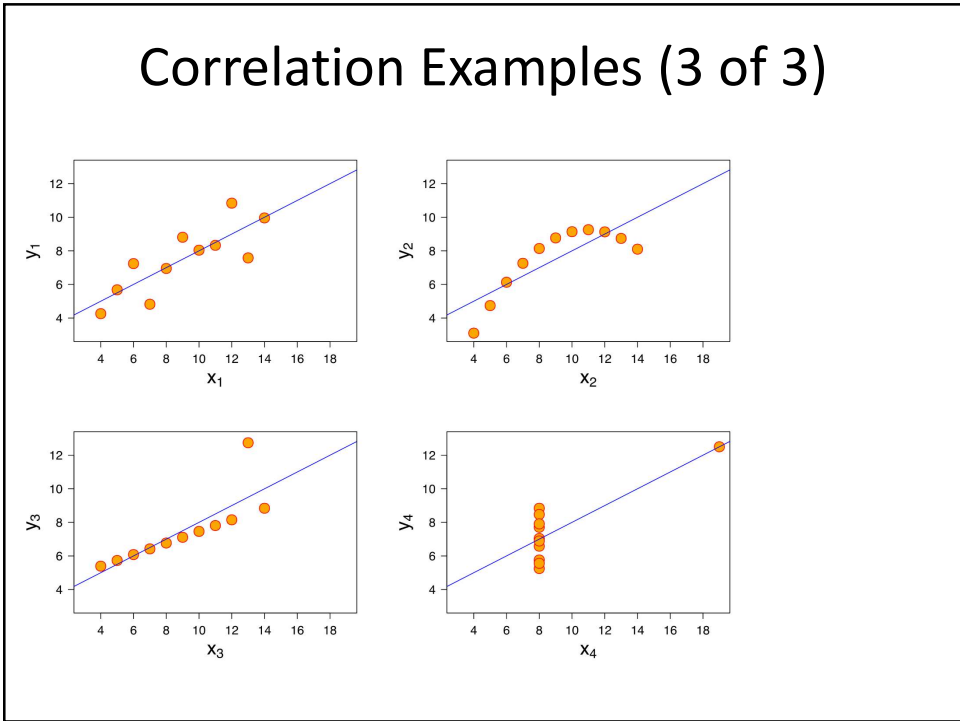


$r = +1$

46

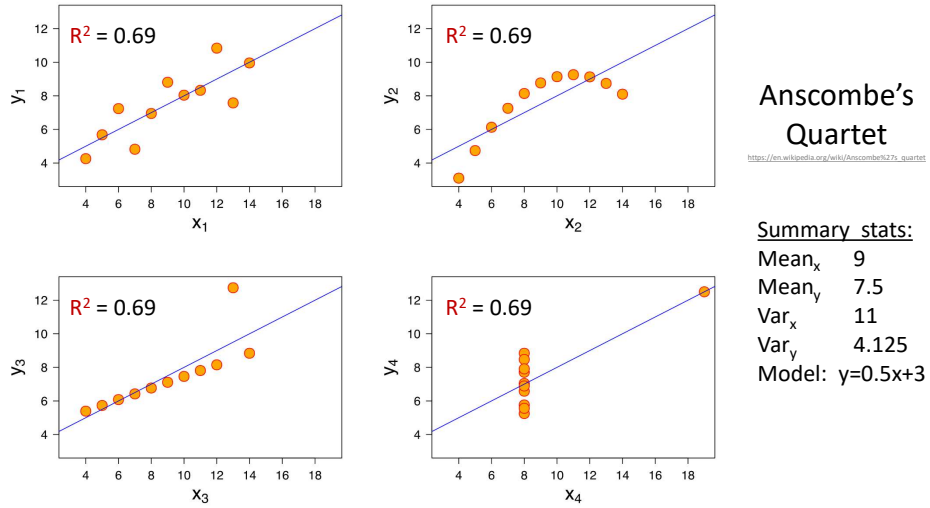


47



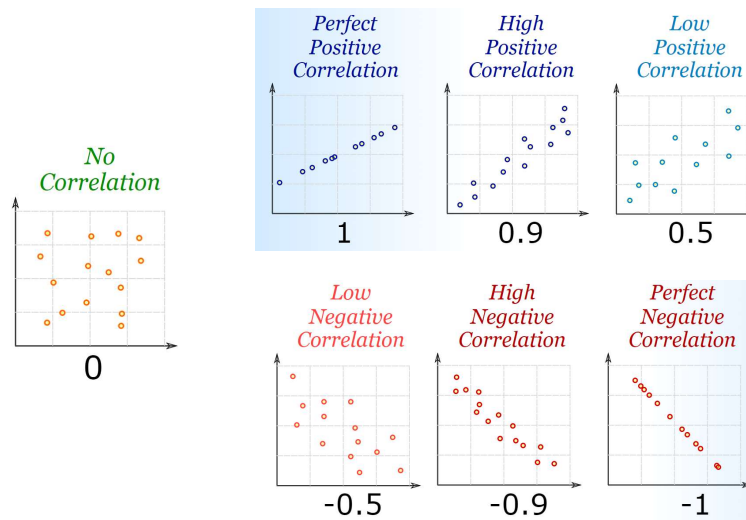
48

Correlation Examples (3 of 3)



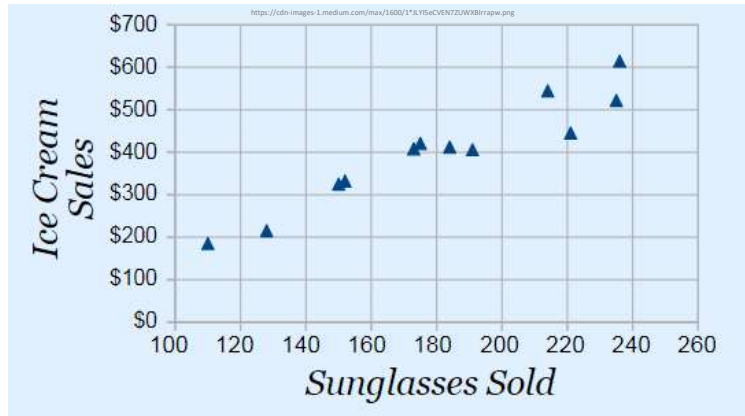
49

Correlation Summary



50

