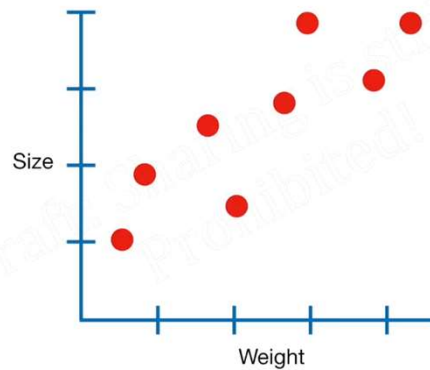


## REGULARIZATION

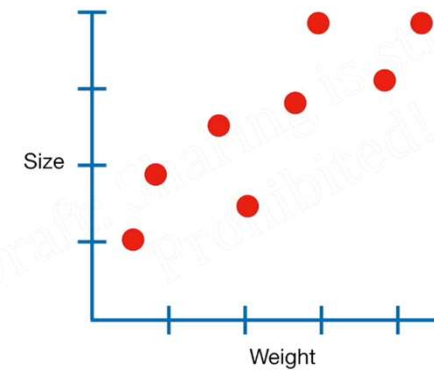
- Ridge Regression
- Lasso regression
- Elastic-Net

## RIDGE REGRESSION

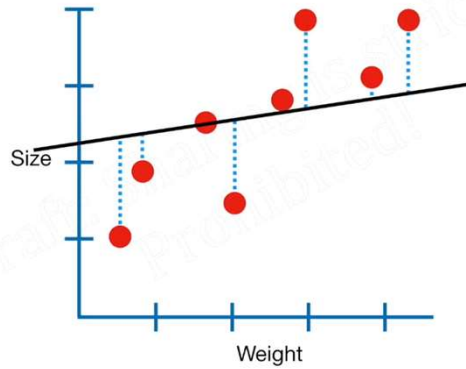
Let's start by collecting **Weight** and **Size** measurements from a bunch of mice...



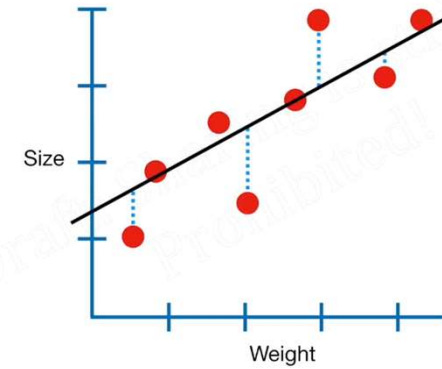
Since these data look relatively linear, we will use **Linear Regression**, AKA **Least Squares**, to model the relationship between **Weight** and **Size**.



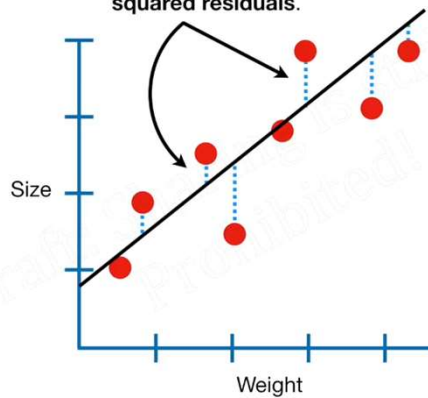
So we'll fit a line to the data using  
**Least Squares.**



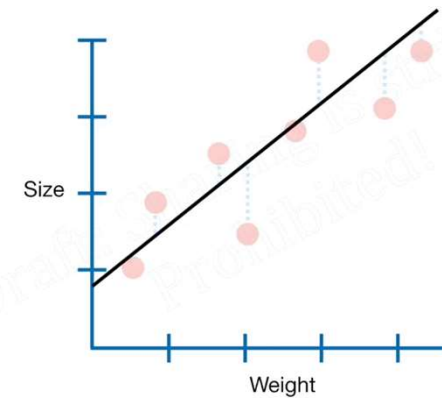
So we'll fit a line to the data using  
**Least Squares.**



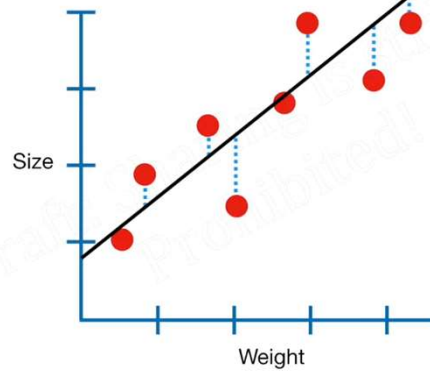
In other words, we find the line that  
results in the **minimum sum of  
squared residuals.**



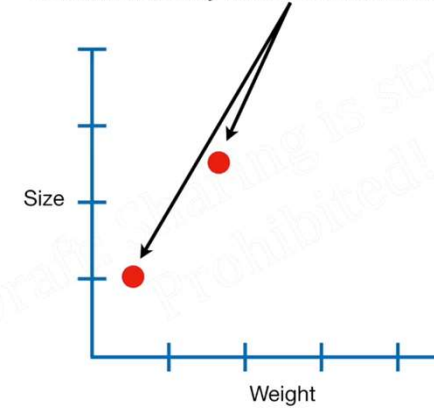
Ultimately, we end up with  
this equation for the line: **Size = 0.9 + 0.75 × Weight**



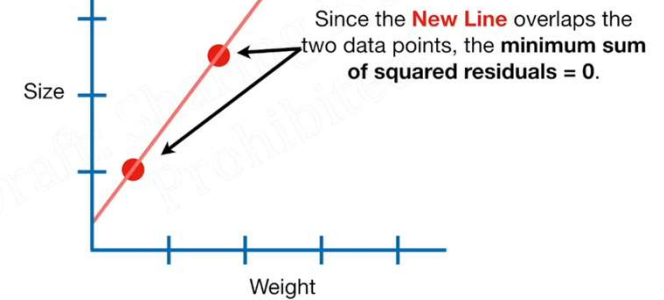
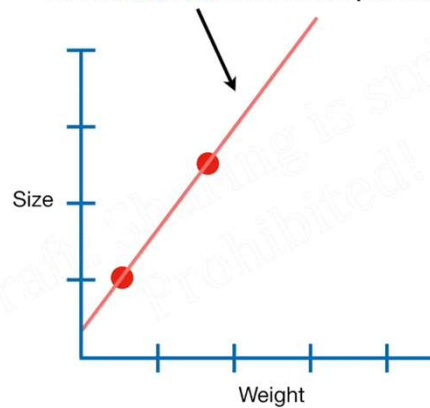
When we have a lot of measurements, we can be fairly confident that the **Least Squares** line accurately reflects the relationship between **Size** and **Weight**.

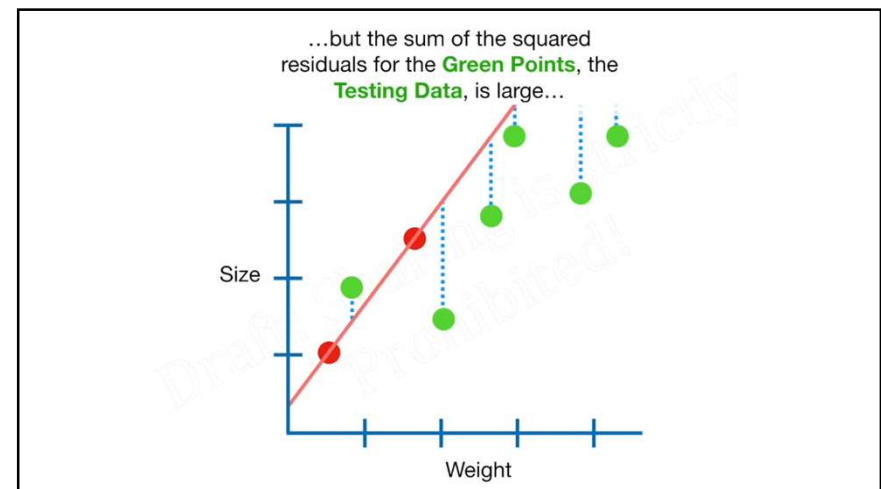
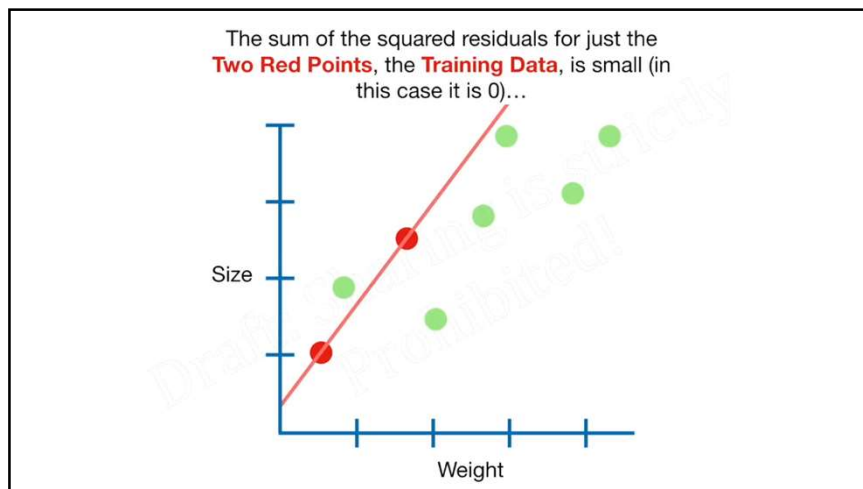
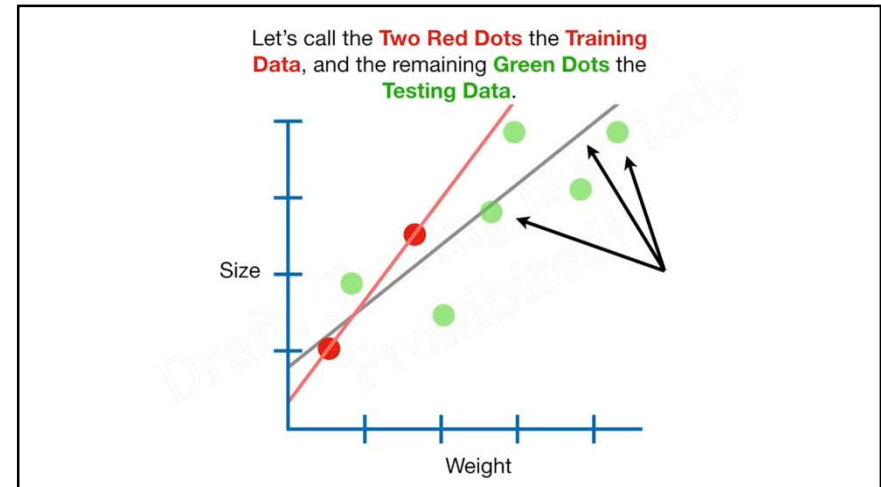
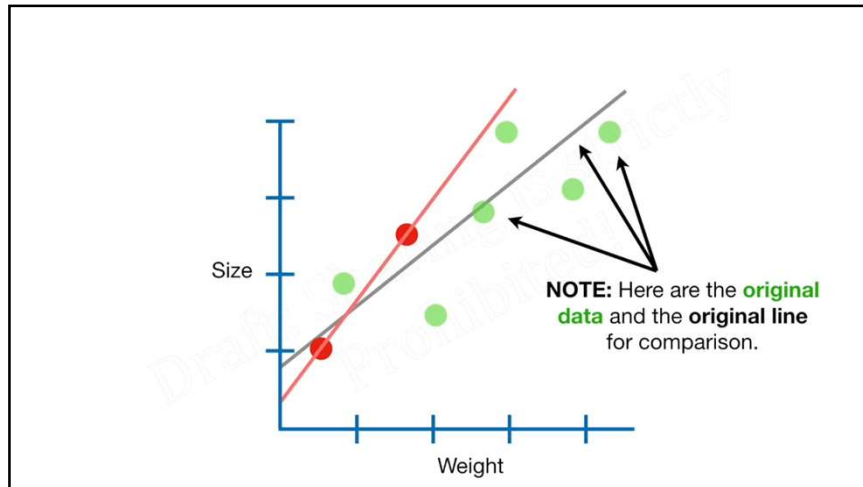


But what if we only have two measurements?

