# StatQuest!!!
# Linear Regression



Awesomeness

Time spent watching StatQuest

# Study Guide!!!

## The Problem

| | Weight | Size |
|---|---|---|
| 1 | 6 | 8 |
| 2 | 2.5 | 2 |
| 3 | 13 | 14 |
| 4 | 10 | 7 |

We measured the **Weight** and **Size** of 4 mice…

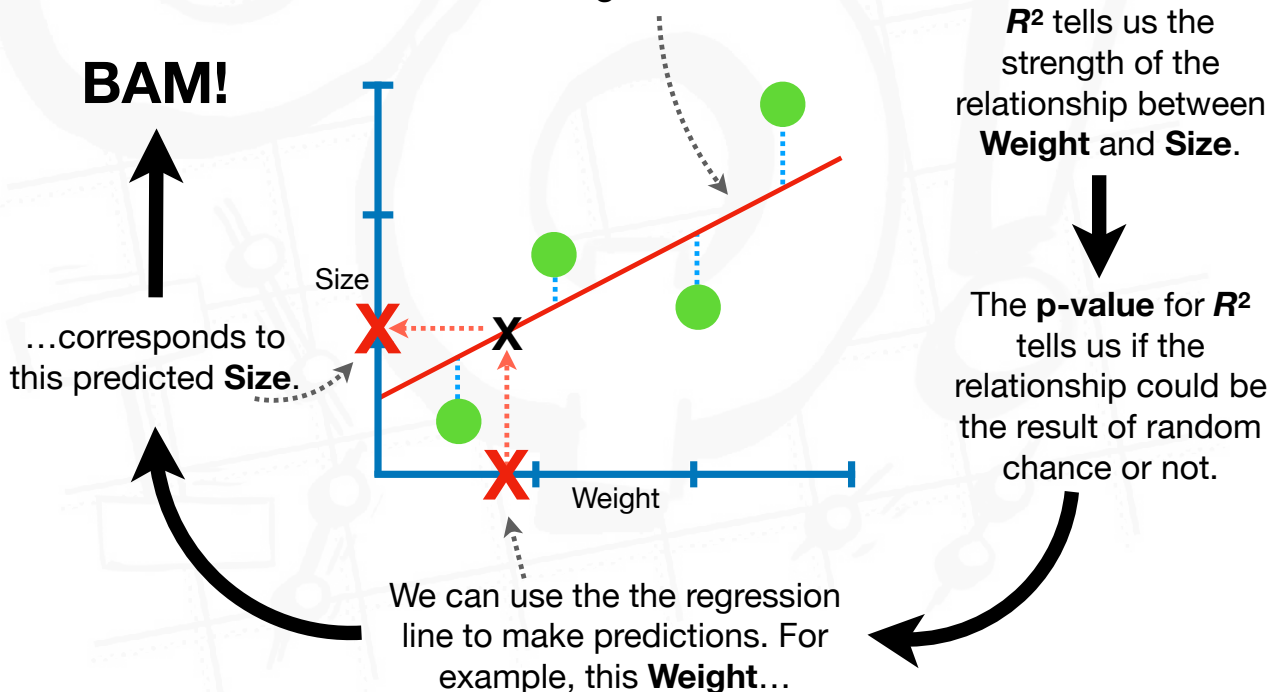…and we can plot the data…

…and it looks like **Weight** and **Size** are related. How can we test this hypothesis? Is the relationship useful?
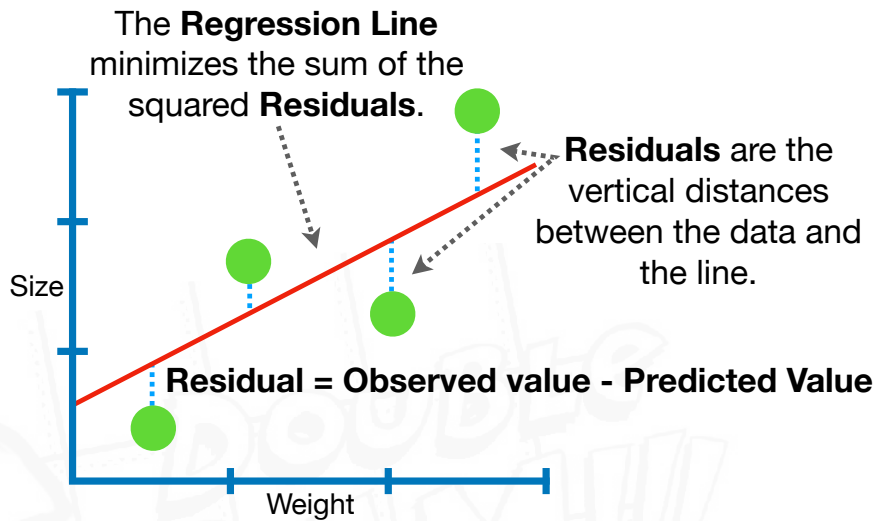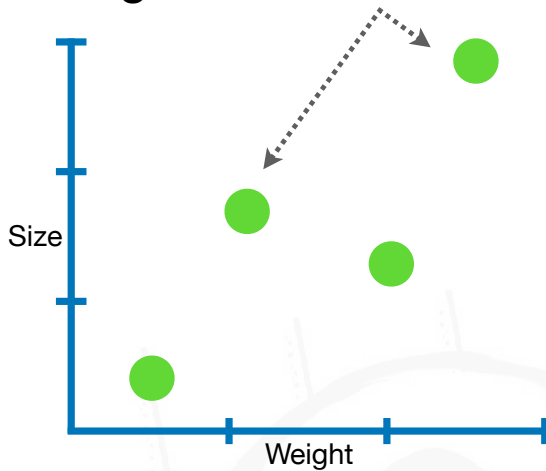
Size

Weight

## The Solution - Linear Regression

The slope of this regression line tells us that as **Weight** increases, so does **Size**. This suggests there is a relationship between **Weight** and **Size**.

$R^2$ tells us the strength of the relationship between **Weight** and **Size**.

