



# Ridge Regression and Lasso Regression

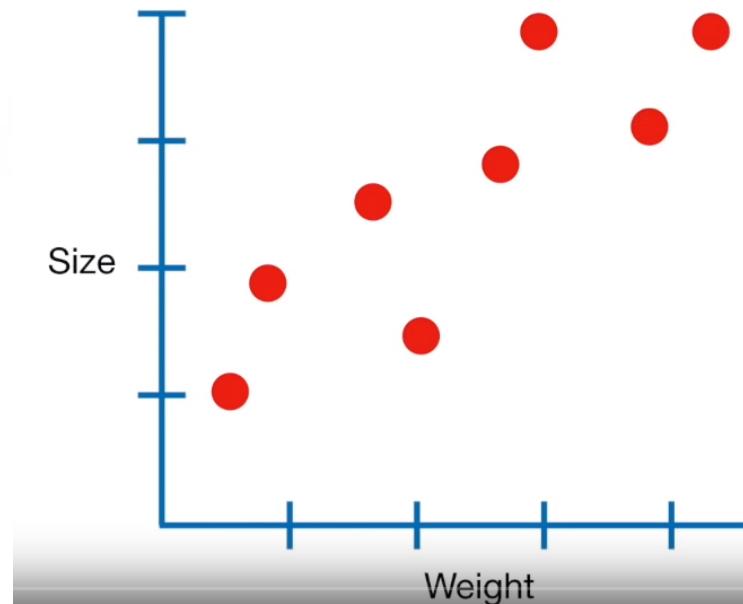
By,

Divesh R. Kubal

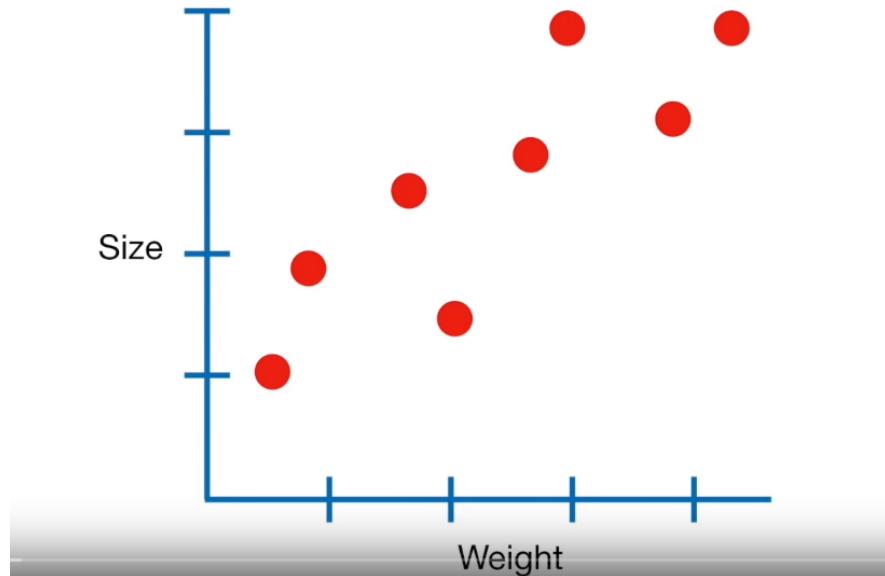
Data Scientist, eClerx Services  
Center of Excellence – Machine Learning

# Ridge Regression

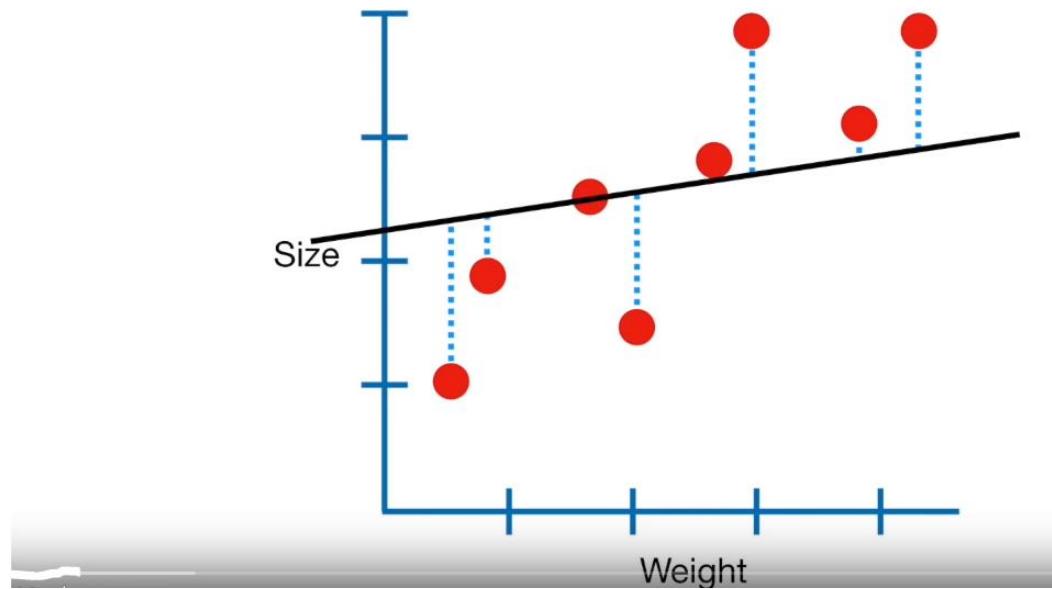
Let's start by collecting **Weight** and **Size** measurements from a bunch of mice...



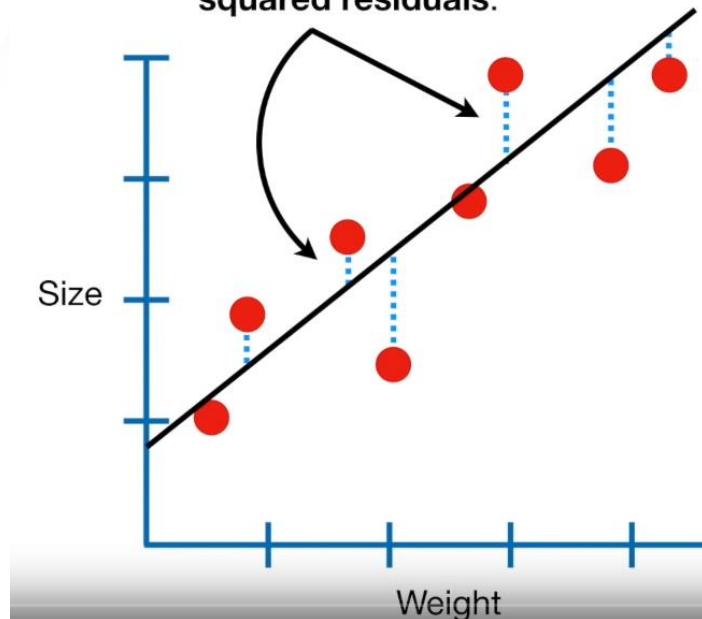
Since these data look relatively linear, we will use **Linear Regression**, AKA **Least Squares**, to model the relationship between **Weight** and **Size**.



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



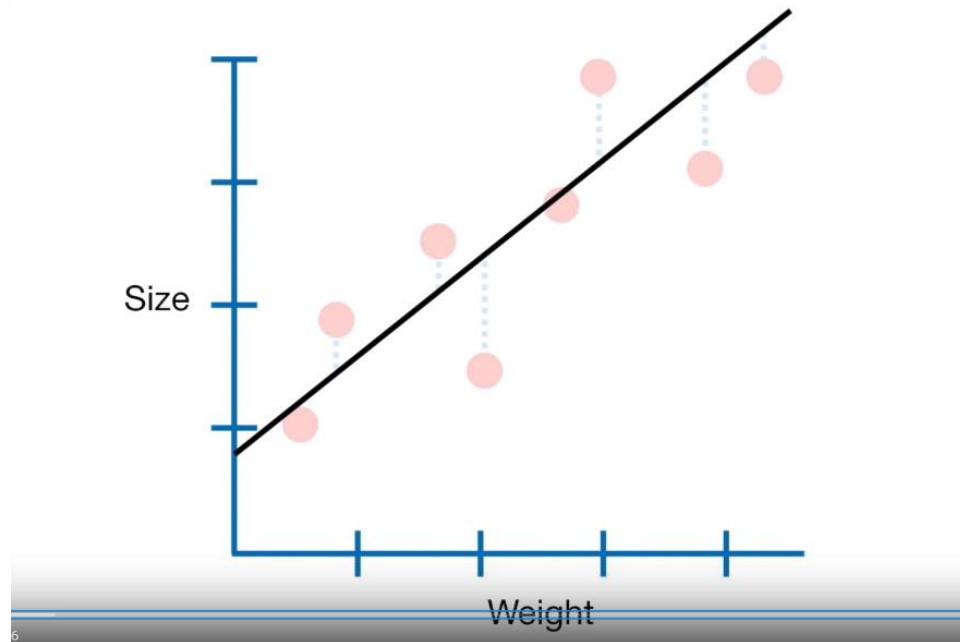
In other words, we find the line that results in the **minimum sum of squared residuals**.





Ultimately, we end up with  
this equation for the line:

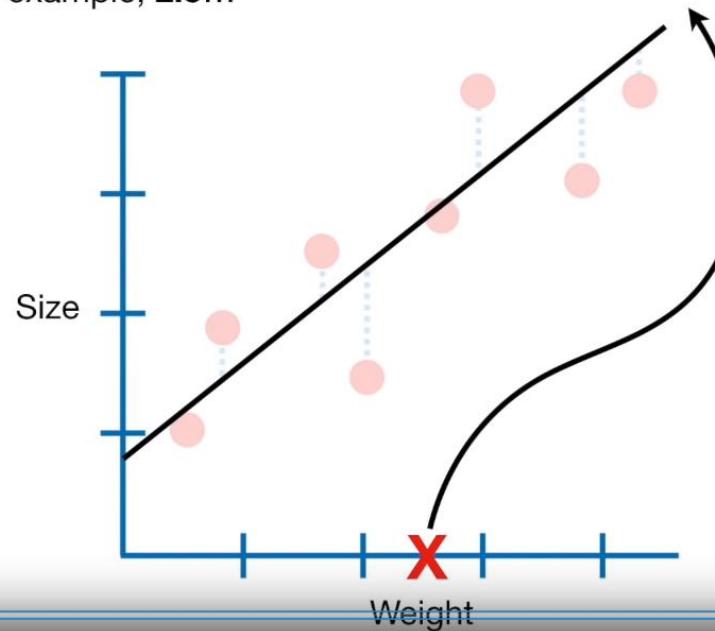
$$\text{Size} = 0.9 + 0.75 \times \text{Weight}$$





We can plug in a value for  
**Weight**, for example, **2.5**...

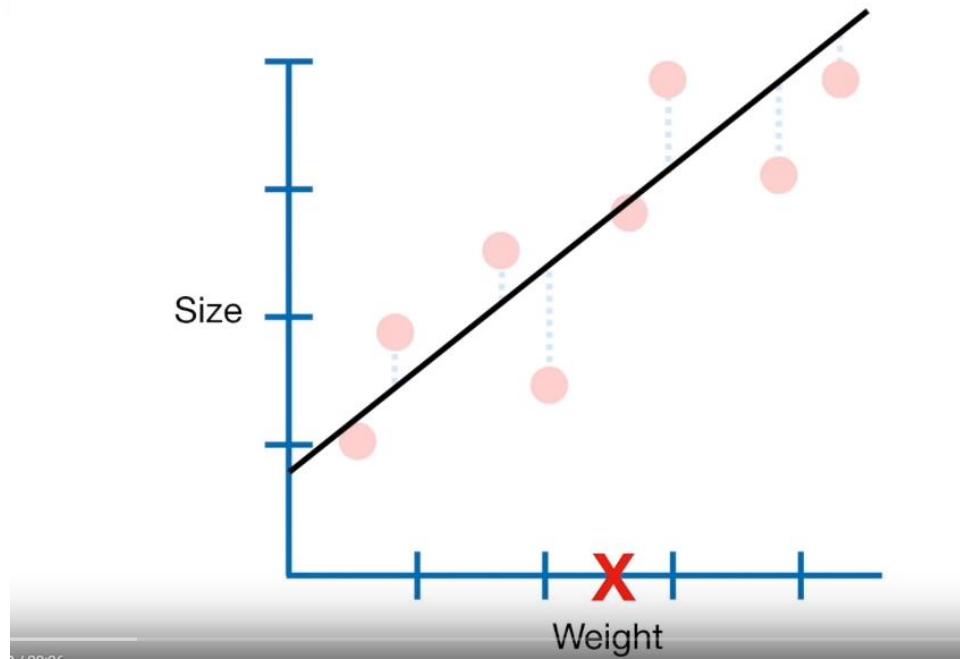
$$\text{Size} = 0.9 + 0.75 \times \text{Weight}$$



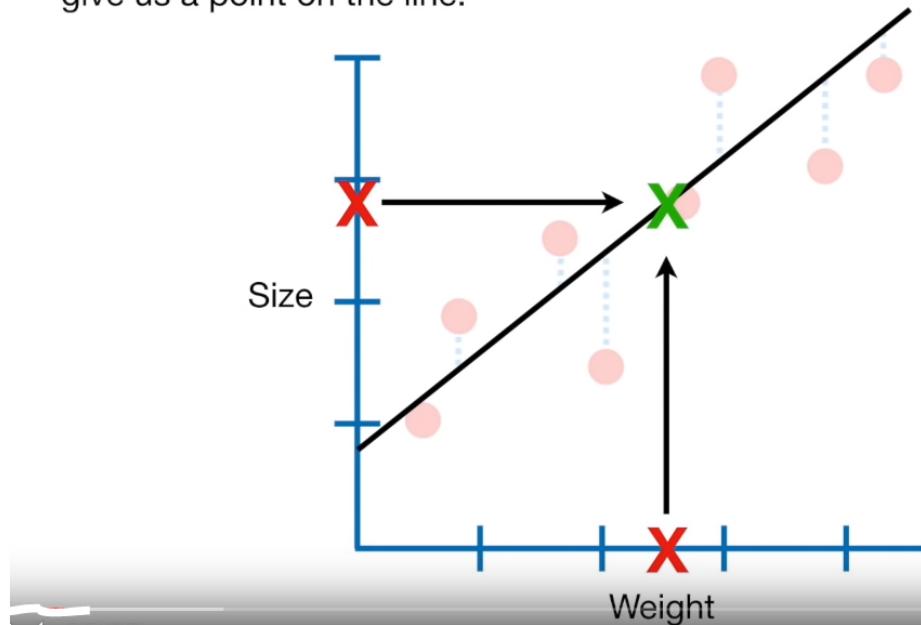


...and do the math...

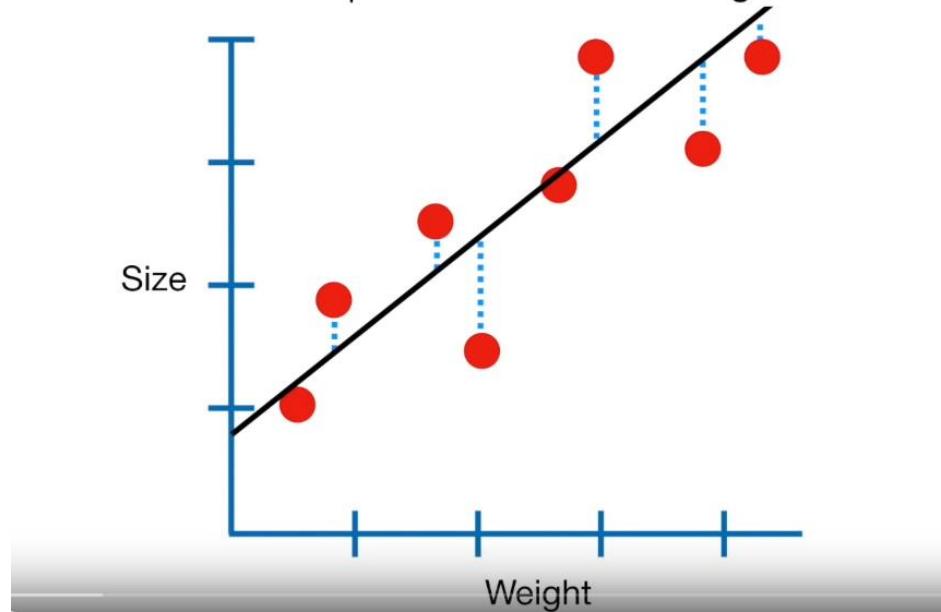
$$\text{Size} = 0.9 + 1.88$$



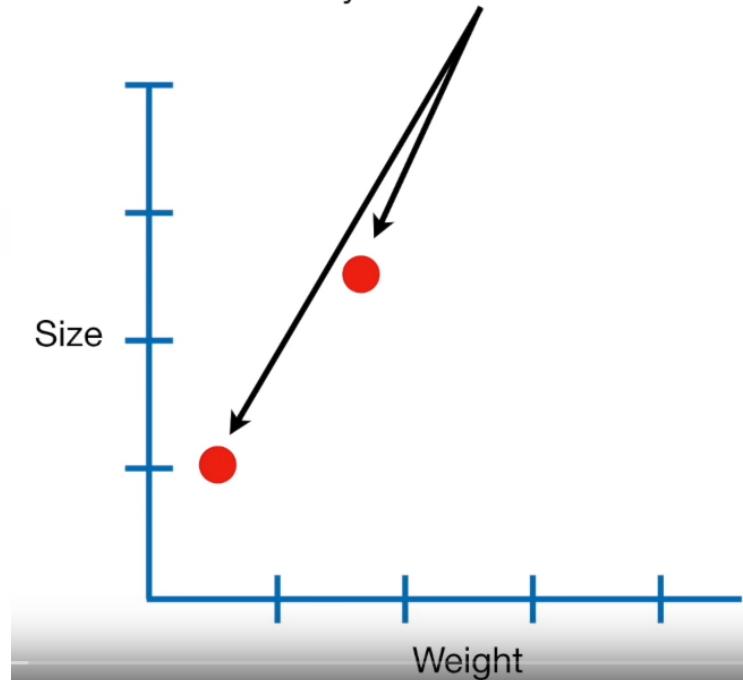
Together, the value for **Weight**,  
**2.5**, and the value for **Size**, **2.8**,  
give us a point on the line.