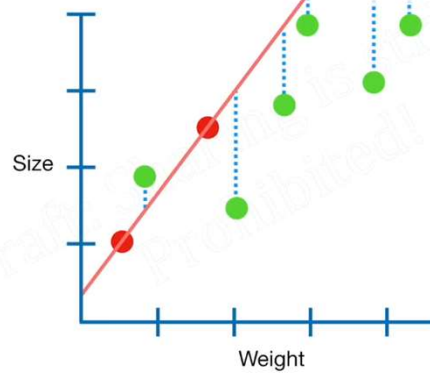


We fit a **New Line** with **Least Squares**...

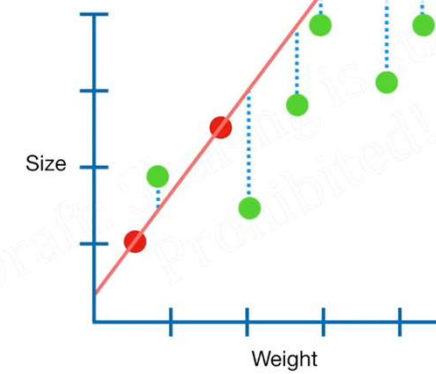




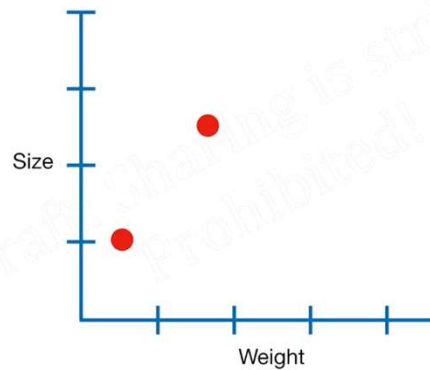
...and that means that the **New Line** has **High Variance**.



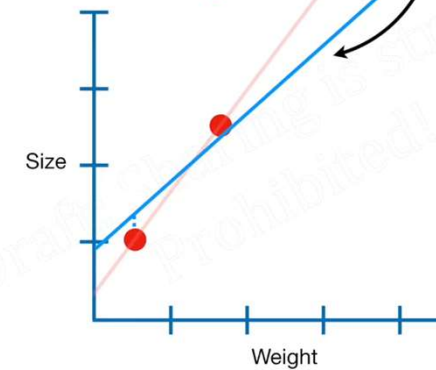
In machine learning lingo, we'd say that the **New Line** is **Over Fit** to the **Training Data**.



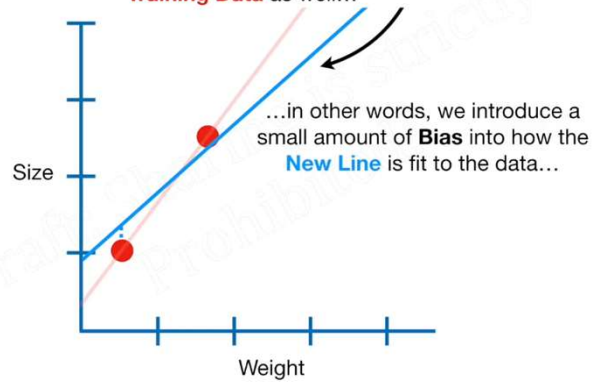
Now let's go back to just the **Training Data**...



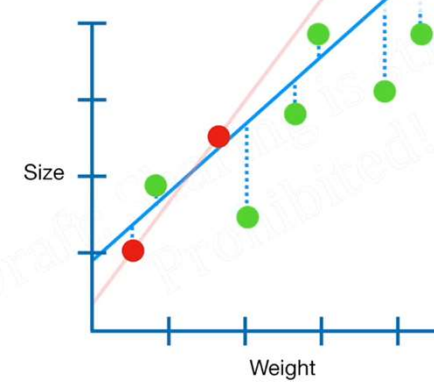
The main idea behind **Ridge Regression** is to find a **New Line** that doesn't fit the **Training Data** as well...



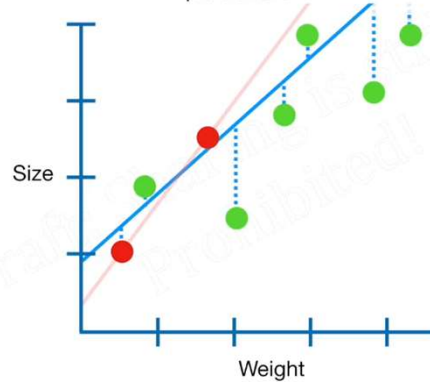
The main idea behind **Ridge Regression** is to find a **New Line** that doesn't fit the **Training Data** as well...



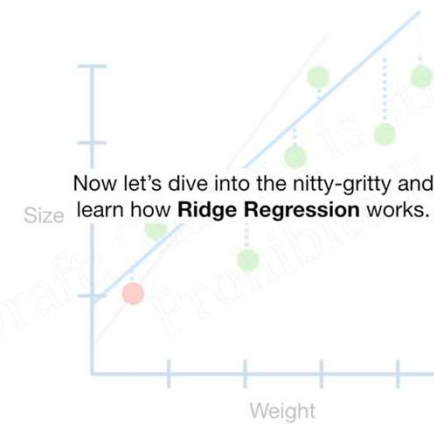
...but in return for that small amount of **Bias**, we get a significant drop in **Variance**.

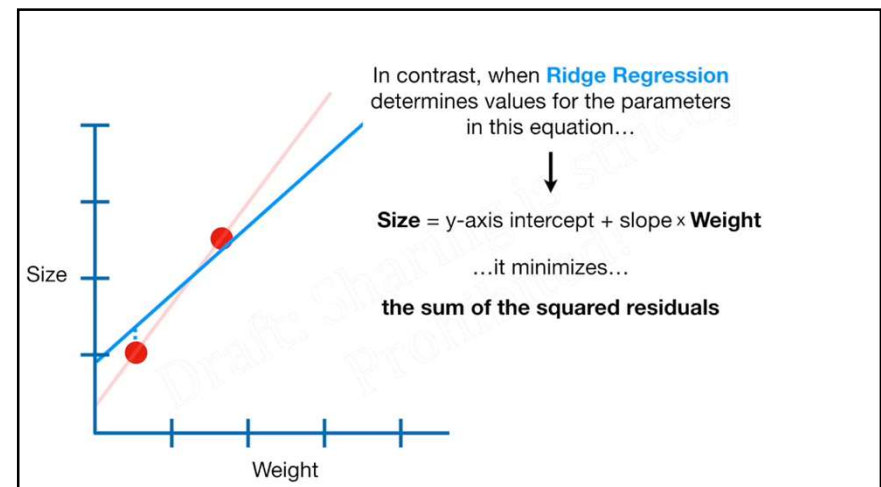
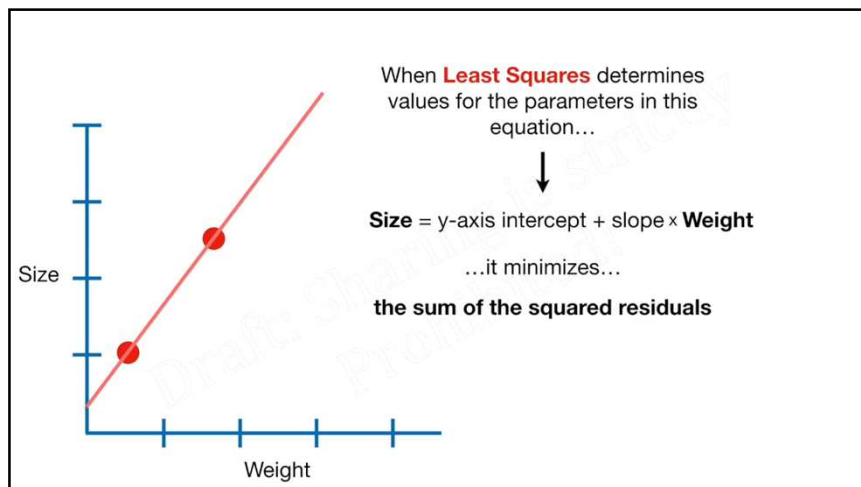
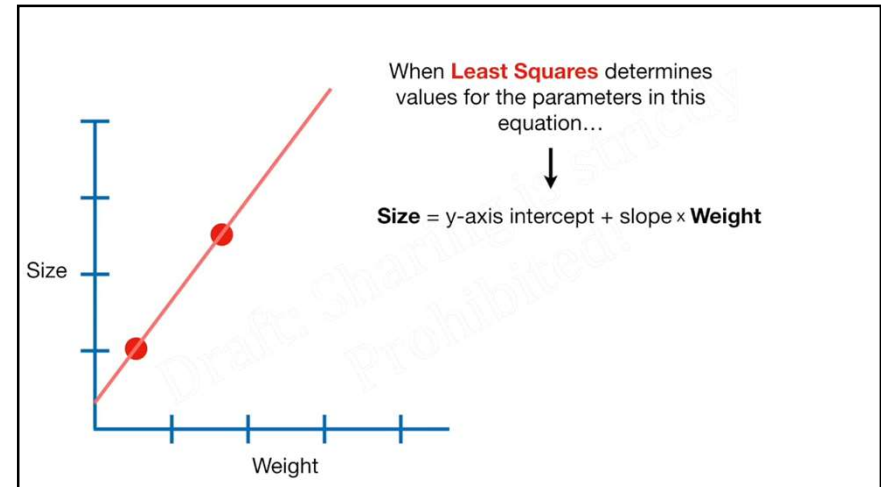
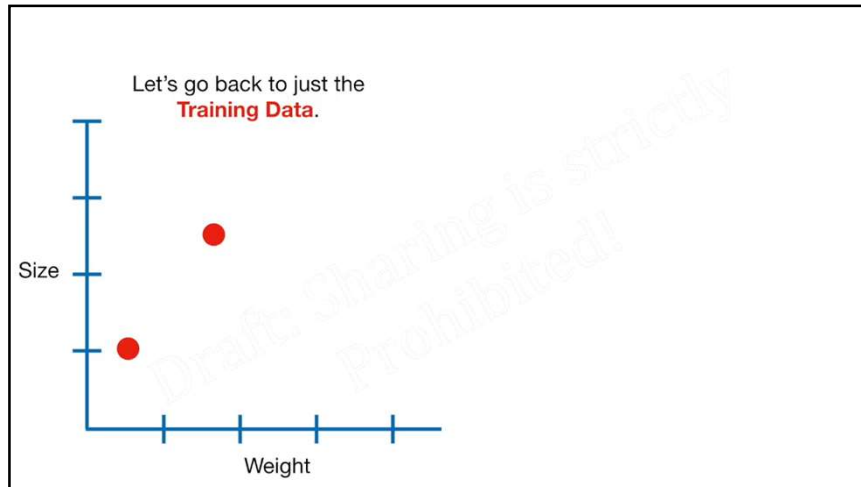


