# Regularization

BCSE0105 MACHINE LEARNING

# What is Regularization?
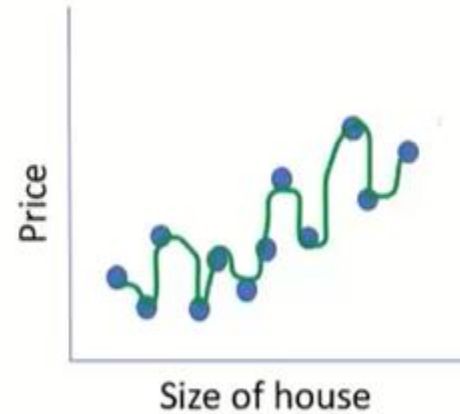
House Price Prediction

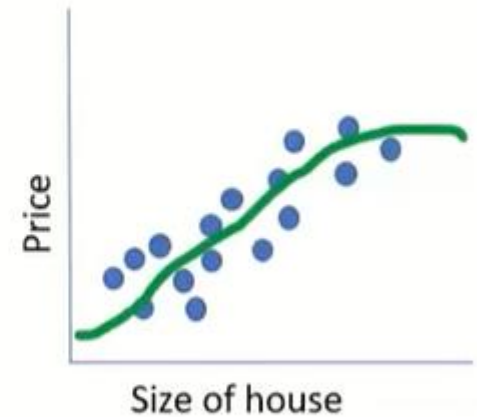| Underfitting | Overfitting | Good Fit |
|---|---|---|

$$Price = \beta_0 + \beta_1 * size$$

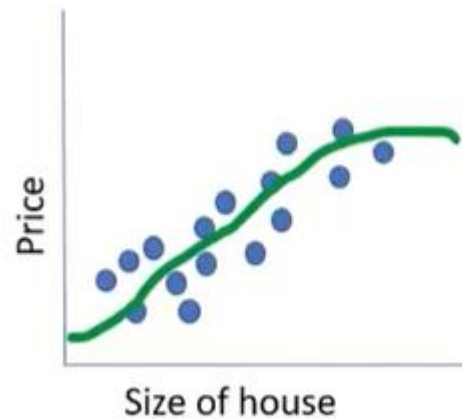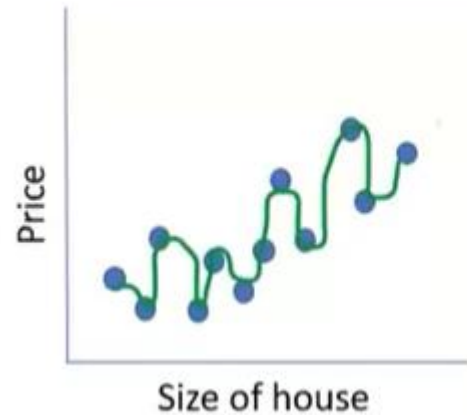$$Price = \beta_0 + \beta_1 * size + \beta_2 * size^2 + \beta_3 * size\,3 + \beta_4 * size\,4$$

$$Price = \beta_0 + \beta_1 * size + \beta_2 * size^2$$

# What is Regularization?

- Constraining a model to make it simpler and reduce the risk of overfitting is called **regularization.**

- A training constraint whose objective is to reduce overfitting and thus improve the model's ability to generalize

- The amount of regularization to apply during learning can be controlled by a hyperparameter.

# Working of Regularization



Price = $\beta_0 + \beta_1 * size + \beta_2 * size^2 + \beta_3 * size\ 3 + \beta_4 * size\ 4$

reduce $\beta_3$ and $\beta_4$ close to zero

Price = $\beta_0 + \beta_1 * size + \beta_2 * size^2$

# How to reduce /shrink coefficients

- Use cost function

$$MSE = \frac{1}{n} \sum_{i=1}^{n} \left( y_{i\,(actual)} - y_{i\,(predicted)} \right)^2$$

This is higher order polynomial equation and our aim is to reduce MSE. For this we have different **regularization techniques.**

# Regularization Techniques

- Used to reduce overfitting
- Used to generalize the model well

# Ridge Regularization

$$MSE = \frac{1}{n} \sum_{i=1}^{n} \left(y_{i\,(actual)} - y_{i\,(predicted)}\right)^2$$

$$Loss = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^{P} \beta_j{}^2$$

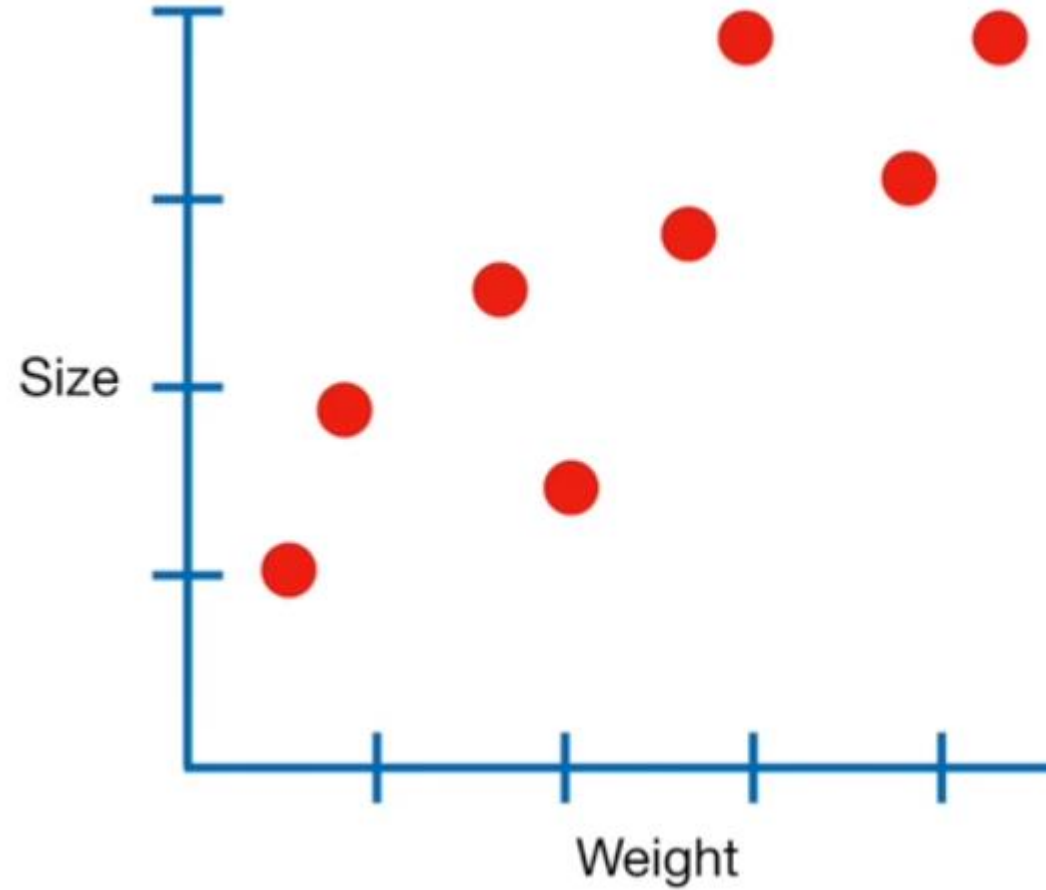Penalty term regularizes the coefficients

$\lambda$ = Tuning parameter

# Lasso Regularization

$$MSE = \frac{1}{n} \sum_{i=1}^{n} \left( y_{i\ (actual)} - y_{i\ (predicted)} \right)^2$$

$$Loss = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^{P} |\beta_j|$$

Penalty term regularizes the coefficients

$\lambda$ = Tuning parameter

# Which Technique To Use?

## Ridge

Lot of features In the dataset and all features have small coefficients.

## Lasso

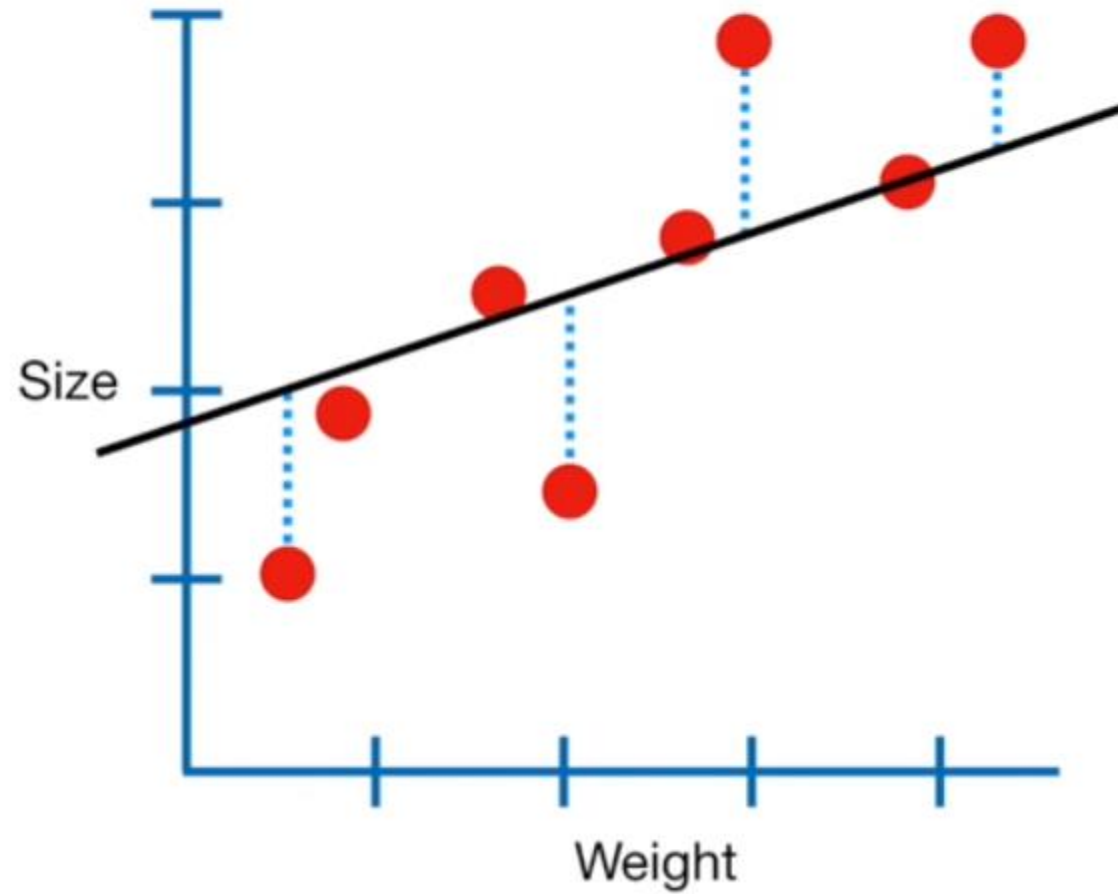Small number of features and few features have high coefficient value.

# Understanding the Idea behind regularization

So we'll fit a line to the data using **Least Squares**.

Size

Weight

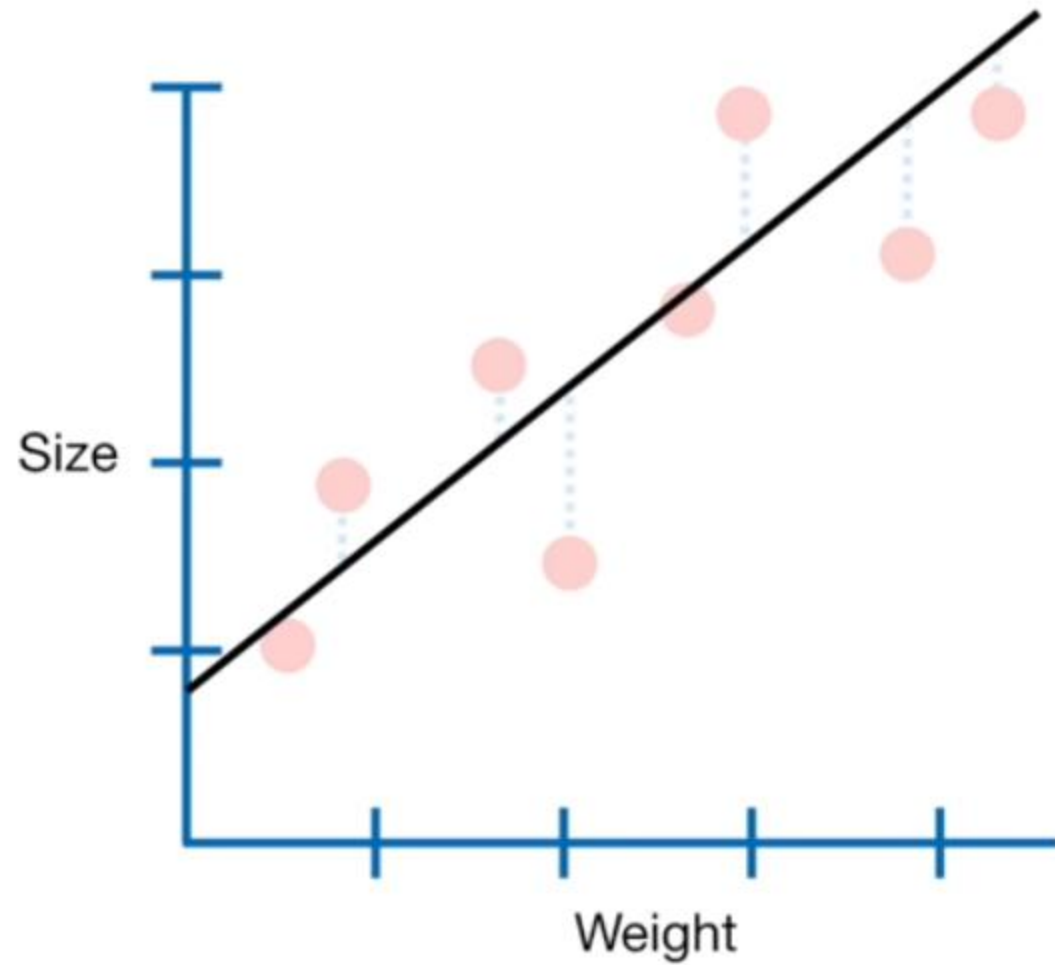In other words, we find the line that results in the **minimum sum of squared residuals**.

