Ordinary Least Squares (OLS)

Gradient Descent (GD)

Maximum Likelihood Estimation (MLE)

# OLS: Ordinary Least Square

| Symbols | Meaning |
|---------|---------|
| $x$ | Independent variable data from observation |
| $\bar{x}$ | Mean of $x$ |
| $y$ | Dependent variable data from observation |
| $\bar{y}$ | Mean of $y$ |
| $\hat{y}$ | Estimate of $y$ by the regression model |
| $n$ | Number of observations |

**Steps:**

1. Get the difference (error): $(y\text{-}\hat{y})$

2. Square the difference: $(y\text{-}\hat{y})^2$

3. Take the sum for all data: $\sum (y-\hat{y})^2$

This is total error. Our objective is to keep this as minimum as possible.

$$Y = f(x) = 4(x - 3)^2 + 5$$

$$SSE = f(?) = \sum (y - \hat{y})^2 = \sum (y - mx - c)^2$$

$$SSE = f(?) = \sum (y - \hat{y})^2 = \sum (y - \theta_1 x - \theta_0)^2$$

$$SSE = f(?) = \sum (y - \hat{y})^2 = \sum (y - \beta_1 x - \beta_0)^2$$

$$SSE = f(?) = \sum (y - \hat{y})^2 = \sum (y - ax - b)^2$$

$$SSE = f(?) = \sum_i^n (y_i - \hat{y}_i)^2 = \sum_i^n (y_i - ax_i - b)^2$$

$$Y = f(x) = 4(x-3)^2 + 5$$

$$SSE = f(m, c) = \sum (y - \hat{y})^2 = \sum (y - mx - c)^2$$
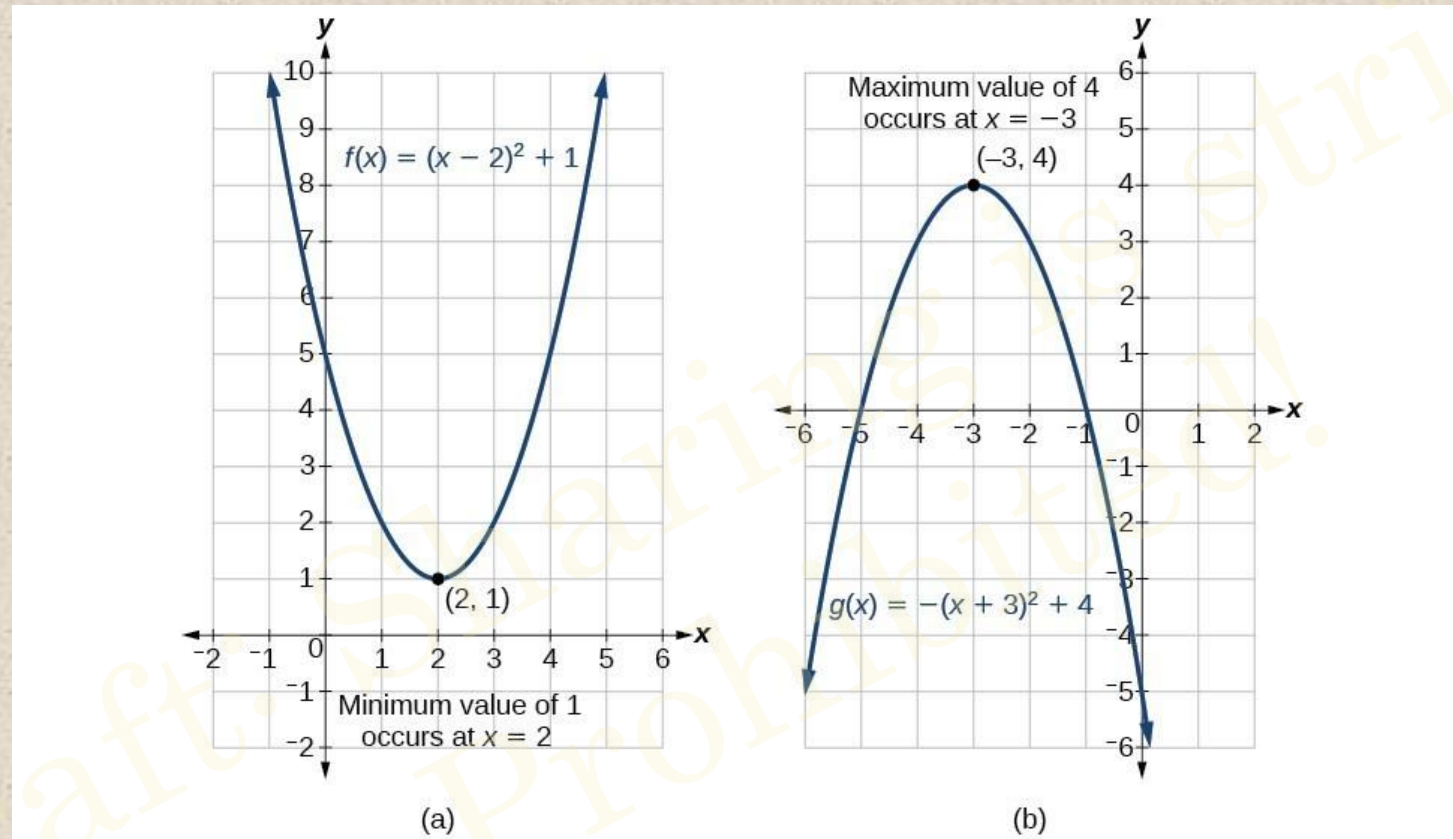
$$SSE = f(\theta_0, \theta_1) = \sum (y - \hat{y})^2 = \sum (y - \theta_1 x - \theta_0)^2$$