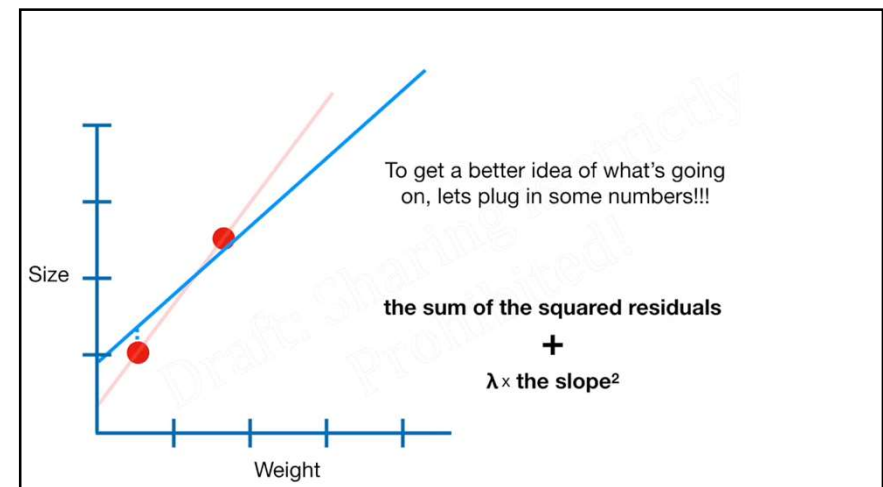
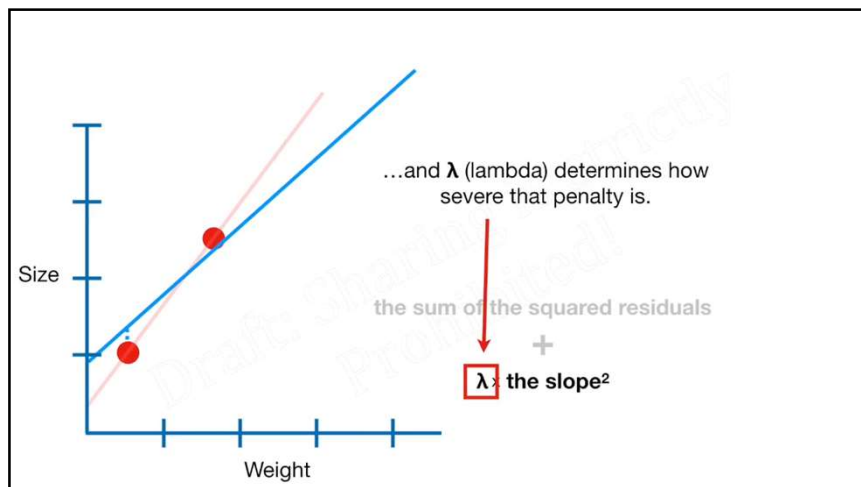
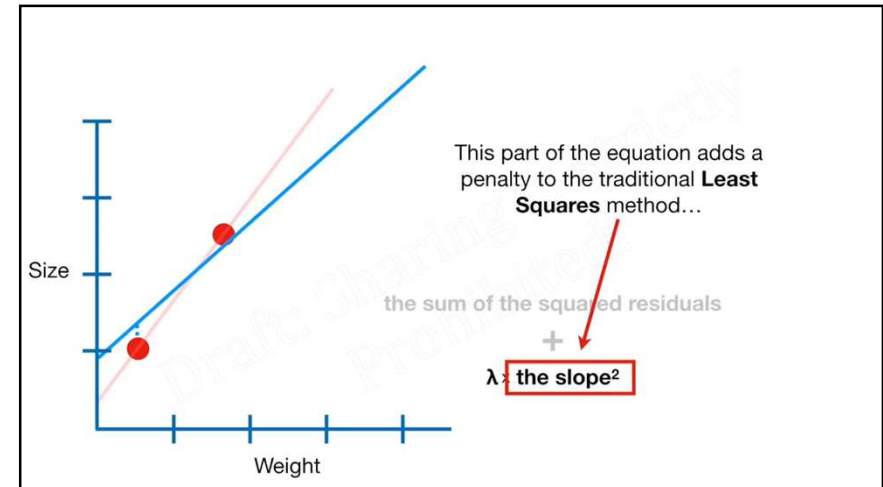
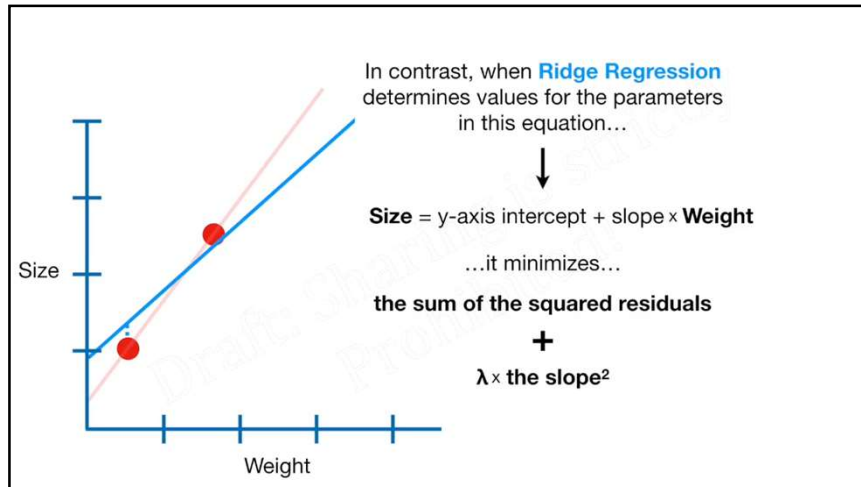
In other words, by starting with a slightly worse fit, **Ridge Regression** can provide better long term predictions.



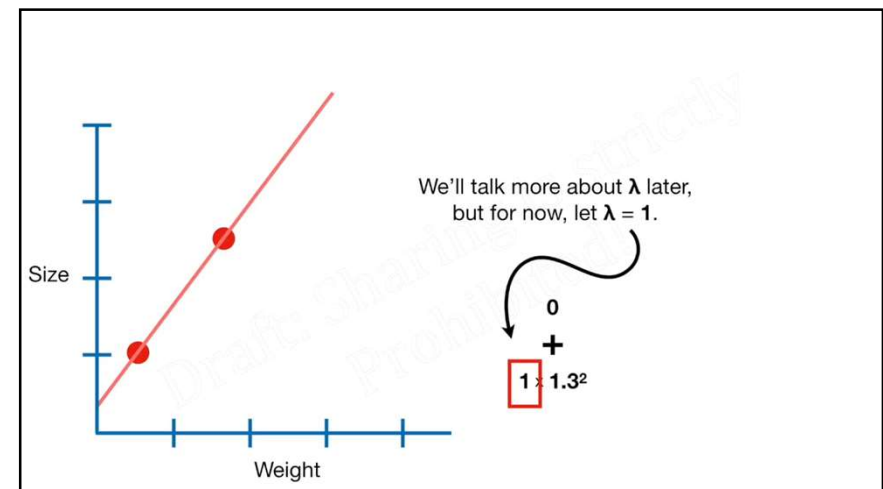
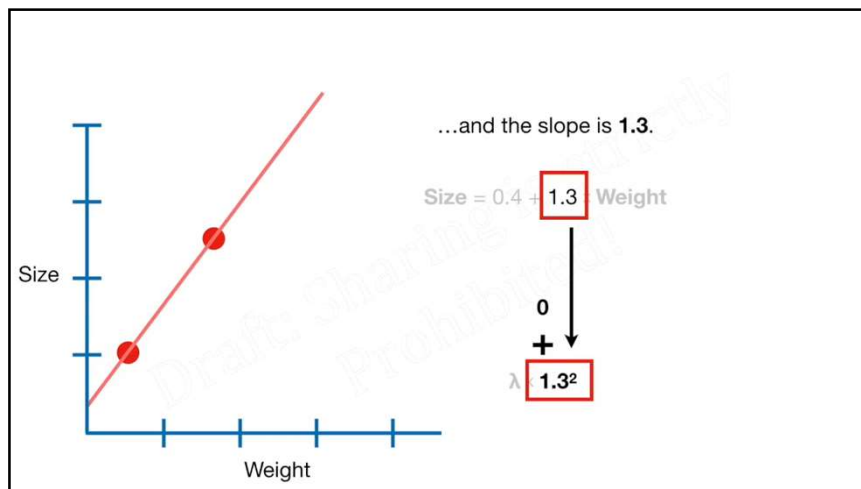
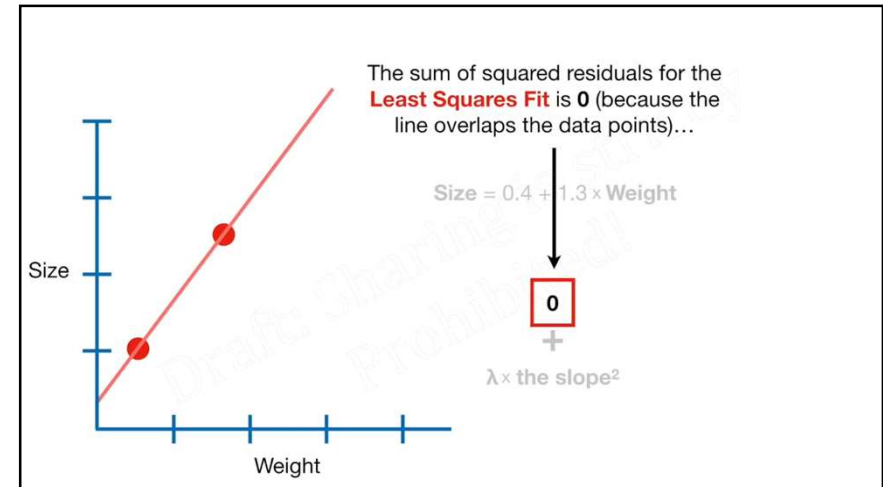
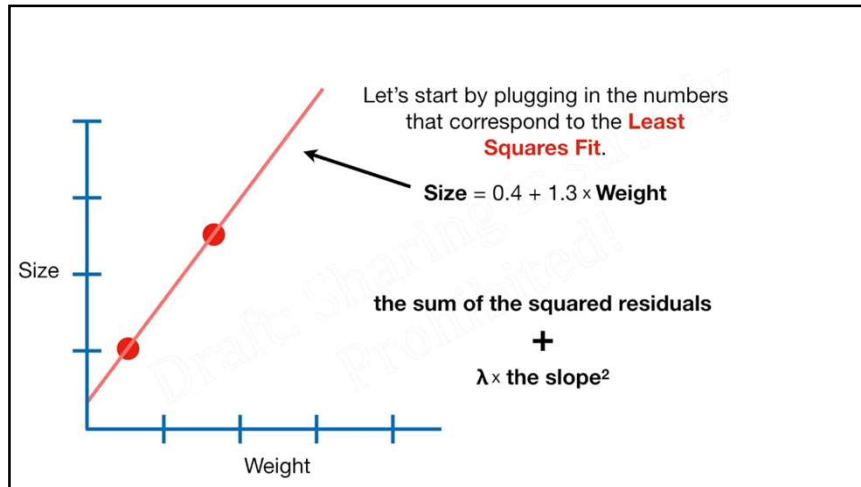
Now let's dive into the nitty-gritty and learn how **Ridge Regression** works.