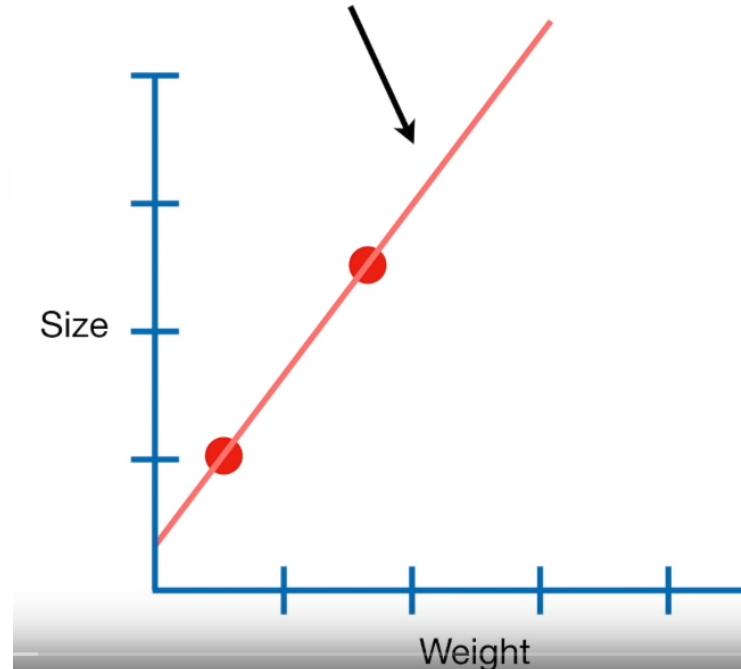
When we have a lot of measurements, we can be fairly confident that the **Least Squares** line accurately reflects the relationship between **Size** and **Weight**.

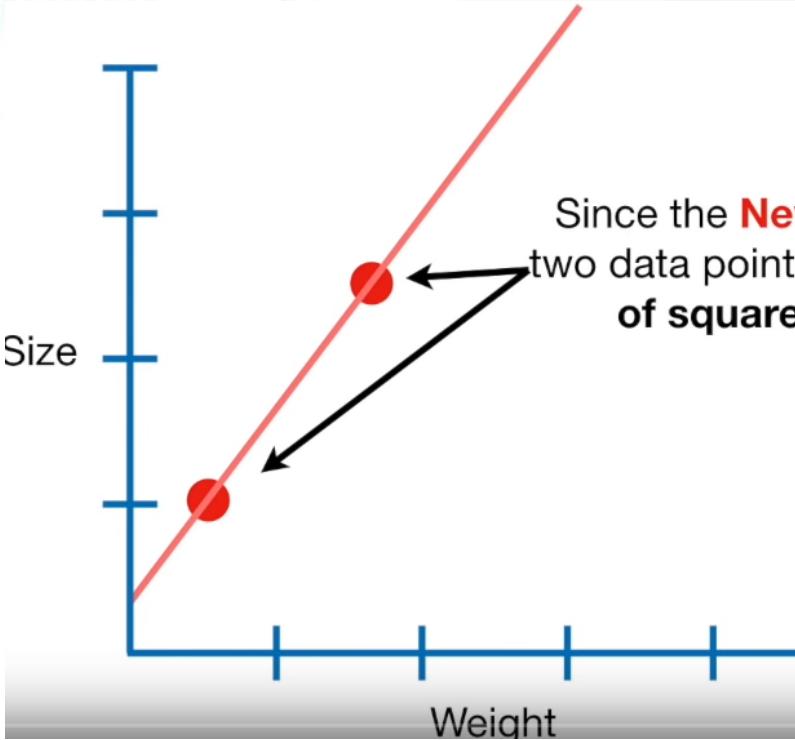


But what if we only have two measurements?

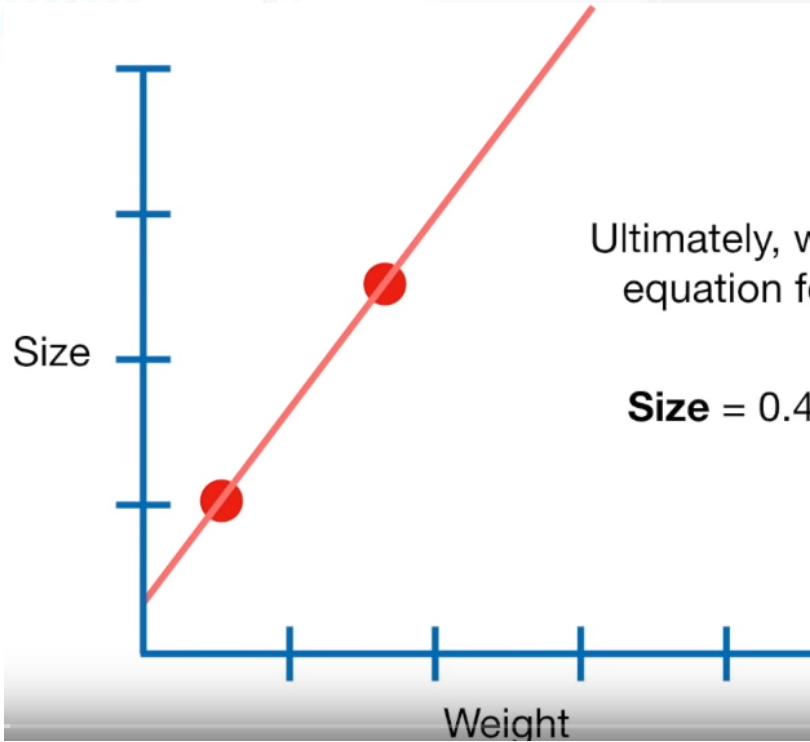


We fit a **New Line** with **Least Squares**...



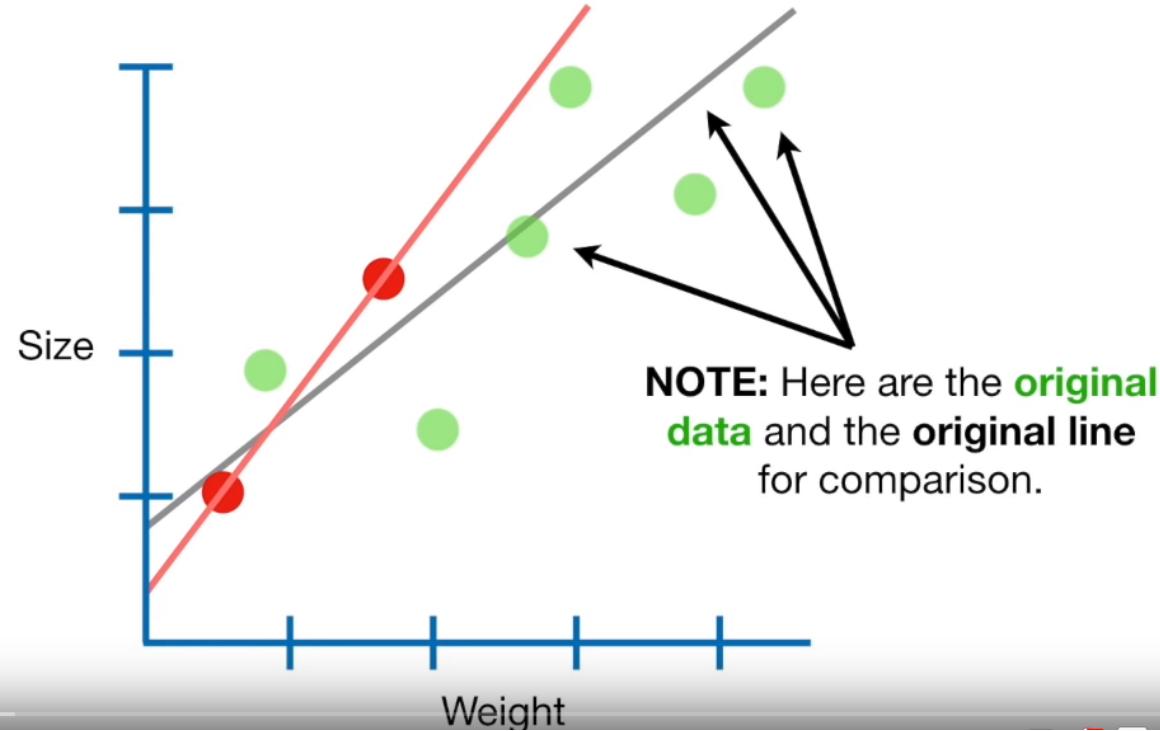


Since the **New Line** overlaps the  
two data points, the **minimum sum  
of squared residuals = 0.**

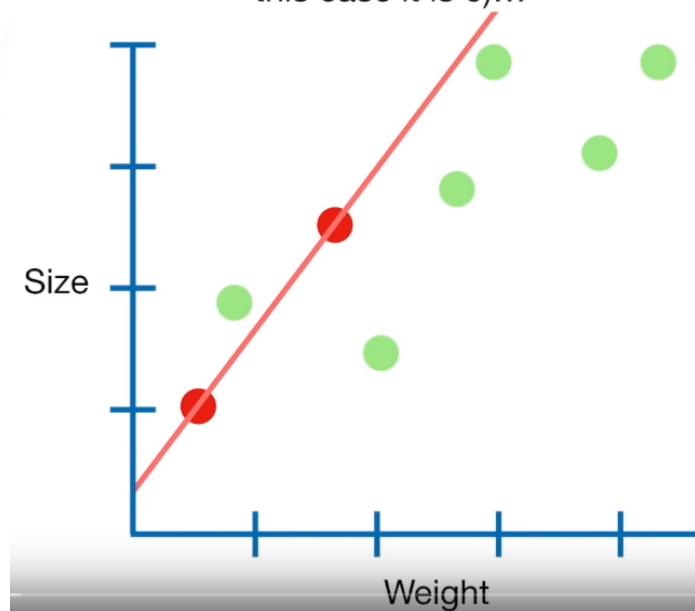


Ultimately, we end up with this equation for the **New Line**:

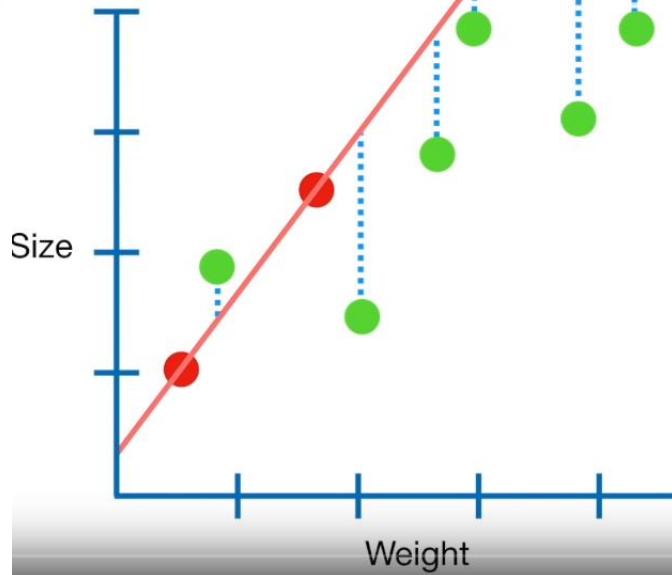
$$\text{Size} = 0.4 + 1.3 \times \text{Weight}$$



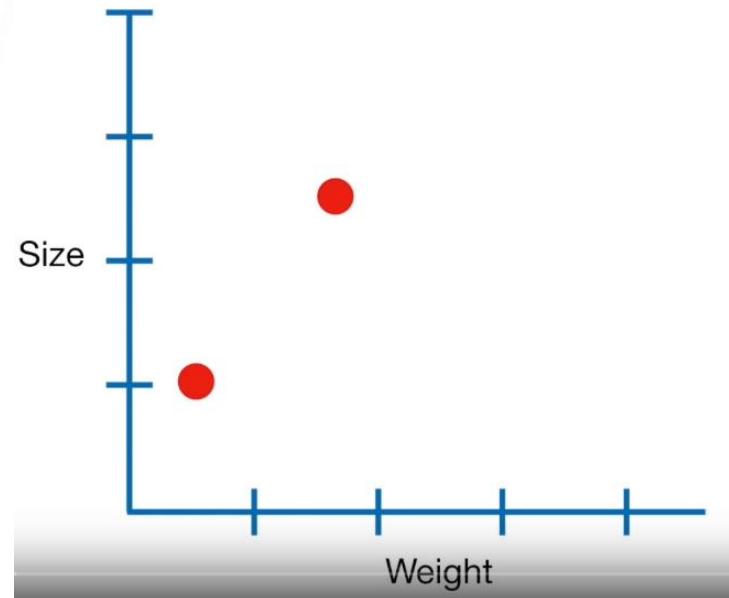
The sum of the squared residuals for just the  
**Two Red Points**, the **Training Data**, is small (in  
this case it is 0)...



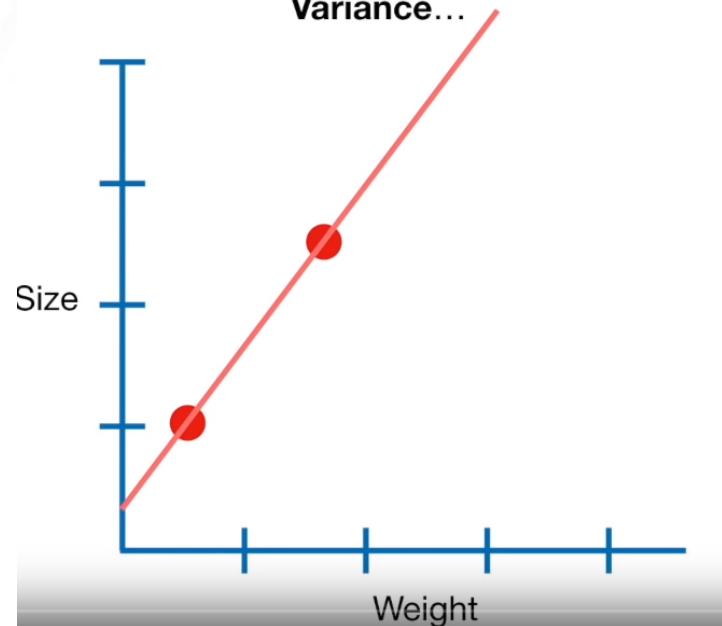
...but the sum of the squared residuals for the **Green Points**, the **Testing Data**, is large...



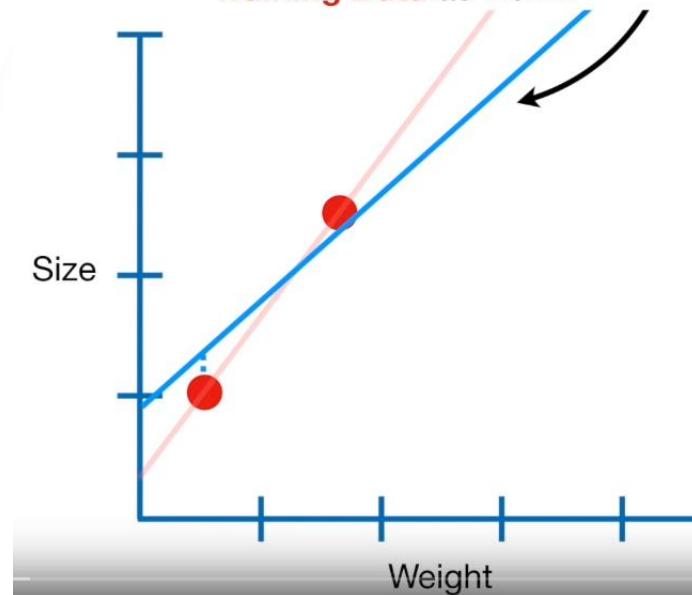
Now let's go back to just the **Training Data**...



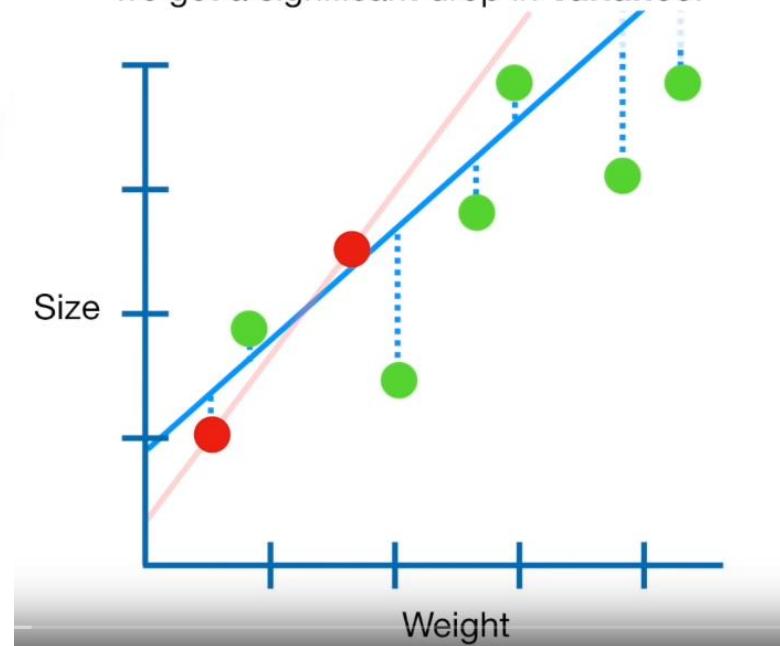
We just saw that **Least Squares** results in a  
**Line** that is **Over Fit** and has **High  
Variance**...



