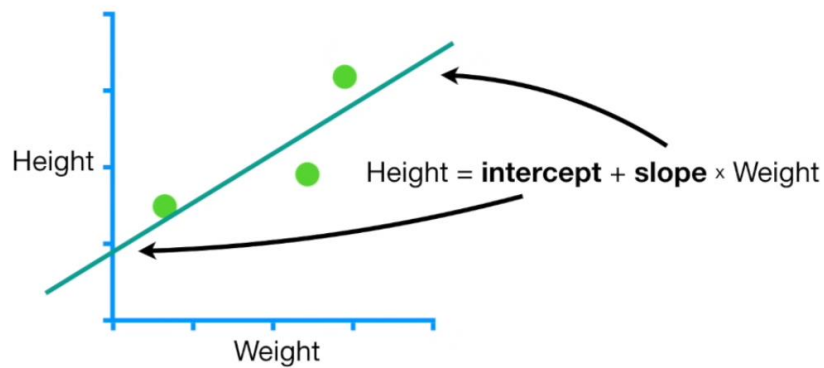


Gradient Descent

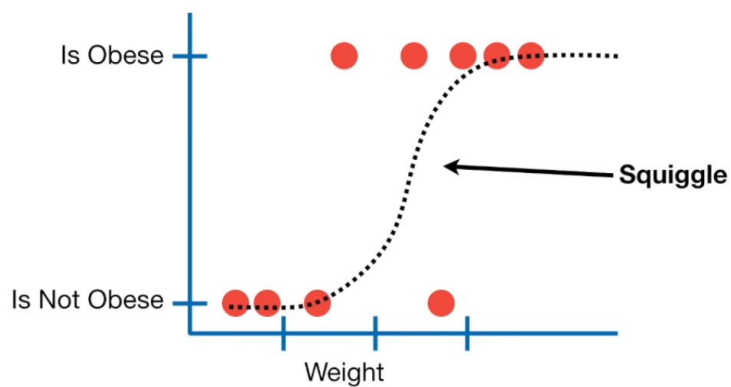
Purpose to find out the best intercept value for best fit line of Regression

Entire game is around slope and intercept of the line(s)

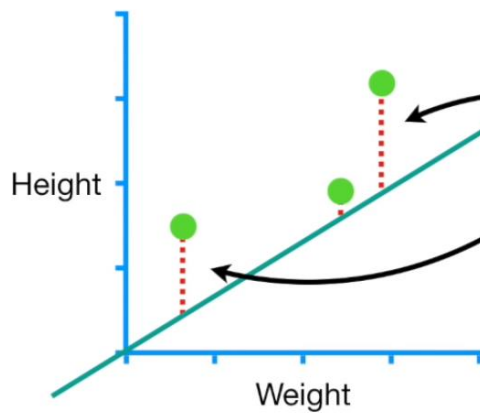
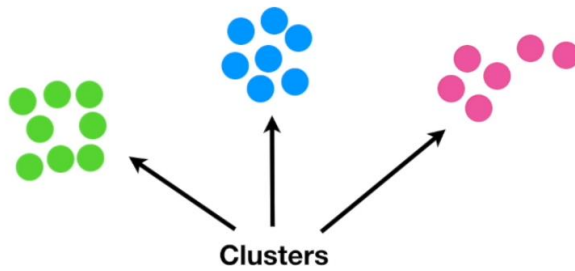
When we fit a line with **Linear Regression**, we optimize the **Intercept** and **Slope**.



When we use **Logistic Regression**, we optimize a squiggle.