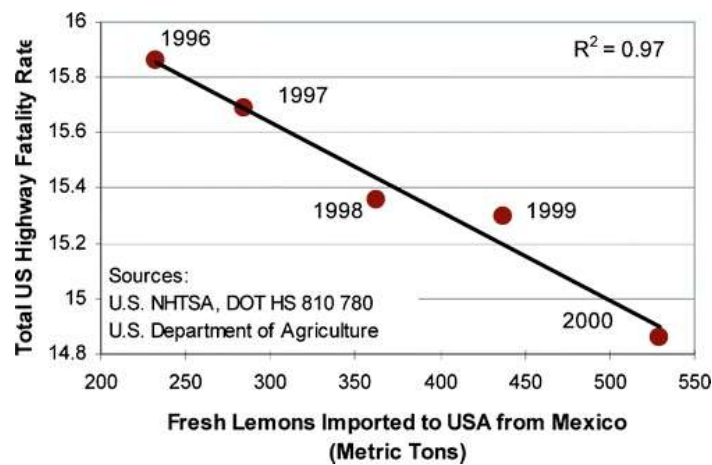
Correlation is not Causation



Buying sunglasses *causes* people to buy ice cream?

51

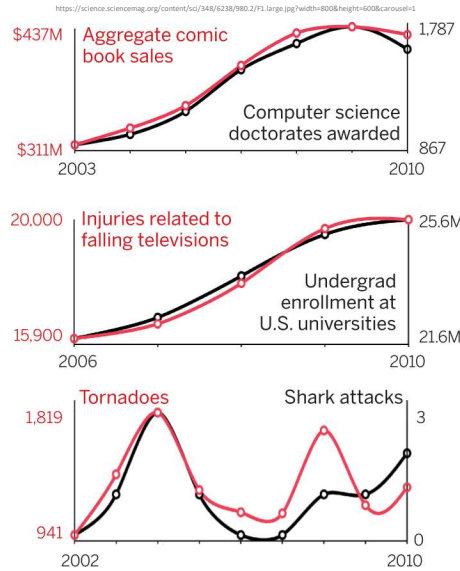
Correlation is not Causation



Importing lemons *causes* fewer highway fatalities?