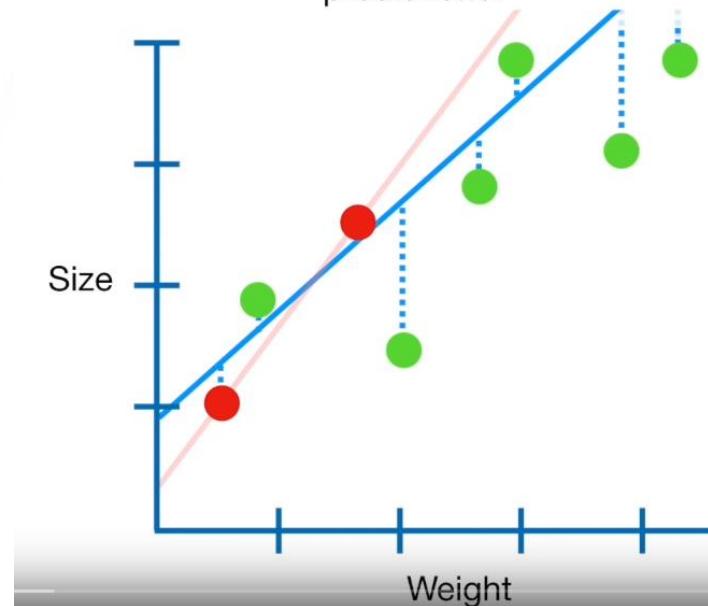
The main idea behind **Ridge Regression**  
is to find a **New Line** that doesn't fit the  
**Training Data** as well...



...but in return for that small amount of **Bias**,  
we get a significant drop in **Variance**.

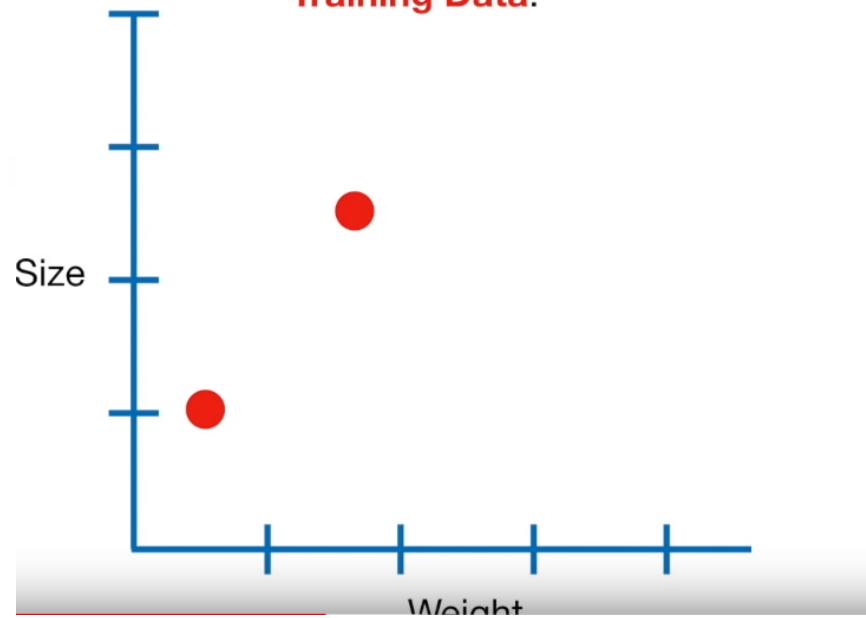


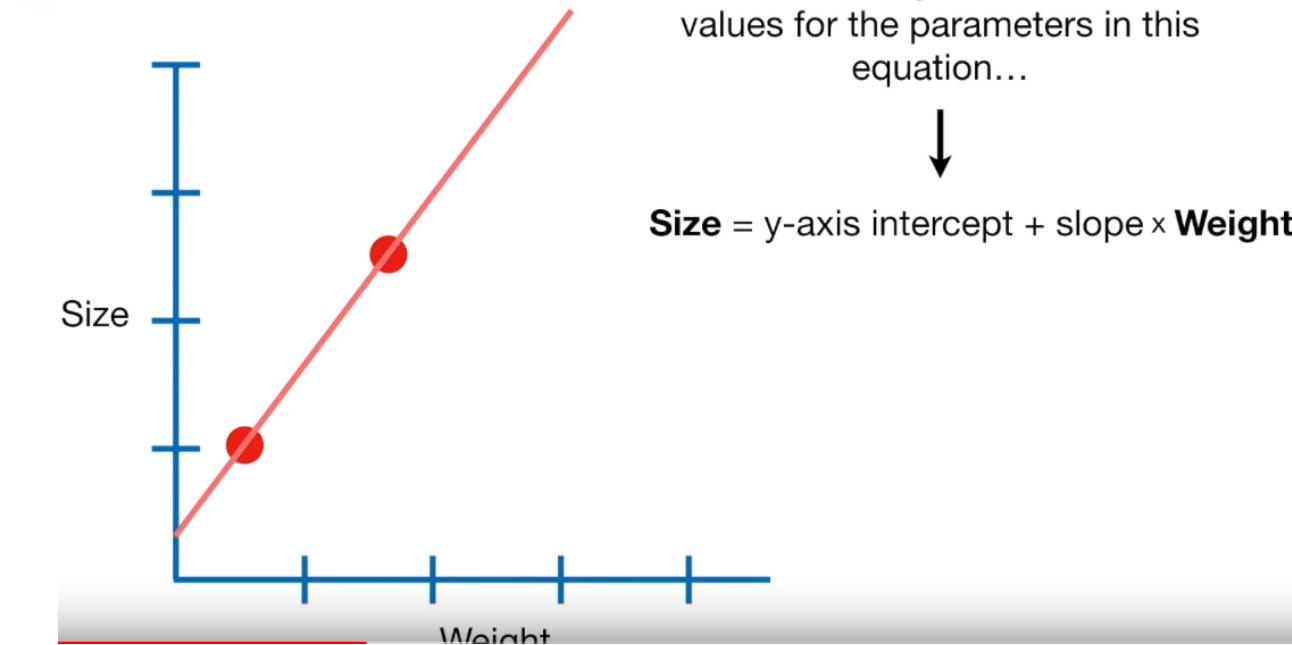
In other words, by starting with a slightly worse fit,  
**Ridge Regression** can provide better long term  
predictions.

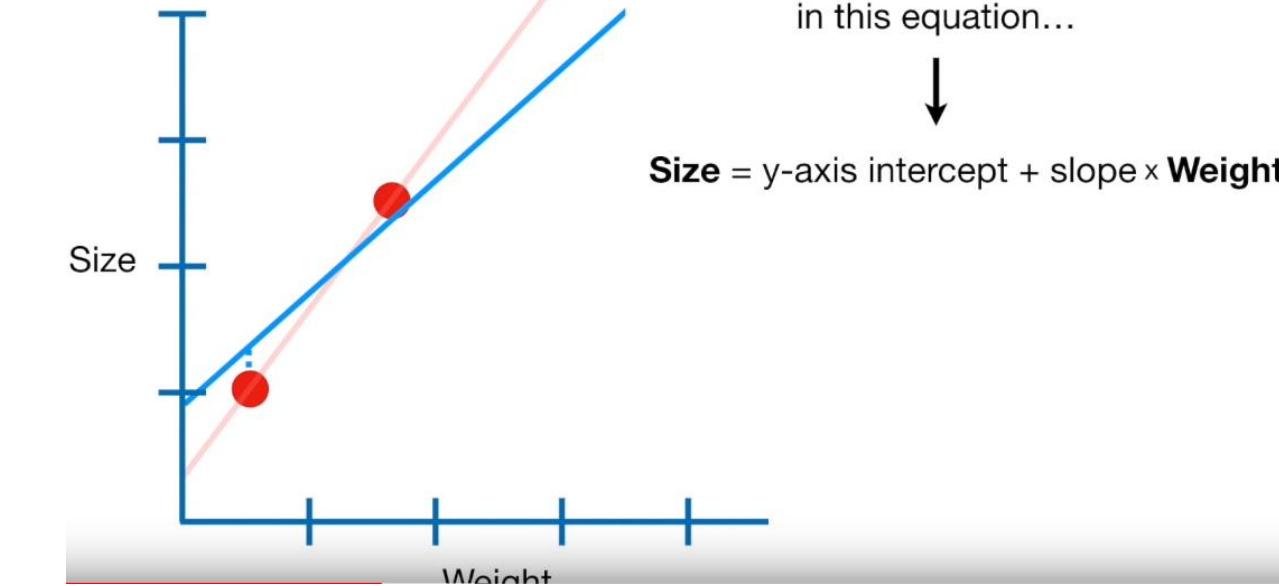


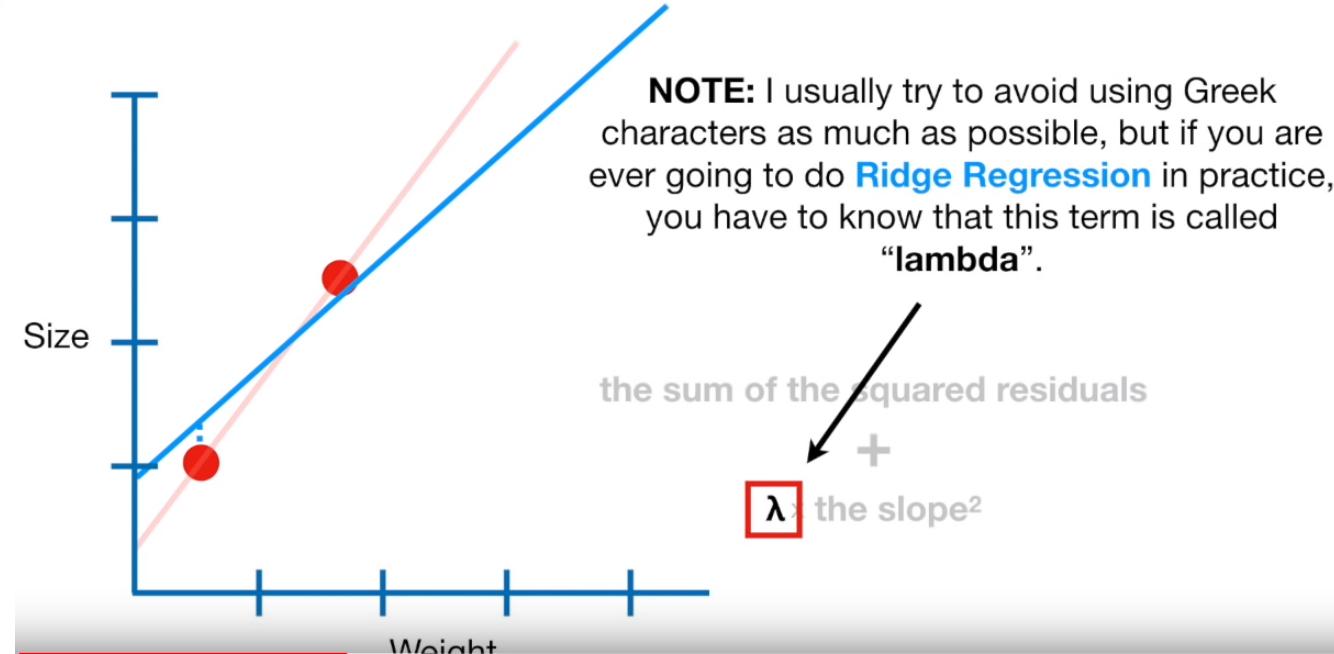


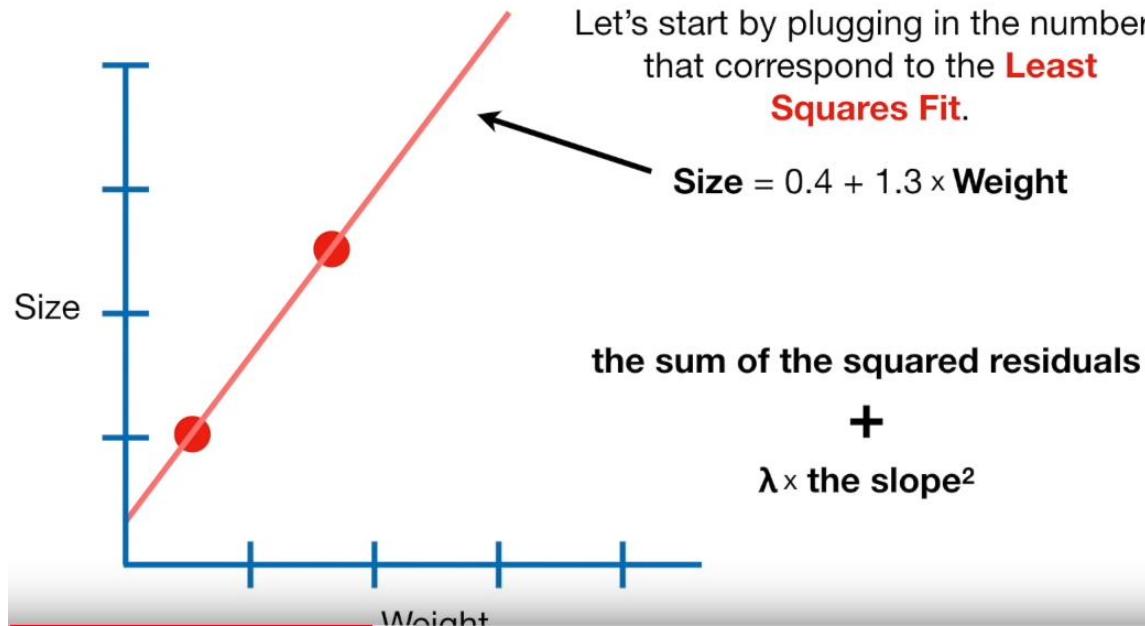
Let's go back to just the  
**Training Data.**