The **p-value** for $R^2$ tells us if the relationship could be the result of random chance or not.

**BAM!**

…corresponds to this predicted **Size**.

Size

Weight

We can use the the regression line to make predictions. For example, this **Weight**…

## NOTES:

# Fitting a line to the data

The **Regression Line** minimizes the sum of the squared **Residuals**.

**Residuals** are the vertical distances between the data and the line.

Size

Weight

Size

Weight

**Residual = Observed value - Predicted Value**
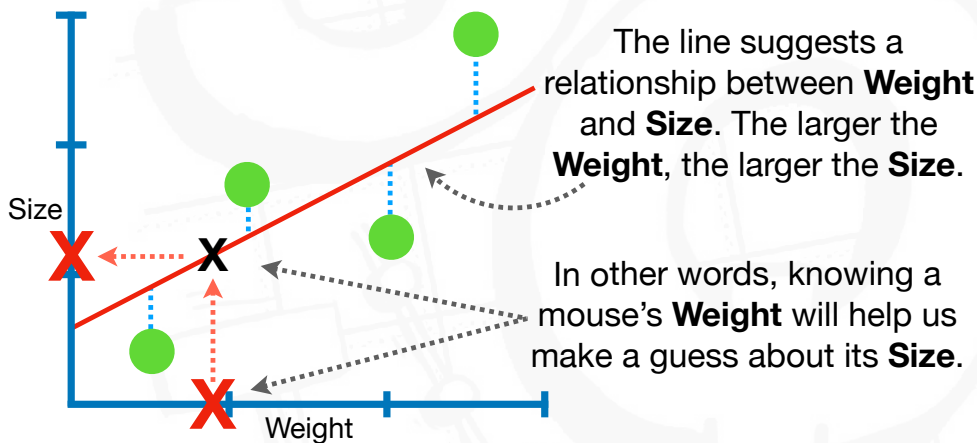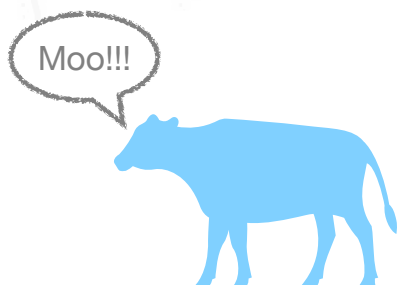
**Residuals** are squared for two reasons: 1) So that negative values do not cancel out positive values when we add them together, 2) Unlike the absolute value, the derivative of a squared number exists for all values, making the math easier.