$$SSE = f(\beta_0, \beta_1) = \sum (y - \hat{y})^2 = \sum (y - \beta_1 x - \beta_0)^2$$
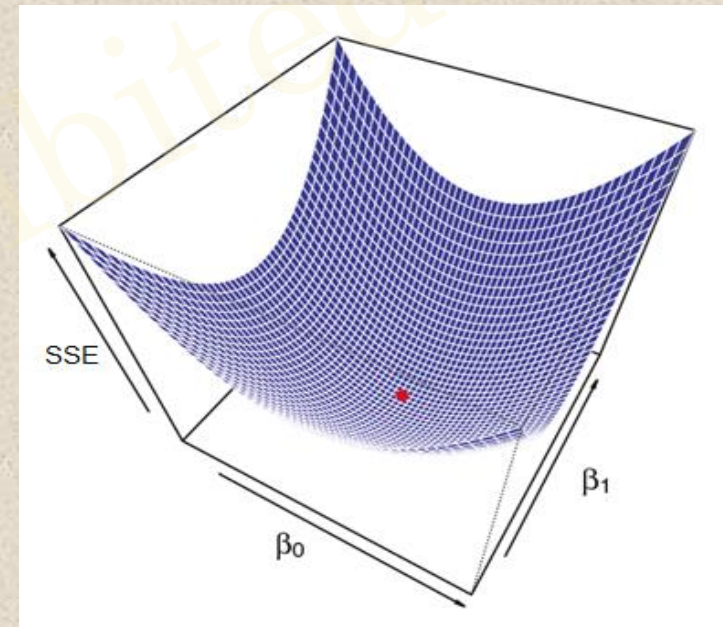
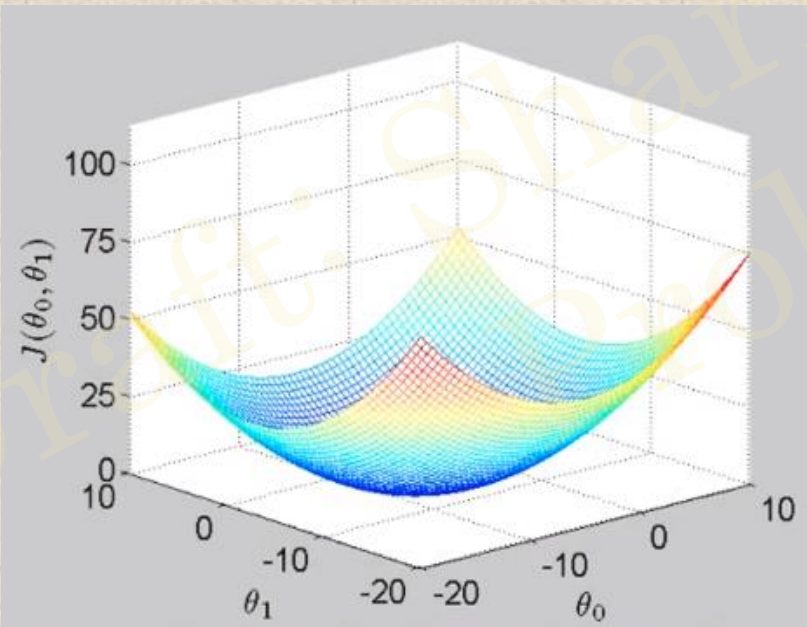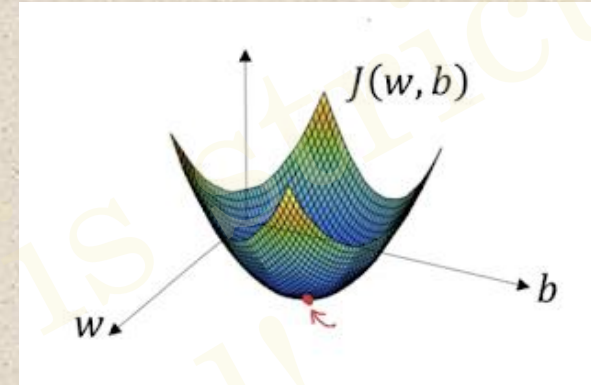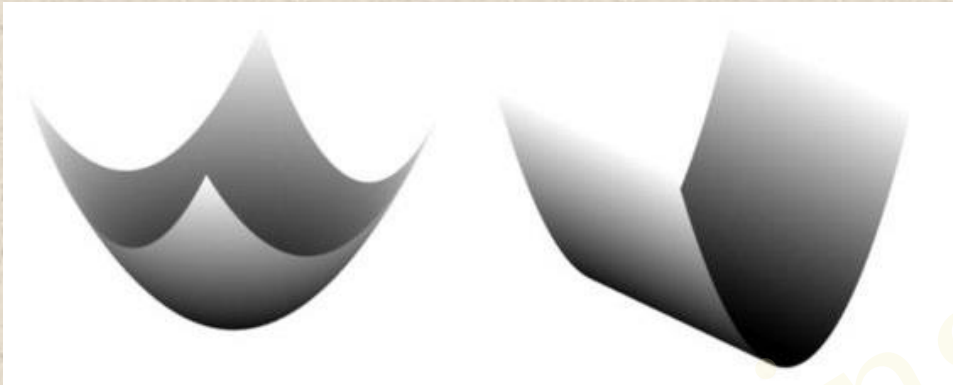$$SSE = f(a, b) = \sum (y - \hat{y})^2 = \sum (y - ax - b)^2$$