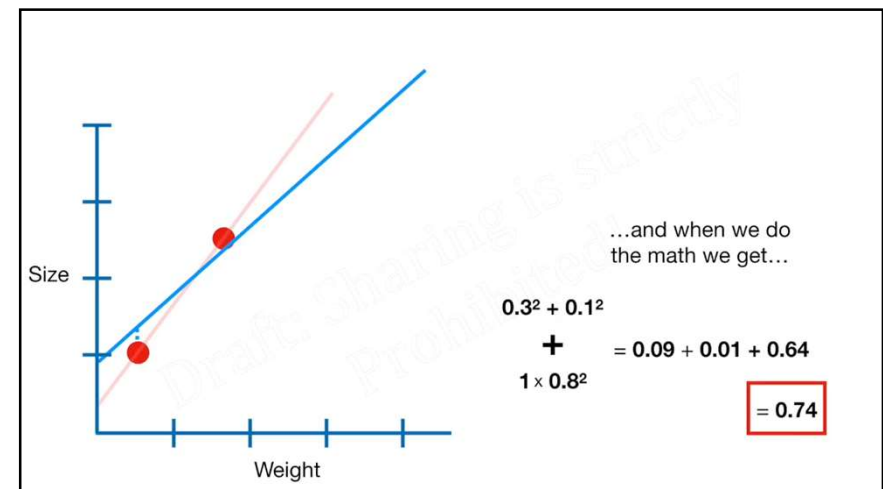
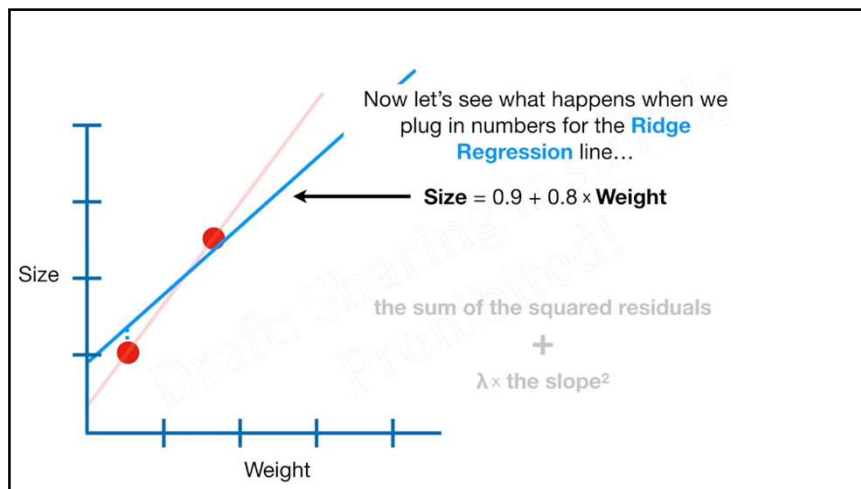
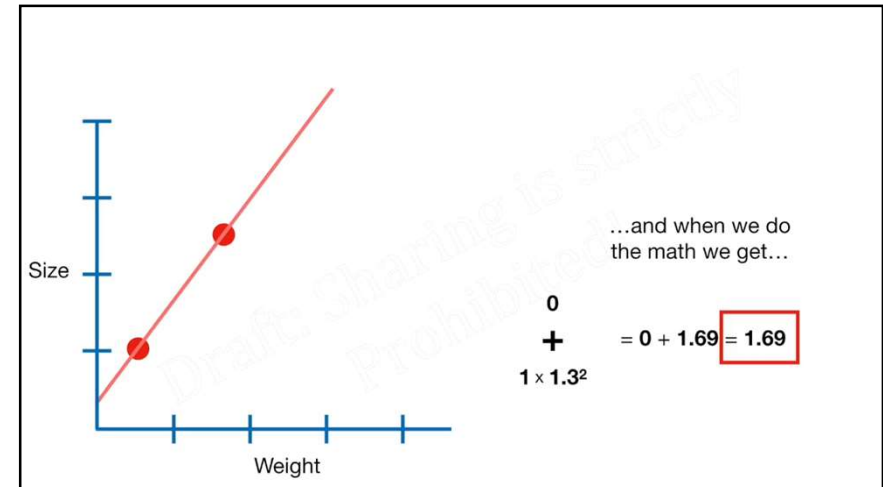
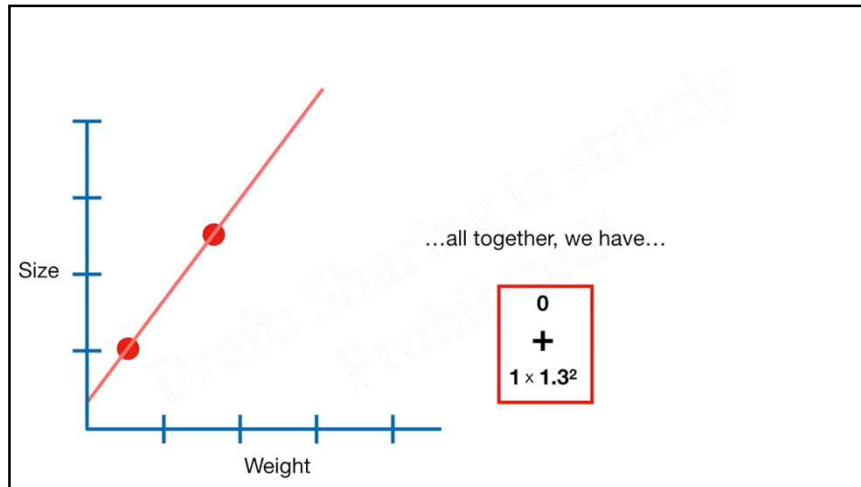


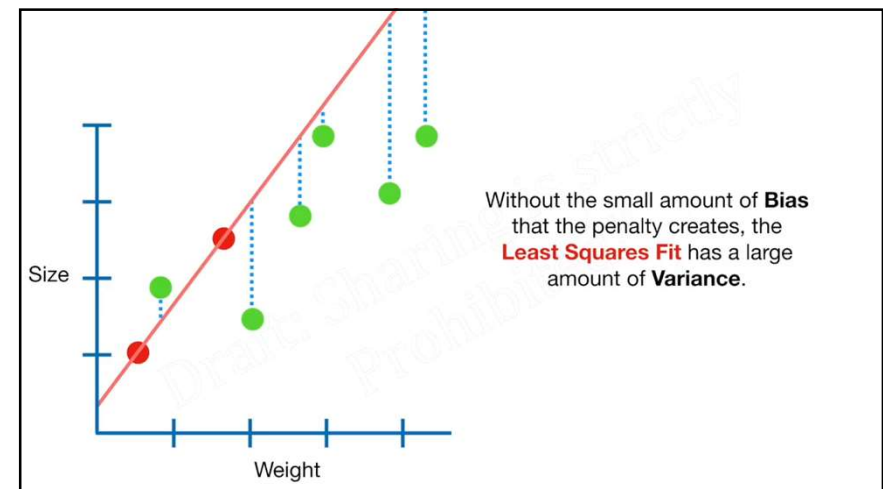
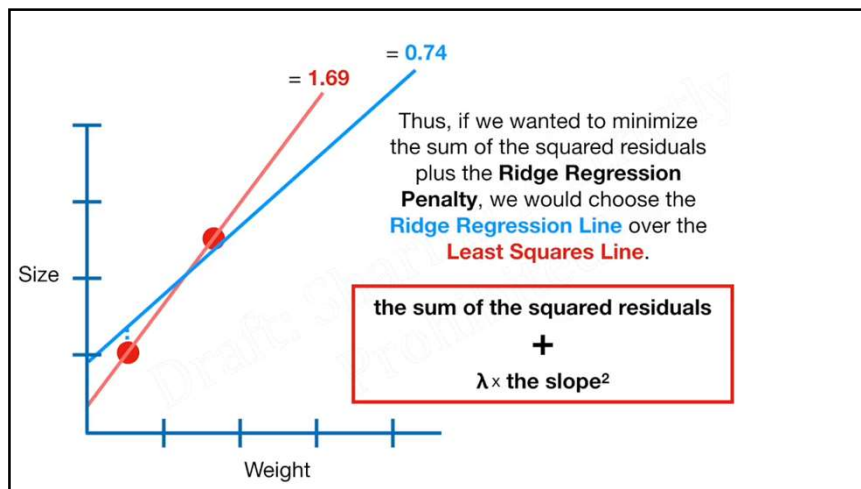
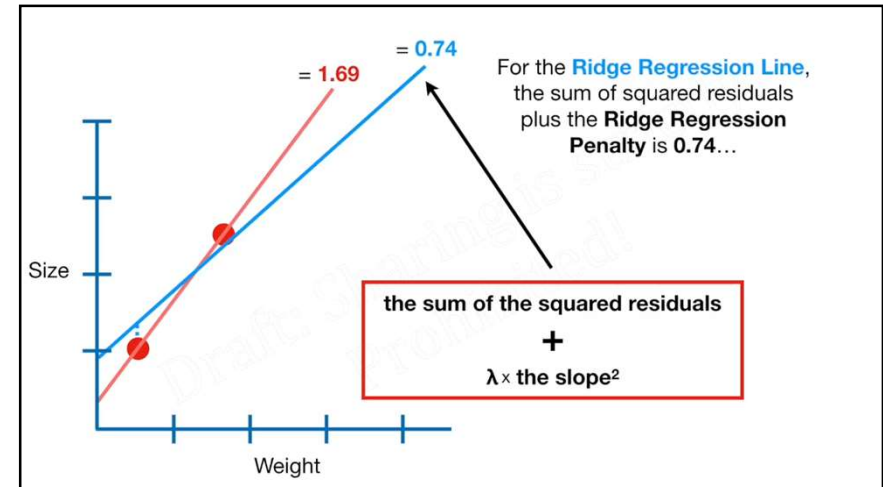
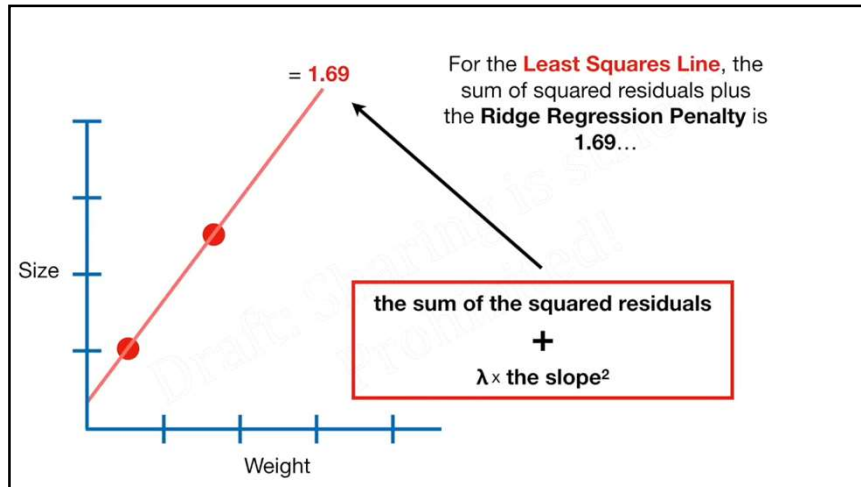