# Interpreting the Regression Line

Size

Weight

The line suggests a relationship between **Weight** and **Size**. The larger the **Weight**, the larger the **Size**.

In other words, knowing a mouse's **Weight** will help us make a guess about its **Size**.

We quantify how good that guess will be with $R^2$, which tells us the accuracy of the guess, and its **p-value**, which tells us if we should believe the guess in the first place.
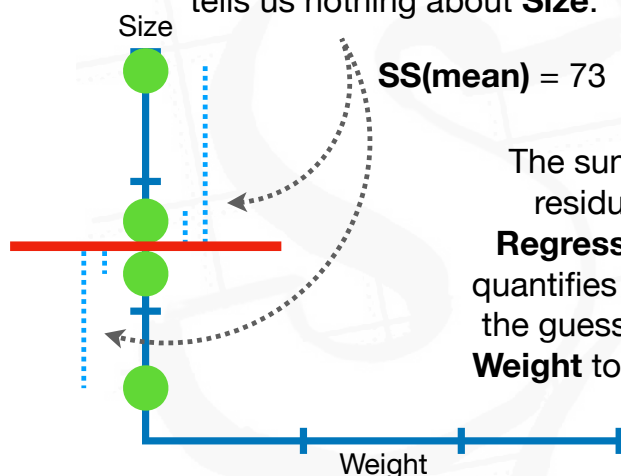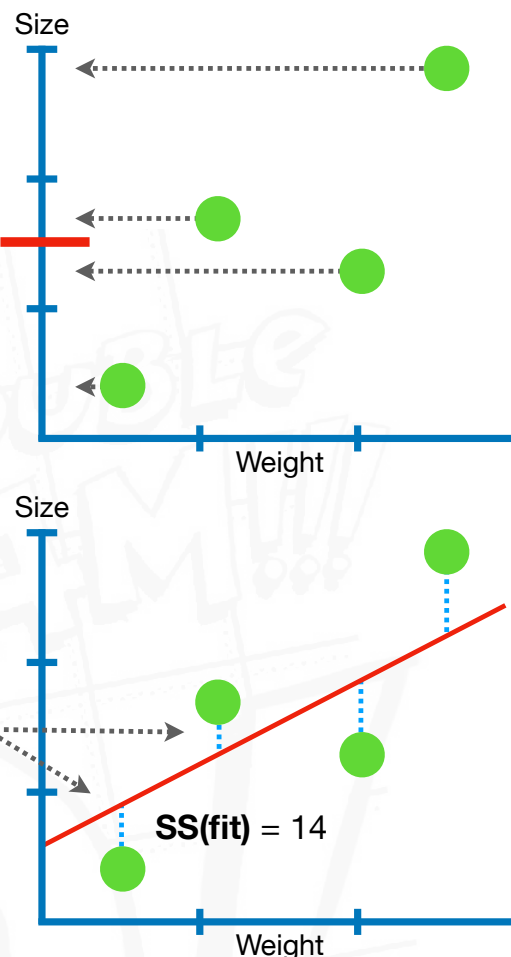
# NOTES:

What do you think of Linear Regression?

Moo!!!

# Calculating $R^2$

If there was no relationship between **Weight** and **Size**, and knowing **Weight** would not help you predict **Size**…

…then the best prediction of **Size** would be the average **Size**.

Size

Weight

The sum of the squared residuals around the mean for **Size**, **SS(mean)**, quantifies how bad (or good) the guess is when we assume that **Weight** tells us nothing about **Size**.

**SS(mean)** = 73

Size

Weight

The sum of the squared residuals around the **Regression Line**, **SS(fit)**, quantifies how bad (or good) the guess is when we allow **Weight** to tell us about **Size**.