Ultimately, we end up with this equation for the line:
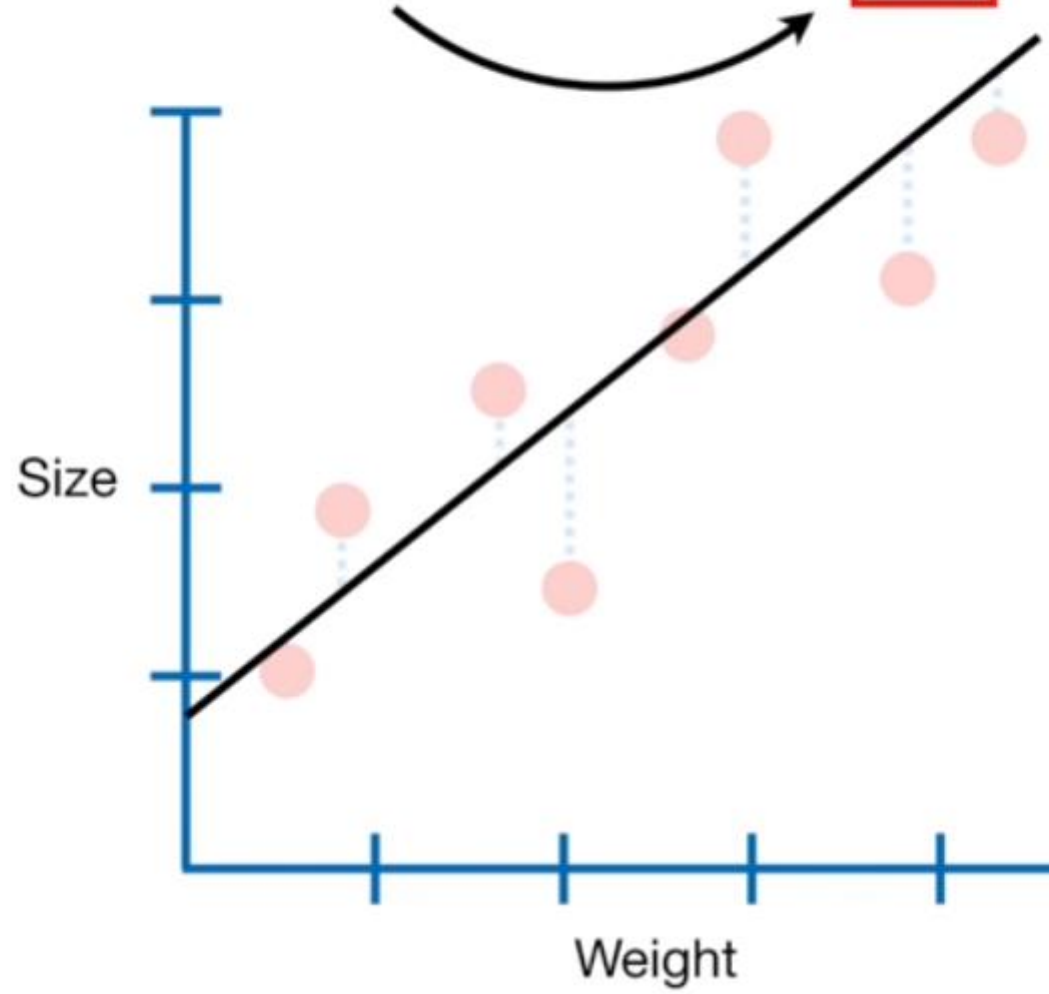
**Size** = 0.9 + 0.75 × **Weight**

...a **y-axis intercept**...

$$\mathbf{Size} = \boxed{0.9} + 0.75 \times \mathbf{Weight}$$

...and a **slope.**     **Size** = 0.9 + 0.75 × **Weight**

We can plug in a value for **Weight**, for example, **2.5**...
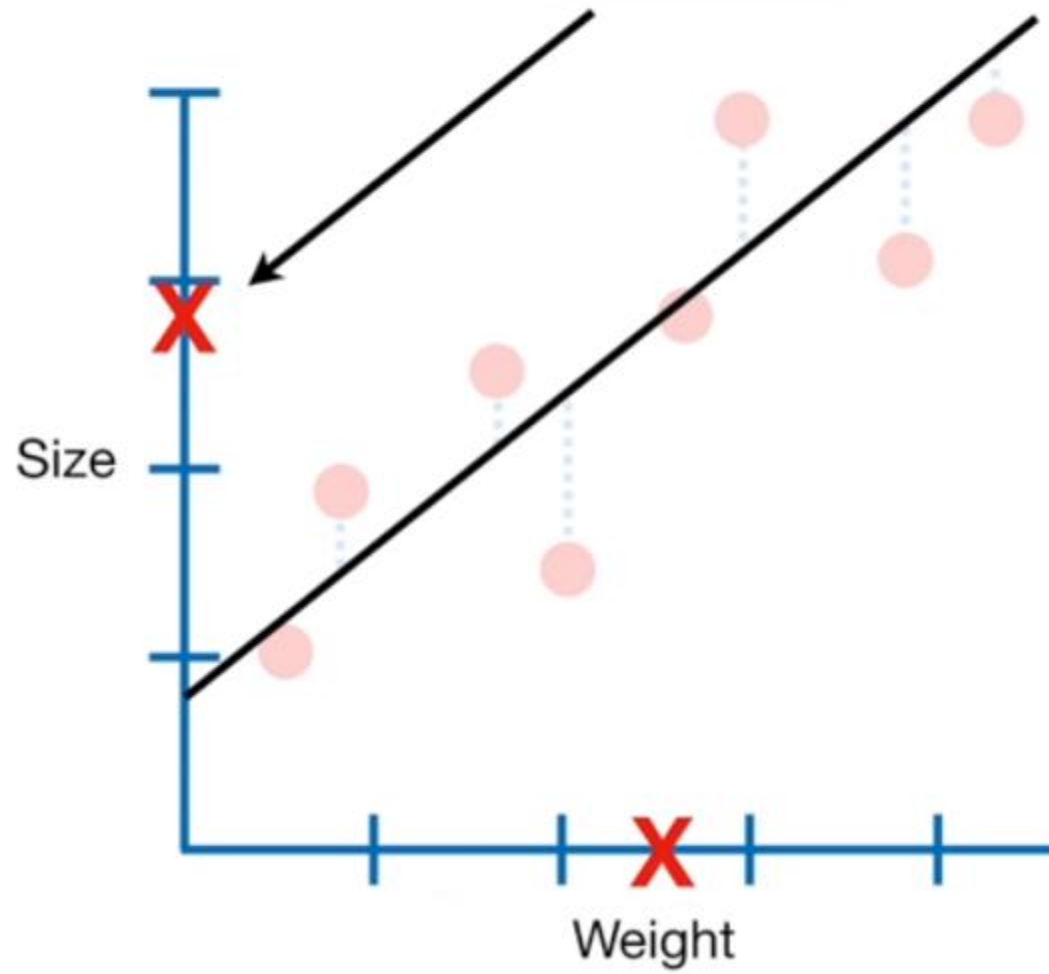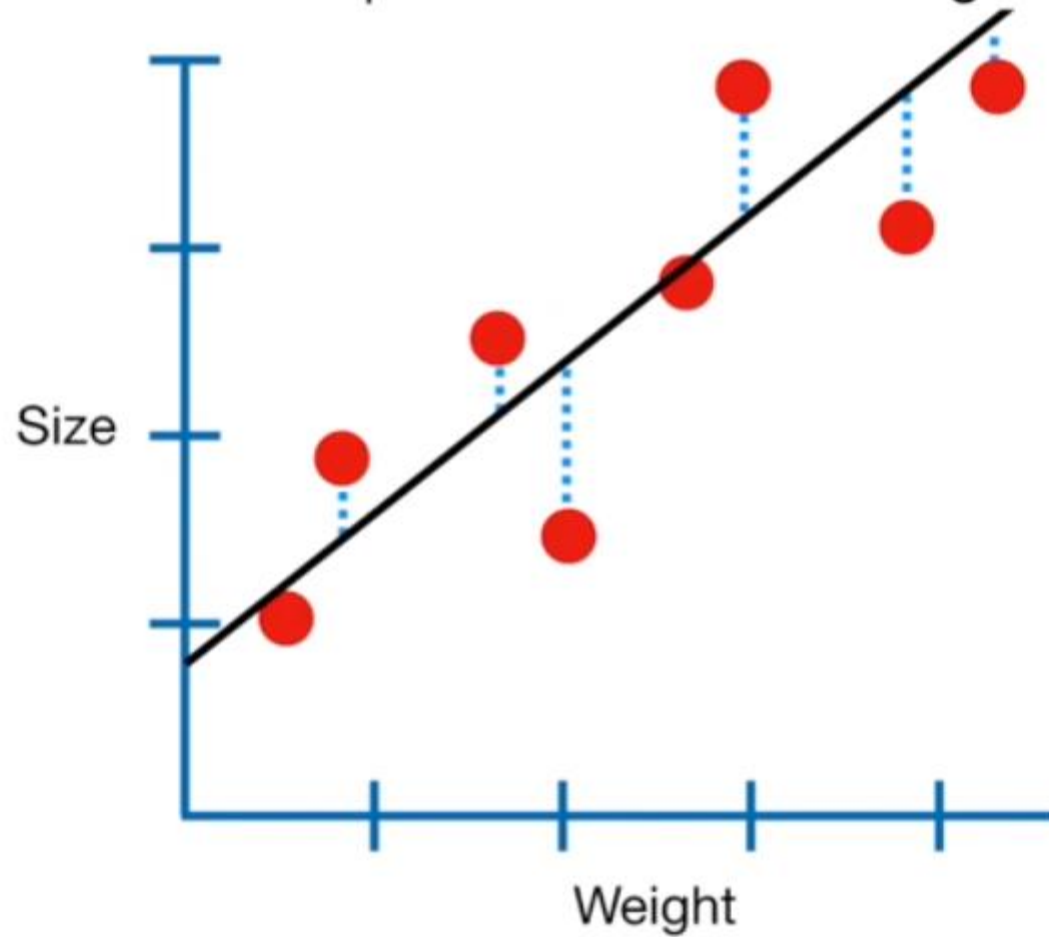
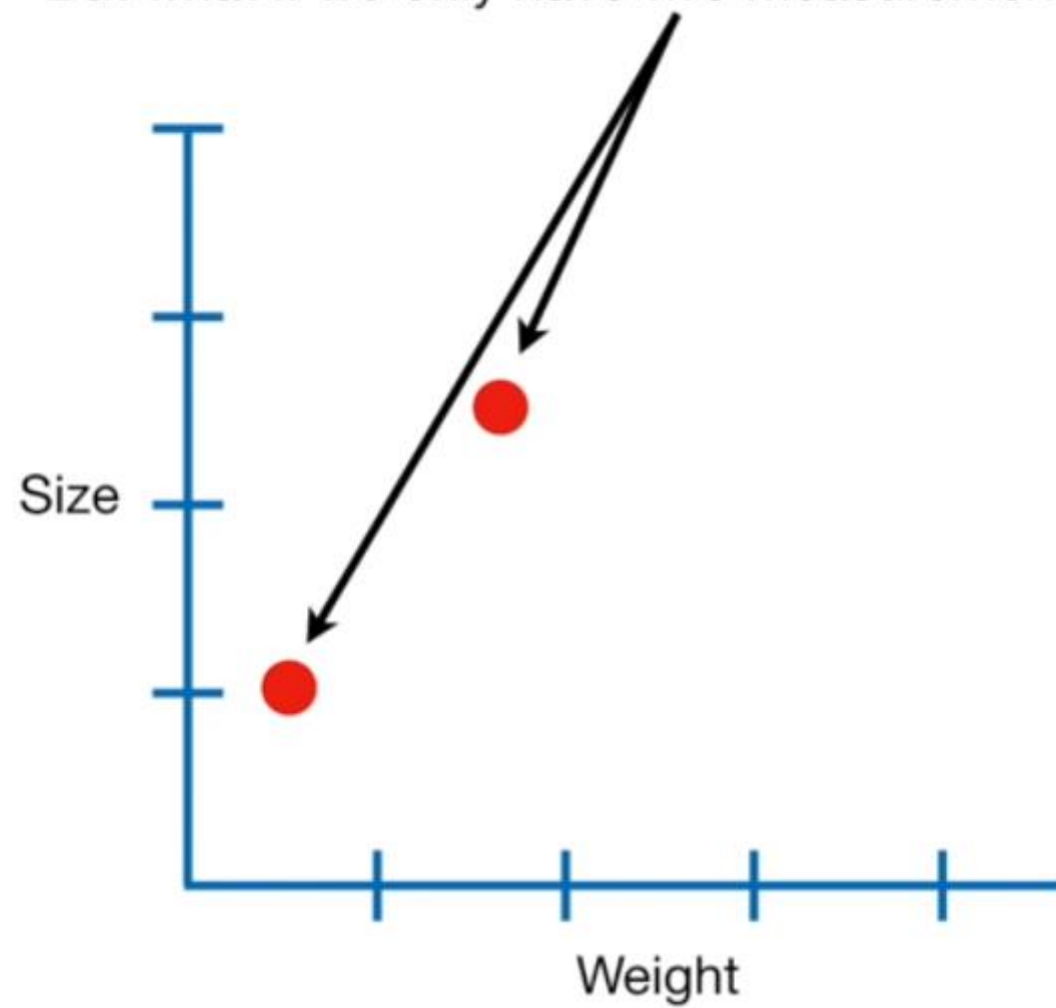**Size** $= 0.9 + 0.75 \times$ **Weight**