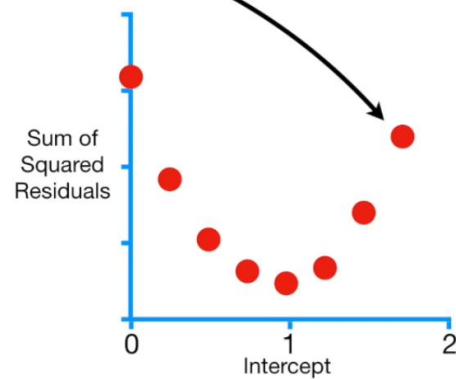
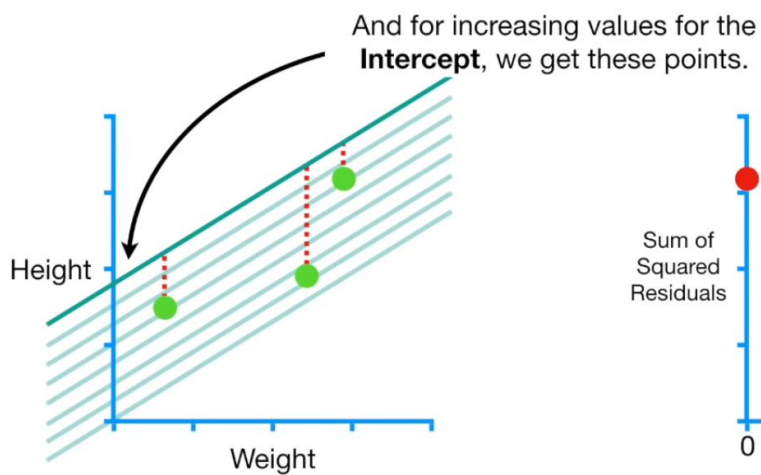
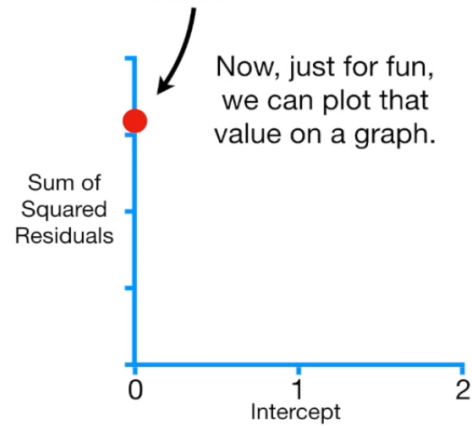
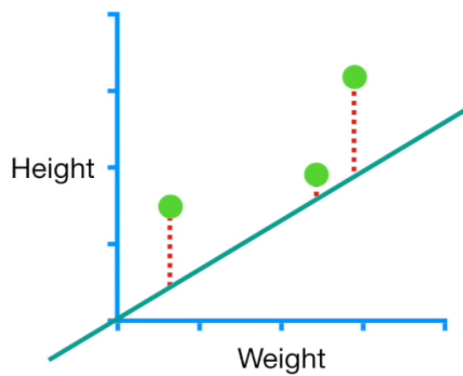
And when we use **t-SNE**, we optimize clusters.



In this example, we will evaluate how well this line fits the data with the **Sum of the Squared Residuals**.

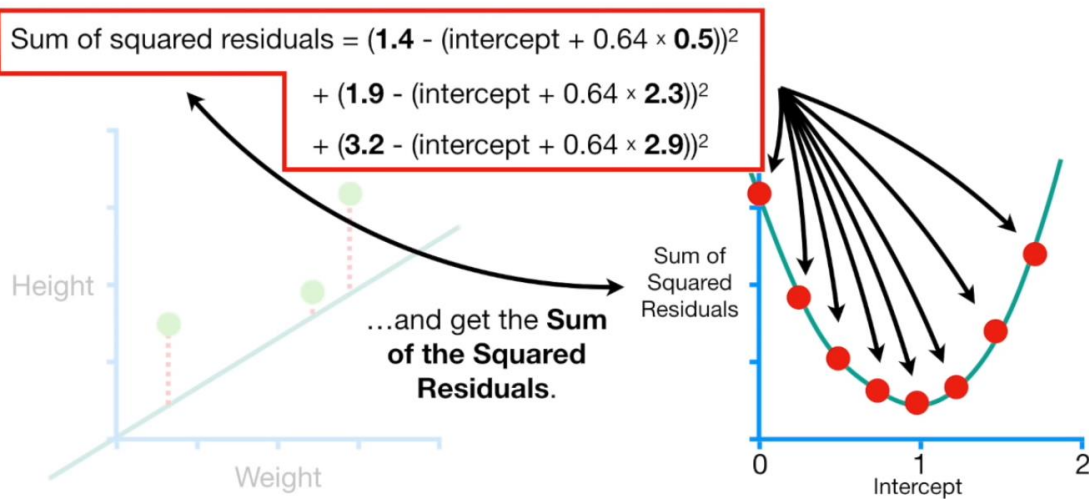
NOTE: In Machine Learning lingo, The Sum of the Squared Residuals is a type of **Loss Function**.

$$\text{Sum of squared residuals} = 1.1^2 + 0.4^2 + 1.3^2 = 3.1$$



Here we have several fitting lines, but we choose the best fit line (least residual value).

Here at intercept 1, we are getting least residual value, which mean line passing through the intercept with the respective slope is best fit line.



$$\frac{d}{d \text{ intercept}} (1.4 - (\text{intercept} + 0.64 \times 0.5))^2 = 2(1.4 - (\text{intercept} + 0.64 \times 0.5)) \times -1$$

$$\frac{d}{d \text{ intercept}} 1.4 - (\text{intercept} + 0.64 \times 0.5)$$

$$\frac{d}{d \text{ intercept}} \cancel{1.4} + (-1)\cancel{\text{intercept}} - 0.64 \times \cancel{0.5} = -1$$

These parts don't contain a term for the **Intercept**, so they go away.