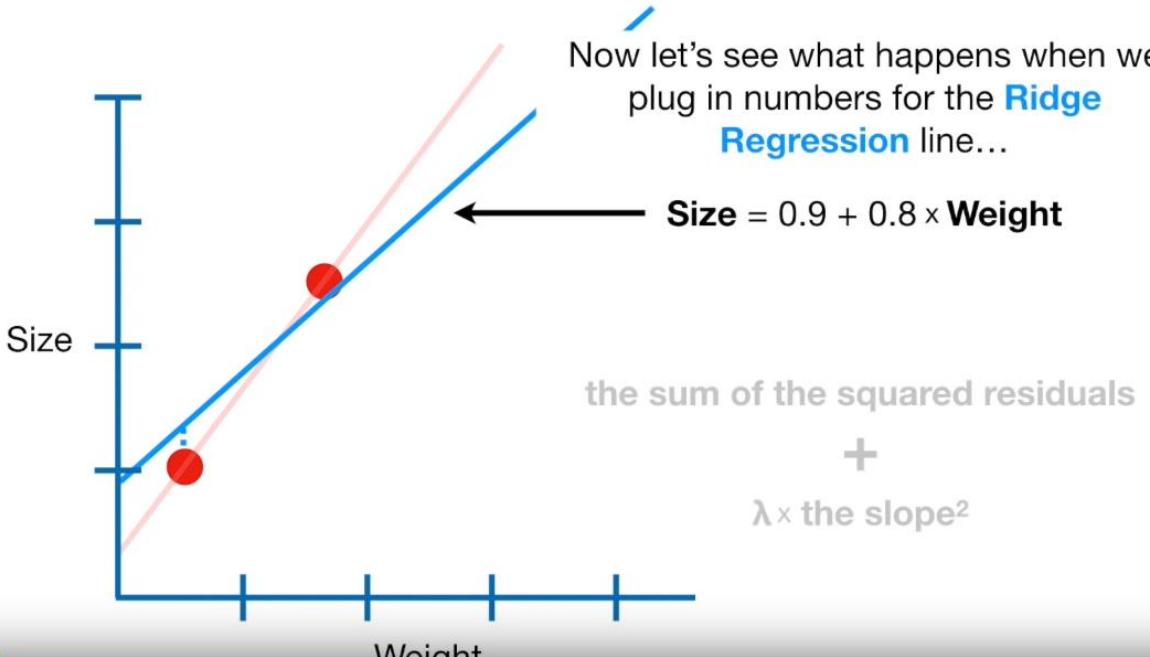


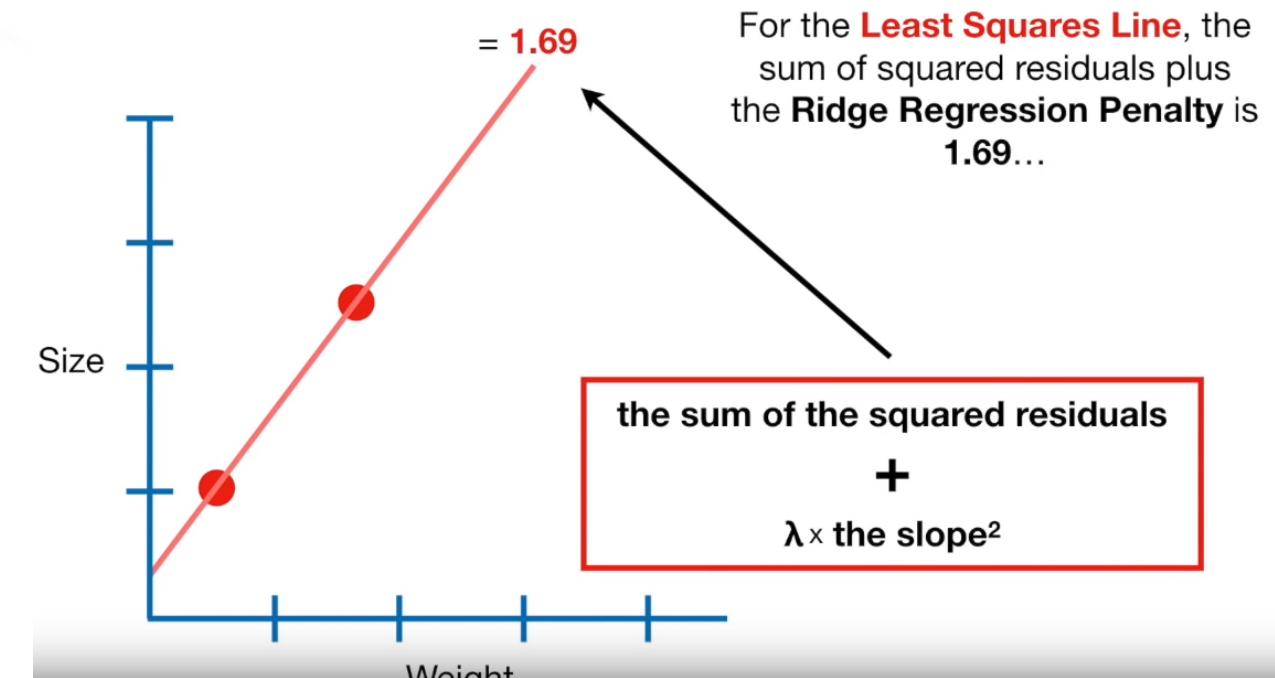


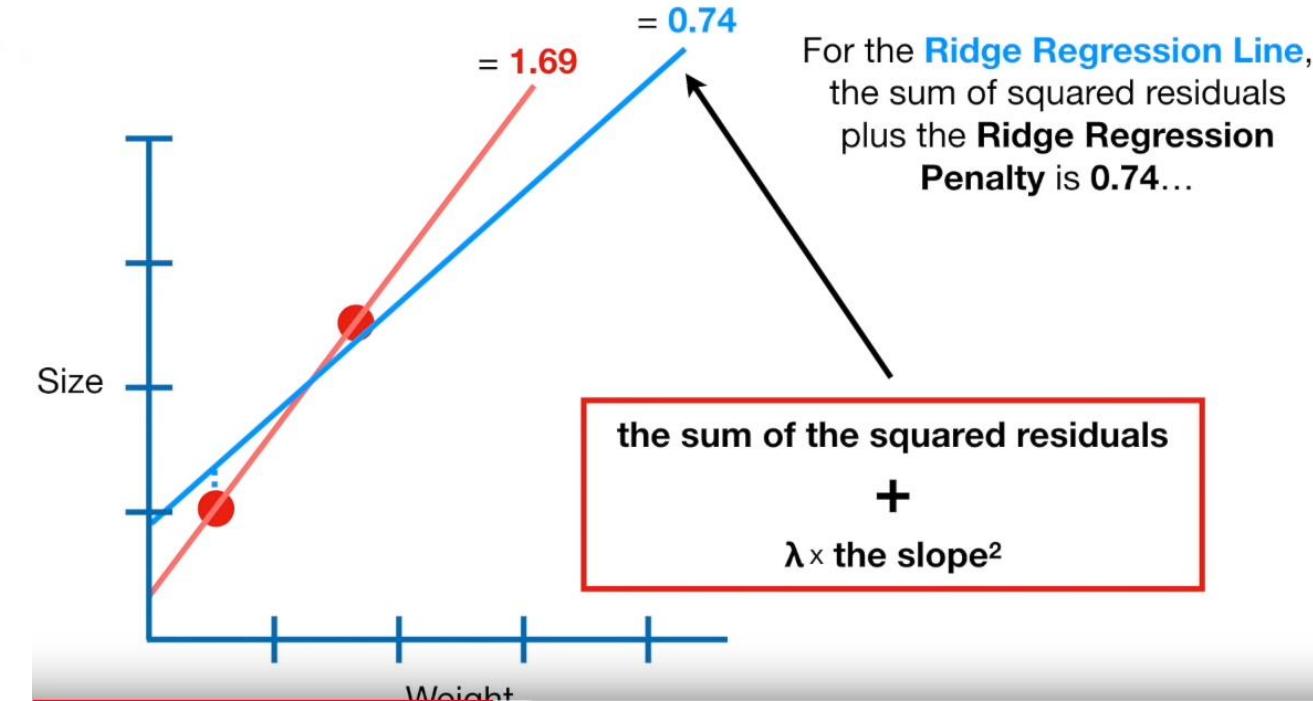


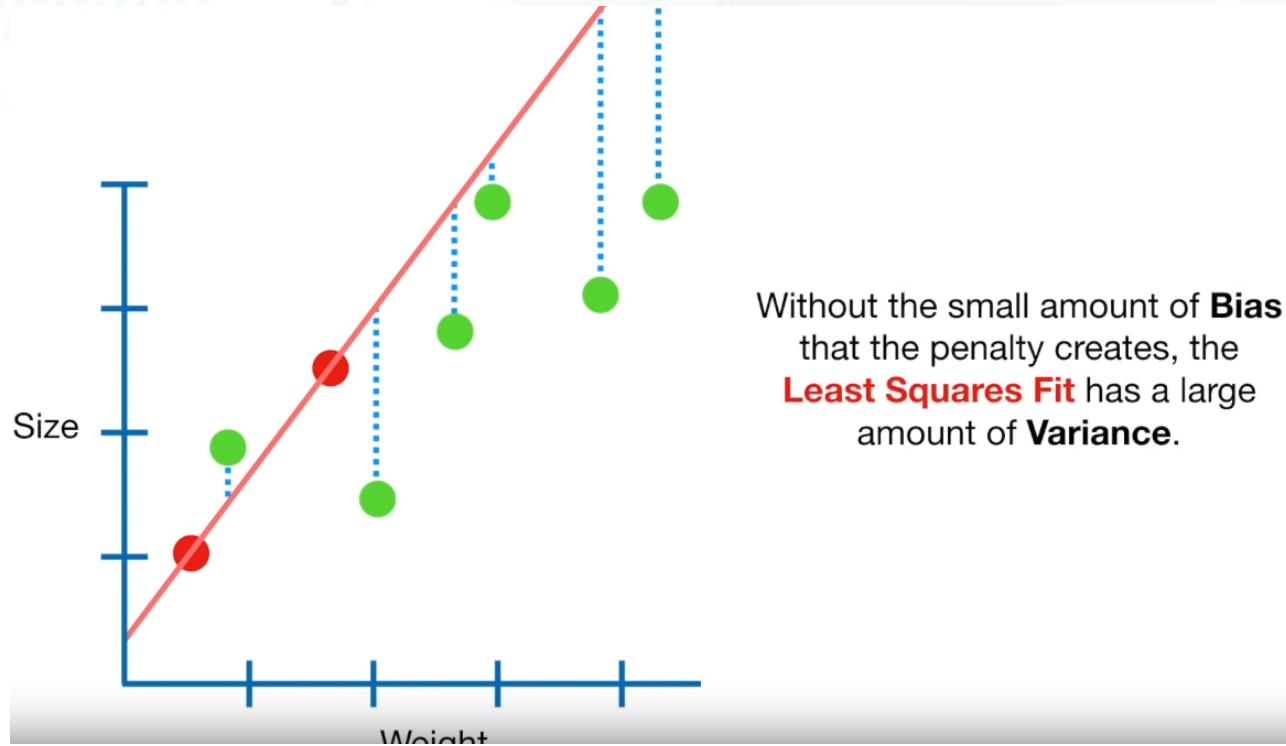




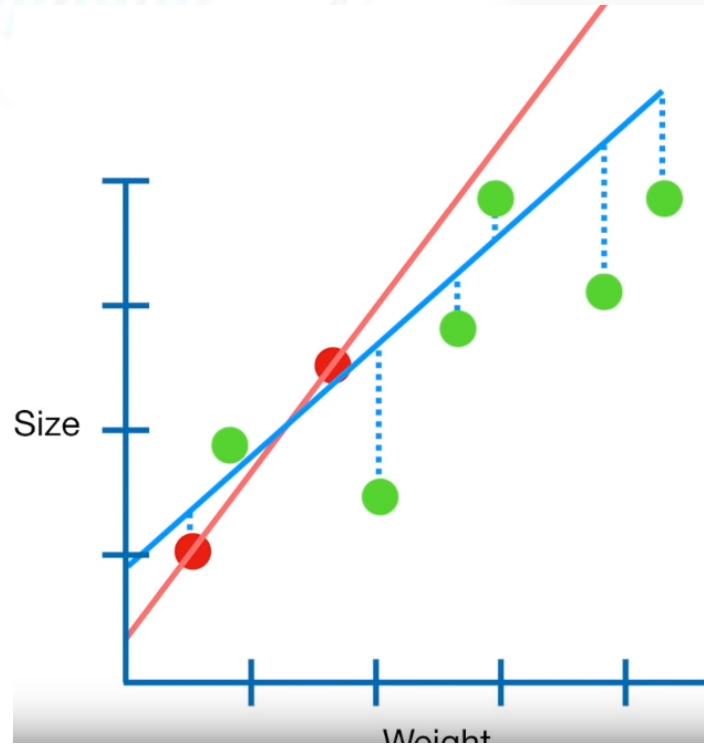




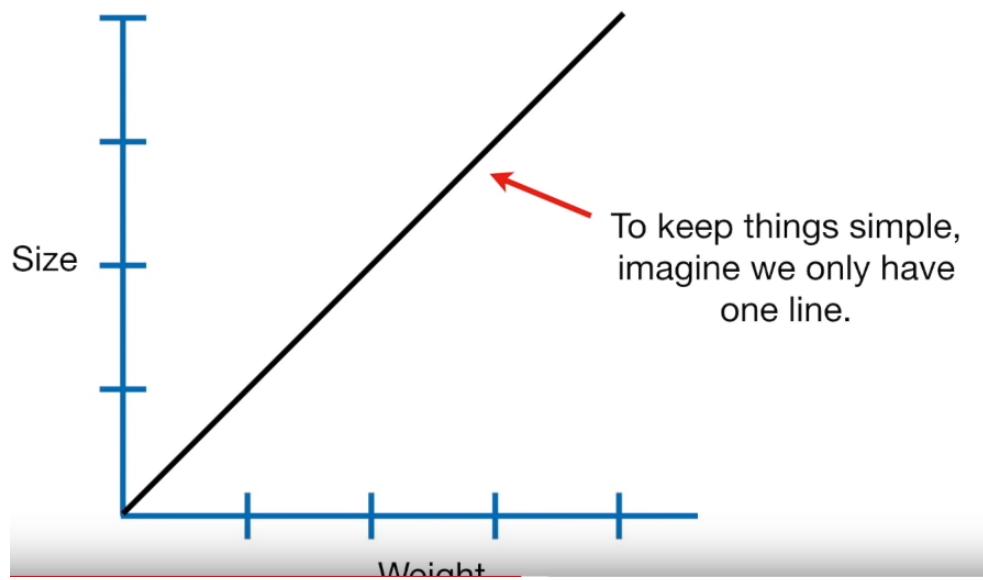


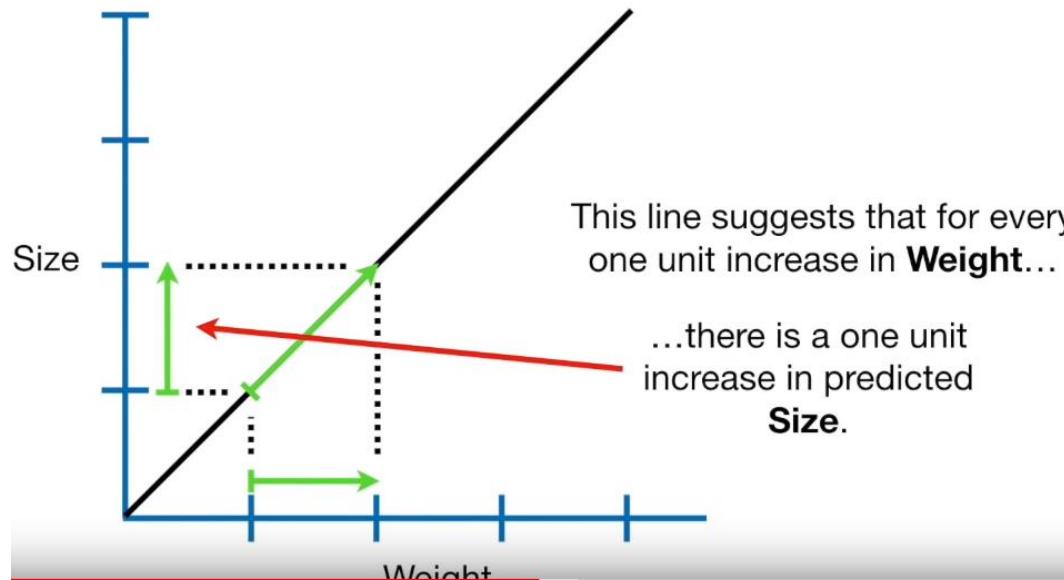


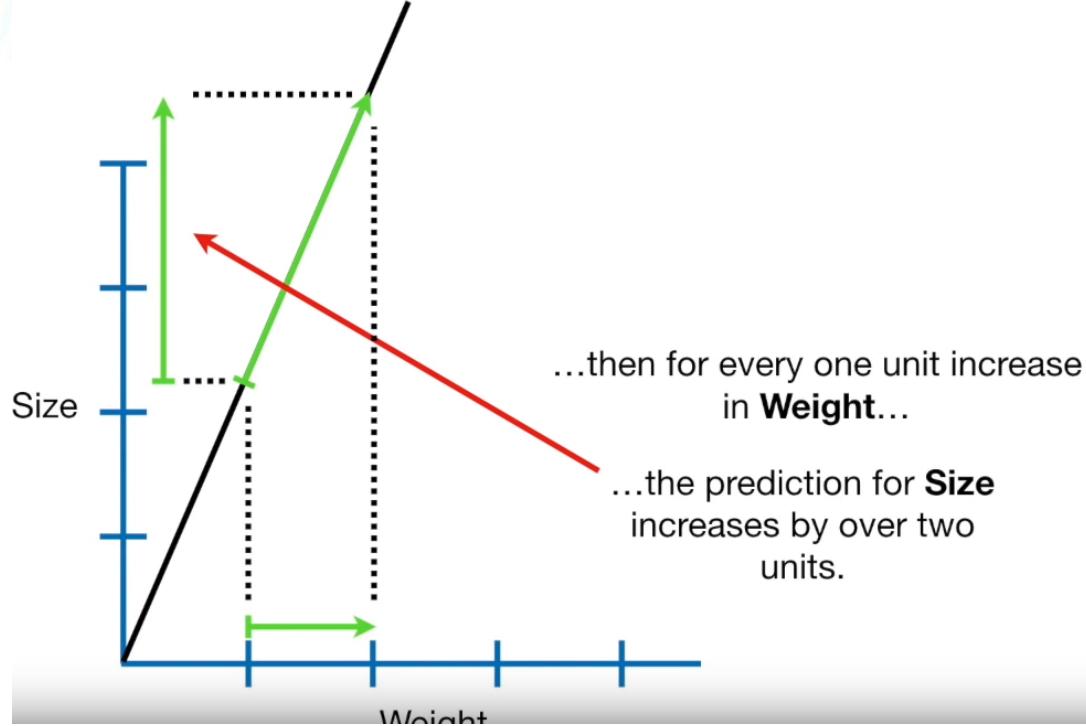
Without the small amount of **Bias** that the penalty creates, the **Least Squares Fit** has a large amount of **Variance**.

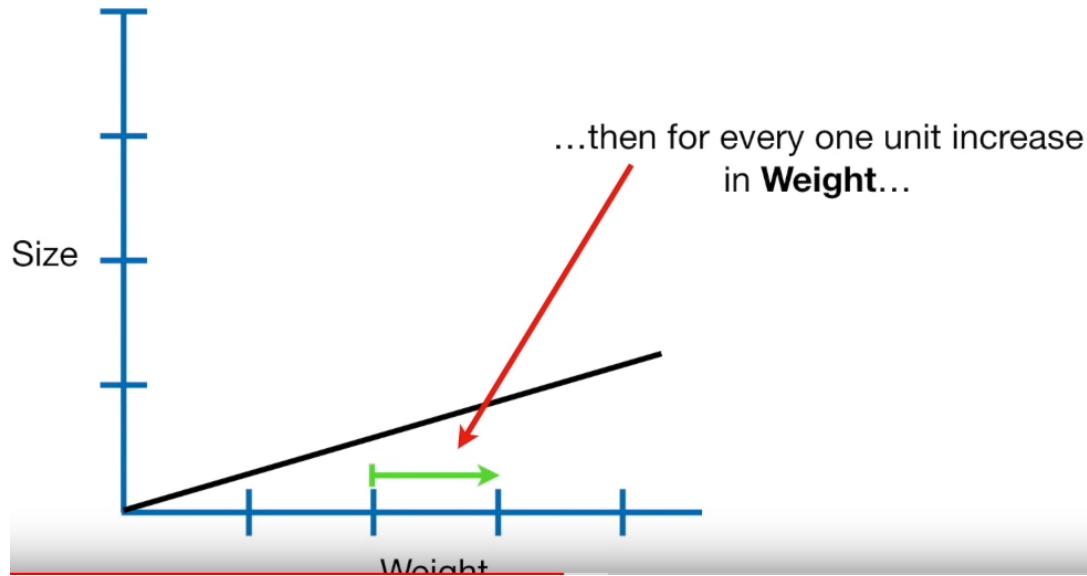


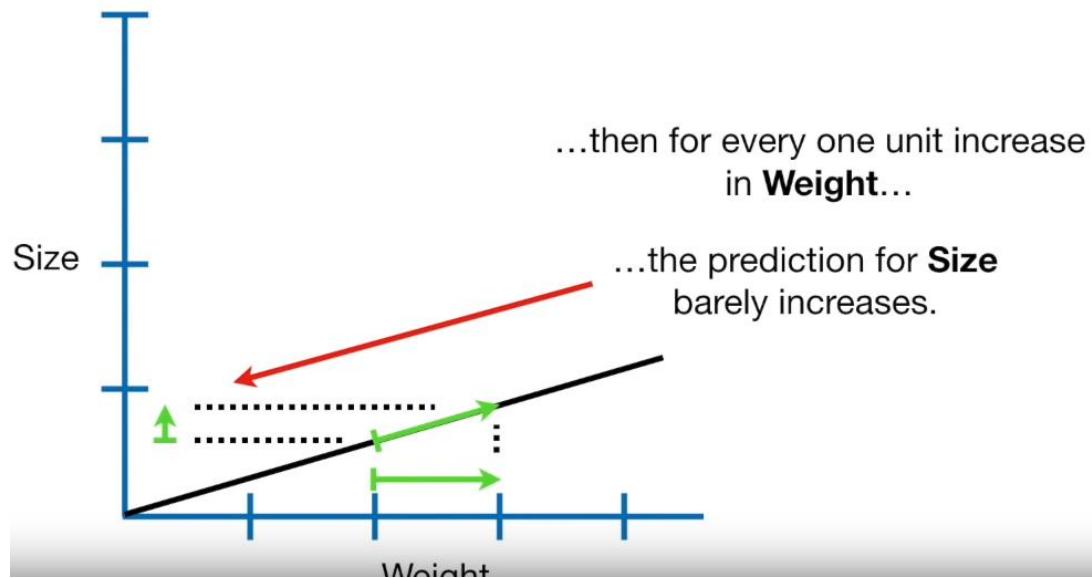
In contrast, the **Ridge Regression Line**, which has the small amount of **Bias** due to the penalty, has less **Variance**.