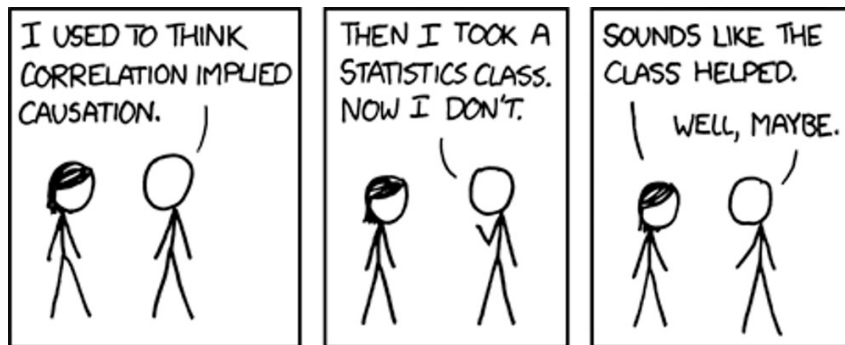
52

Correlation is not Causation



53

Correlation is not Causation



<https://xkcd.com/552/>

54

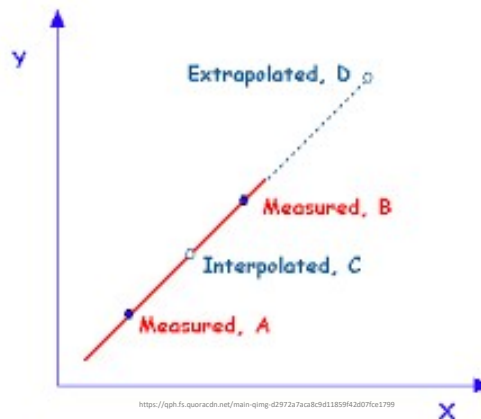
Outline

- Introduction (done)
- Simple Linear Regression (done)
- Measures of Variation (done)
- Misc (next)

55

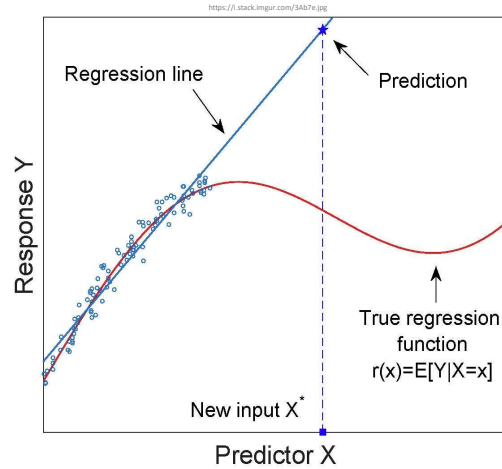
Extrapolation versus Interpolation

- Prediction
 - Interpolation – within measured X-range
 - Extrapolation – outside measured X-range