We use **Chain Rule** for taking derivatives to check the slope of the line.

Negative sign is nothing but minus sign of the part of predicted value (line), and 1 is derivative of respective part is 1.

$$\frac{d}{d \text{ intercept}} (1.4 - (\text{intercept} + 0.64 \times 0.5))^2 = 2(1.4 - (\text{intercept} + 0.64 \times 0.5)) \times -1$$

$$= -2(1.4 - (\text{intercept} + 0.64 \times 0.5))$$

...and this...

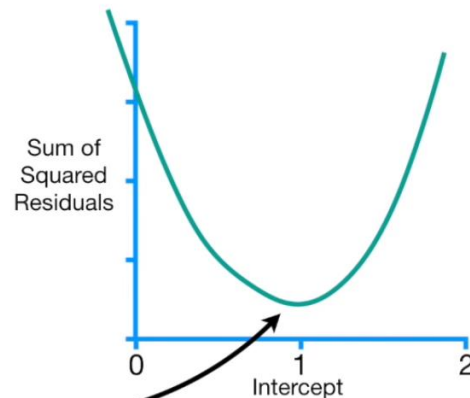
...is the derivative
of the first part...

$$\frac{d}{d \text{ intercept}} \text{ Sum of squared residuals} = \frac{d}{d \text{ intercept}} (1.4 - (\text{intercept} + 0.64 \times 0.5))^2$$

$$+ \frac{d}{d \text{ intercept}} (1.9 - (\text{intercept} + 0.64 \times 2.3))^2$$

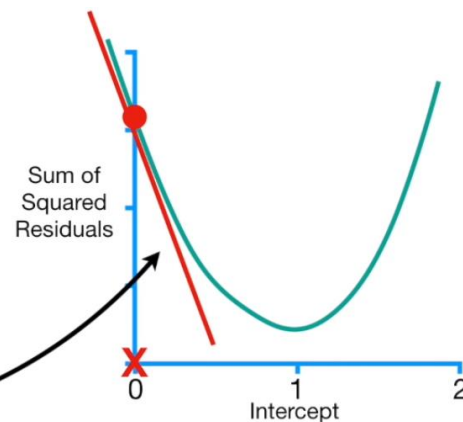
$$+ \frac{d}{d \text{ intercept}} (3.2 - (\text{intercept} + 0.64 \times 2.9))^2$$

$$\begin{aligned} \frac{d}{d \text{ intercept}} \text{ Sum of squared residuals} = & -2(1.4 - (\text{intercept} + 0.64 \times 0.5)) \\ & + -2(1.9 - (\text{intercept} + 0.64 \times 2.3)) \\ & + -2(3.2 - (\text{intercept} + 0.64 \times 2.9)) \end{aligned}$$



Now that we have the derivative,
Gradient Descent will use it to find
where the Sum of Squared
Residuals is lowest.

$$\begin{aligned} \frac{d}{d \text{ intercept}} \text{ Sum of squared residuals} = & -2(1.4 - (0 + 0.64 \times 0.5)) \\ & + -2(1.9 - (0 + 0.64 \times 2.3)) \\ & + -2(3.2 - (0 + 0.64 \times 2.9)) \\ = & -5.7 \end{aligned}$$



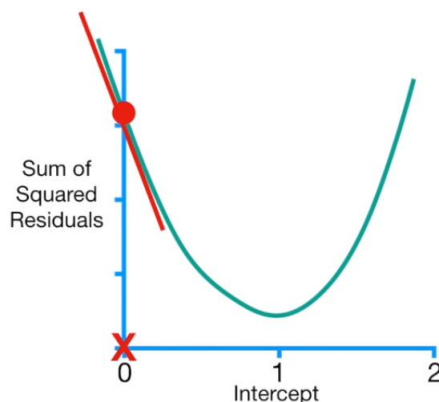
So when the **Intercept** = 0,
the slope of the curve = **-5.7**.

When slope is far from zero, we should take **big steps** (not much big, but moderate) and when we are close zero then we should take **baby steps**, so that we could get a very close value of the intercept.

$$\begin{aligned} \frac{d}{d \text{ intercept}} \text{ Sum of squared residuals} &= \\ &= -2(1.4 - (0 + 0.64 \times 0.5)) \\ &+ -2(1.9 - (0 + 0.64 \times 2.3)) \\ &+ -2(3.2 - (0 + 0.64 \times 2.9)) \\ &= -5.7 \end{aligned}$$

$$\text{Step Size} = -5.7 \times 0.1$$

NOTE: We'll talk more about **Learning Rates** later.