Size

Weight

Together, the value for **Weight**, **2.5**, and the value for **Size**, **2.8**, give us a point on the line.
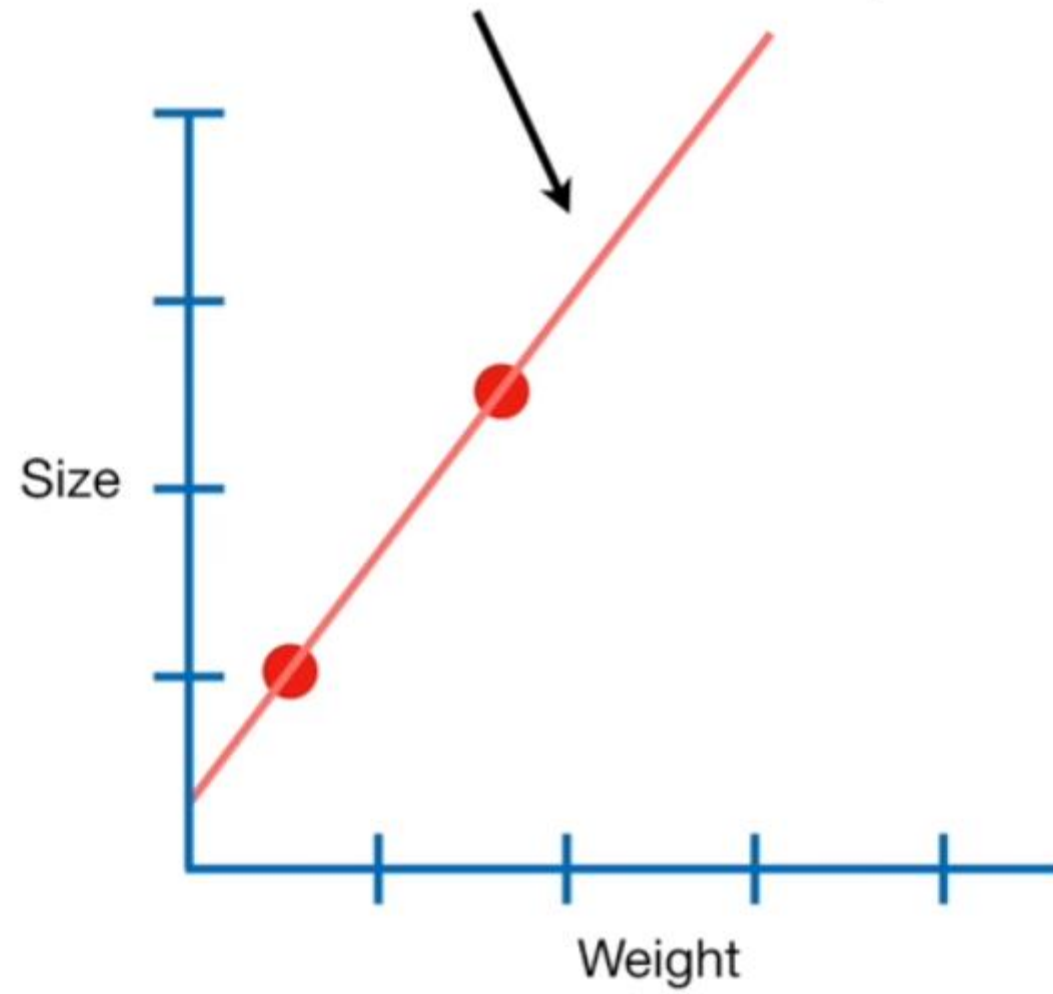
When we have a lot of measurements, we can be fairly confident that the **Least Squares** line accurately reflects the relationship between **Size** and **Weight**.
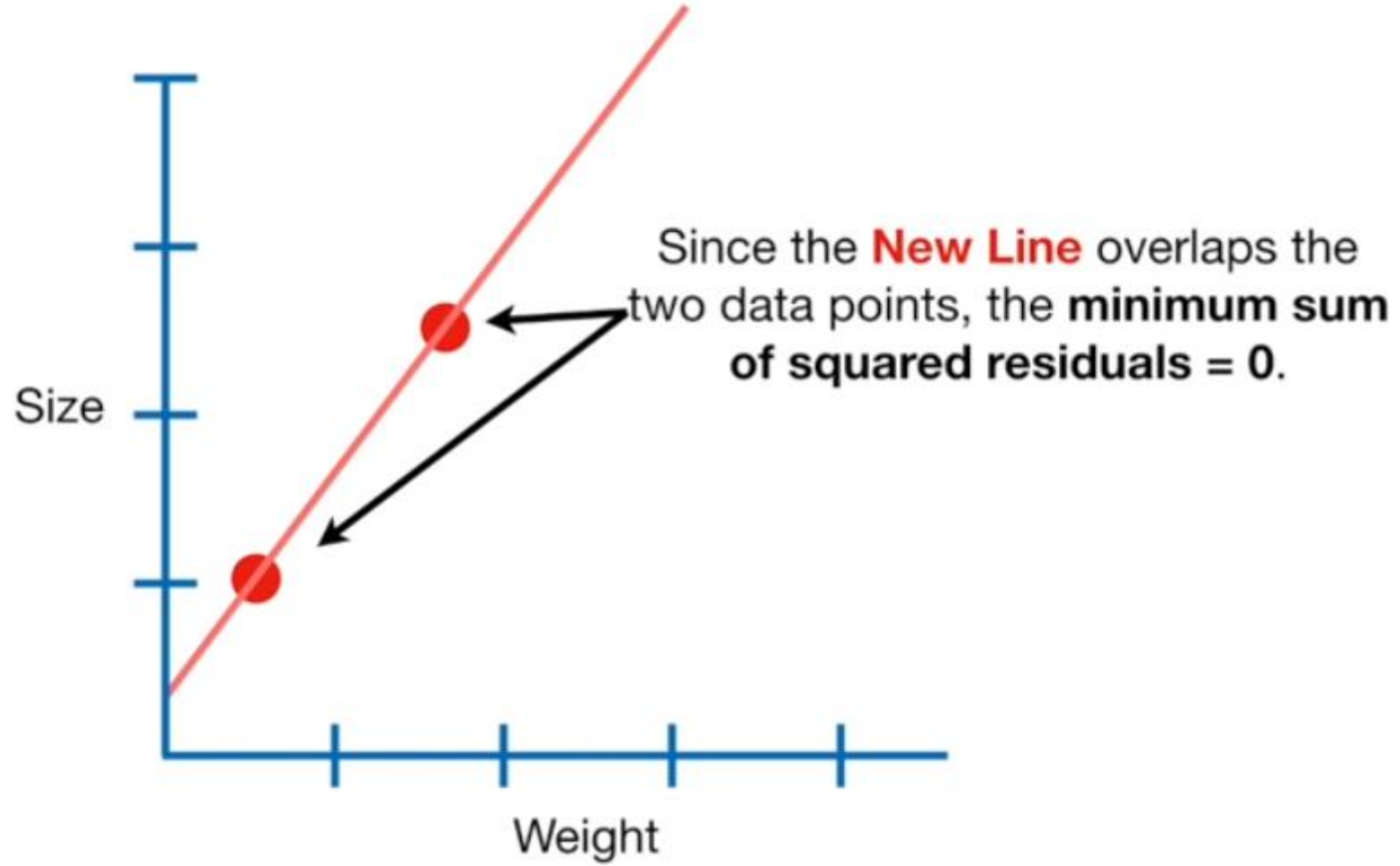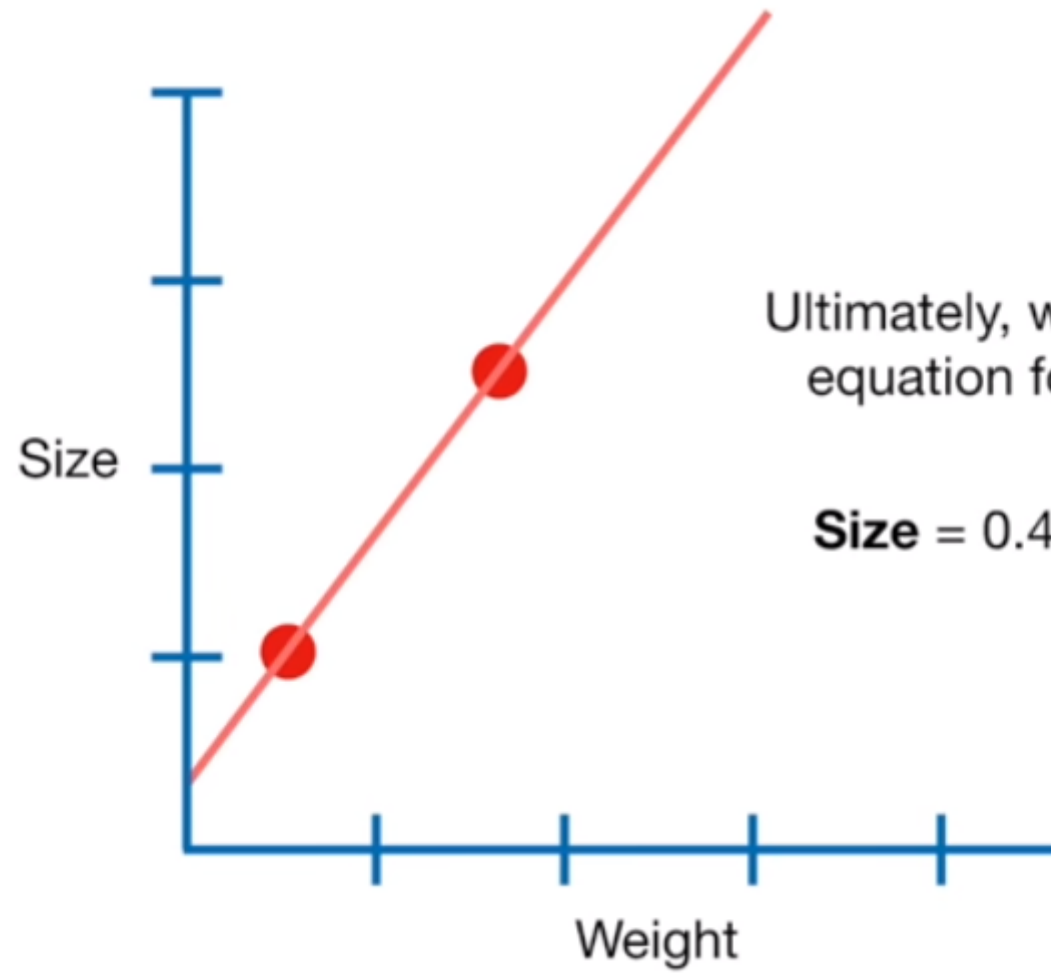
We fit a **New Line** with **Least Squares**…