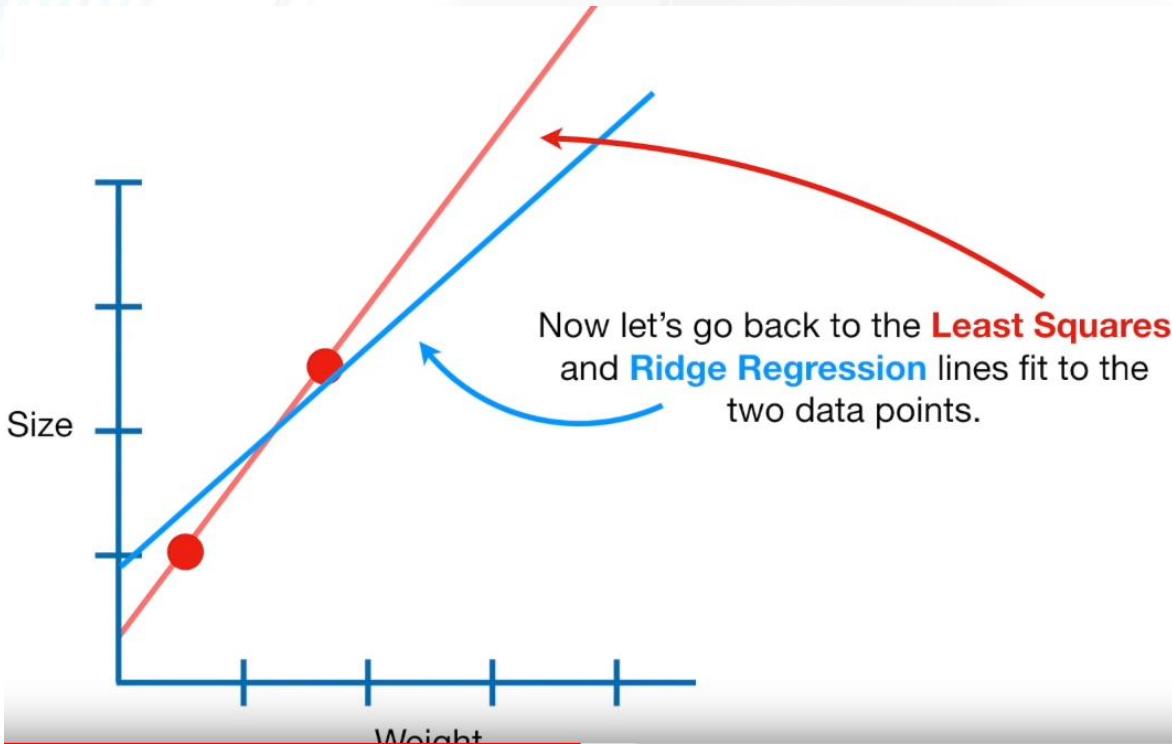


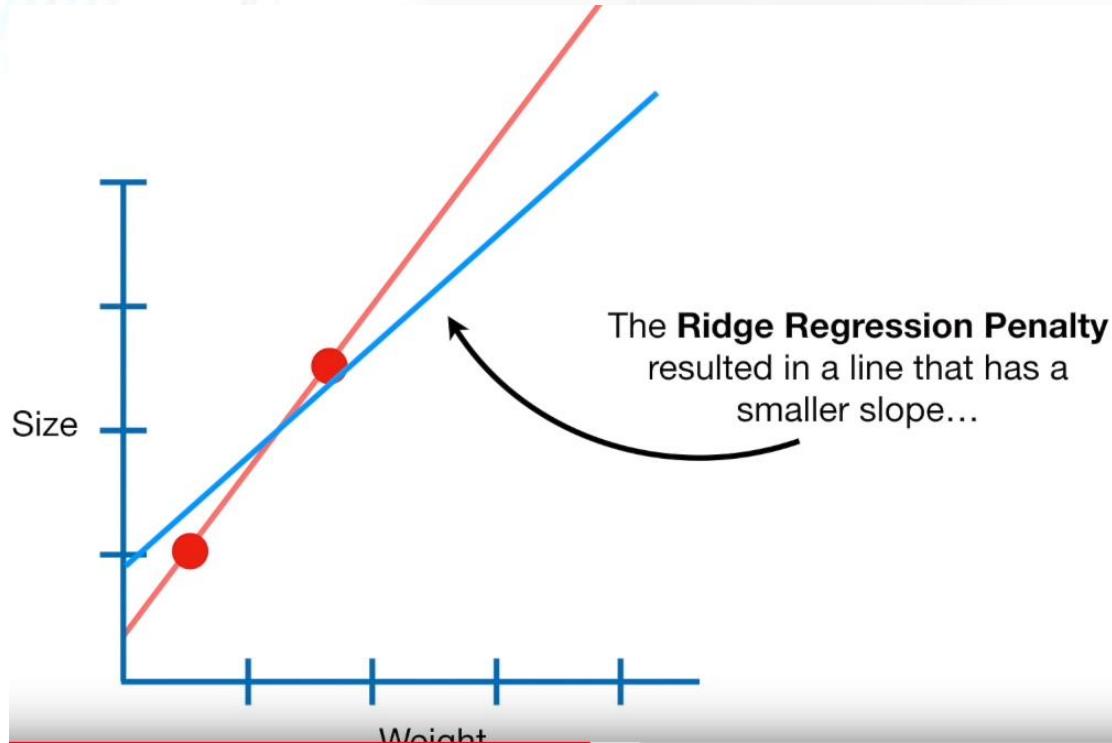


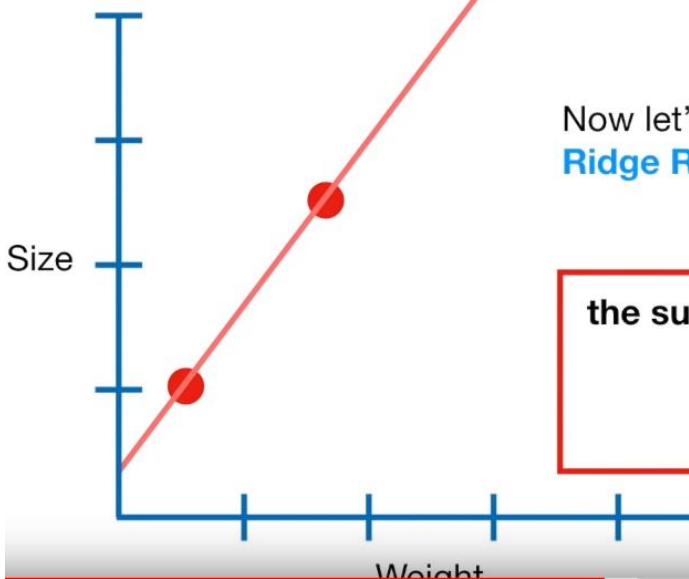






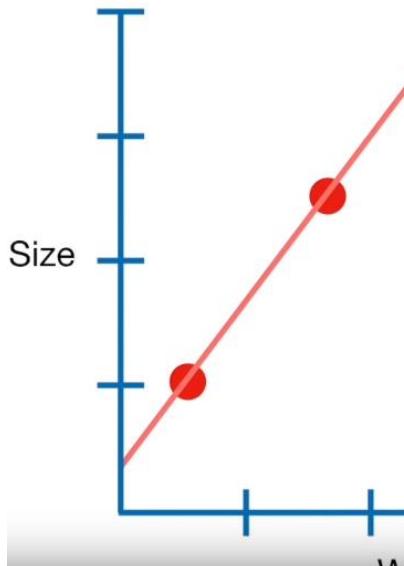






Now let's go back to the equation that  
**Ridge Regression** tries to minimize...

the sum of the squared residuals  
+  
 $\lambda \times \text{the slope}^2$



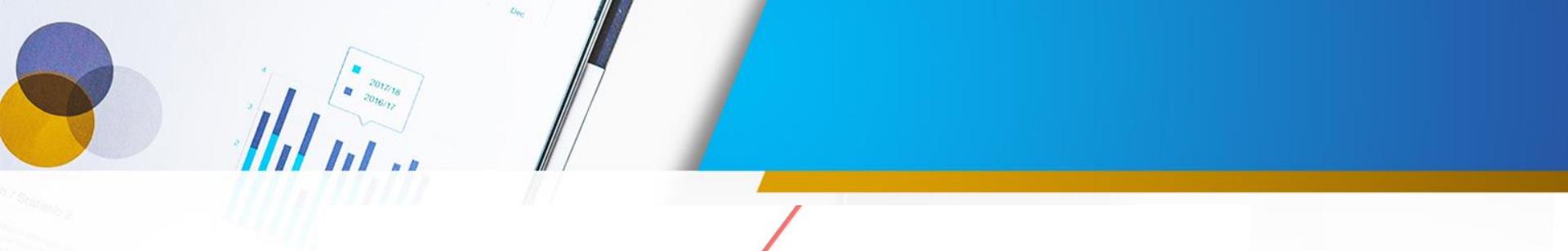
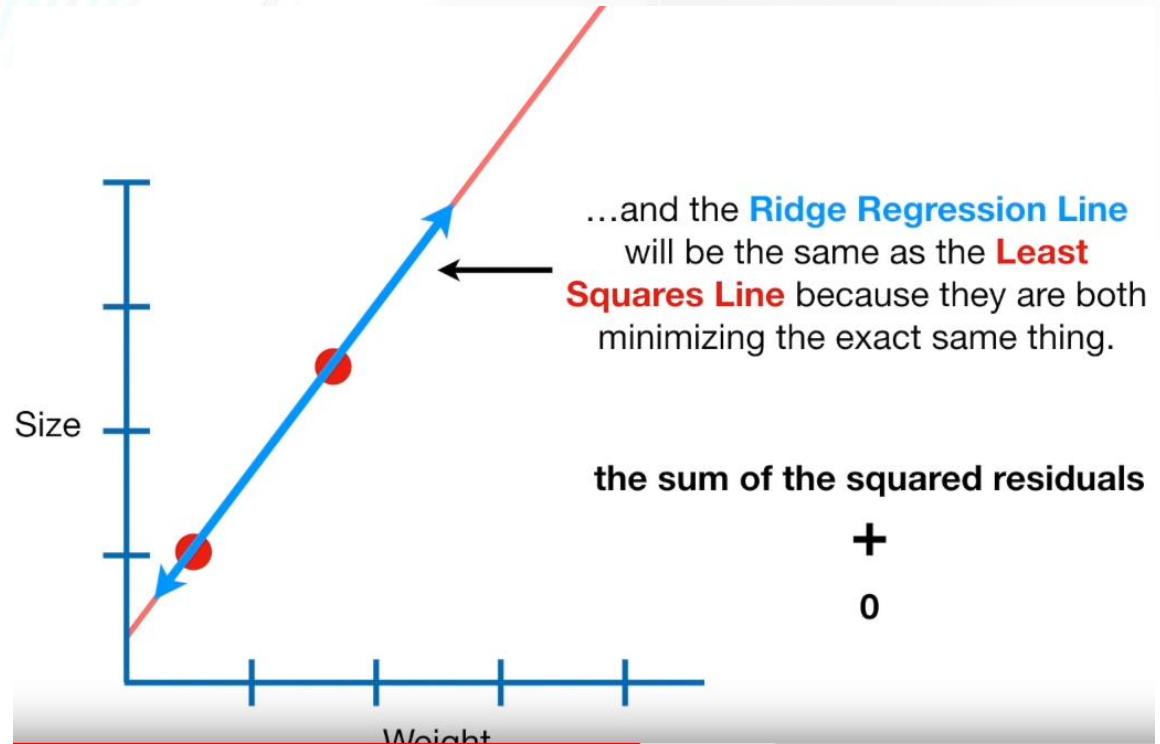
...and that means that the **Ridge Regression Line** will only minimize the sum of squared residuals...