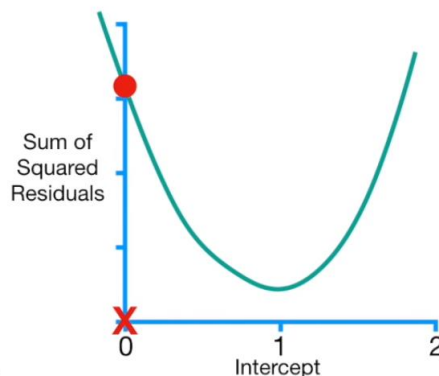


$$\begin{aligned} \frac{d}{d \text{ intercept}} \text{ Sum of squared residuals} &= \\ &= -2(1.4 - (0 + 0.64 \times 0.5)) \\ &+ -2(1.9 - (0 + 0.64 \times 2.3)) \\ &+ -2(3.2 - (0 + 0.64 \times 2.9)) \\ &= -5.7 \end{aligned}$$

$$\text{Step Size} = -5.7 \times 0.1 = -0.57$$

$$\text{New Intercept} = \text{Old Intercept} - \text{Step Size}$$

...minus the **Step Size**.



Step-size we kept by 0.1. We change it respectively along the calculations.

Negative sign indicates the slope going down, which means line is going towards down.

Our aim is to go close to zero.

$$\frac{d}{d \text{ intercept}} \text{ Sum of squared residuals} =$$

$$-2(1.4 - (0 + 0.64 \times 0.5))$$

$$+ -2(1.9 - (0 + 0.64 \times 2.3))$$

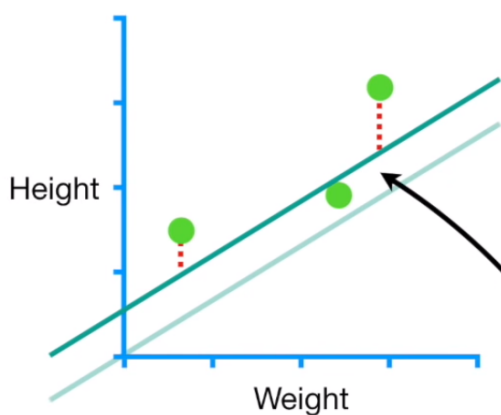
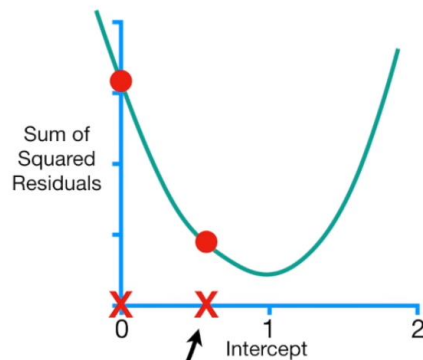
$$+ -2(3.2 - (0 + 0.64 \times 2.9))$$

$$= -5.7$$

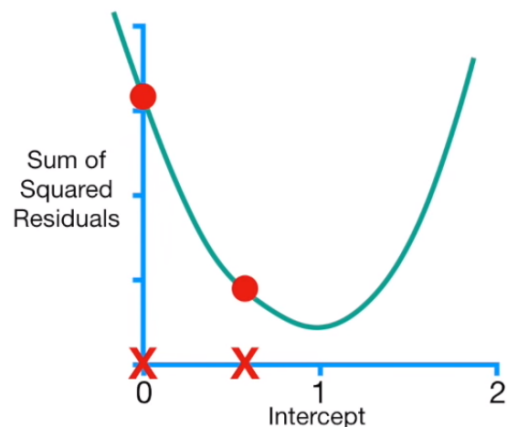
$$\text{Step Size} = -5.7 \times 0.1 = -0.57$$