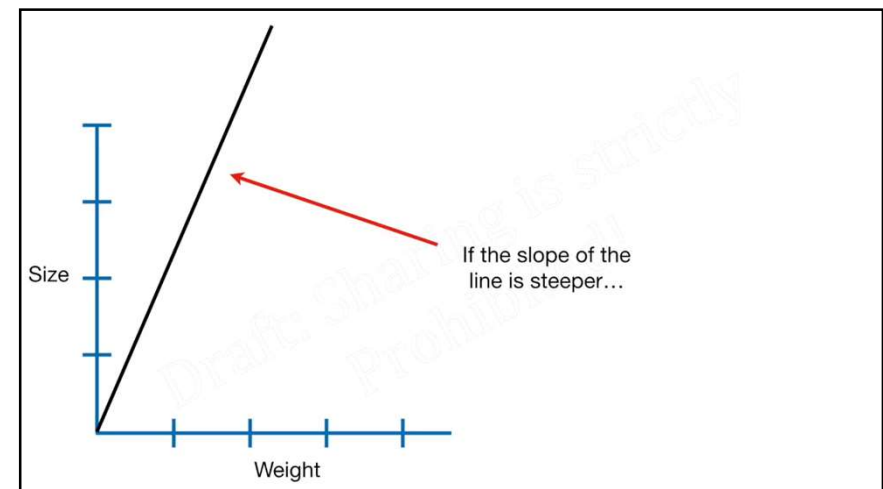
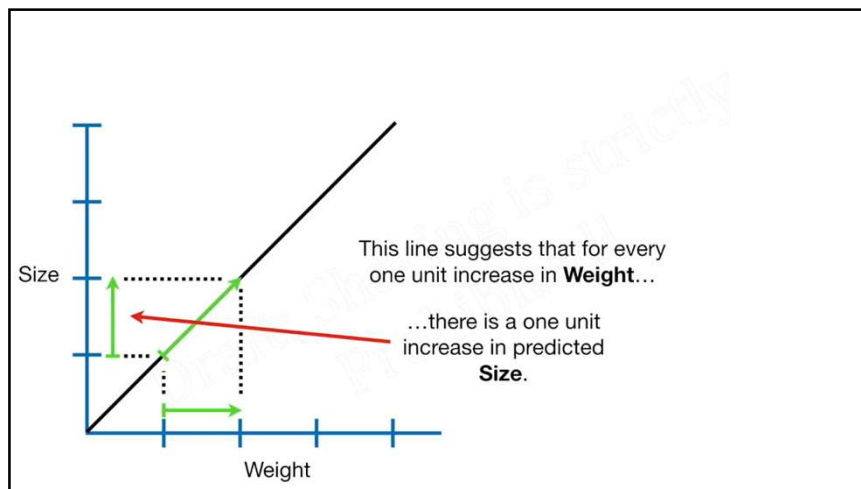
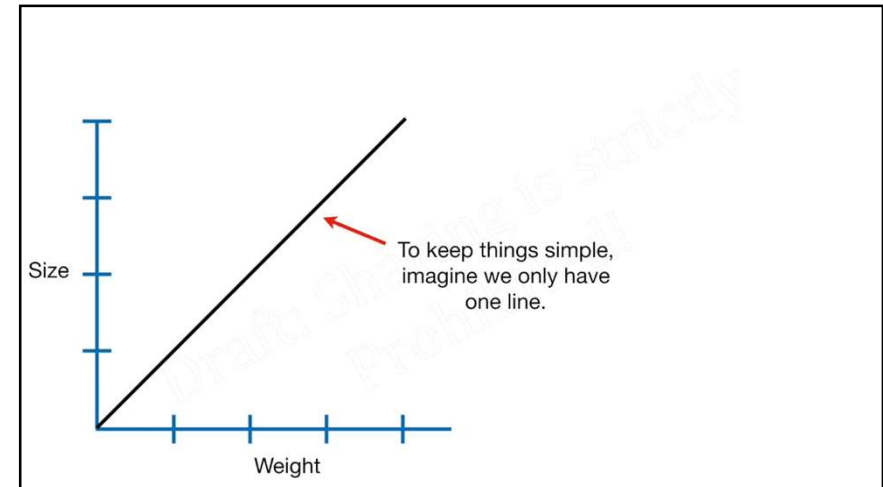
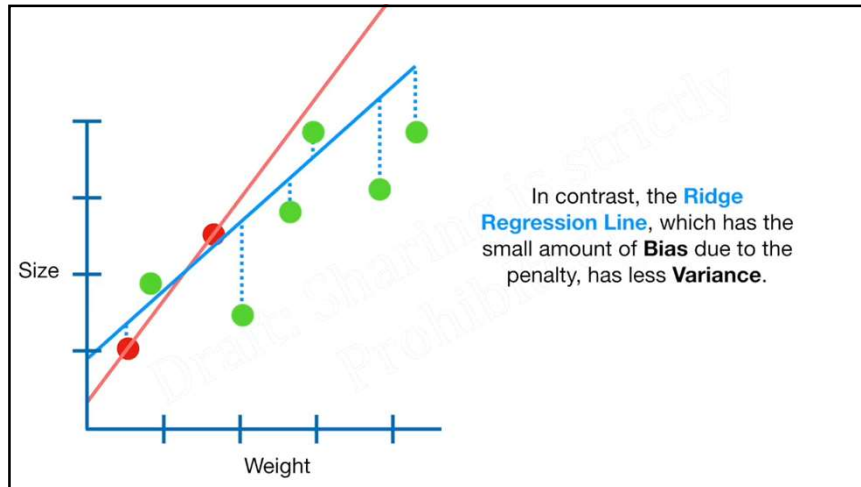


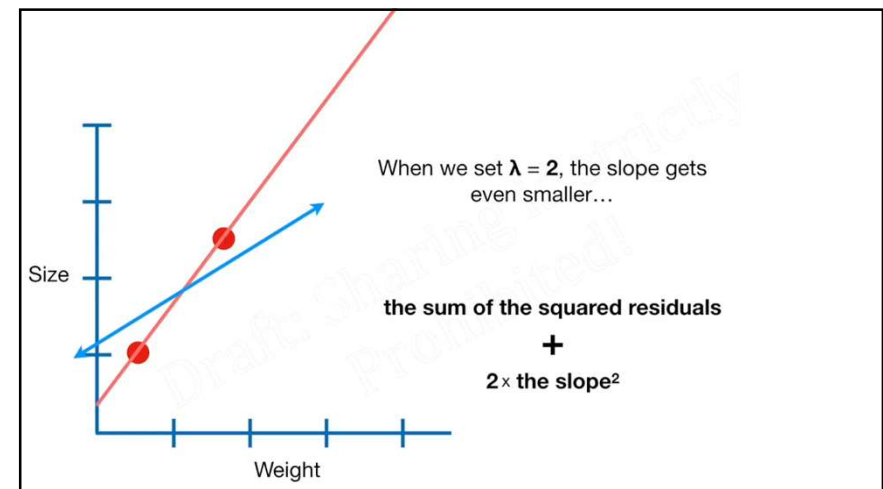
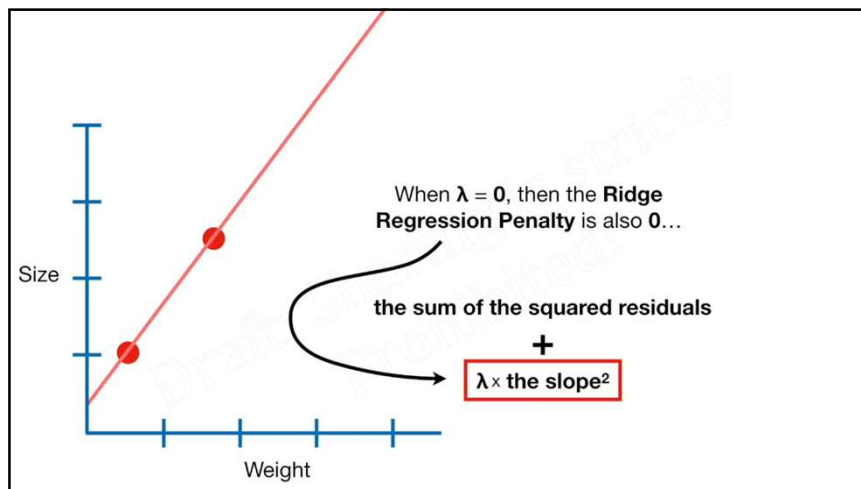
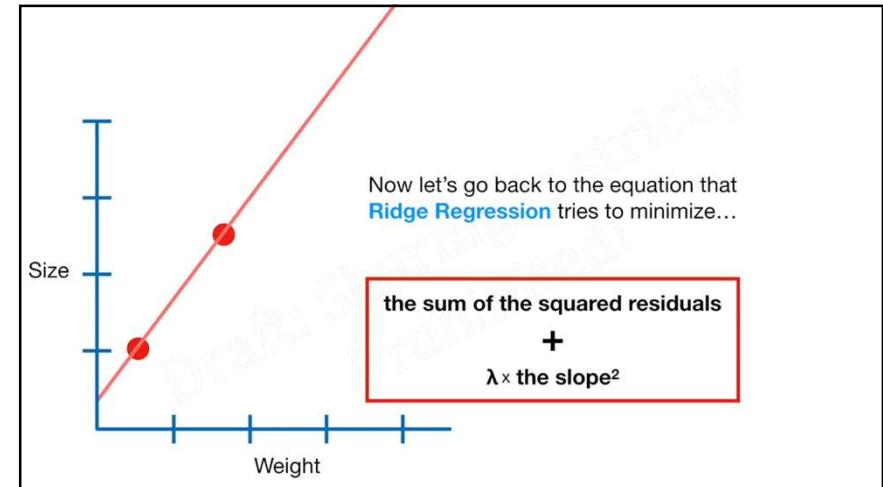
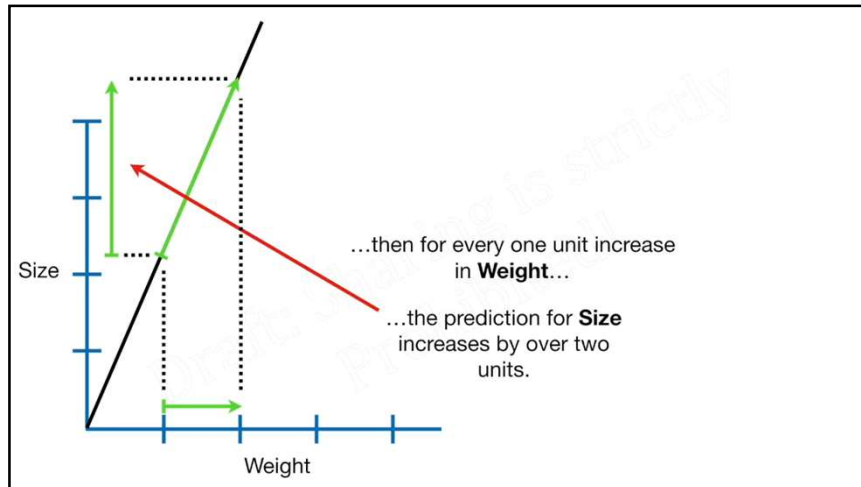