56

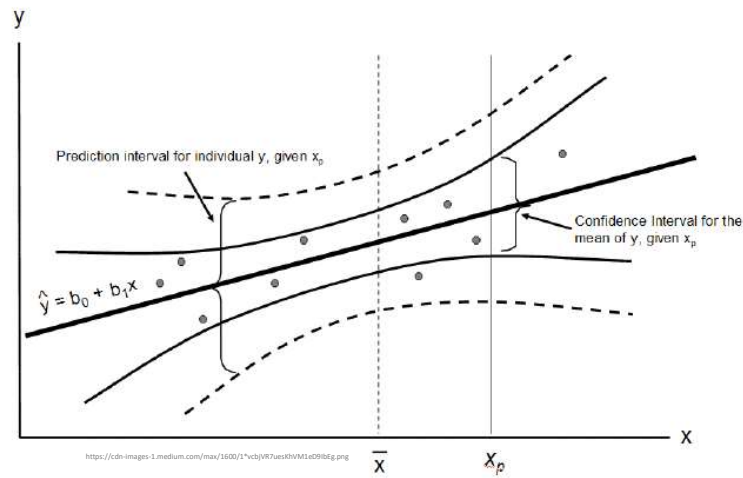
Be Careful When Extrapolating



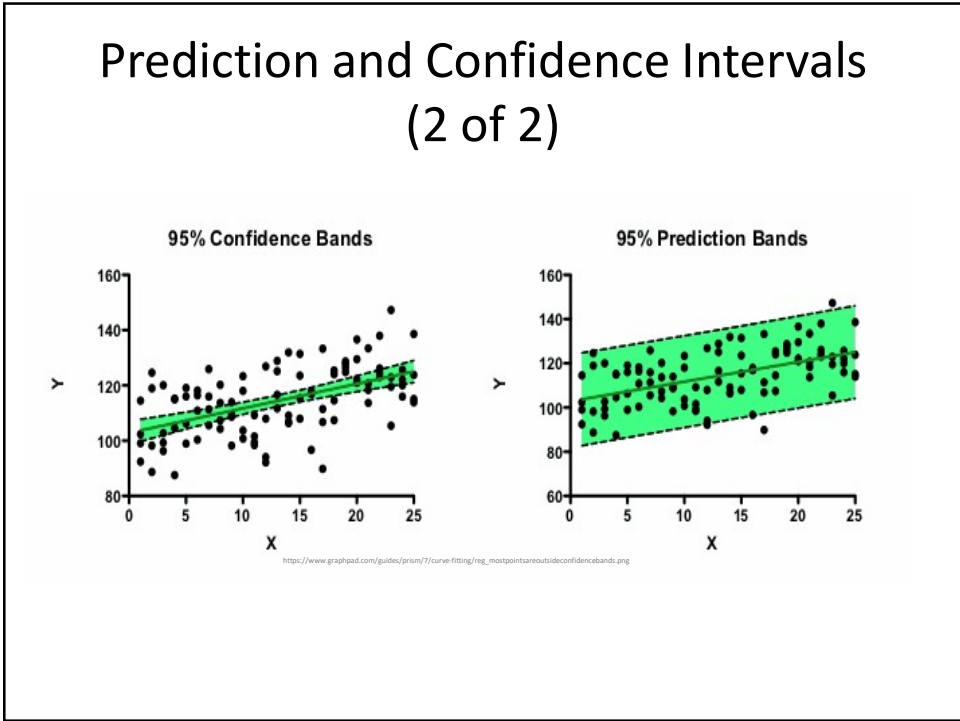
If **extrapolate**, make sure have reason to assume model continues

57

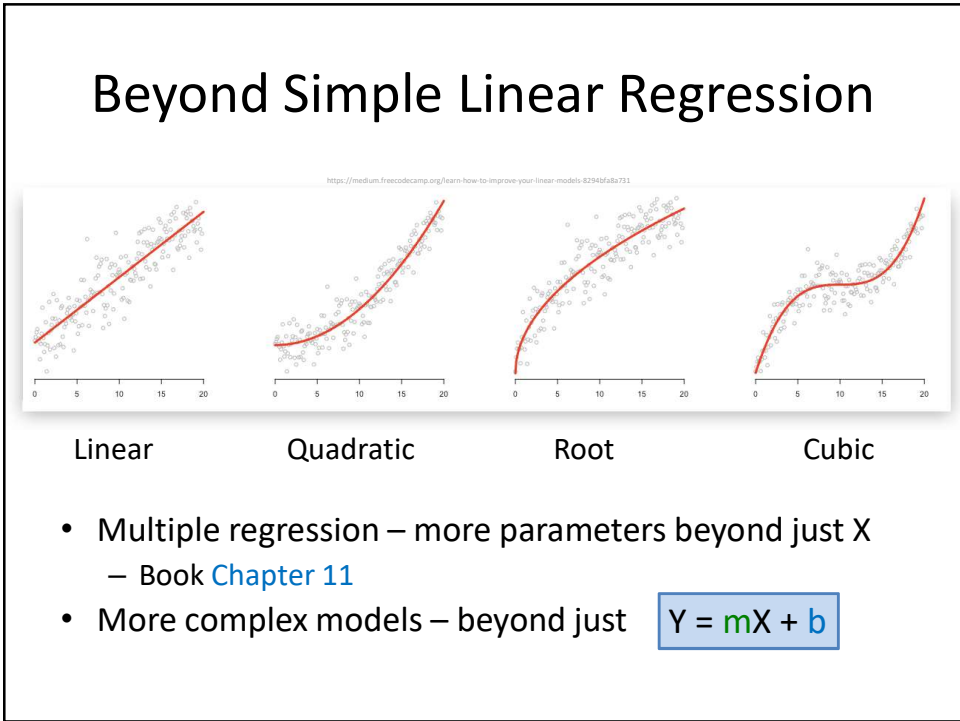
Prediction and Confidence Intervals (1 of 2)