$$\text{New Intercept} = 0 - (-0.57) = \boxed{0.57}$$

...and the the **New Intercept = 0.57.**



...we can see how much the residuals shrink when the **Intercept = 0.57.**



$$\frac{d}{d \text{ intercept}} \text{ Sum of squared residuals} =$$

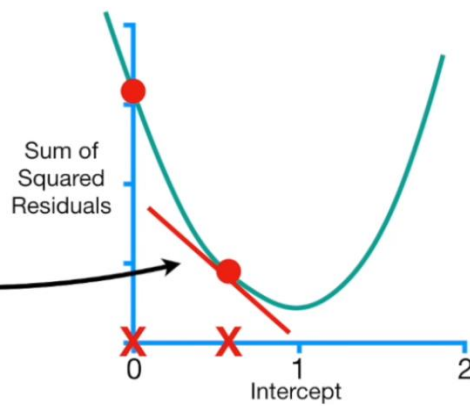
$$-2(1.4 - (0.57 + 0.64 \times 0.5))$$

$$+ -2(1.9 - (0.57 + 0.64 \times 2.3))$$

$$+ -2(3.2 - (0.57 + 0.64 \times 2.9))$$

$$= \boxed{-2.3}$$

...and that tells us the slope of the curve = **-2.3.**

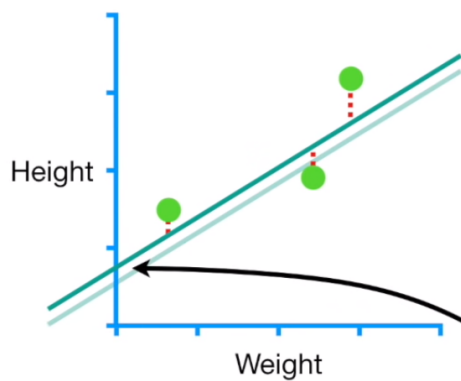
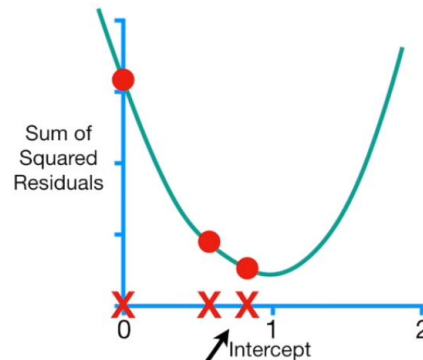


$$\begin{aligned} \frac{d}{d \text{ intercept}} \text{ Sum of squared residuals} &= \\ &= -2(1.4 - (0.57 + 0.64 \times 0.5)) \\ &\quad + -2(1.9 - (0.57 + 0.64 \times 2.3)) \\ &\quad + -2(3.2 - (0.57 + 0.64 \times 2.9)) \\ &= -2.3 \end{aligned}$$

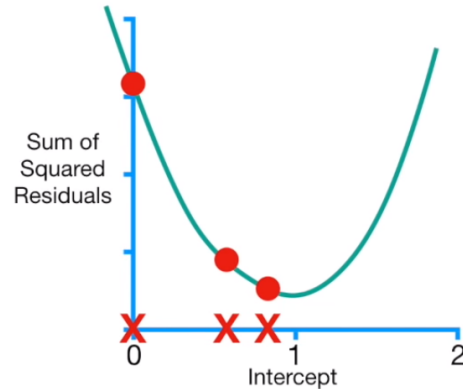
$$\text{Step Size} = -2.3 \times 0.1 = -0.23$$

$$\text{New Intercept} = 0.57 - (-0.23) = \boxed{0.8}$$

...and the **New Intercept = 0.8**

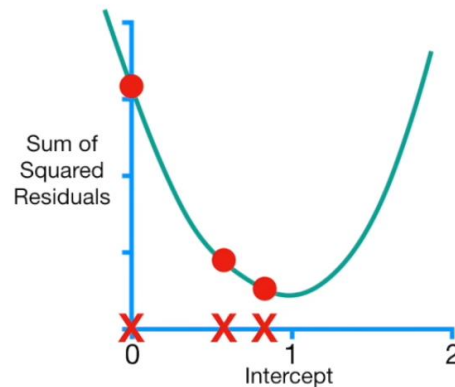


...to when the **Intercept = 0.8**



$$\begin{aligned} \frac{d}{d \text{ intercept}} \text{ Sum of squared residuals} &= \\ &= -2(1.4 - (0.8 + 0.64 \times 0.5)) \\ &\quad + -2(1.9 - (0.8 + 0.64 \times 2.3)) \\ &\quad + -2(3.2 - (0.8 + 0.64 \times 2.9)) \\ &= \boxed{-0.9} \end{aligned}$$

...and we get **-0.9**.