Size

Weight

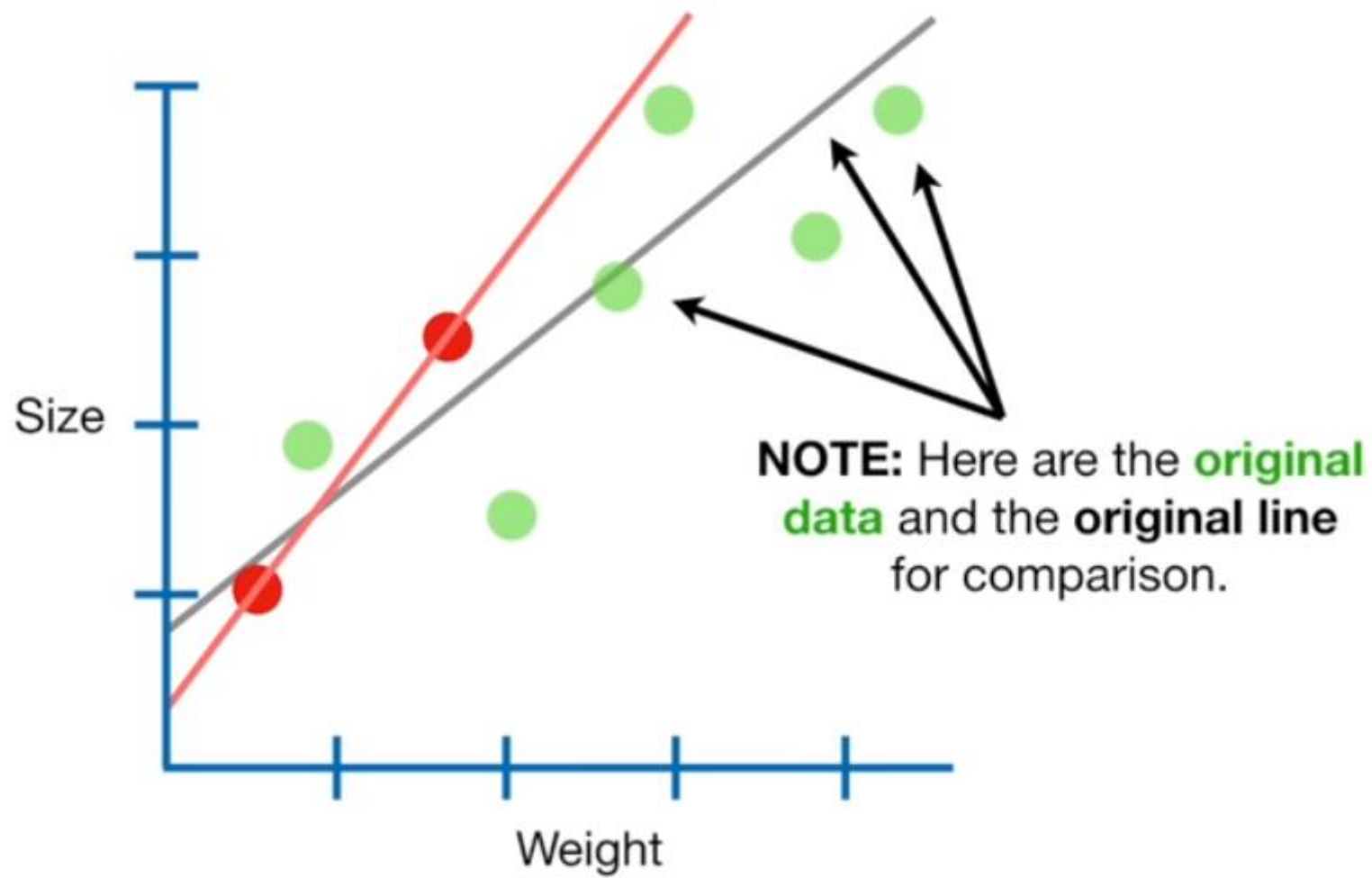Since the **New Line** overlaps the two data points, the **minimum sum of squared residuals = 0**.

Ultimately, we end up with this equation for the **New Line**:

**Size** $= 0.4 + 1.3 \times$ **Weight**

**NOTE:** Here are the **original data** and the **original line** for comparison.
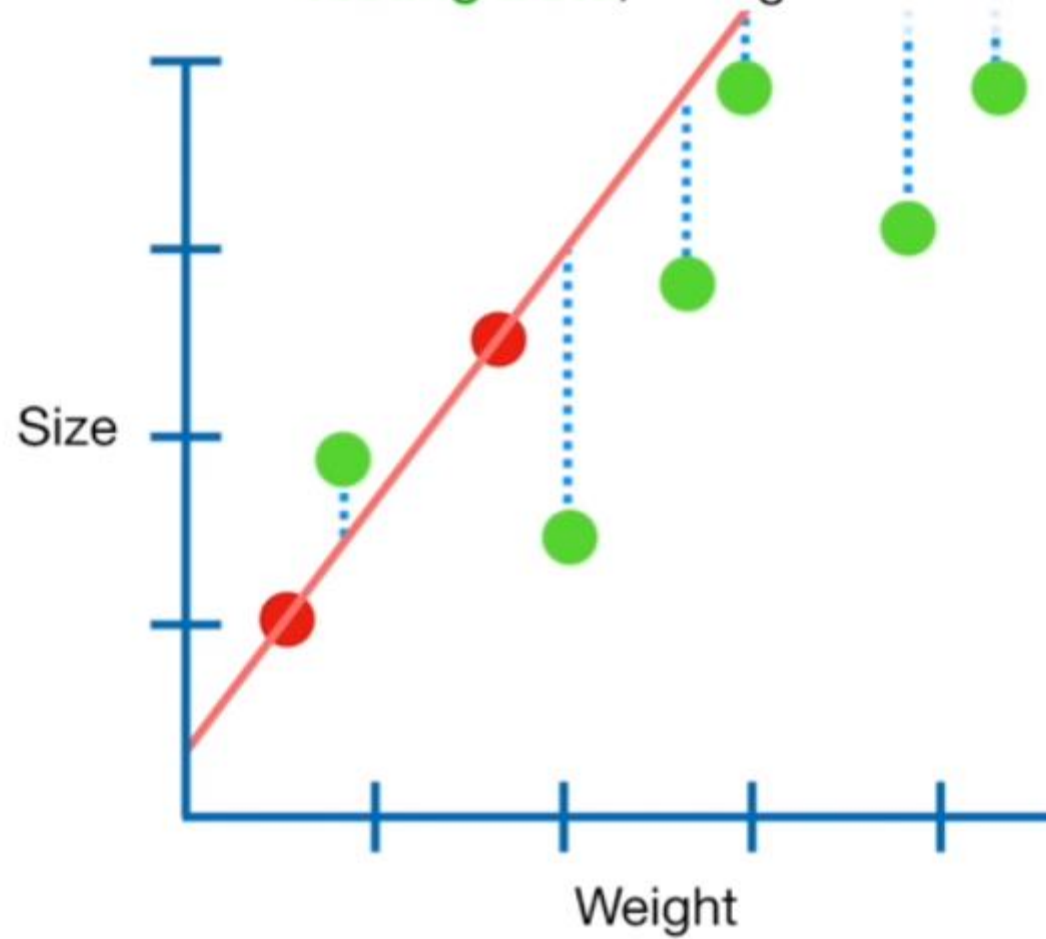
Let's call the **Two Red Dots** the **Training Data**, and the remaining **Green Dots** the **Testing Data**.

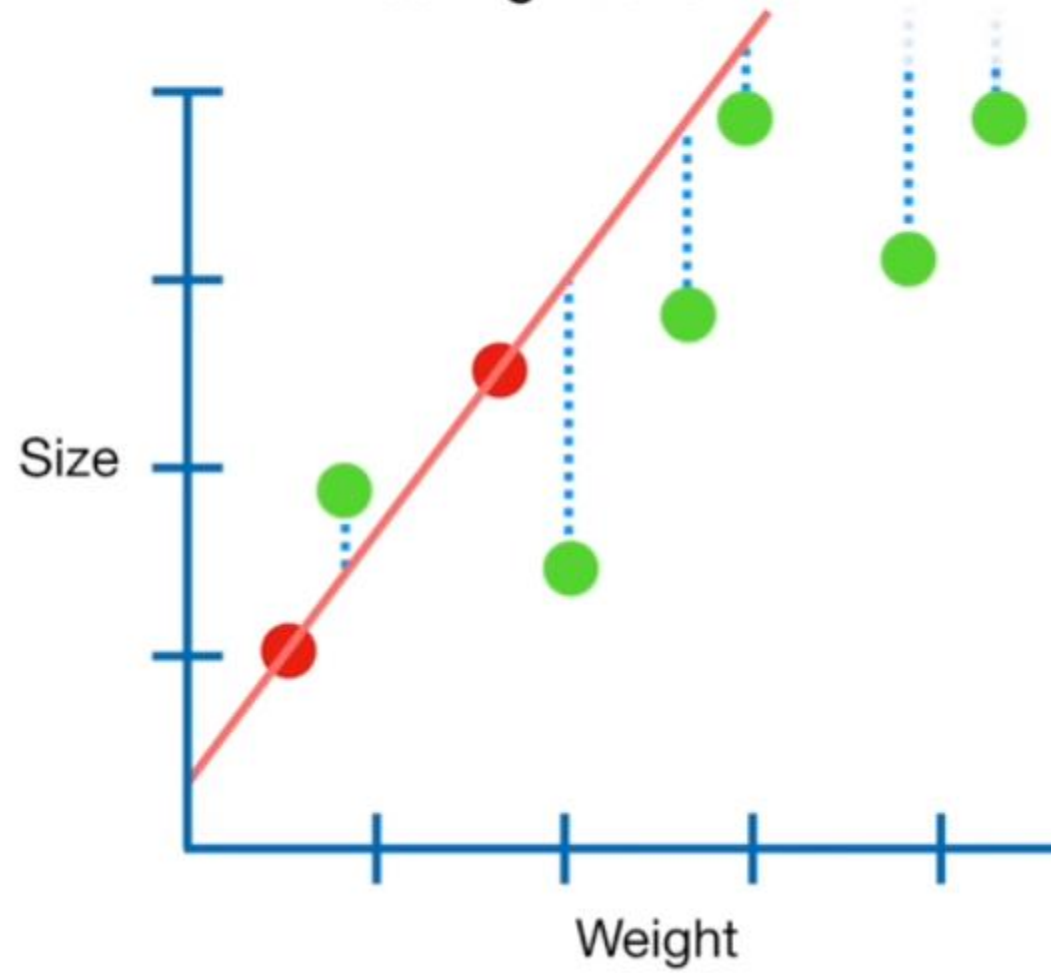The sum of the squared residuals for just the **Two Red Points**, the **Training Data**, is small (in this case it is 0)…
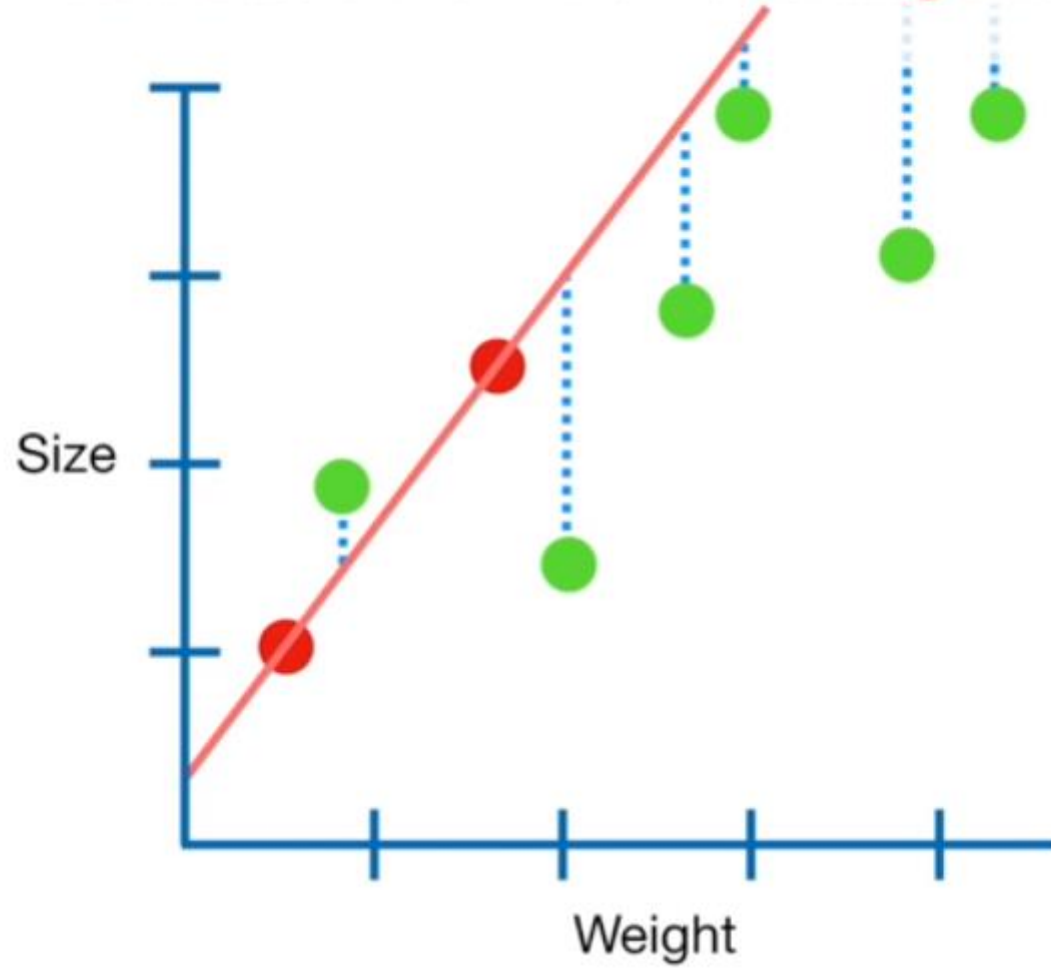
We just saw that **Least Squares** results in a
<span style="color:red">**Line**</span> that is **Over Fit** and has **High
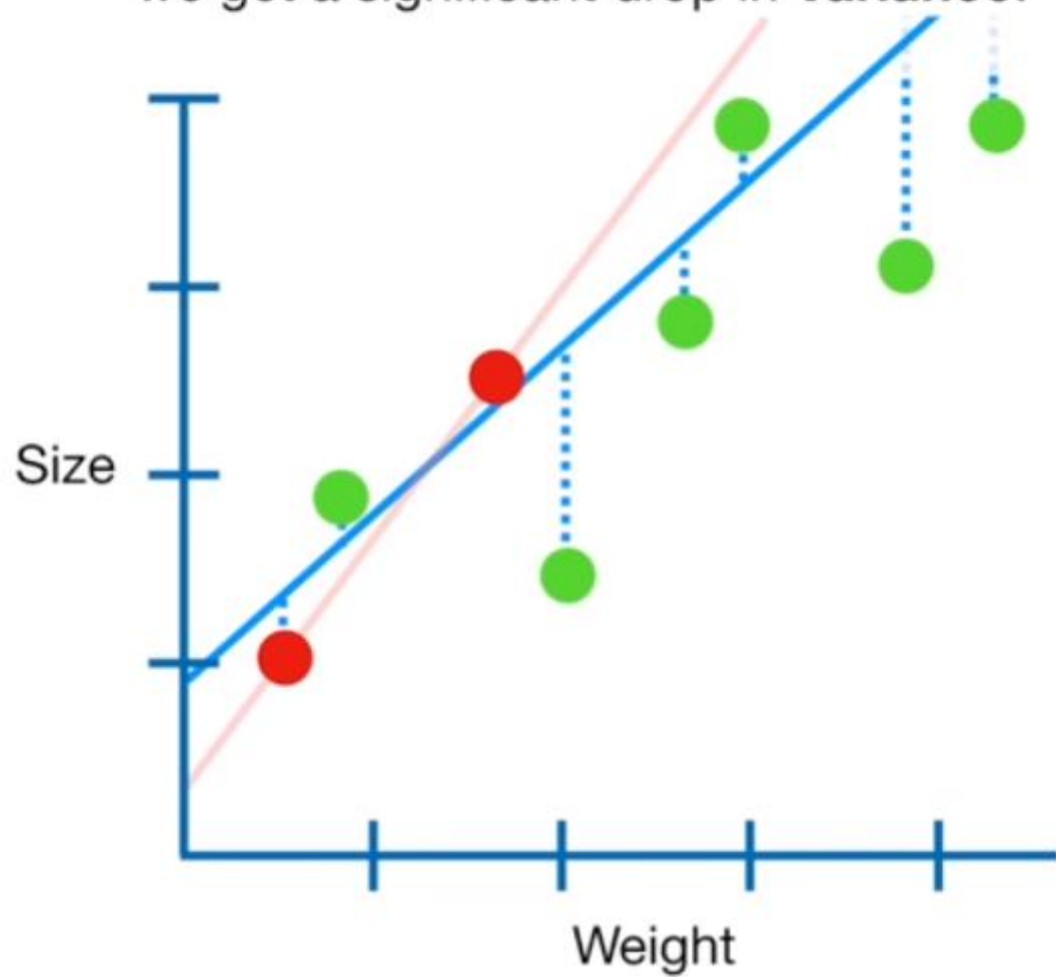Variance**…