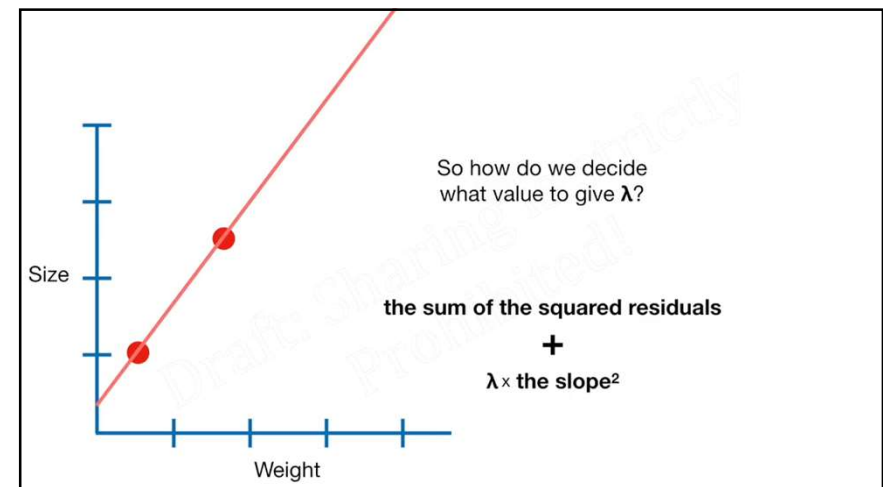
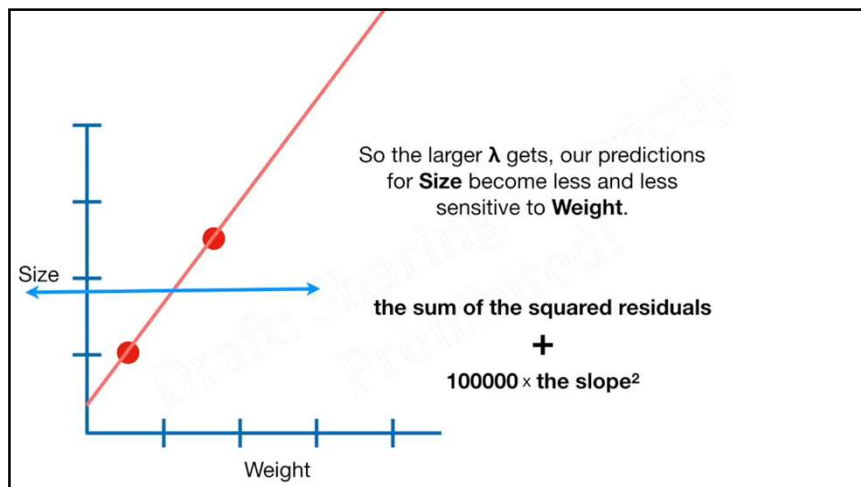
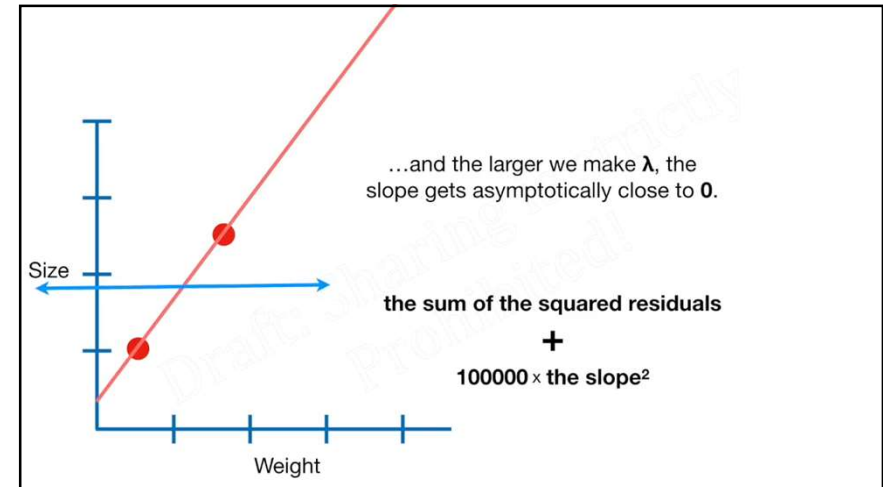
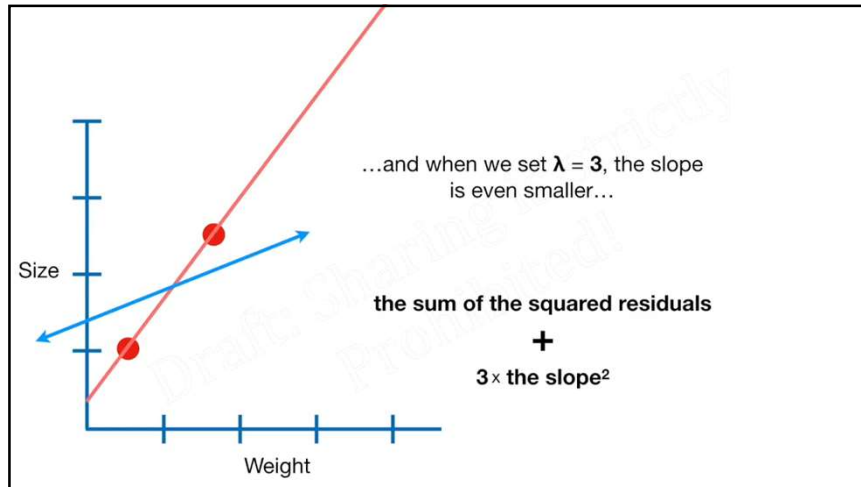


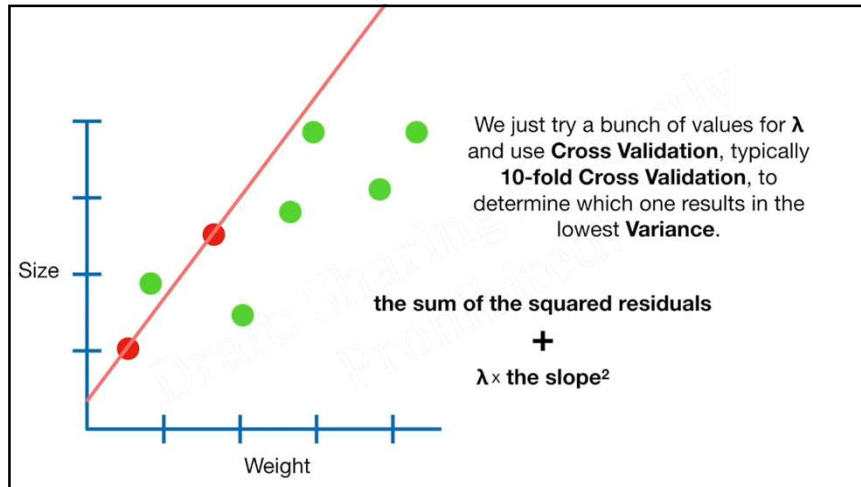




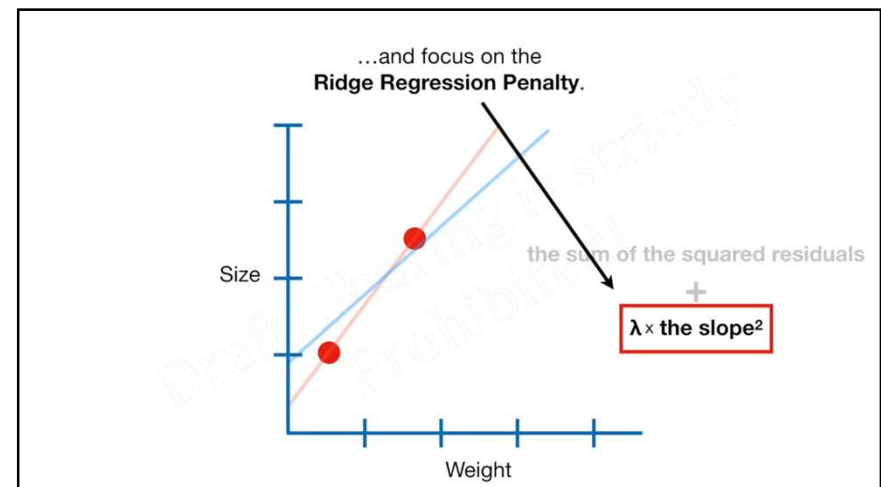
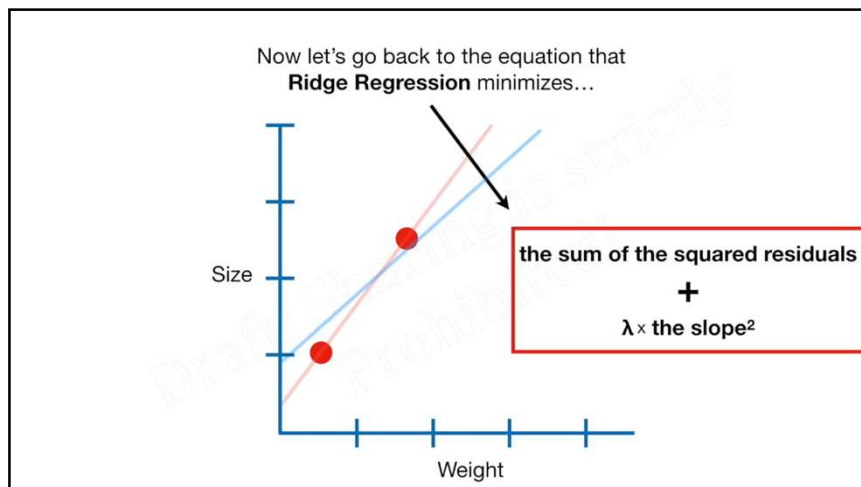


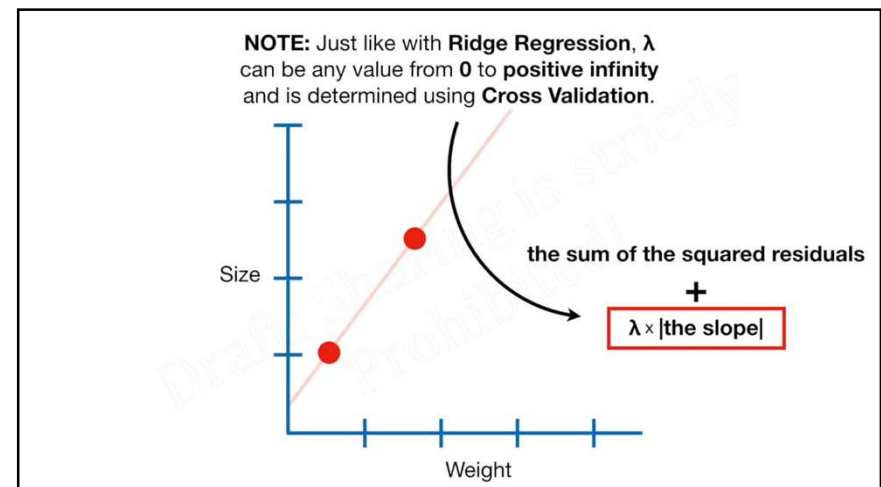
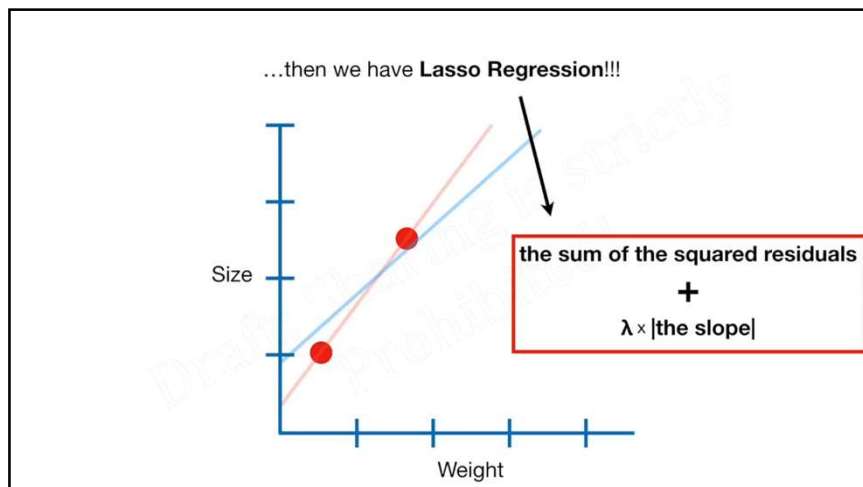
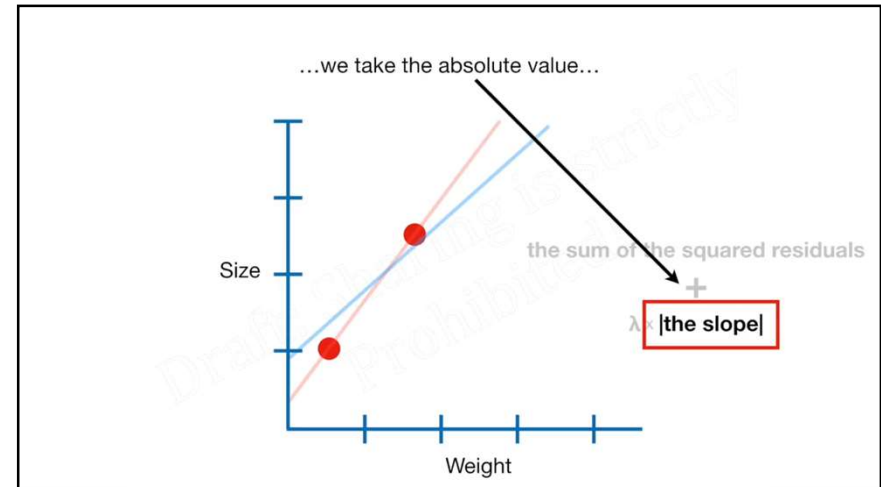
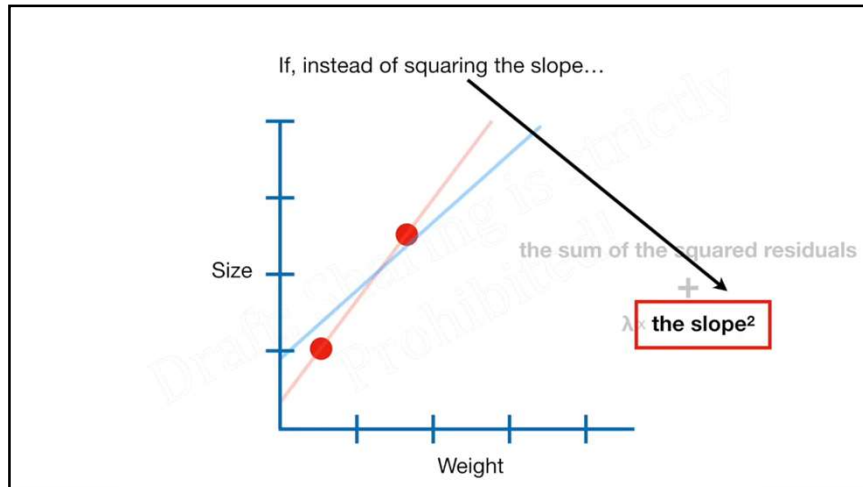