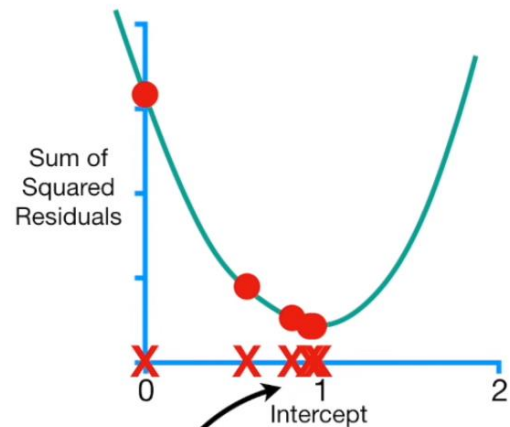
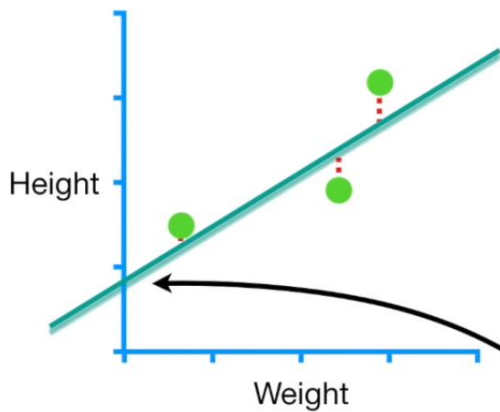
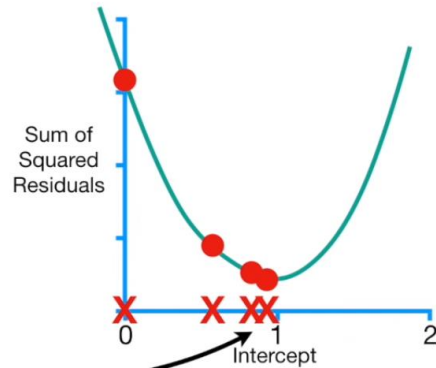


$$\begin{aligned} \frac{d}{d \text{ intercept}} \text{ Sum of squared residuals} &= \\ &= -2(1.4 - (0.8 + 0.64 \times 0.5)) \\ &+ -2(1.9 - (0.8 + 0.64 \times 2.3)) \\ &+ -2(3.2 - (0.8 + 0.64 \times 2.9)) \\ &= -0.9 \end{aligned}$$

$$\text{Step Size} = -0.9 \times 0.1 = -0.09$$

$$\text{New Intercept} = 0.8 - (-0.09) = \boxed{0.89}$$

...and the **New Intercept = 0.89**



Then we take another step and the **New Intercept = 0.92...**

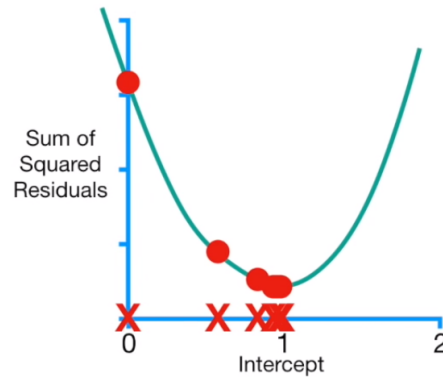
Big step and **Baby step** is automatically handle by Gradient Descent, as shown in above figure.

As we approach towards zero residuals, we get a proper slope. That's all the magic done by Gradient Descent.

After 6 steps, the **Gradient Descent** estimate for the **Intercept** is **0.95**.

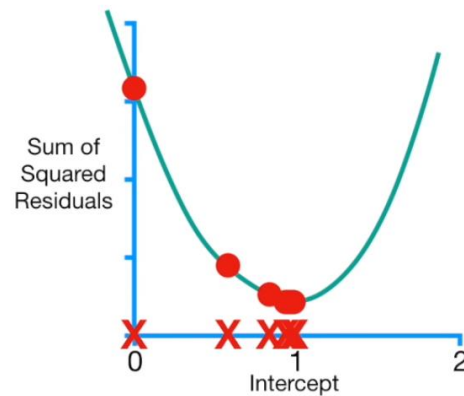
NOTE: The **Least Squares** estimate for the intercept is also **0.95**.

So we know that **Gradient Descent** has done its job, but without comparing its solution to a gold standard, how does **Gradient Descent** know to stop taking steps?



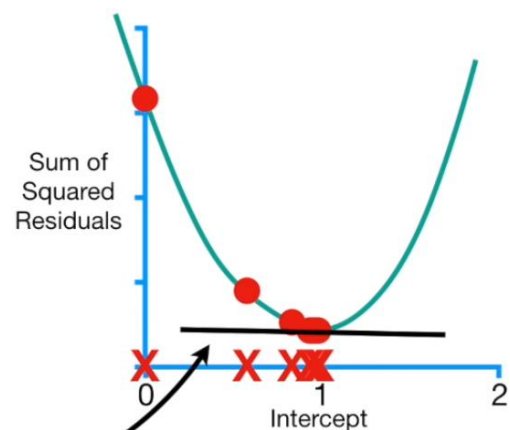
Gradient Descent stops when the **Step Size** is **Very Close To 0**.

$$\text{Step Size} = \text{Slope} \times \text{Learning Rate}$$



The **Step Size** will be **Very Close to 0** when the **Slope** is very close to **0**.