Size

Weight

In contrast, when **Ridge Regression** determines values for the parameters in this equation…

↓

**Size** = y-axis intercept + slope × **Weight**

…it minimizes…

**the sum of the squared residuals**

**+**

$\lambda \times$ **the slope²**

...and **λ** (lambda) determines how severe that penalty is.
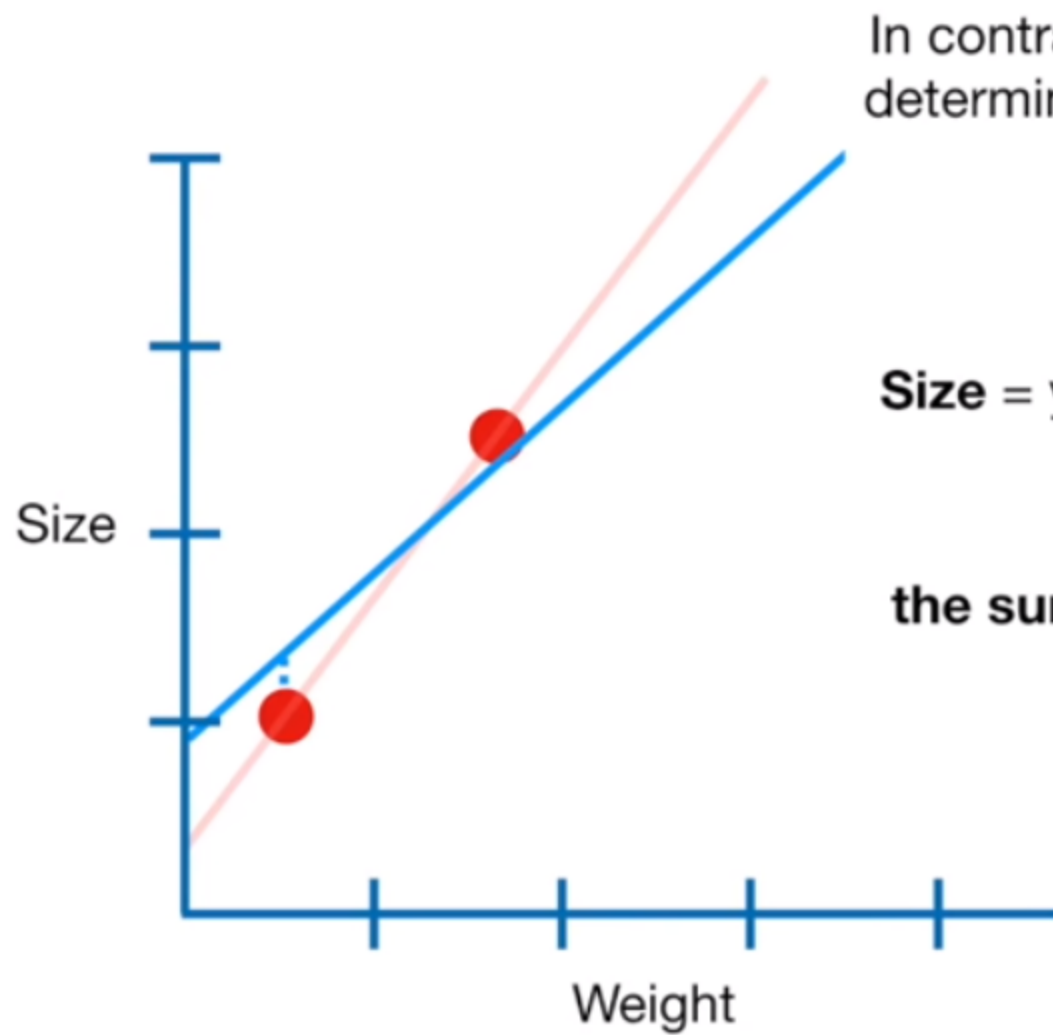
the sum of the squared residuals

+

**λ** × **the slope²**

To get a better idea of what's going on, lets plug in some numbers!!!

**the sum of the squared residuals**

**+**

$\lambda \times$ **the slope$^2$**

Let's start by plugging in the numbers that correspond to the **Least Squares Fit**.

**Size** = 0.4 + 1.3 × **Weight**

**the sum of the squared residuals**

**+**

$\lambda \times$ **the slope²**

The sum of squared residuals for the **Least Squares Fit** is **0** (because the line overlaps the data points)…

$$Size = 0.4 + 1.3 \times Weight$$

**the sum of the squared residuals**

$+$

$\lambda \times$ the slope$^2$

Size

Weight

...and the slope is **1.3**.
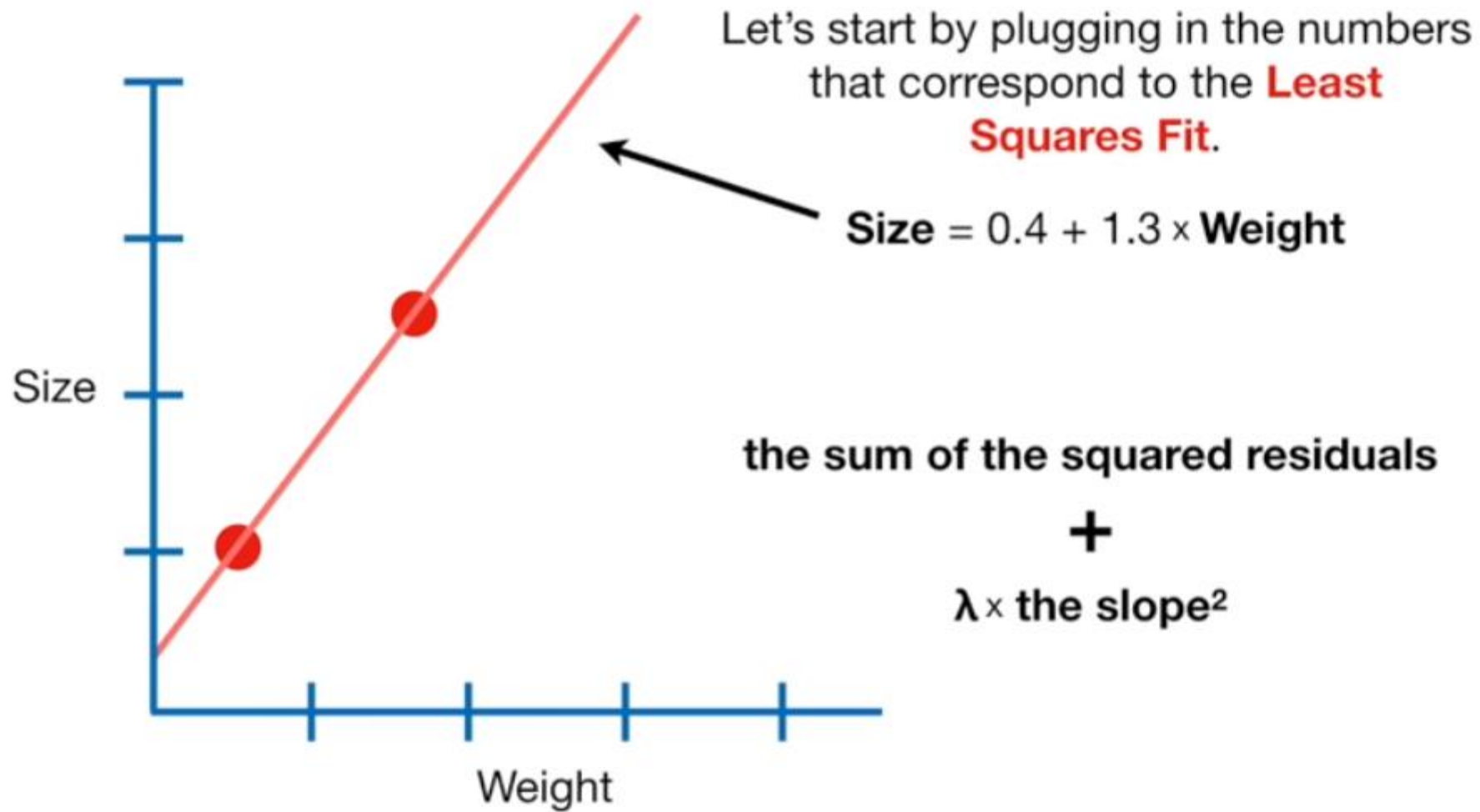
Size = 0.4 + **1.3** Weight

0

+

λ × **1.3²**

# Lambda =1



...all together, we have...

$$0 + 1 \times 1.3^2$$

...and when we do the math we get...

$0$

$+$

$1 \times 1.3^2$

$= 0 + 1.69 = 1.69$

Now let's see what happens when we plug in numbers for the **Ridge Regression** line...

**Size** = 0.9 + 0.8 × **Weight**

the sum of the squared residuals

+

$\lambda \times$ the slope$^2$

The slope is **0.8**…

Size = 0.9 + **0.8** × Weight

$0.3^2 + 0.1^2$

**+**

λ × **the slope²**

...and just like before, we'll let $\lambda = 1$.

$$0.3^2 + 0.1^2$$

$$+$$

$$1 \times 0.8^2$$

...and when we do the math we get...

$0.3^2 + 0.1^2$

$+$

$1 \times 0.8^2$

$= 0.09 + 0.01 + 0.64$

$= 0.74$

= 0.74

= 1.69

For the **Ridge Regression Line**, the sum of squared residuals plus the **Ridge Regression Penalty** is **0.74**…
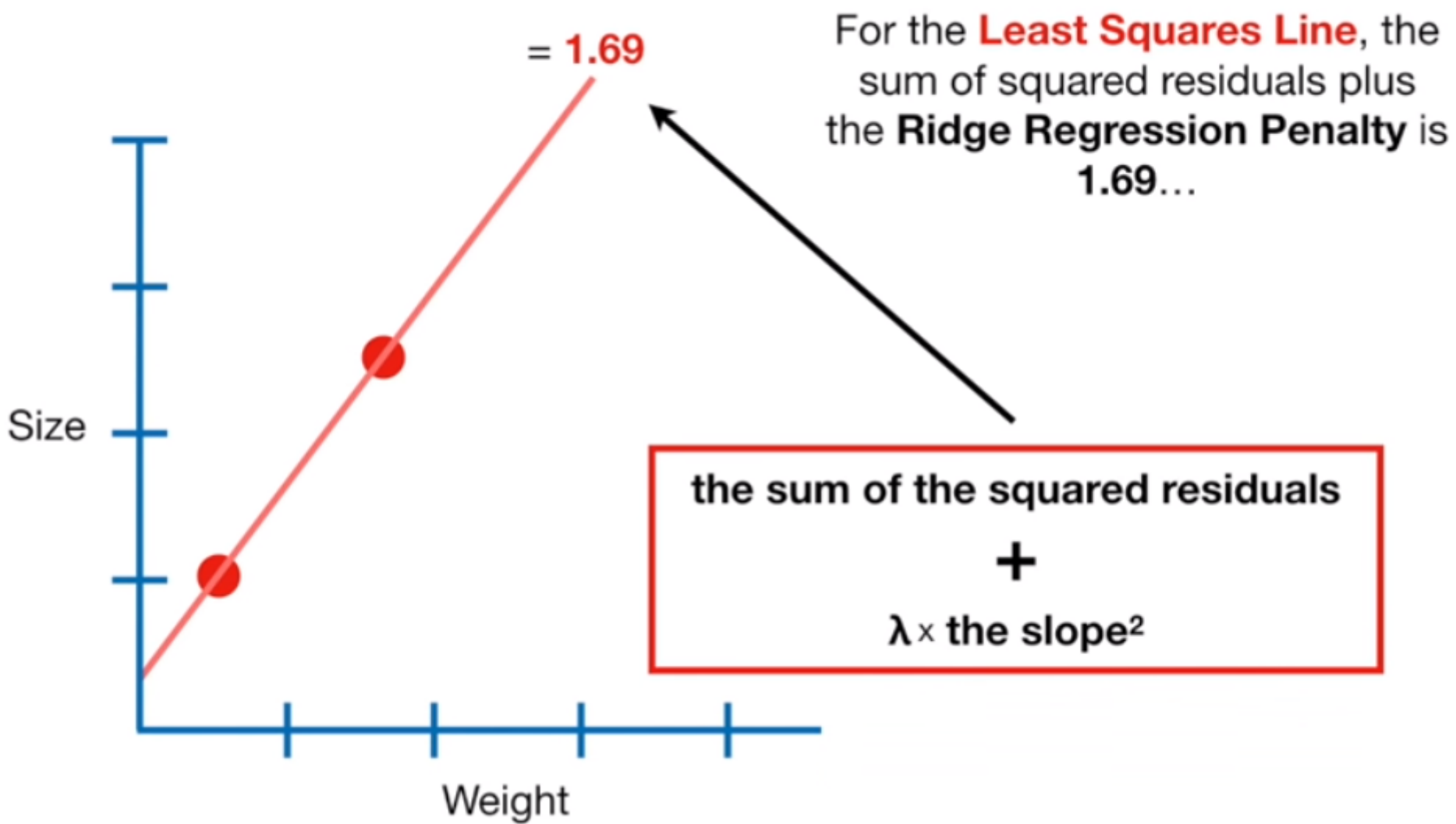
Size