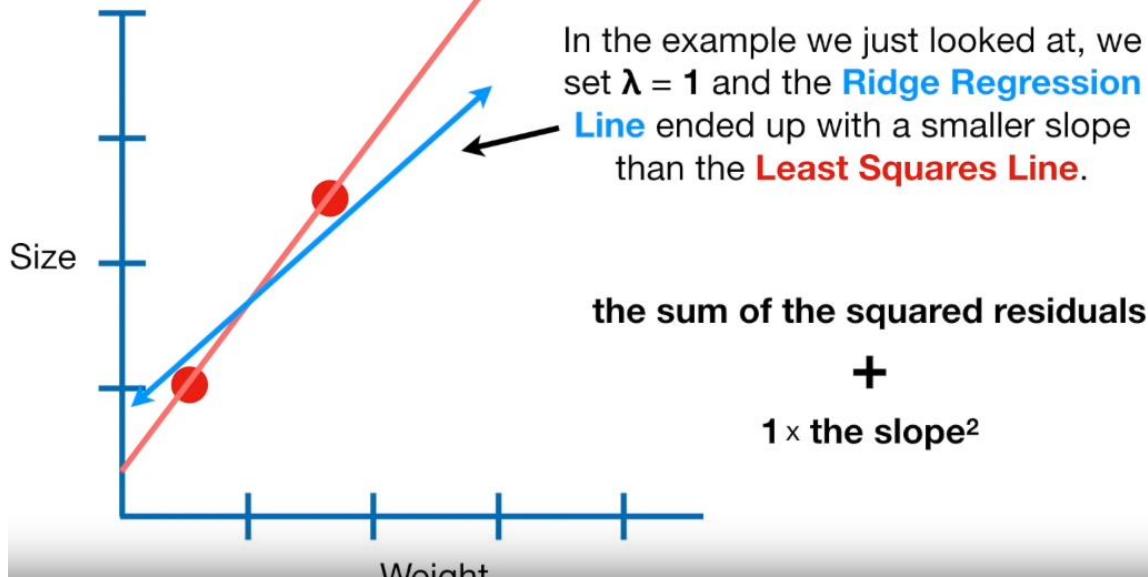


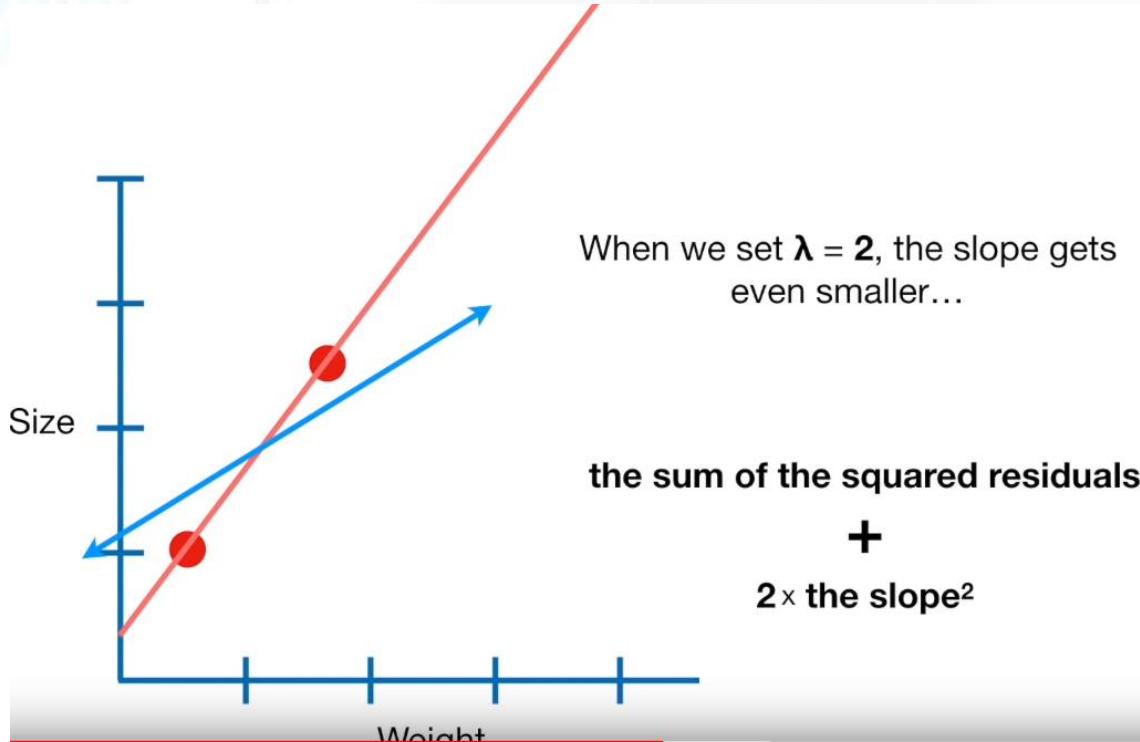
**the sum of the squared residuals**

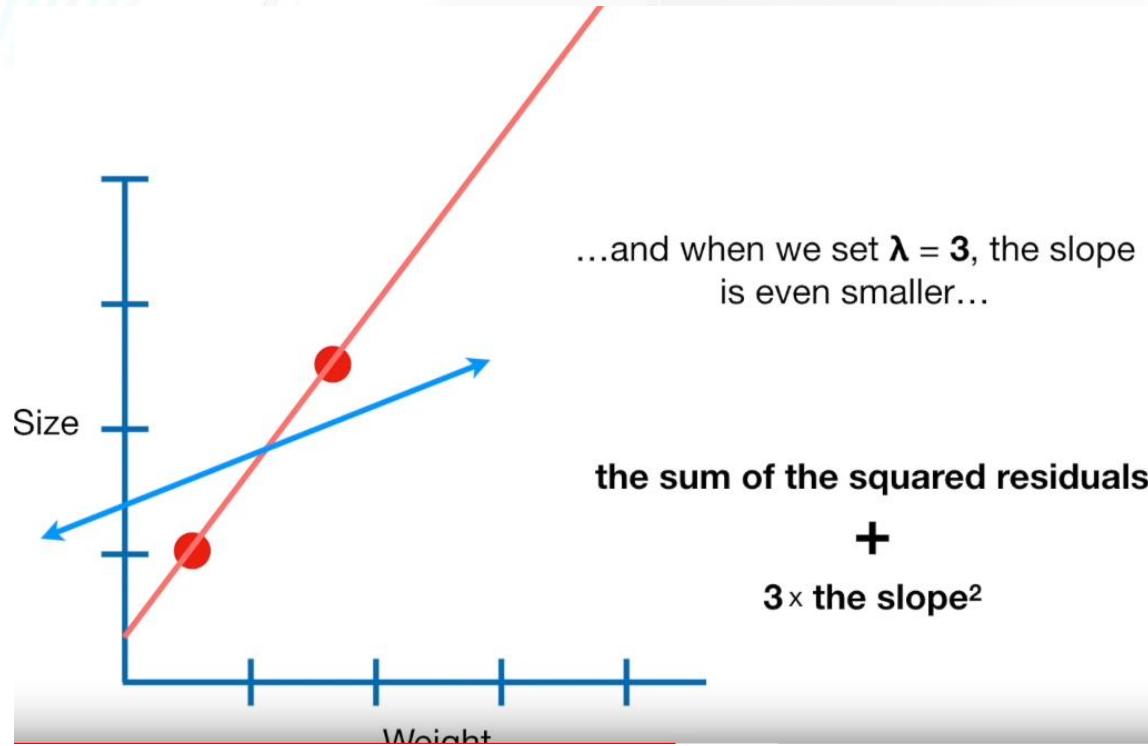
+

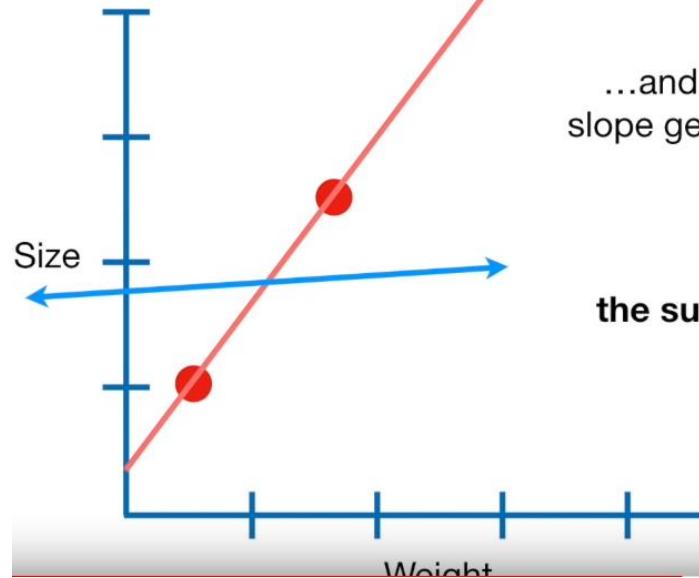
0









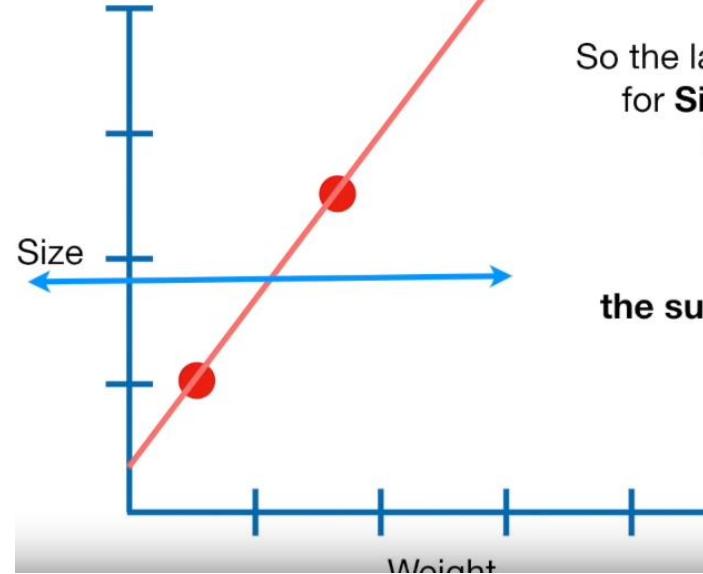


...and the larger we make  $\lambda$ , the slope gets asymptotically close to 0.

the sum of the squared residuals

+

$100 \times \text{the slope}^2$

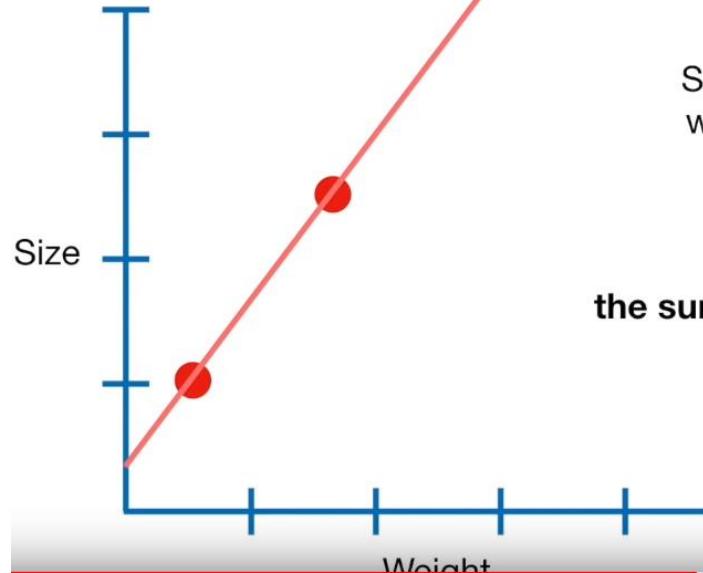


So the larger  $\lambda$  gets, our predictions for **Size** become less and less sensitive to **Weight**.

the sum of the squared residuals

+

100000 × the slope<sup>2</sup>

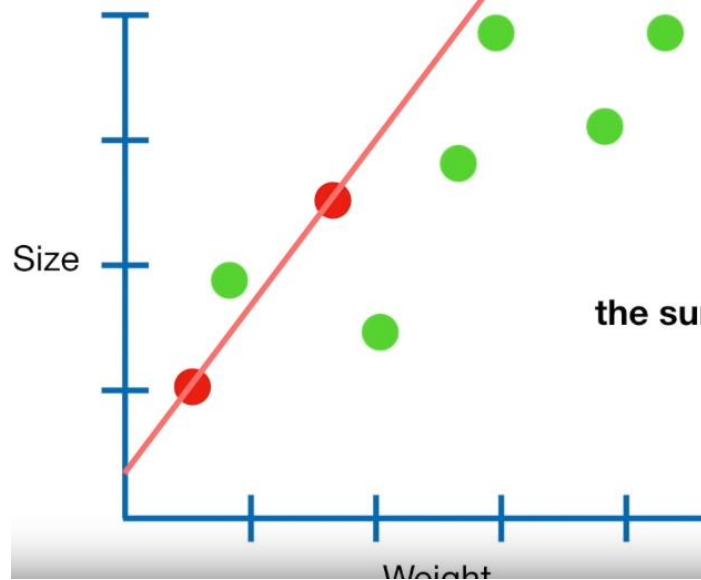


So how do we decide  
what value to give  $\lambda$ ?

**the sum of the squared residuals**

+

$\lambda \times \text{the slope}^2$



We just try a bunch of values for  $\lambda$  and use **Cross Validation**, typically **10-fold Cross Validation**, to determine which one results in the lowest **Variance**.

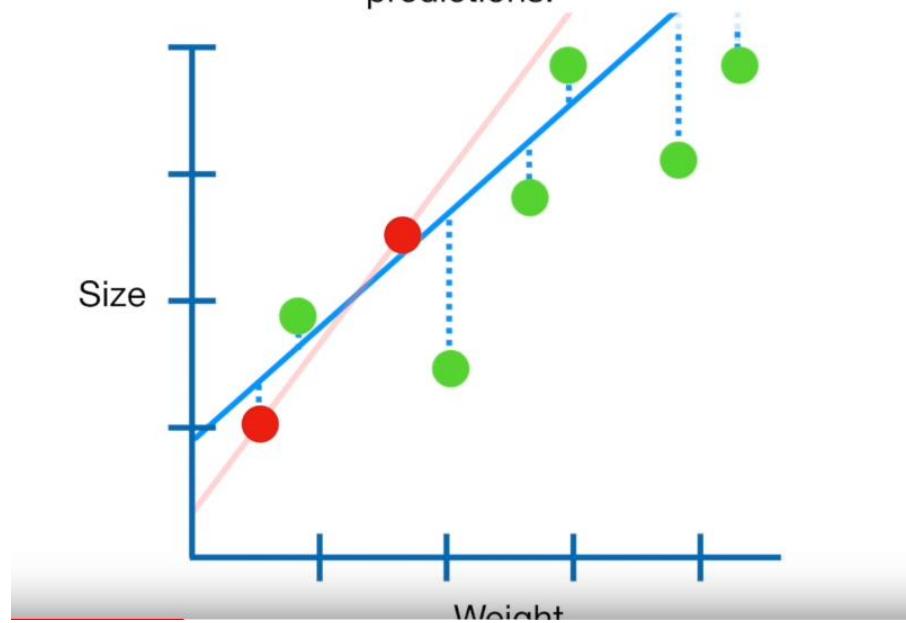
**the sum of the squared residuals**

+

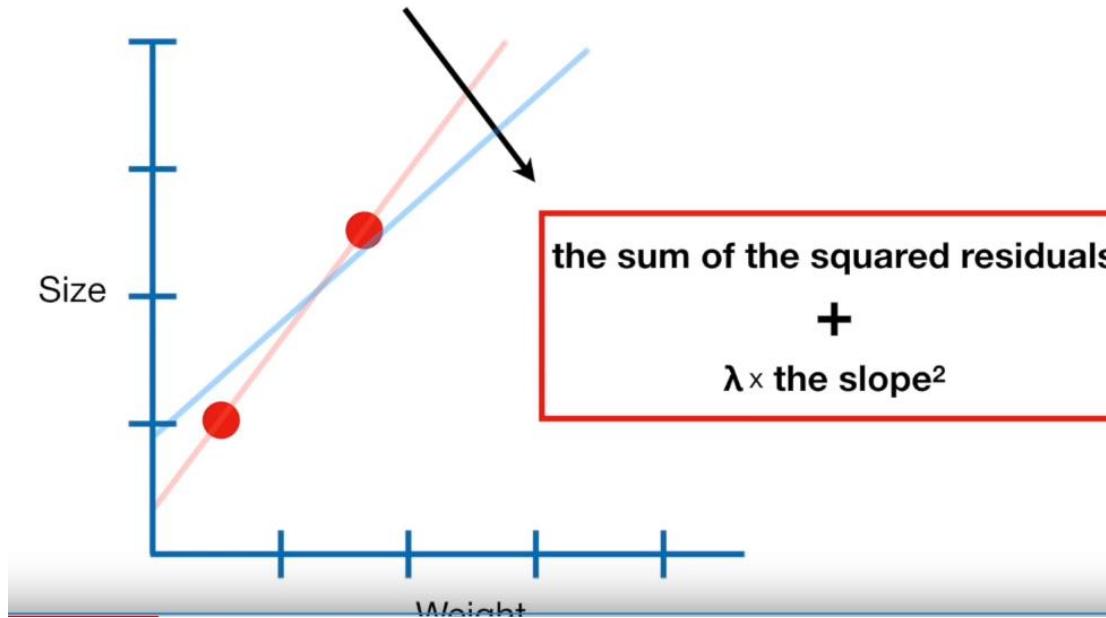
$\lambda \times \text{the slope}^2$

# Lasso Regression

The main idea was that by starting with a slightly worse fit, **Ridge Regression** provided better long term predictions.

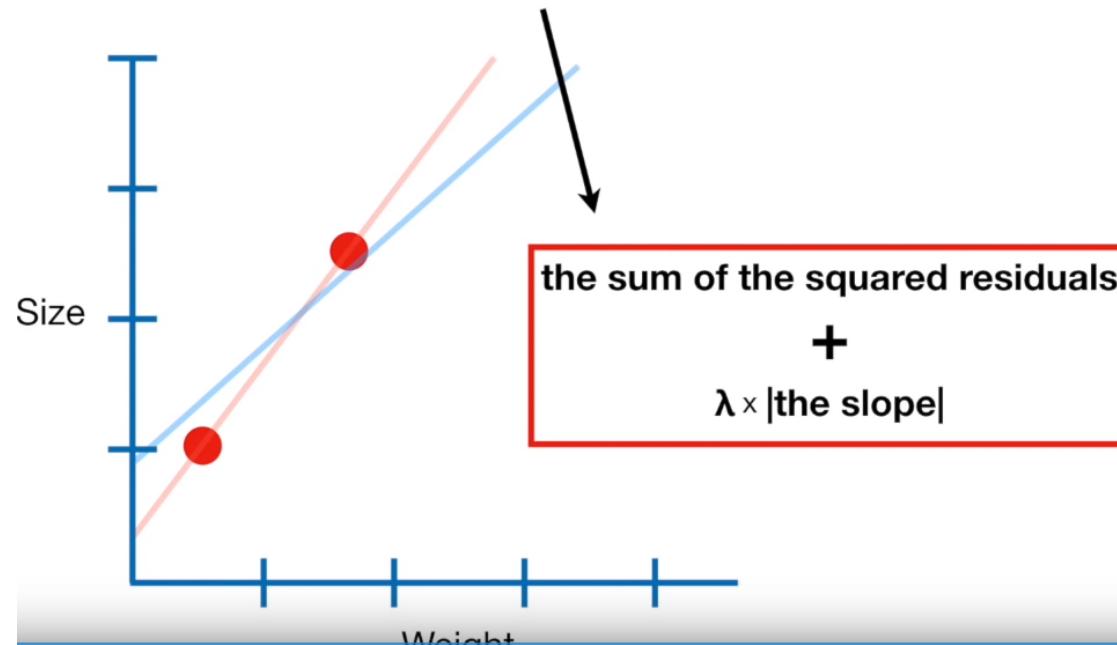


Now let's go back to the equation that  
**Ridge Regression** minimizes...

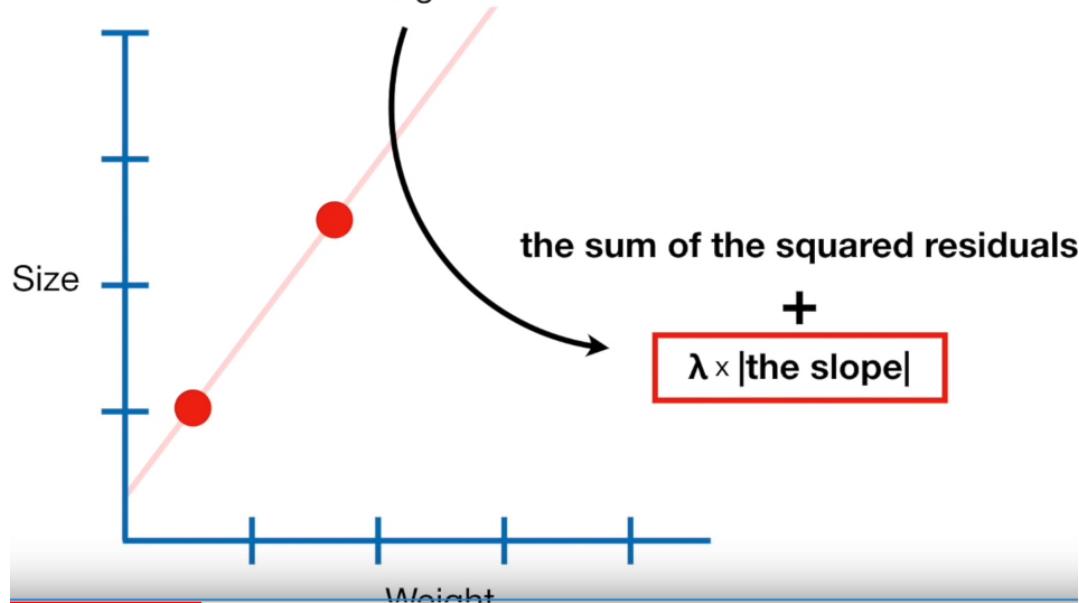


# Lasso Regression

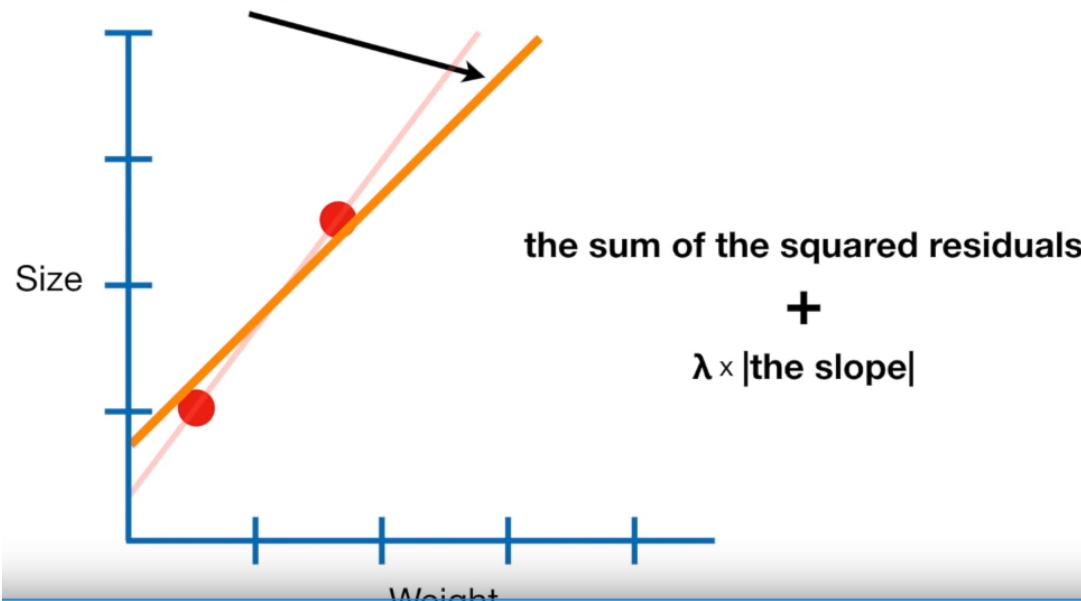
...then we have **Lasso Regression!!!**



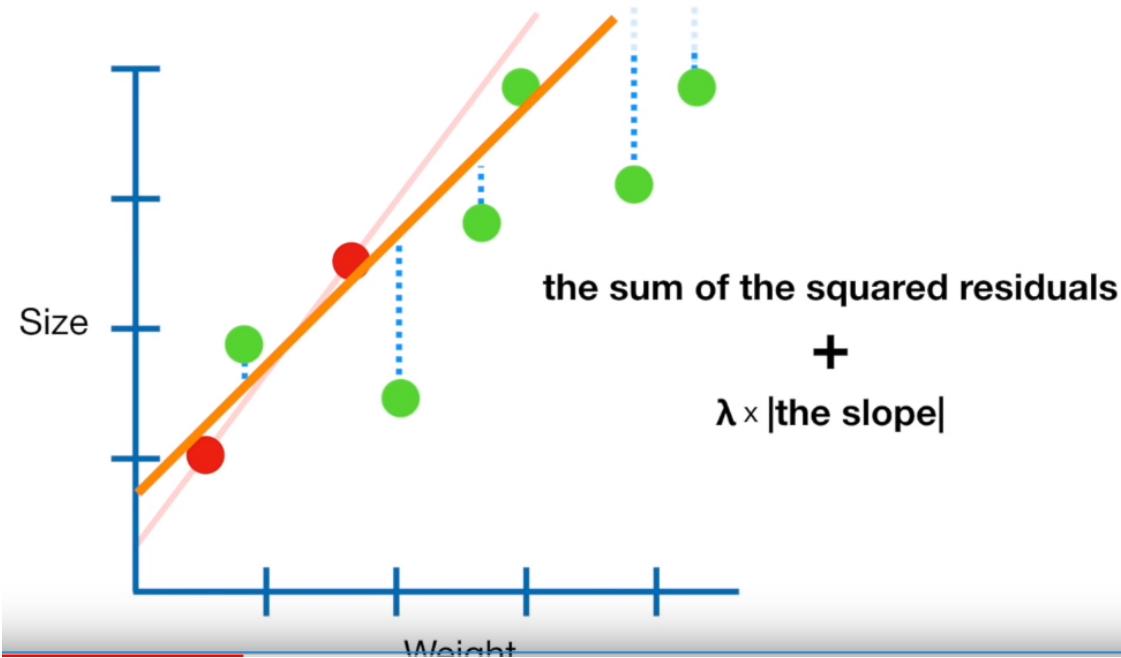
**NOTE:** Just like with **Ridge Regression**,  $\lambda$  can be any value from **0** to **positive infinity** and is determined using **Cross Validation**.