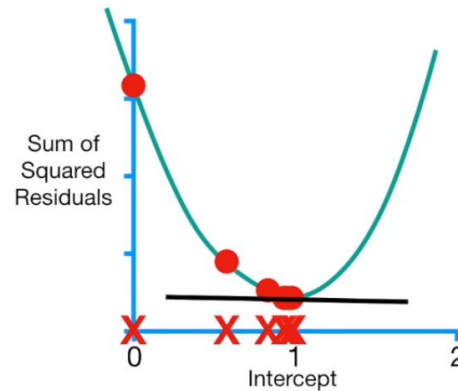
$$\text{Step Size} = \text{Slope} \times \text{Learning Rate}$$



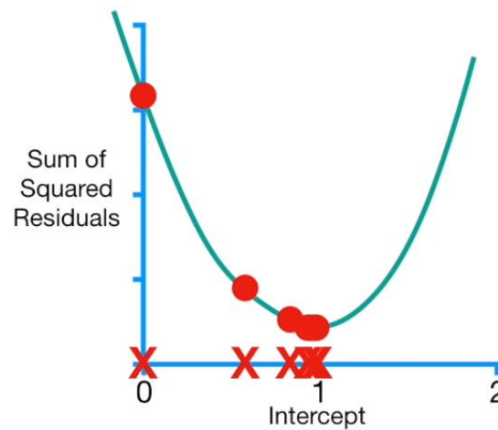
Then we would plug in
0.009 for the **Slope** and **0.1**
 for the **Learning Rate**..

↙ ↘

Step Size = 0.009 × 0.1



So, even if the **Step Size** is
 large, if there have been
 more than the **Maximum**
Number of Steps, Gradient
Descent will stop.



$$\frac{d}{d \text{ intercept}} \text{ Sum of squared residuals} =$$

$$-2(1.4 - (\text{intercept} + \text{slope} \times 0.5))$$

$$+ -2(1.9 - (\text{intercept} + \text{slope} \times 2.3))$$

$$+ -2(3.2 - (\text{intercept} + \text{slope} \times 2.9))$$

NOTE: When you have two or
 more derivatives of the same
 function, they are called a
Gradient.

$$\frac{d}{d \text{ slope}} \text{ Sum of squared residuals} =$$

$$-2 \times 0.5(1.4 - (\text{intercept} + \text{slope} \times 0.5))$$

$$+ -2 \times 2.9(3.2 - (\text{intercept} + \text{slope} \times 2.9))^2$$

$$+ -2 \times 2.3(1.9 - (\text{intercept} + \text{slope} \times 2.3))^2$$

$$\frac{d}{d \text{ intercept}} \text{ Sum of squared residuals} =$$

$$-2(1.4 - (\text{intercept} + \text{slope} \times 0.5))$$

$$+ -2(1.9 - (\text{intercept} + \text{slope} \times 2.3))$$

$$+ -2(3.2 - (\text{intercept} + \text{slope} \times 2.9))$$

We will use this **Gradient** to **descend** to lowest point in the **Loss Function**, which, in this case, is the Sum of the Squared Residuals...

...thus, this is why this algorithm is called **Gradient Descent**!

$$\frac{d}{d \text{ slope}} \text{ Sum of squared residuals} =$$

$$-2 \times 0.5(1.4 - (\text{intercept} + \text{slope} \times 0.5))$$

$$+ -2 \times 2.9(3.2 - (\text{intercept} + \text{slope} \times 2.9))^2$$

$$+ -2 \times 2.3(1.9 - (\text{intercept} + \text{slope} \times 2.3))^2$$

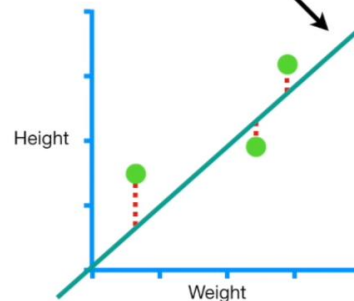
$$\frac{d}{d \text{ intercept}} \text{ Sum of squared residuals} =$$

$$-2(1.4 - (\text{intercept} + \text{slope} \times 0.5))$$

$$+ -2(1.9 - (\text{intercept} + \text{slope} \times 2.3))$$

$$+ -2(3.2 - (\text{intercept} + \text{slope} \times 2.9))$$

Thus, this line, with **Intercept = 0** and **Slope = 1**, is where we will start.