Weight

**the sum of the squared residuals**

**+**

$\lambda \times$ **the slope²**
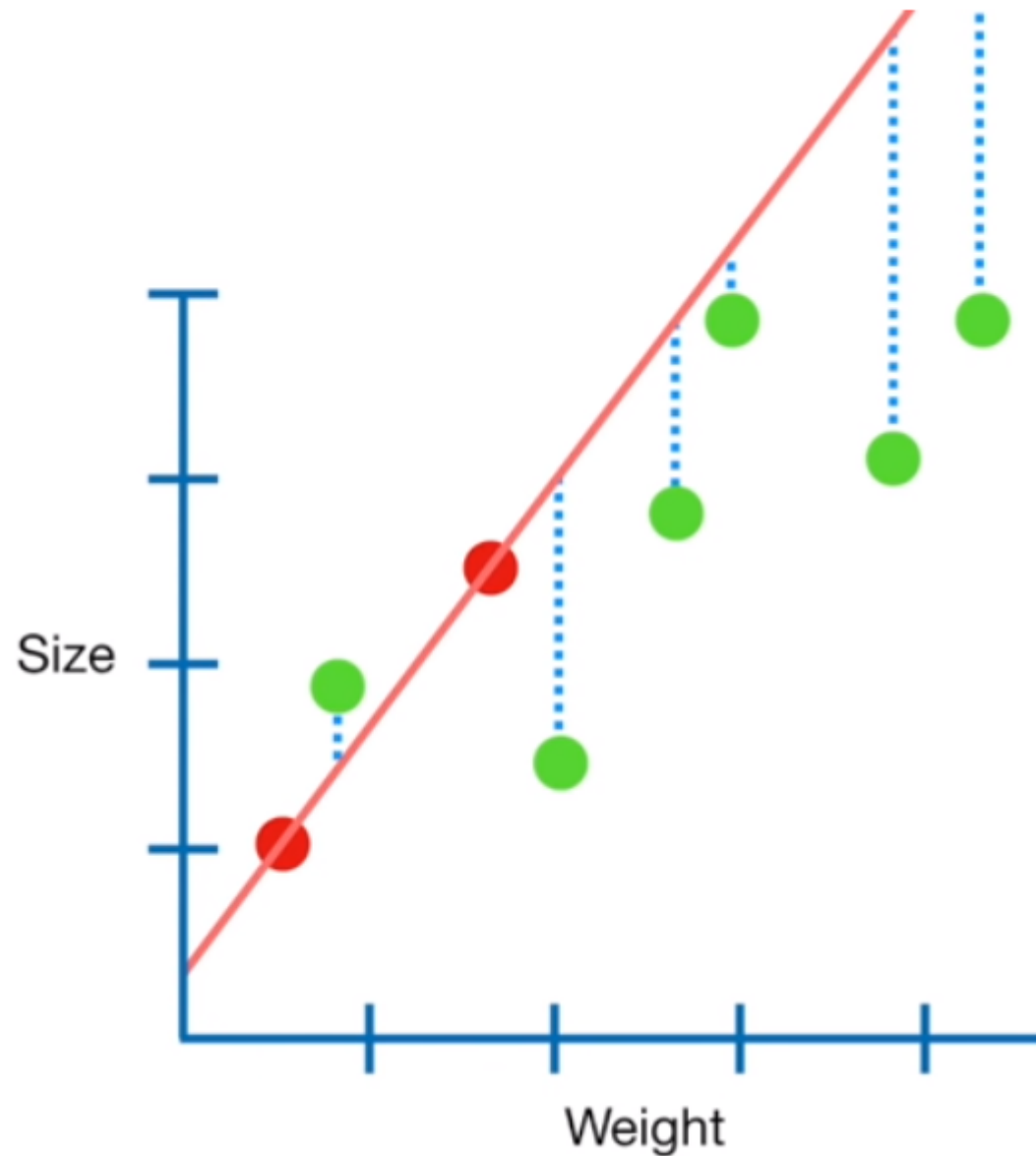
= 0.74

= 1.69

Size
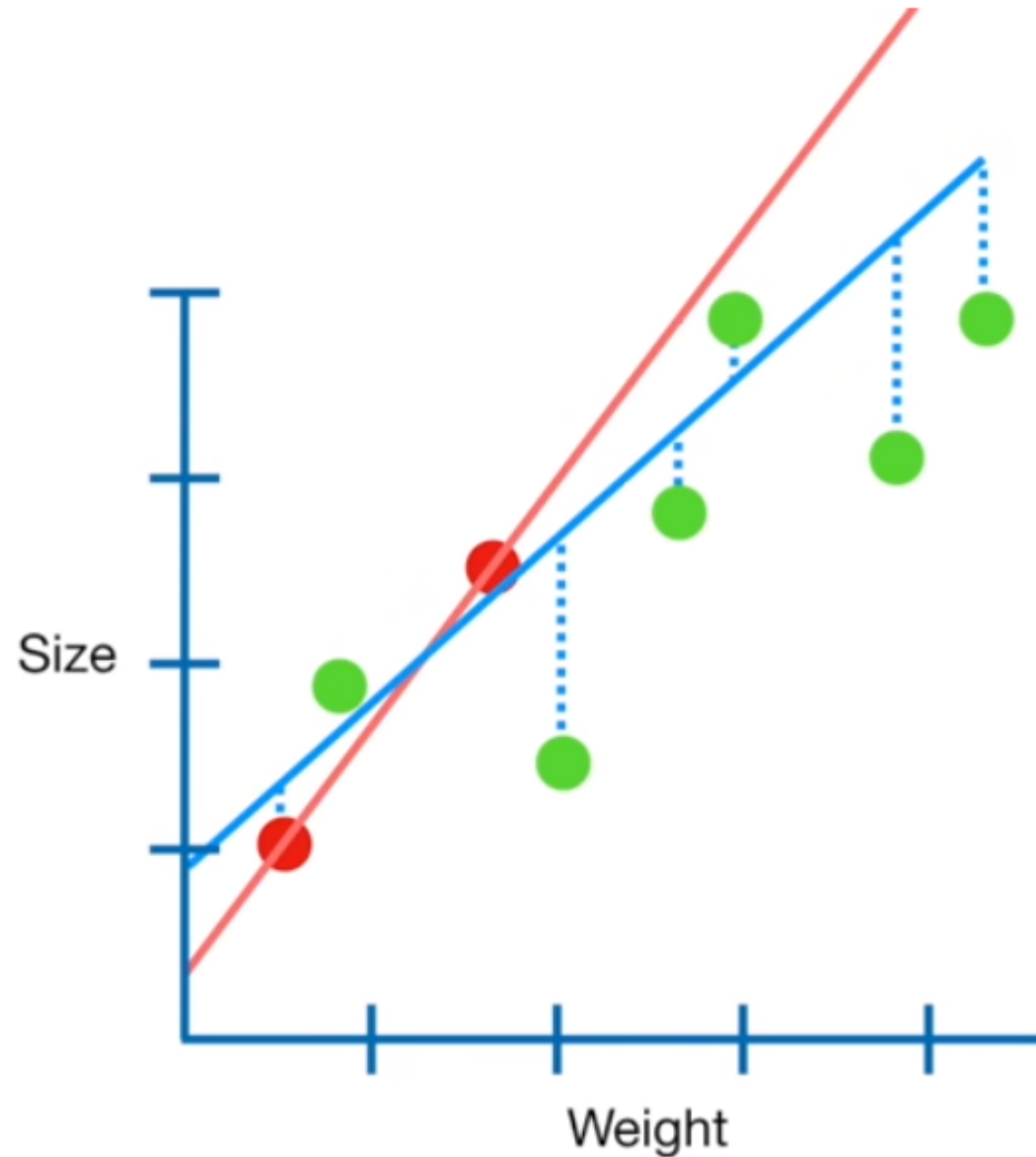
Weight

Thus, if we wanted to minimize the sum of the squared residuals plus the **Ridge Regression Penalty**, we would choose the **Ridge Regression Line** over the **Least Squares Line**.

**the sum of the squared residuals**
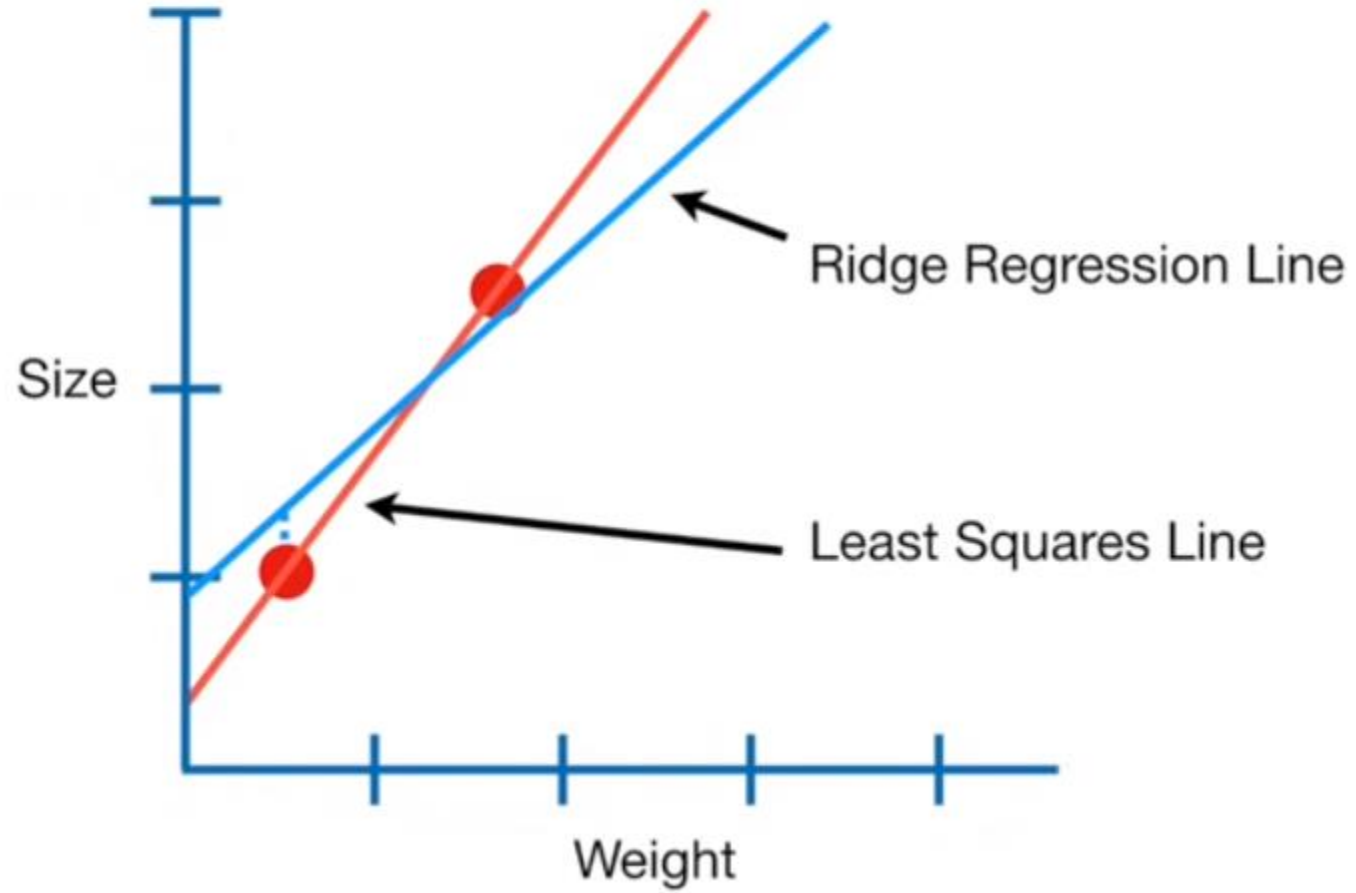
**+**

$\lambda \times$ **the slope$^2$**

Without the small amount of **Bias** that the penalty creates, the **Least Squares Fit** has a large amount of **Variance**.