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



...but less **Variance** than **Least Squares**.



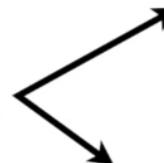


**the sum of the squared residuals**

**+**

**$\lambda \times \text{the slope}^2$**

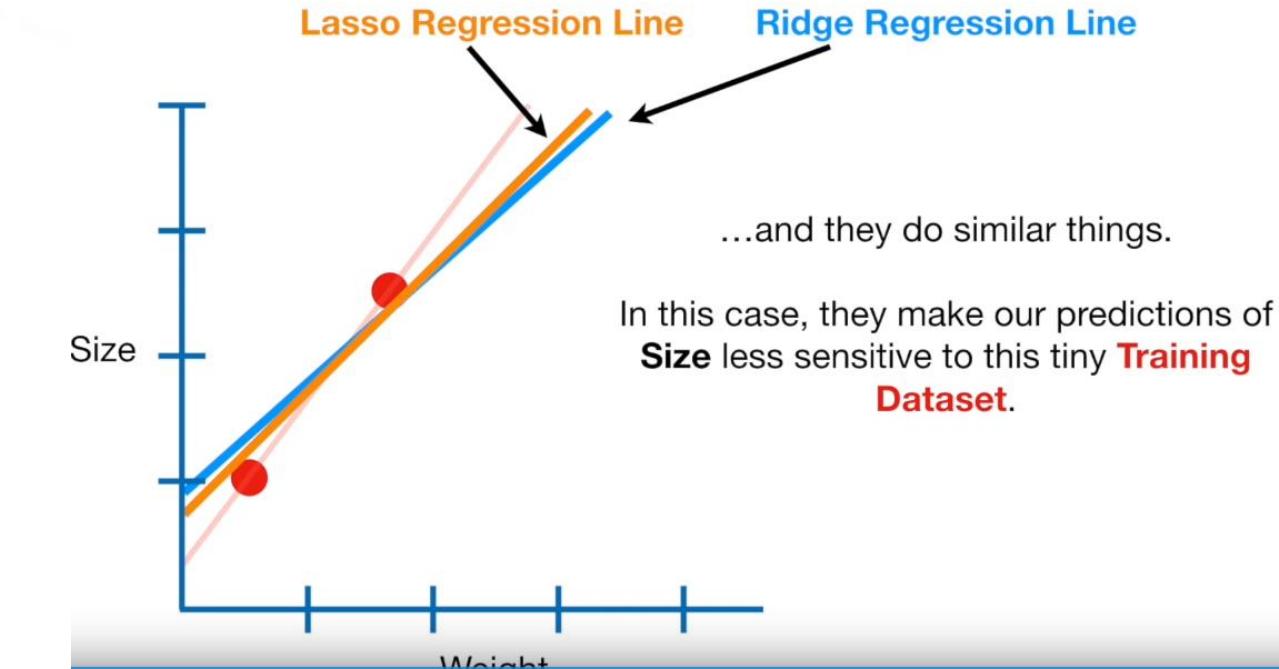
...look very similar....

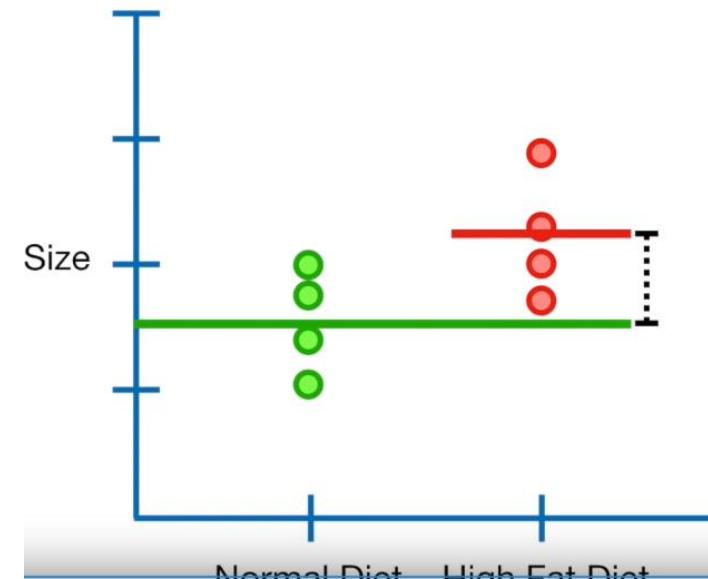


**the sum of the squared residuals**

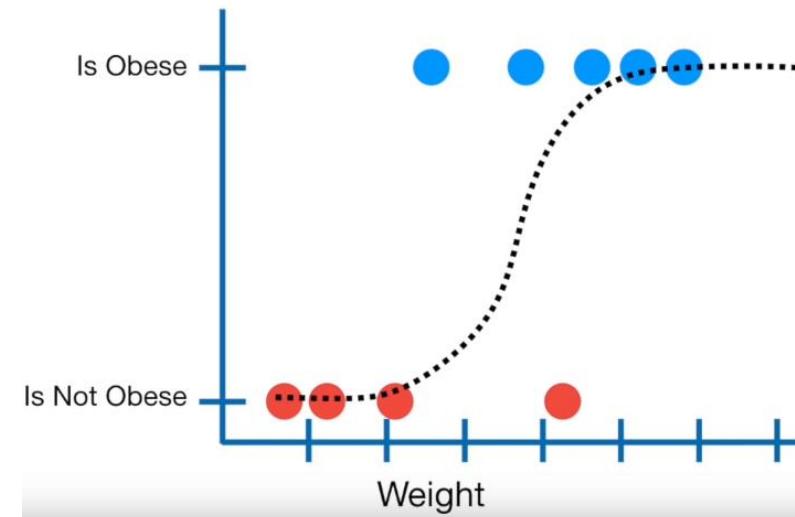
**+**

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

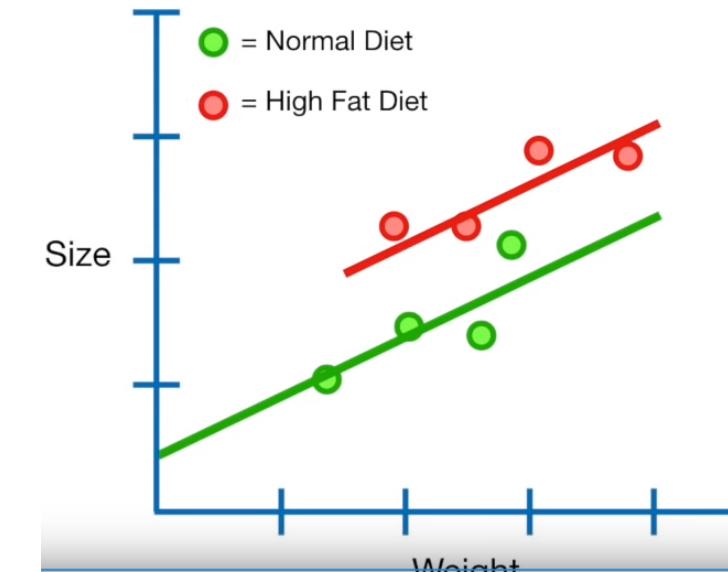




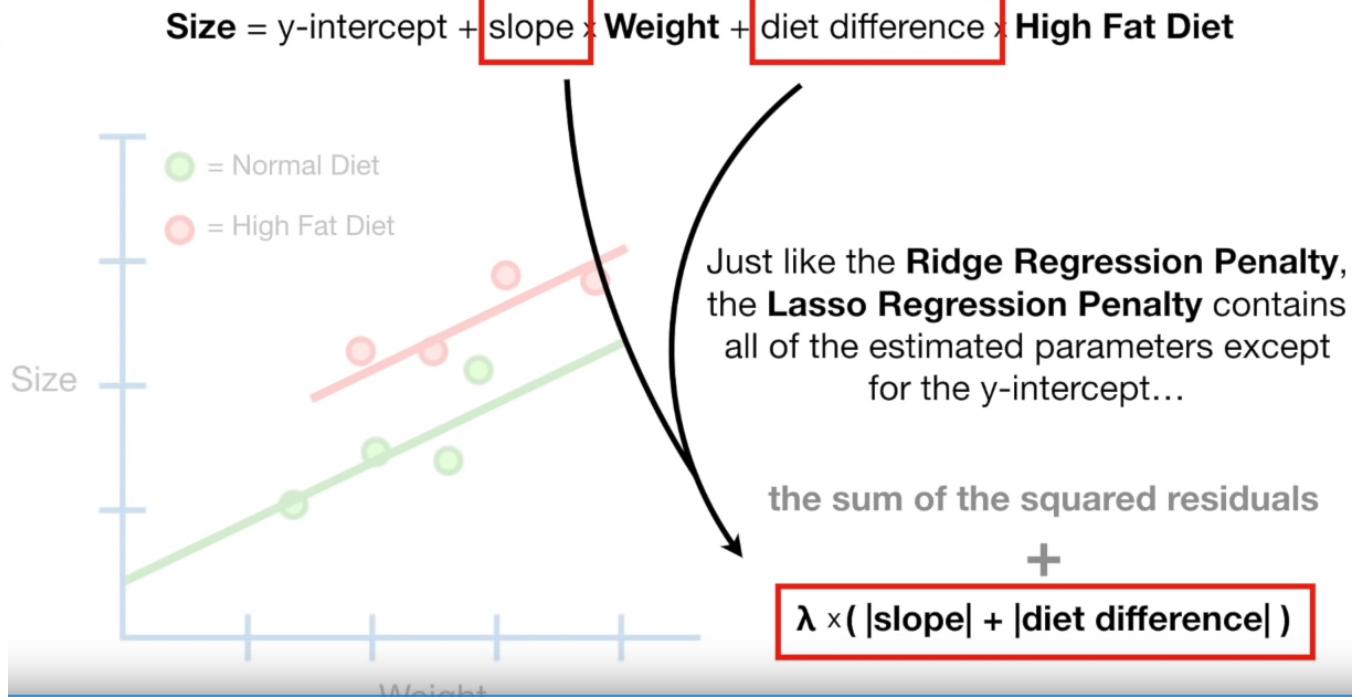
Both **Ridge** and **Lasso Regression** can be applied in the same contexts, like this situation, where we are using two **Different Diets** to predict **Size**...

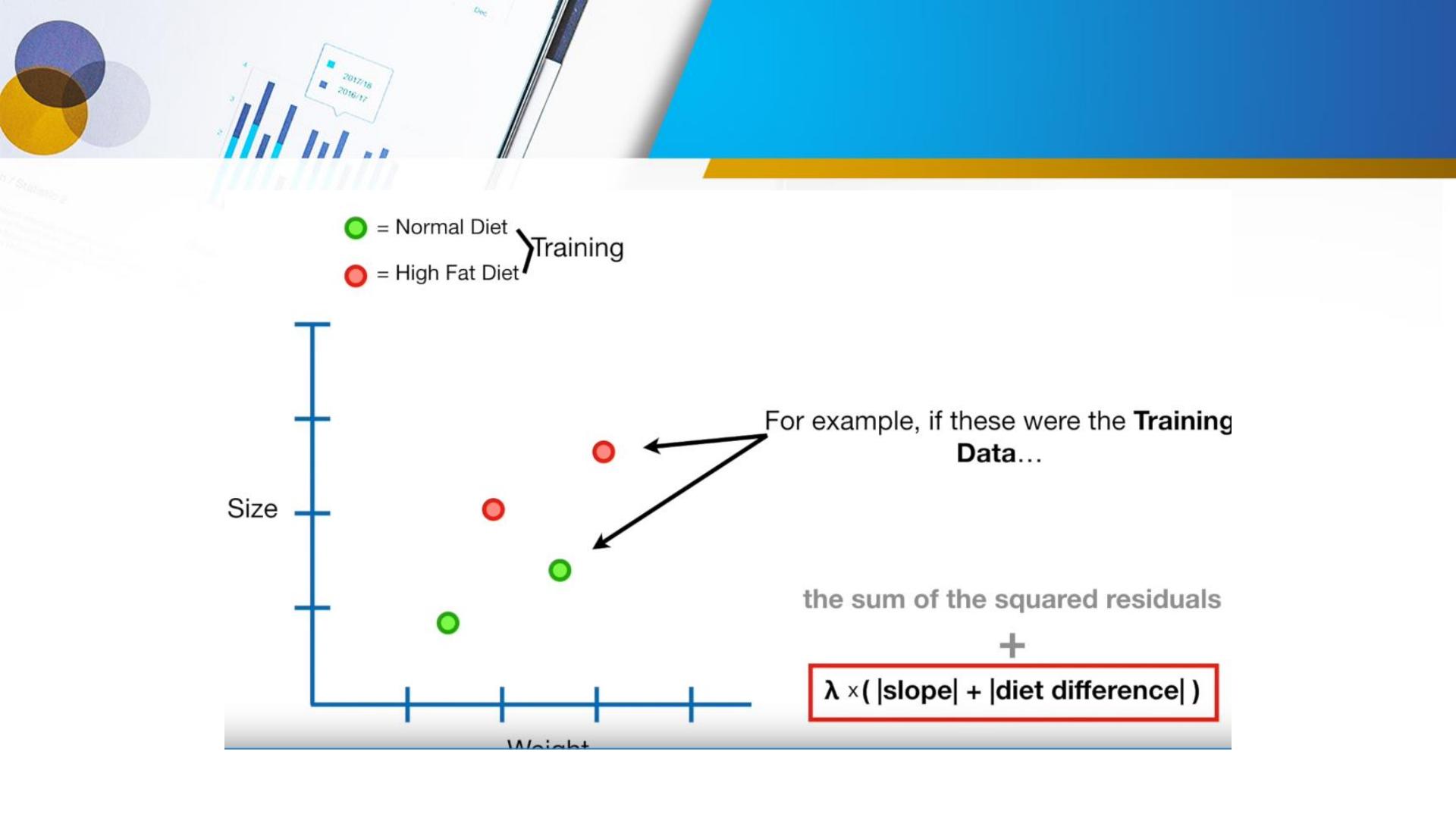


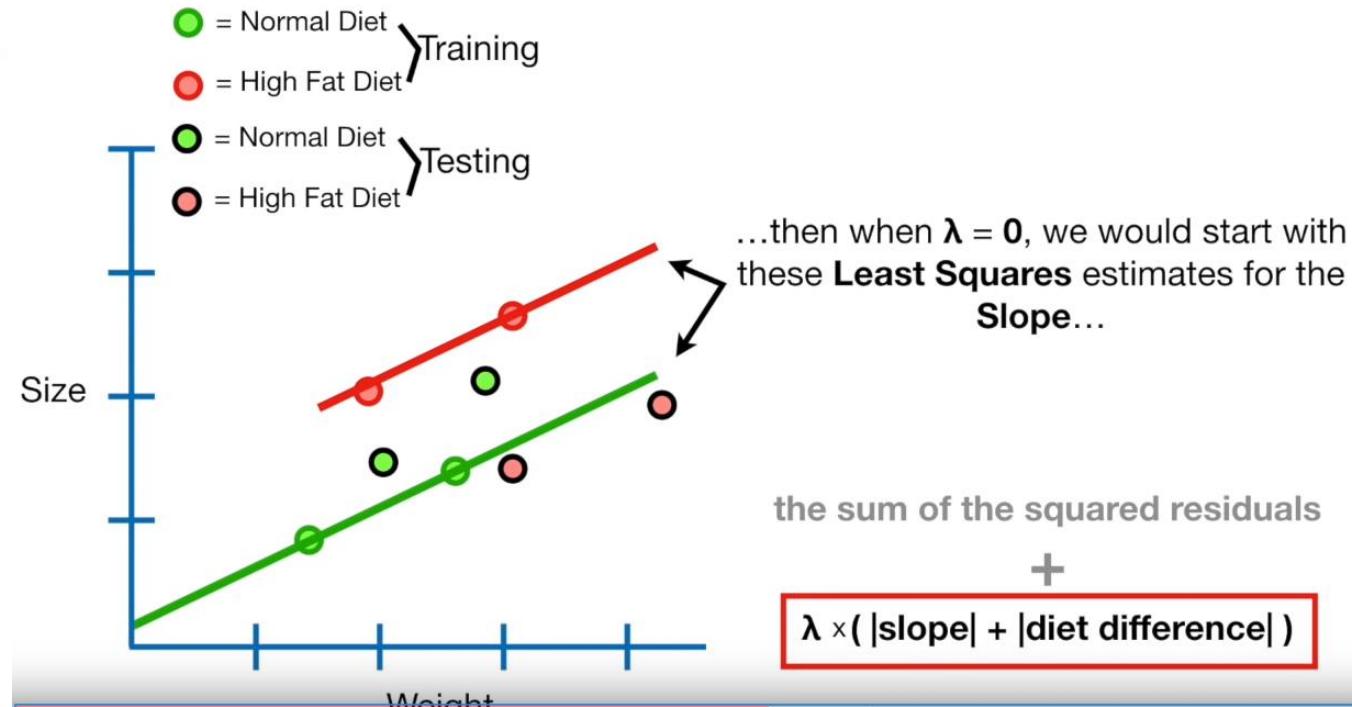
...or in a Logistic Regression  
setting where we use **Weight** to  
predict **Obesity**...