Size

**SS(fit)** = 14

Weight

---

$R^2$ compares the guesses using only the mean, **SS(mean)**, to the guesses made with the regression line, **SS(fit)**

By dividing **SS(mean)** and **SS(fit)** by the number of observations, $n$, we get the variances.

$$R^2 = \frac{SS(mean) - SS(fit)}{SS(mean)} = \frac{\frac{SS(mean)}{n} - \frac{SS(fit)}{n}}{\frac{SS(mean)}{n}} = \frac{Var(mean) - Var(fit)}{Var(mean)}$$

Plugging in **SS(mean) = 73** and **SS(fit) = 14** from the data…

$$R^2 = \frac{73 - 14}{73} = \frac{\frac{73}{4} - \frac{14}{4}}{\frac{73}{4}} = \frac{18 - 3.5}{18} = 0.81$$

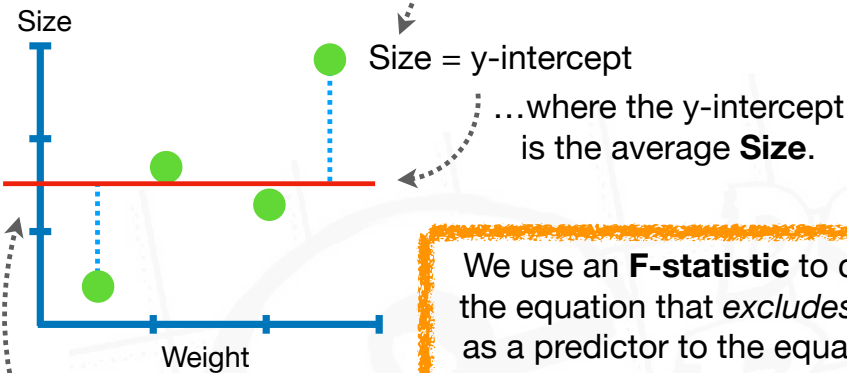When $R^2 = 0$, the average **Size** (or average y-axis variable, whatever that happens to be) is the best guess we can make. When $R^2 = 1$, the **Regression Line** makes perfect guesses.

This tells us that there is **81%** less variation around the **Regression Line** than around the mean value for **Size**.

Alternatively, we can say that **Weight** "explains" **81%** of the variation in **Size**.

# Calculating the p-value for $R^2$

When we ignore **Weight** when predicting **Size**, the equation for making predictions is simply…

Size

Size = y-intercept

…where the y-intercept is the average **Size**.

Weight

$p_{mean}$ = 1 because the y-intercept is the only parameter for this line.

**The Numerator** is the average variance explained per extra parameter that the **Regression Line** uses.

$$\frac{SS(mean) - SS(fit)}{p_{fit} - p_{mean}}$$

This is the extra number of parameters used for the **Regression Line**.

When we use **Weight** predict **Size**, the equation for making predictions is…

Size = y-intercept + (slope × Weight)

Size

Weight

$p_{fit}$ = 2 because the y-intercept and slope are parameters for this line.

**The Denominator** is the variation in size not explained by the **Regression Line**.

$$SS(fit) / (n - p_{fit})$$

Dividing **SS(fit)** by $n$-$p_{fit}$ compensates for the extra variables in the **Regression Line**.

We use an **F-statistic** to compare the equation that *excludes* **Weight** as a predictor to the equation that *includes* **Weight** as a predictor.
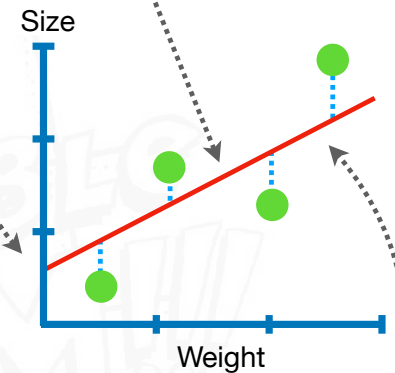
$$F = \frac{\frac{SS(mean) - SS(fit)}{p_{fit} - p_{mean}}}{SS(fit) / (n - p_{fit})}$$

**NOTE: SS(mean)** and **SS(fit)** are defined in **Calculating $R^2$**.

The top is proportional to the amount of variance explained by the **Regression Line**. (see **Calculating $R^2$** for details)

Plugging in: $n$ = 4

**SS(mean) = 73**     $p_{mean}$ = 1

**SS(fit) = 14**     $p_{fit}$ = 2

Degrees of Freedom:
**DF1** = $p_{mean}$ - $p_{fit}$ = 1
**DF2** = $n$ - $p_{fit}$ = 2

$$F = \frac{\frac{73 - 14}{2 - 1}}{14 / (4 - 2)} = 8.32$$

The **F-statistic**, **8.32**, tells us that the **p-value**, **0.1**, is the area under the curve from **8.32** to infinity.