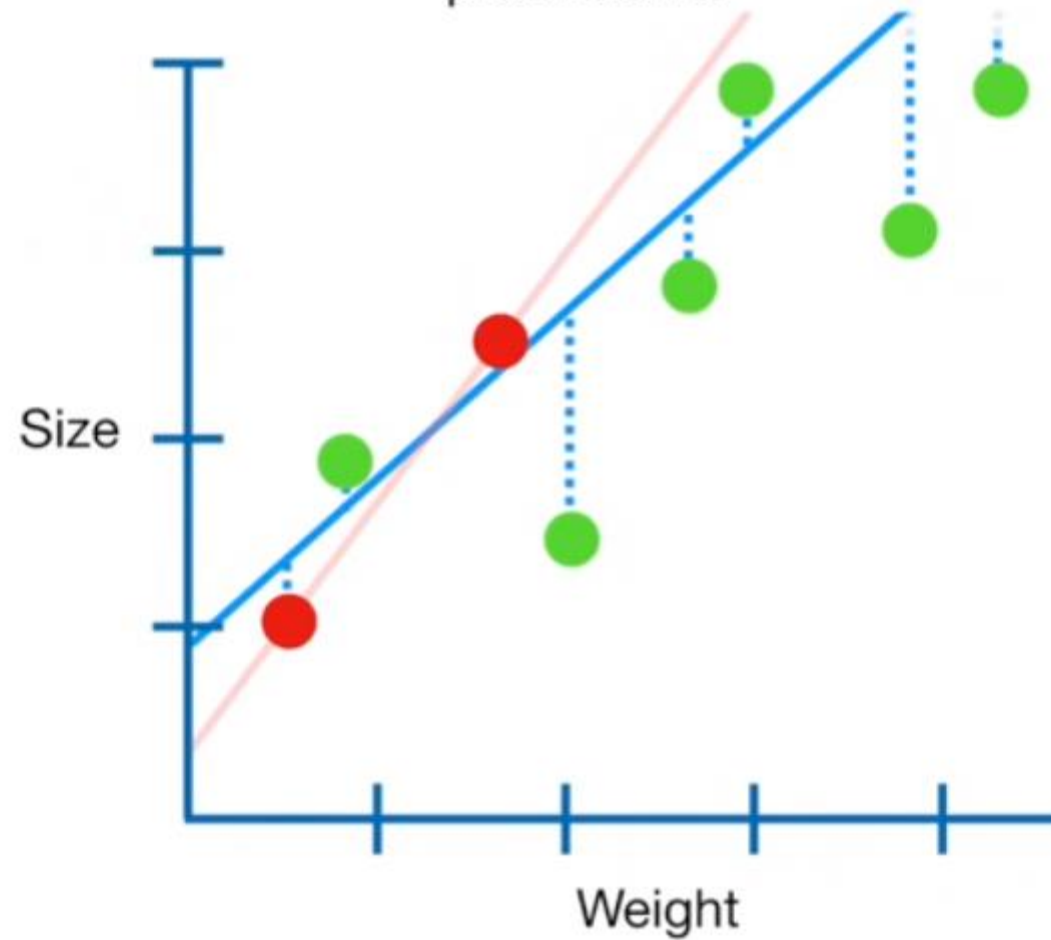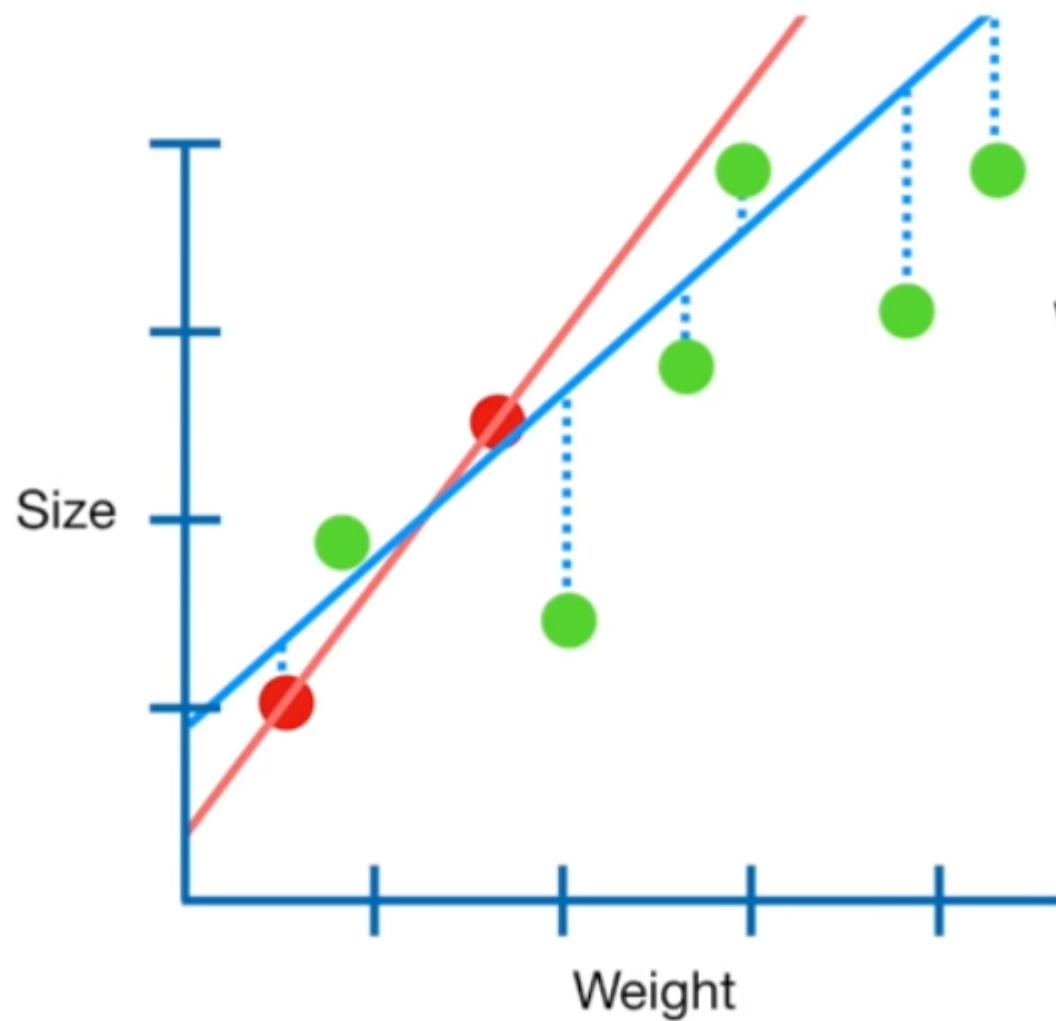In contrast, the **Ridge Regression Line**, which has the small amount of **Bias** due to the penalty, has less **Variance**.
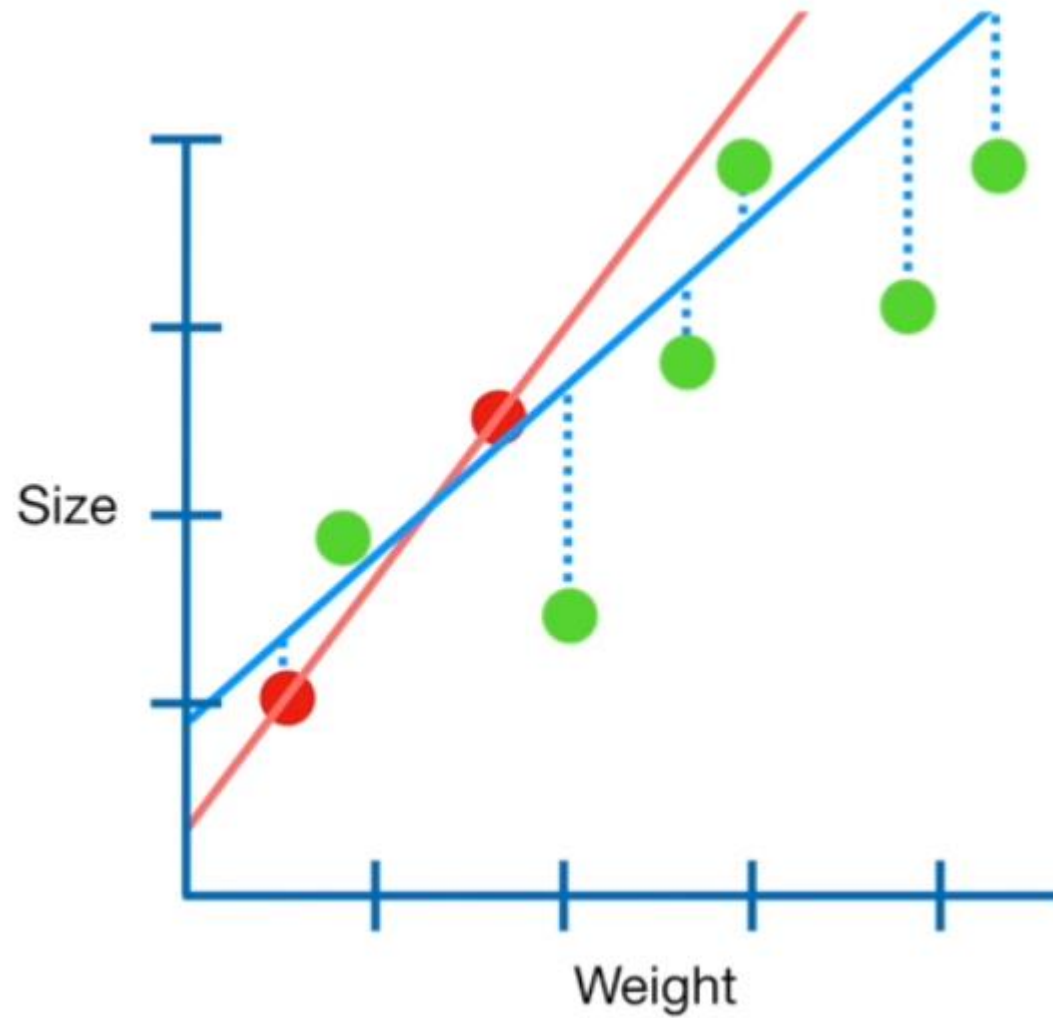
The main idea was that by starting with a slightly worse fit, **Ridge Regression** provided better long term predictions.

When the sample sizes are relatively small, then **Ridge Regression** can improve predictions made from new data (i.e. reduce **Variance**) by making the predictions less sensitive to the **Training Data**.

The **Ridge Regression Penalty** itself is **λ** times the sum of all squared parameters, except for the y-intercept…
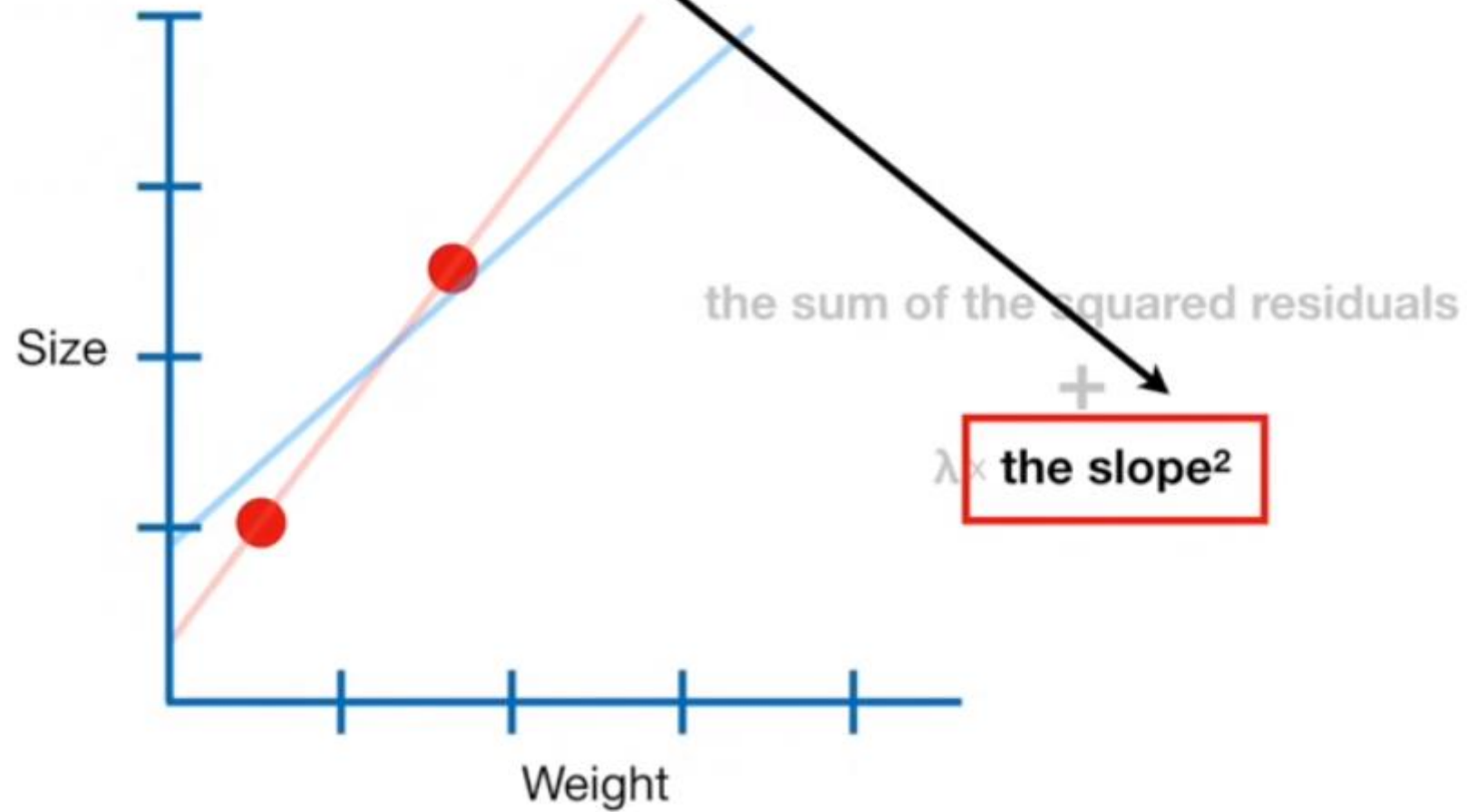
the sum of the squared residuals