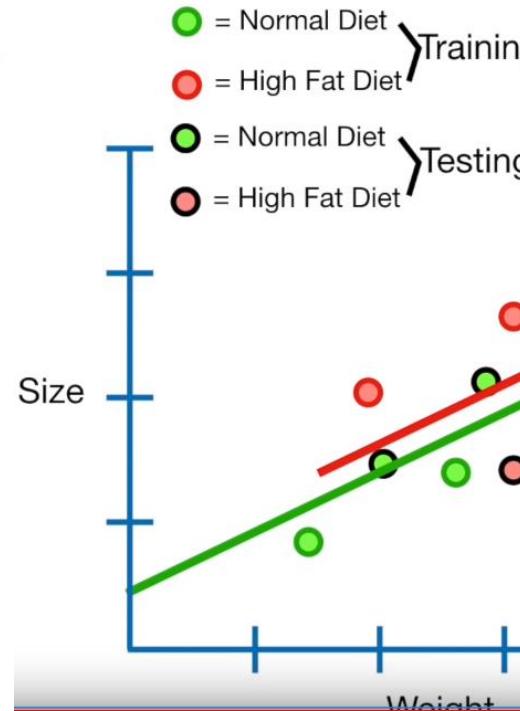


...and both **Ridge** and **Lasso** Regression can be applied to complicated models that combine different types of data.









...but as we increase the value for  $\lambda$ ,  
**Ridge and Lasso Regression** may  
shrink **Diet Difference** a lot more than  
they shrink the **Slope**.

the sum of the squared residuals

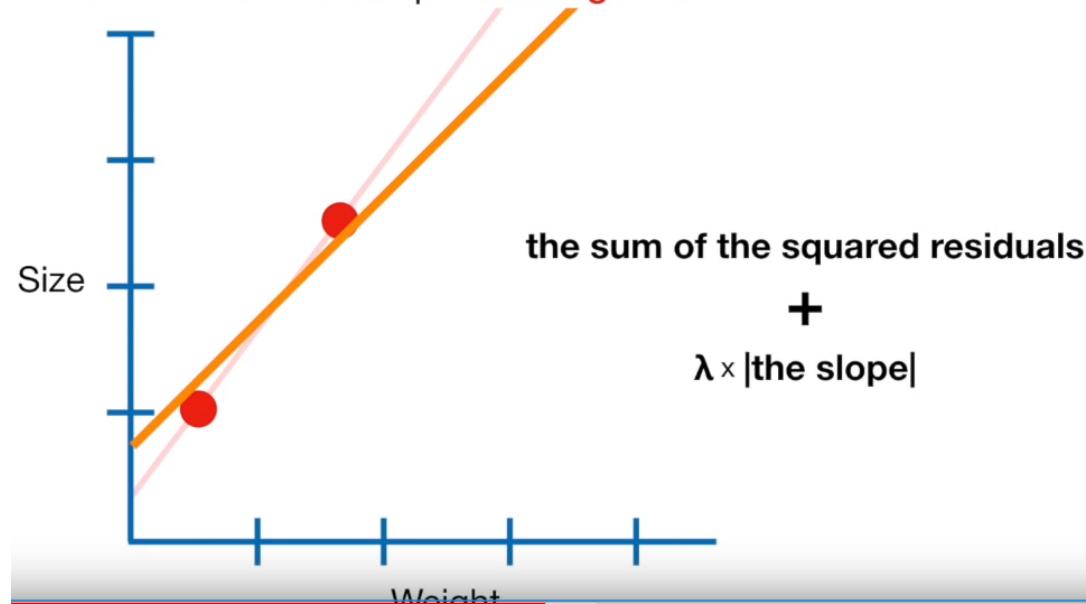
$$+ \boxed{\lambda \times (|\text{slope}| + |\text{diet difference}|)}$$



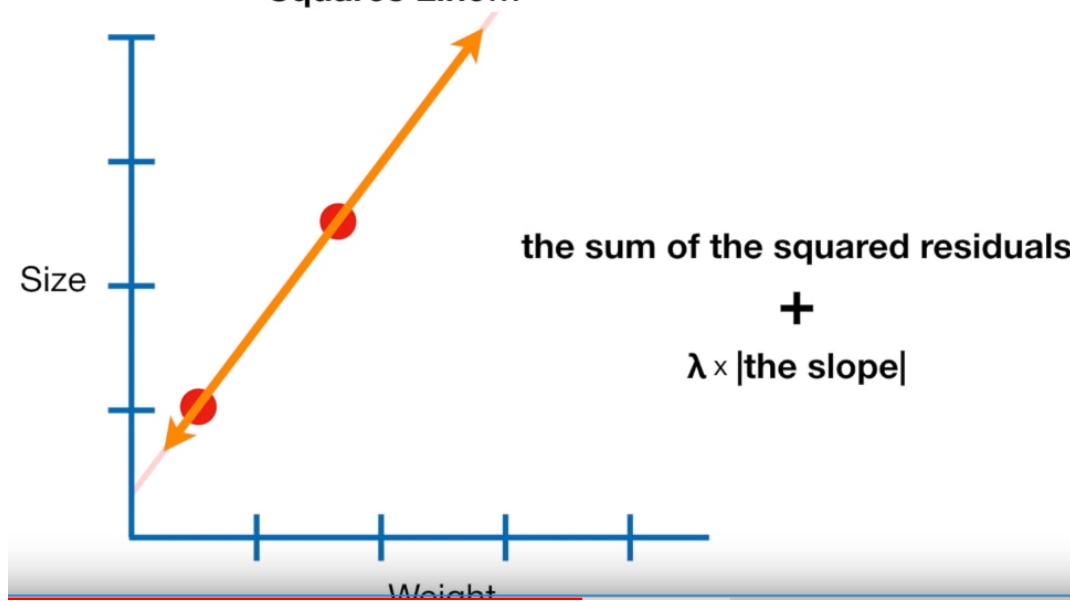
OK, we've seen how **Ridge** and  
**Lasso Regression** are similar.

Now let's talk about the big  
difference between them.

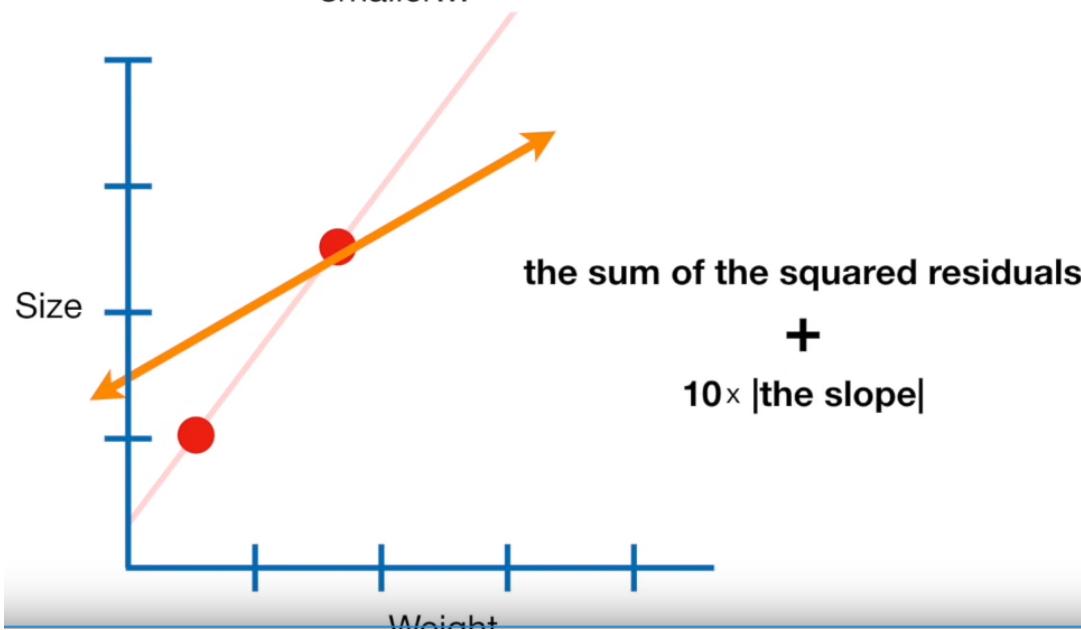
To see what makes **Lasso Regression** different from **Ridge Regression**, let's go back to the two sample **Training Data**.



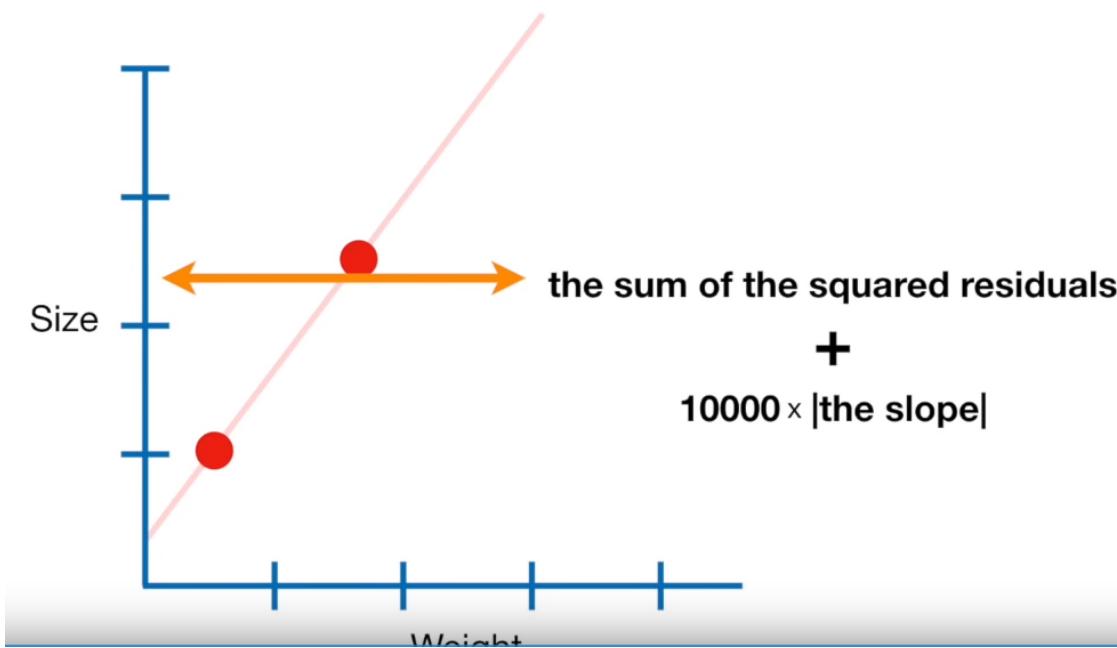
When  $\lambda = 0$ , then the **Lasso Regression Line** will be the same as the **Least Squares Line**...



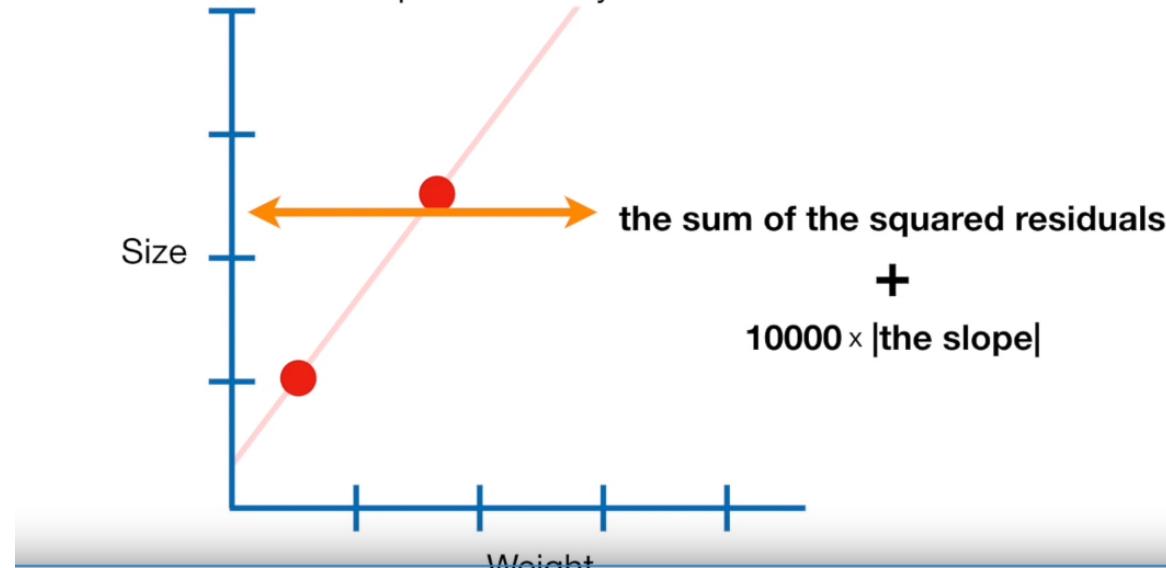
As  $\lambda$  increases in value, the slope gets smaller...



...until the slope = **0**.



The big difference between **Ridge** and **Lasso Regression** is that **Ridge Regression** can only shrink the slope asymptotically close to **0** while **Lasso Regression** can shrink the slope all the way to **0**.





(i)

$$\begin{aligned} \text{Size} = & \text{y-intercept} + \text{slope} \times \mathbf{Weight} + \text{diet difference} \times \mathbf{High\ Fat\ Diet} \\ & + \text{astrological offset} \times \mathbf{Sign} + \text{airspeed scalar} \times \mathbf{Airspeed\ of\ Swallow} \end{aligned}$$


To appreciate this difference,  
let's look at a big, huge, crazy  
equation...



Size = y-intercept + slope × **Weight** + diet difference × **High Fat Diet**  
+ astrological offset × **Sign** + airspeed scalar × **Airspeed of Swallow**

When we apply **Ridge Regression** to this equation, we find the minimal sum of the squared residuals plus the **Ridge Regression Penalty**...



$$\lambda \times (\text{slope}^2 + \text{diet difference}^2 + \text{astrological offset}^2 + \text{airspeed scalar}^2)$$

**Size** = y-intercept + slope  $\times$  **Weight** + diet difference  $\times$  **High Fat Diet**

+ astrological offset  $\times$  **Sign** + airspeed scalar  $\times$  **Airspeed of Swallow**

...and the larger we make  $\lambda$ ...

...and these parameters might shrink a lot, but they will never be equal to 0.

$\lambda$

( slope<sup>2</sup> + diet difference<sup>2</sup> + astrological offset<sup>2</sup> + airspeed scalar<sup>2</sup> )



Size = y-intercept + slope × Weight + diet difference × High Fat Diet  
+ astrological offset × Sign + airspeed scalar × Airspeed of Swallow

In contrast, with **Lasso Regression**...



$\lambda \times ( |\text{slope}| + |\text{diet difference}| + |\text{astrological offset}| + |\text{airspeed scalar}| )$



**Size** = y-intercept + slope  $\times$  **Weight** + diet difference  $\times$  **High Fat Diet**

- astrological offset  $\times$  **Sign** + airspeed scalar  $\times$  **Airspeed of Swallow**

...when we increase the value for  $\lambda$ ...

...and these parameters will go all the way to 0...

$\lambda$

( $|\text{slope}| + |\text{diet difference}| + |\text{astrological offset}| + |\text{airspeed scalar}|$ )



**Size** = y-intercept + slope  $\times$  **Weight** + diet difference  $\times$  **High Fat Diet**



...and we're left with a way to predict **Size** that only includes **Weight** and **Diet**...


$$\text{Size} = \text{y-intercept} + \text{slope} \times \mathbf{Weight} + \text{diet difference} \times \mathbf{High\ Fat\ Diet}$$

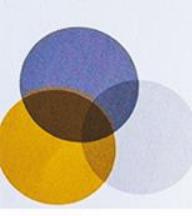
~~+ astrological offset  $\times$  **Sign** + airspeed scalar  $\times$  **Airspeed of Swallow**~~

Since **Lasso Regression** can exclude useless variables from equations, it is a little better than **Ridge Regression** at reducing the **Variance** in models that contain a lot of useless variables.


$$\text{Size} = \text{y-intercept} + \text{slope} \times \mathbf{Weight} + \text{diet difference} \times \mathbf{High\ Fat\ Diet}$$

~~+ astrological offset  $\times$  **Sign** + airspeed scalar  $\times$  **Airspeed of Swallow**~~

In contrast, **Ridge Regression** tends to do a little better  
when most variables are useful.



Ridge Regression is very similar to...



the sum of the squared residuals

+

$\lambda \times \text{the slope}^2$

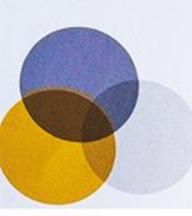
...Lasso Regression



the sum of the squared residuals

+

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



The superficial difference is that  
**Ridge Regression** squares the  
variables...

the sum of the squared residuals

+

$\lambda \times \text{the slope}^2$



**Size** = y-intercept + slope × **Weight** + diet difference × **High Fat Diet**

But the big difference is that **Lasso Regression** can exclude useless variables from equations.

This makes the final equation simpler and easier to interpret.