58

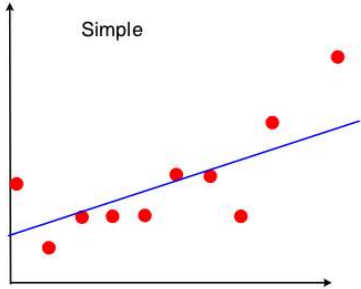


59



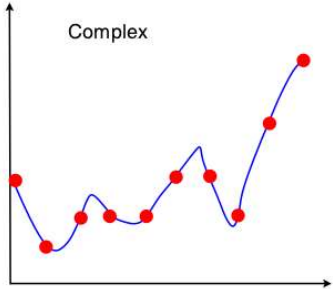
60

More Complex Models



Simple

$$y = 12x + 9$$



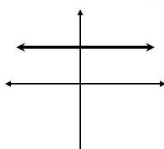
Complex

$$y = 18x^4 + 13x^3 - 9x^2 + 3x + 20$$

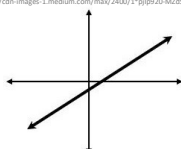
- Higher order polynomial model has less error
- A “perfect” fit (no error)
- How does a polynomial do this?

61

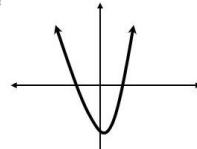
Graphs of Polynomial Functions



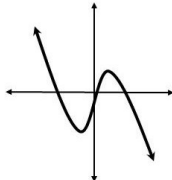
Constant Function
(degree = 0)



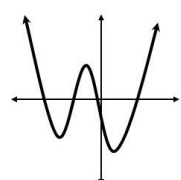
Linear Function
(degree = 1)



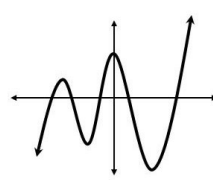
Quadratic Function
(degree = 2)



Cubic Function
(deg. = 3)



Quartic Function
(deg. = 4)

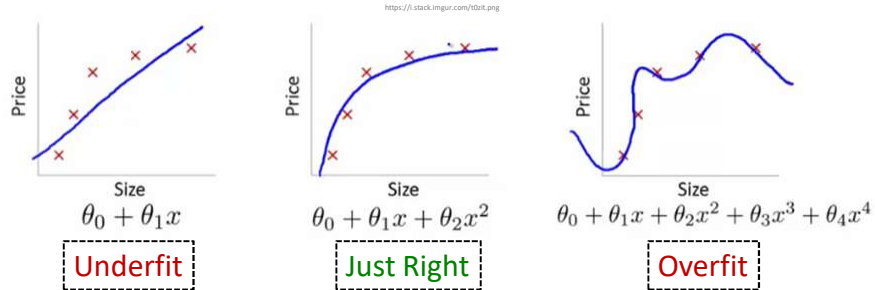


Quintic Function
(deg. = 5)

Higher degree, more potential “wiggles”
But **should** you use?

62

Underfit and Overfit



- **Overfit** analysis matches data too closely with more parameters than can be justified
 - **Underfit** analysis does not adequately match data since parameters are missing
- Both model do not predict well (i.e., for non-observed values)
- **Just right** – fit data well “enough” with as few parameters as possible