**+**

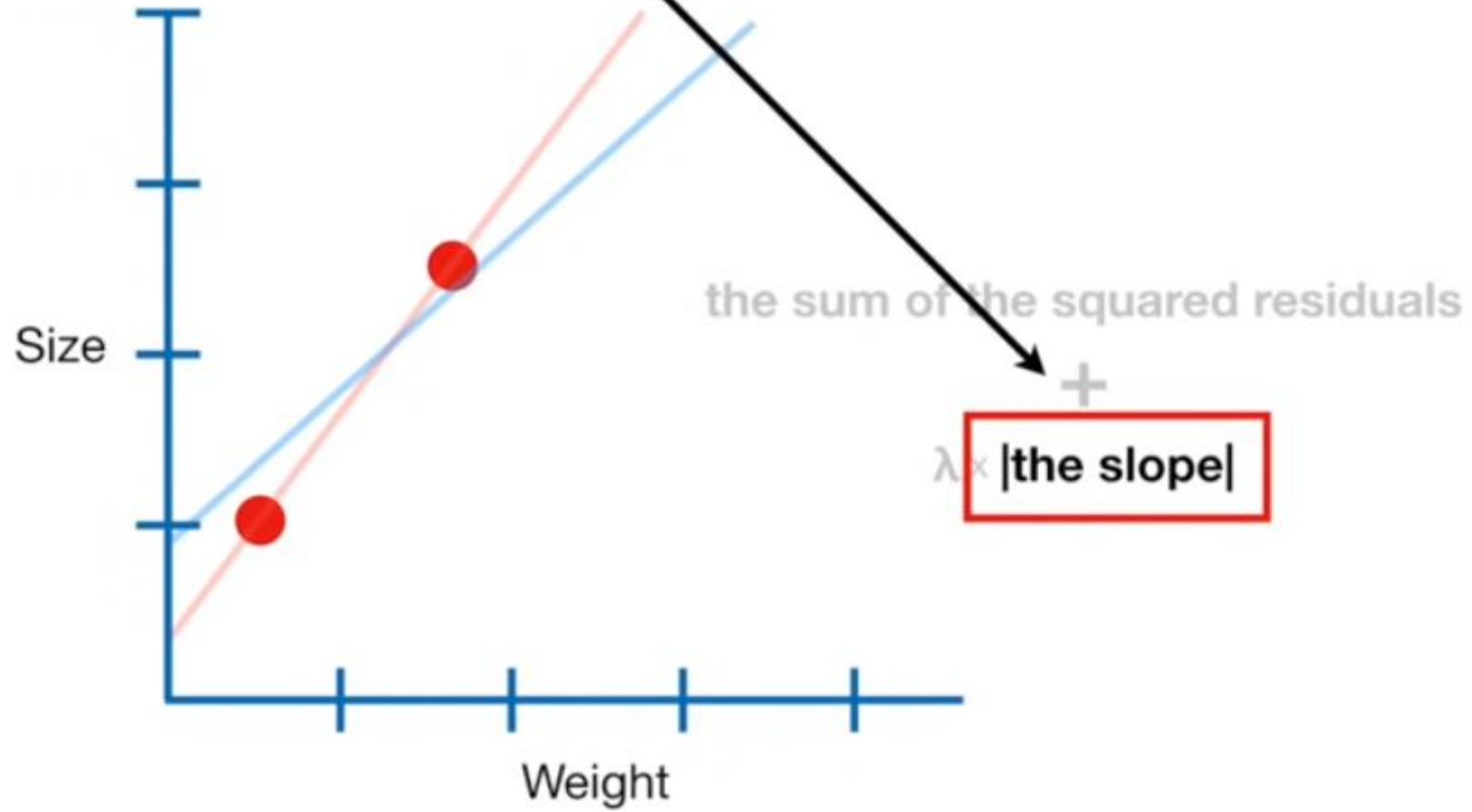$$\lambda \times \text{Slope}^2$$

$\lambda$ can be any value from **0** to
**positive infinity**.

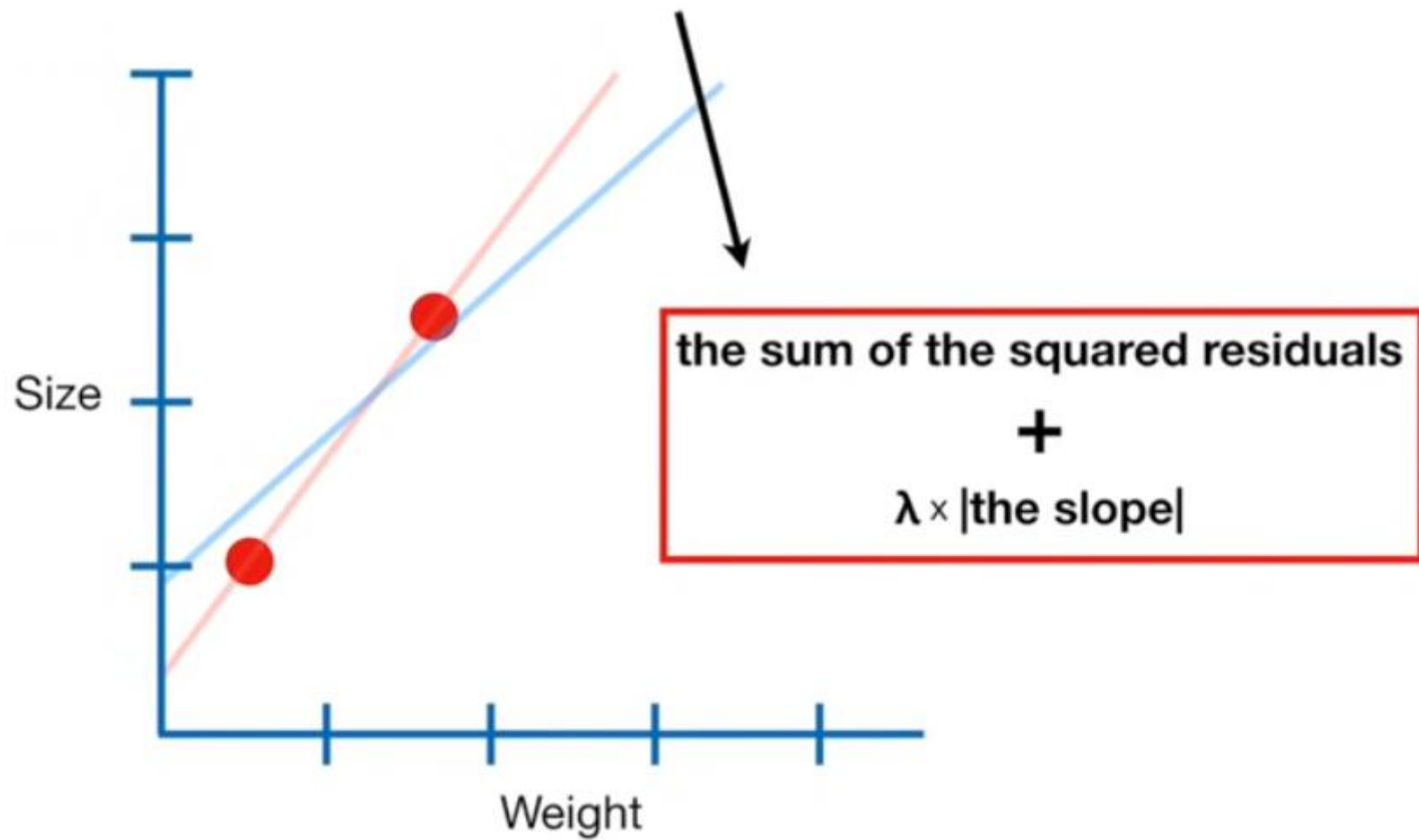**Lasso Regression** is very, very similar to **Ridge Regression**, but it has some very, very important differences...
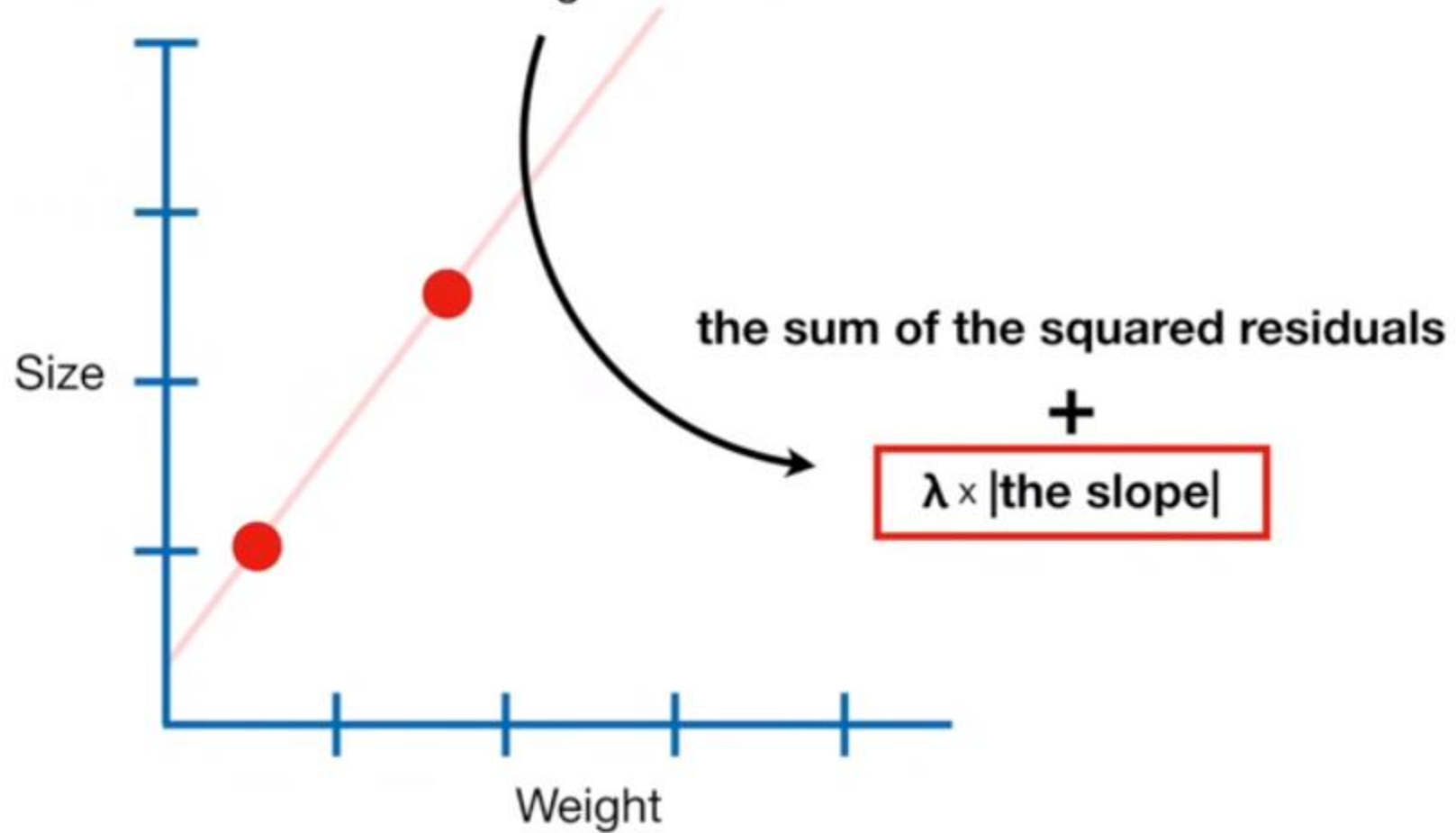
If, instead of squaring the slope…

Size

Weight

the sum of the squared residuals

$+$

$\lambda \times$ **the slope²**

...we take the absolute value...

the sum of the squared residuals

$+$

$\lambda \times$ |the slope|

...then we have **Lasso Regression**!!!



the sum of the squared residuals

**+**

$\lambda \times |\text{the slope}|$

Like **Ridge Regression**, **Lasso Regression** (the Orange Line) results in a line with a little bit of **Bias**...



the sum of the squared residuals

**+**

$\lambda \times |\text{the slope}|$

Size

Weight

the sum of the squared residuals

**+**

$\lambda \times$ the slope$^2$

...look very similar....

the sum of the squared residuals

**+**

$\lambda \times |$the slope$|$