F-distribution with DF1=1 and DF2=2

*p*-value = 0.1 = **red area**

8.32

# Using more than one variable to make predictions:

When we use more than one variable to make predictions, we compare a **Simple Model** to a **Fancy Model**.
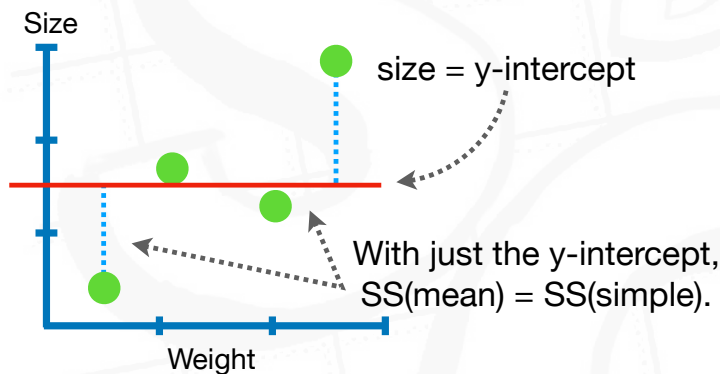
$$R^2 = \frac{SS(mean) - SS(fancy)}{SS(mean)}$$

$$F = \frac{\dfrac{SS(simple) - SS(fancy)}{p_{fancy} - p_{simple}}}{SS(fancy) / (n - p_{fit})}$$

---
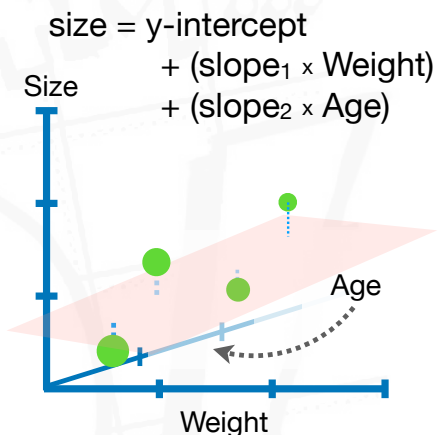
### Example 1

In this example, the **Simple Model** has **1** parameter, the y-intercept, so $p_{simple} = 1$…

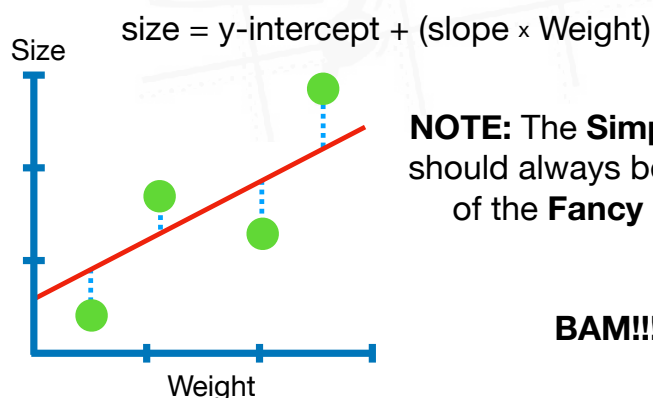…while this **Fancy Model** has **3** parameters, the y-intercept and **2** slopes, so $p_{fancy} = 3$.

Size

size = y-intercept

Weight

With just the y-intercept, SS(mean) = SS(simple).

Size

size = y-intercept
+ (slope₁ × Weight)
+ (slope₂ × Age)

size = y-intercept + (slope$_1$ × Weight) + (slope$_2$ × Age)

Age

Weight

**NOTE:** $R^2$ is typically only calculated between a fancy model and the simplest model that only contains the y-intercept.

---

### Example 2

In this example, the **Simple Model** has **2** parameters, the y-intercept and a slope, so $p_{simple} = 2$…

…while this **Fancy Model** has **3** parameters, the y-intercept and **2** slopes, so $p_{fancy} = 3$.

Size

size = y-intercept + (slope × Weight)

**NOTE:** The **Simple Model** should always be a subset of the **Fancy Model**.

**BAM!!!**

Size

size = y-intercept + (slope$_1$ × Weight) + (slope$_2$ × Age)

Age

Weight