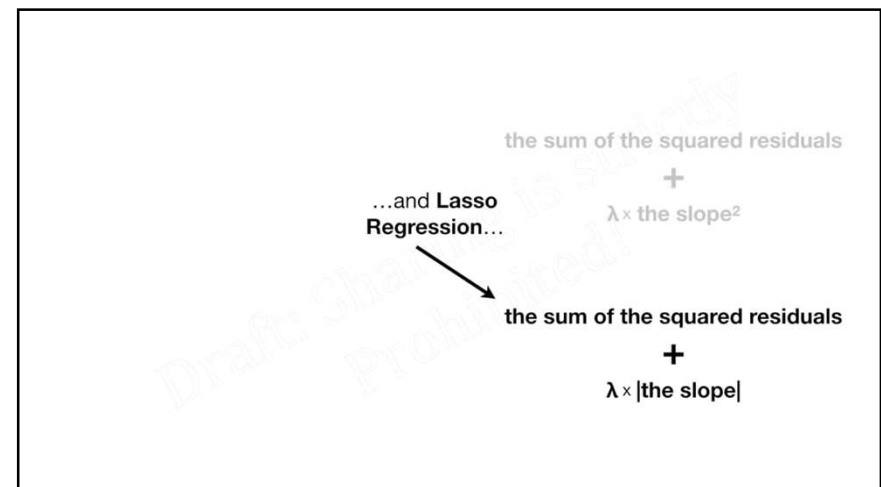
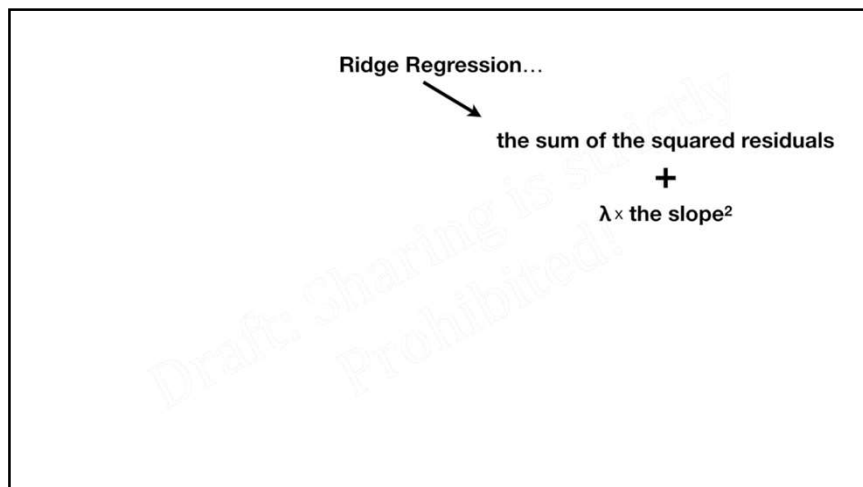
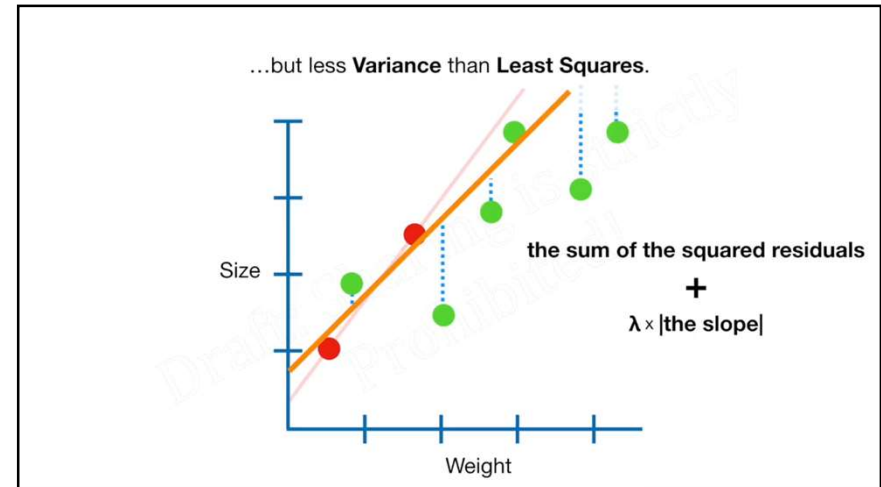
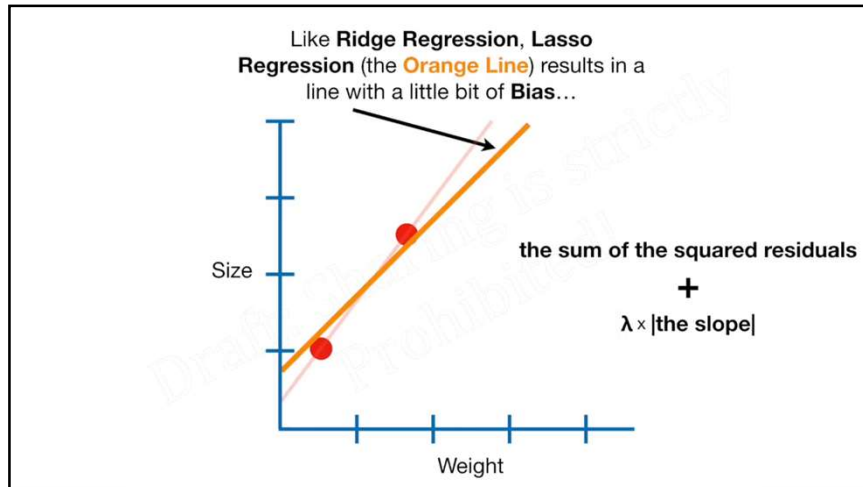


## LASSO REGRESSION





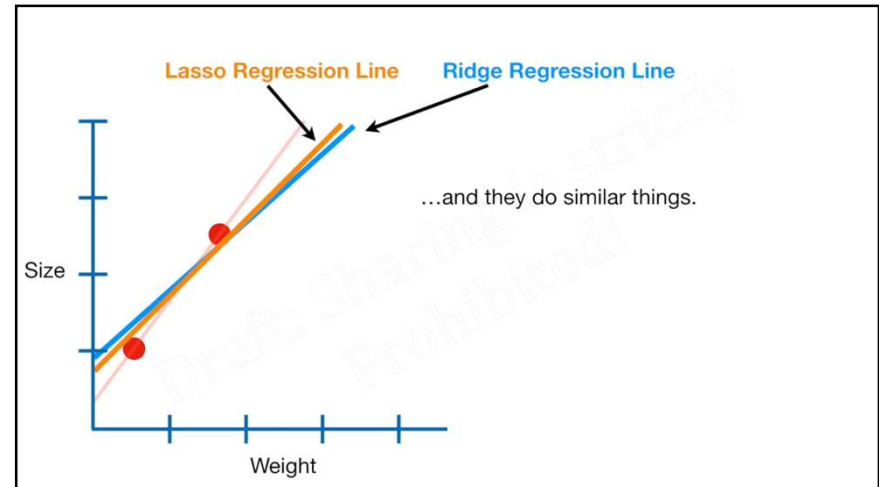




...look very similar....

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

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

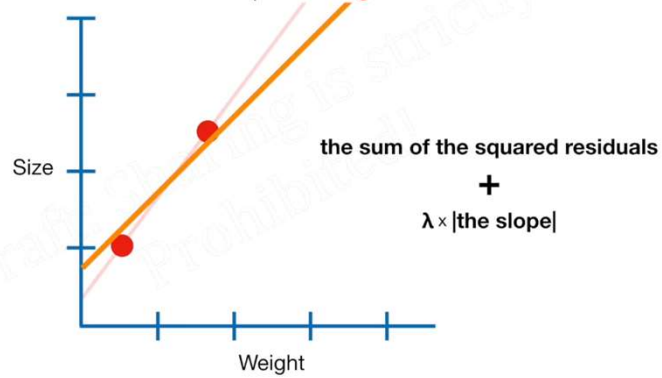


OK, we've seen how **Ridge** and **Lasso Regression** are similar.

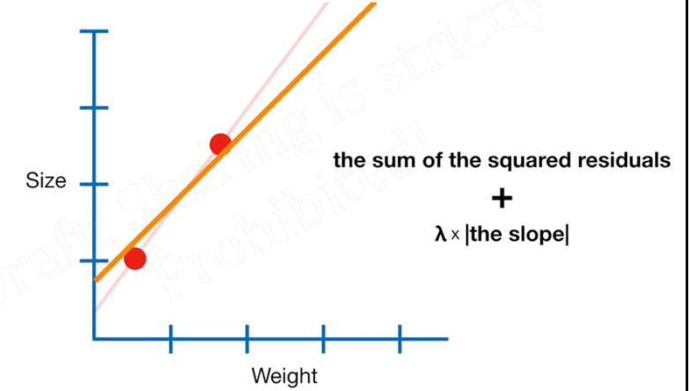
OK, we've seen how **Ridge** and **Lasso Regression** are similar.

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

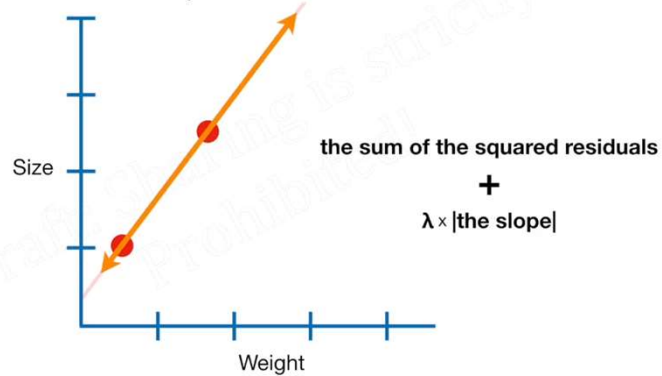
To see what makes **Lasso Regression** different from **Ridge Regression**, let's go back to the two sample **Training Data**.