$$\frac{d}{d \text{ slope}} \text{ Sum of squared residuals} =$$

$$-2 \times 0.5(1.4 - (\text{intercept} + \text{slope} \times 0.5))$$

$$+ -2 \times 2.9(3.2 - (\text{intercept} + \text{slope} \times 2.9))^2$$

$$+ -2 \times 2.3(1.9 - (\text{intercept} + \text{slope} \times 2.3))^2$$

$$\begin{aligned} \frac{d}{d \text{ intercept}} \text{ Sum of squared residuals} = & -2(1.4 - (0 + 1 \times 0.5)) \\ & + -2(1.9 - (0 + 1 \times 2.3)) \\ & + -2(3.2 - (0 + 1 \times 2.9)) = -1.6 \end{aligned}$$

$$\text{Step Size}_{\text{Intercept}} = -1.6 \times 0.01 = -0.016$$

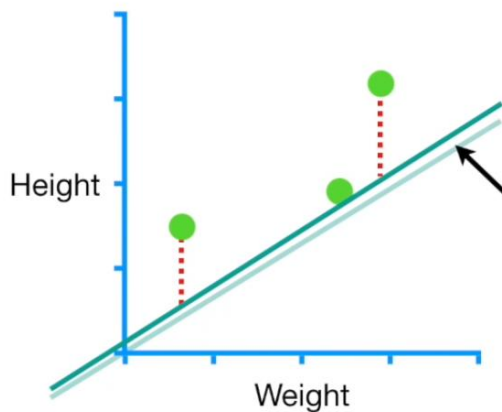
$$\text{New Intercept} = 0 - (-0.016) = 0.016$$

...and we end up
with a **New Intercept**
and a **New Slope**.

$$\begin{aligned} \frac{d}{d \text{ slope}} \text{ Sum of squared residuals} = & -2 \times 0.5(1.4 - (0 + 1 \times 0.5)) \\ & + -2 \times 2.9(3.2 - (0 + 1 \times 2.9))^2 \\ & + -2 \times 2.3(1.9 - (0 + 1 \times 2.3))^2 = -0.8 \end{aligned}$$

$$\text{Step Size}_{\text{Slope}} = -0.8 \times 0.01 = -0.008$$

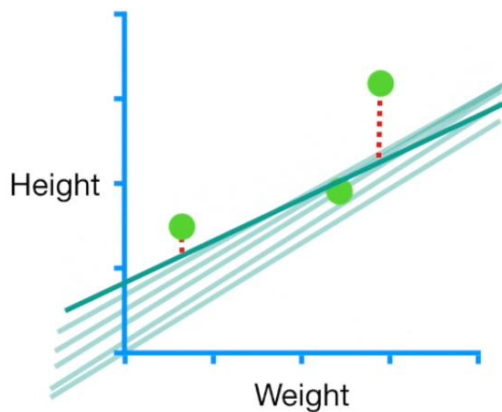
$$\text{New Slope} = 1 - (-0.008) = 1.008$$



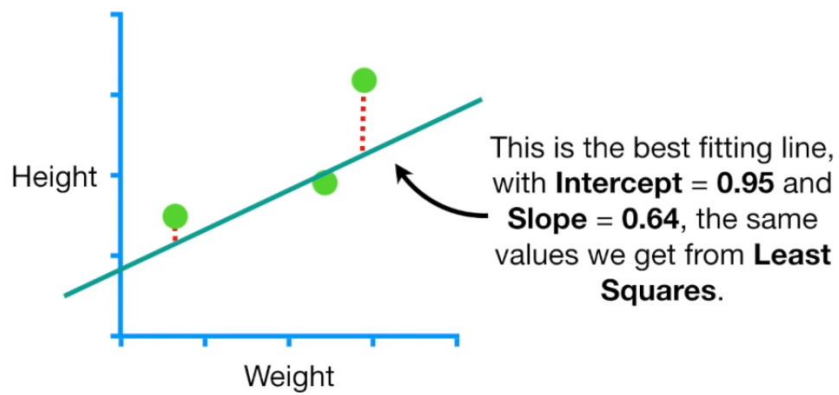
$$\text{New Intercept} = 0 - (-0.016) = 0.016$$

...and this is the new line
(with **Slope = 1.008** and
Intercept = 0.016) after
the first step.

$$\text{New Slope} = 1 - (-0.008) = 1.008$$



Now we just repeat what we did
until all of the **Steps Sizes** are very
small or we reach the **Maximum**
Number of Steps.



Here the points are very less, but what if the points are in millions of values. Hence we approach to **Stochastic Gradient Descent**. Because, for less points the mathematics is simpler but for huge values it's not. It is great if we have tones of data and lots of parameters. The mathematics is all same like regular gradient descent, but stochastic descent takes either a single value from the data or from the cluster for each step.