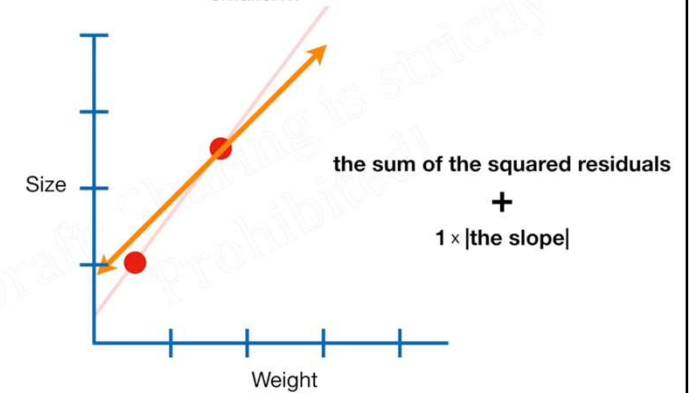
...and let's focus on what happens when we increase the value for  $\lambda$ .

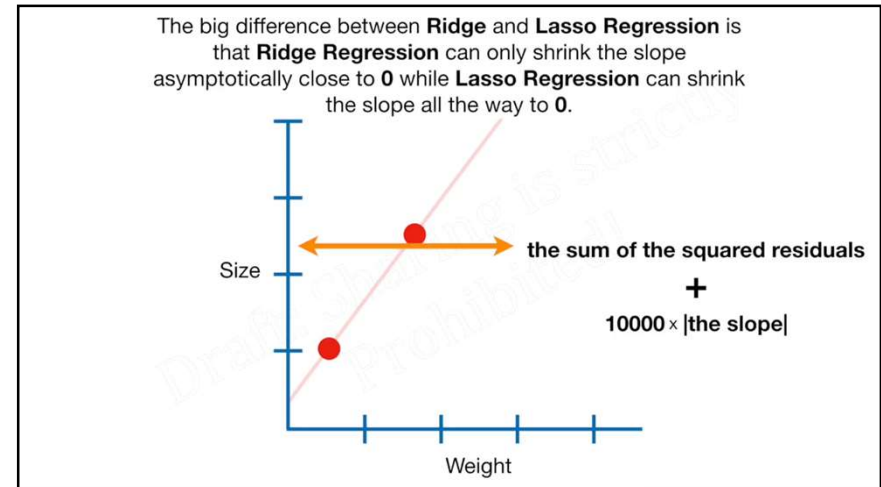
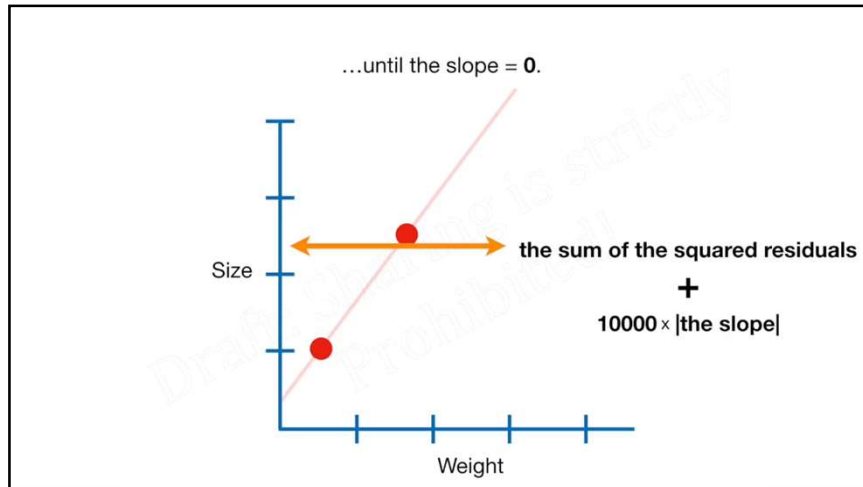


When  $\lambda = 0$ , then the **Lasso Regression Line** will be the same as the **Least Squares Line**...



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





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

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

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

The goal of this equation is to predict **Size**.

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

The terms for **Weight** and **High Fat Diet**, are both reasonable things to use to predict **Size**...

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

...but the **Astrological Sign** and the **Airspeed of a Swallow** (African or European) are terrible ways to predict **Size**.

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

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

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

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

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

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

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

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

...these parameters might shrink a little bit...

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

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

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

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

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

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

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

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

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

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

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

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

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

...then these parameters will shrink a little bit...

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

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

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

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

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

$$\text{Size} = \text{y-intercept} + \text{slope} \times \text{Weight} + \text{diet difference} \times \text{High Fat Diet} \\ + 0 \times \text{Sign} + 0 \times \text{Airspeed of Swallow}$$

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

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

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

$$\text{Size} = \text{y-intercept} + \text{slope} \times \text{Weight} + \text{diet difference} \times \text{High Fat Diet} \\ + \cancel{0 \times \text{Sign}} + \cancel{0 \times \text{Airspeed of Swallow}}$$

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

...and these terms go away...

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

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

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

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

~~$$+ \text{astrological offset} \times \text{Sign} + \text{airspeed scalar} \times \text{Airspeed of Swallow}$$~~

...and excludes all of the silly stuff!!!

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

~~$$+ \text{astrological offset} \times \text{Sign} + \text{airspeed scalar} \times \text{Airspeed of Swallow}$$~~

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.

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

~~$$+ \text{astrological offset} \times \text{Sign} + \text{airspeed scalar} \times \text{Airspeed of Swallow}$$~~

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



Ridge Regression is very similar to... → 
$$\begin{array}{c} \text{the sum of the squared residuals} \\ + \\ \lambda \times \text{the slope}^2 \end{array}$$

Ridge Regression is very similar to... → 
$$\begin{array}{c} \text{the sum of the squared residuals} \\ + \\ \lambda \times \text{the slope}^2 \end{array}$$

...Lasso Regression → 
$$\begin{array}{c} \text{the sum of the squared residuals} \\ + \\ \lambda \times |\text{the slope}| \end{array}$$

The superficial difference is that **Ridge Regression** squares the variables... → 
$$\begin{array}{c} \text{the sum of the squared residuals} \\ + \\ \lambda \times \text{the slope}^2 \end{array}$$

The superficial difference is that **Ridge Regression** squares the variables... → 
$$\begin{array}{c} \text{the sum of the squared residuals} \\ + \\ \lambda \times \text{the slope}^2 \end{array}$$

...and **Lasso Regression** takes the absolute value. → 
$$\begin{array}{c} \text{the sum of the squared residuals} \\ + \\ \lambda \times |\text{the slope}| \end{array}$$

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

~~+ astrological offset × **Sign** + airspeed scalar × **Airspeed of Swallow**~~

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

**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.

## ELASTIC-NET

...but what do we do when we have a model that includes tons more variables?

**Size** = y-intercept + slope<sub>1</sub> × **Weight** + diet difference × **High Fat Diet**

+ slope<sub>2</sub> × **Age** + slope<sub>3</sub> × **Size of Father** + .... + **tons more variables...**

...and when you have millions of parameters, then you will almost certainly need to use some sort of regularization to estimate them.

**Size** = y-intercept + slope<sub>1</sub> × **Weight** + diet difference × **High Fat Diet**  
 + slope<sub>2</sub> × **Age** + slope<sub>3</sub> × **Size of Father** + .... + tons more variables...

However, the variables in those models might be useful or useless. We don't know in advance.

**Size** = y-intercept + slope<sub>1</sub> × **Weight** + diet difference × **High Fat Diet**  
 + slope<sub>2</sub> × **Age** + slope<sub>3</sub> × **Size of Father** + .... + tons more variables...

So how do you choose if you should use **Lasso** or **Ridge Regression**?

**Size** = y-intercept + slope<sub>1</sub> × **Weight** + diet difference × **High Fat Diet**  
 + slope<sub>2</sub> × **Age** + slope<sub>3</sub> × **Size of Father** + .... + tons more variables...

The good news is that you don't have to choose, instead, you use **Elastic-Net Regression!!!**

**Size** = y-intercept + slope<sub>1</sub> × **Weight** + diet difference × **High Fat Diet**  
 + slope<sub>2</sub> × **Age** + slope<sub>3</sub> × **Size of Father** + .... + tons more variables...

**Elastic-Net Regression** sounds super fancy, but if you already know about **Lasso** and **Ridge Regression**, it's super simple.

Just like **Lasso** and **Ridge Regression**, **Elastic-Net Regression** starts with **Least Squares**...

↓  
the sum of the squared residuals

$$+ \lambda_1 \times |\text{variable}_1| + \dots + |\text{variable}_x| + \lambda_2 \times \text{variable}_1^2 + \dots + \text{variable}_x^2$$

...then it combines the **Lasso Regression Penalty**...

↓  
the sum of the squared residuals

$$+ \lambda_1 \times |\text{variable}_1| + \dots + |\text{variable}_x| + \lambda_2 \times \text{variable}_1^2 + \dots + \text{variable}_x^2$$

...with the **Ridge Regression Penalty**.

↓  
the sum of the squared residuals

$$+ \lambda_1 \times |\text{variable}_1| + \dots + |\text{variable}_x| + \lambda_2 \times \text{variable}_1^2 + \dots + \text{variable}_x^2$$

Altogether, **Elastic Net Regression** combines the strengths of **Lasso** and **Ridge Regression**.



the sum of the squared residuals

+

$$\lambda_1 \times |\text{variable}_1| + \dots + |\text{variable}_x| + \lambda_2 \times \text{variable}_1^2 + \dots + \text{variable}_x^2$$

**NOTE:** The **Lasso Regression Penalty** and the **Ridge Regression Penalty** get their own  $\lambda$ s.

$\lambda_1$  for Lasso...

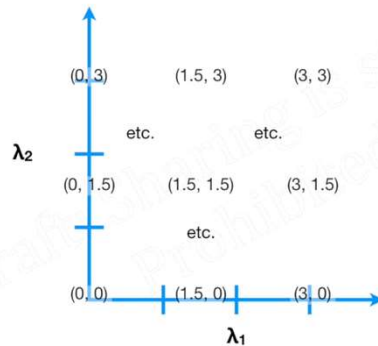
...and  $\lambda_2$  for Ridge.

the sum of the squared residuals

+

$$\lambda_1 |\text{variable}_1| + \dots + |\text{variable}_x| + \lambda_2 \text{variable}_1^2 + \dots + \text{variable}_x^2$$

We use **Cross Validation** on different combinations of  $\lambda_1$  and  $\lambda_2$  to find the best values.



The hybrid **Elastic-Net Regression** is especially good at dealing with situations when there are correlations between parameters.

the sum of the squared residuals

+

$$\lambda_1 \times |\text{variable}_1| + \dots + |\text{variable}_x| + \lambda_2 \times \text{variable}_1^2 + \dots + \text{variable}_x^2$$

This is because on it's own, **Lasso Regression** tends to pick just one of the correlated terms and eliminates the others...

$$\begin{array}{c} \text{the sum of the squared residuals} \\ + \\ \lambda_1 \times |\text{variable}_1| + \dots + |\text{variable}_x| \end{array}$$

...whereas **Ridge Regression** tends to shrink all of the parameters for the correlated variables together.

$$\begin{array}{c} \text{the sum of the squared residuals} \\ + \\ \lambda_2 \times \text{variable}_1^2 + \dots + \text{variable}_x^2 \end{array}$$

By combining **Lasso** and **Ridge Regression**, **Elastic-Net Regression** groups and shrinks the parameters associated with the correlated variables and leaves them in equation or removes them all at once.

$$\begin{array}{c} \text{the sum of the squared residuals} \\ + \\ \lambda_1 \times |\text{variable}_1| + \dots + |\text{variable}_x| + \lambda_2 \times \text{variable}_1^2 + \dots + \text{variable}_x^2 \end{array}$$

# THANK YOU

StatQuest with Josh Starmer  
<https://www.youtube.com/watch?v=Q81RR3yKw30>