$$SSE = f(a, b) = \sum_i^n (y_i - \hat{y}_i)^2 = \sum_i^n (y_i - ax_i - b)^2$$

(a) $f(x) = (x - 2)^2 + 1$, Minimum value of 1 occurs at $x = 2$, point $(2, 1)$

(b) $g(x) = -(x + 3)^2 + 4$, Maximum value of 4 occurs at $x = -3$, point $(-3, 4)$

Differentiate **y**, Set its value to 0, solve the equation to find the value of **x**.

$$SSE = f(a, b) = \sum_i^n (y_i - \hat{y}_i)^2 = \sum_i^n (y_i - ax_i - b)^2$$

Give me $(a, b)$, where the value of SSE is minimum.

Differentiate SSE partially:

a)   With respect to $a$, Set its value to $0$, Solve the equation to find the value of $a$.

b)   With respect to $b$, Set its value to $0$, Solve the equation to find the value of $b$.

Let us denote SSE as S for simplicity:

$$S = \sum (y - \hat{y})^2 = \sum (y - ax - b)^2$$

$$\frac{\partial S}{\partial a} = 0$$

$$\frac{\partial S}{\partial a} = \frac{\partial \left( \sum (y - ax - b)^2 \right)}{\partial a} = 2 \sum \left( (y - ax - b) \cdot (0 - x - 0) \right)$$

$$2 \sum \left( (y - ax - b) \cdot (-x) \right) = 0$$

$$\sum (-xy) + a \sum x^2 + b \sum x = 0$$

$$\sum x = n\overline{x}$$

$$b = \frac{\sum xy - a \sum x^2}{n\overline{x}}$$

$$\frac{\partial S}{\partial b} = 0$$

$$\frac{\partial S}{\partial b} = \frac{\partial \left( \sum (y - ax - b)^2 \right)}{\partial b} = 2 \sum (y - ax - b) \cdot (0 - 0 - 1))$$

$$-2 \sum (y - ax - b) = 0$$

$$-\sum y + a \sum x + b \sum 1 = 0$$

$$\sum 1 = n \quad \sum x = n\overline{x} \quad \sum y = n\overline{y}$$

$$-n\overline{y} + an\overline{x} + nb = 0 \qquad a\overline{x} + b = \overline{y}$$

$$a\overline{x} + \frac{\sum xy}{n\overline{x}} - \frac{a \sum x^2}{n\overline{x}} = \overline{y}$$

$$a \left( \overline{x} - \frac{\sum x^2}{n\overline{x}} \right) + \frac{\sum xy}{n\overline{x}} = \overline{y}$$

$$a \left( n\overline{x}^2 - \sum x^2 \right) + \sum xy = n\overline{xy}$$

$$a = \frac{n\overline{x}\,\overline{y} - \sum xy}{\left( n\overline{x}^2 - \sum x^2 \right)}$$

$$\hat{y} = slope * x + intercept$$

$$slope = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - \left( \sum x \right)^2}$$

$$intercept = \overline{y} - slope \cdot \overline{x}$$

Mother of Dragons

steep slope
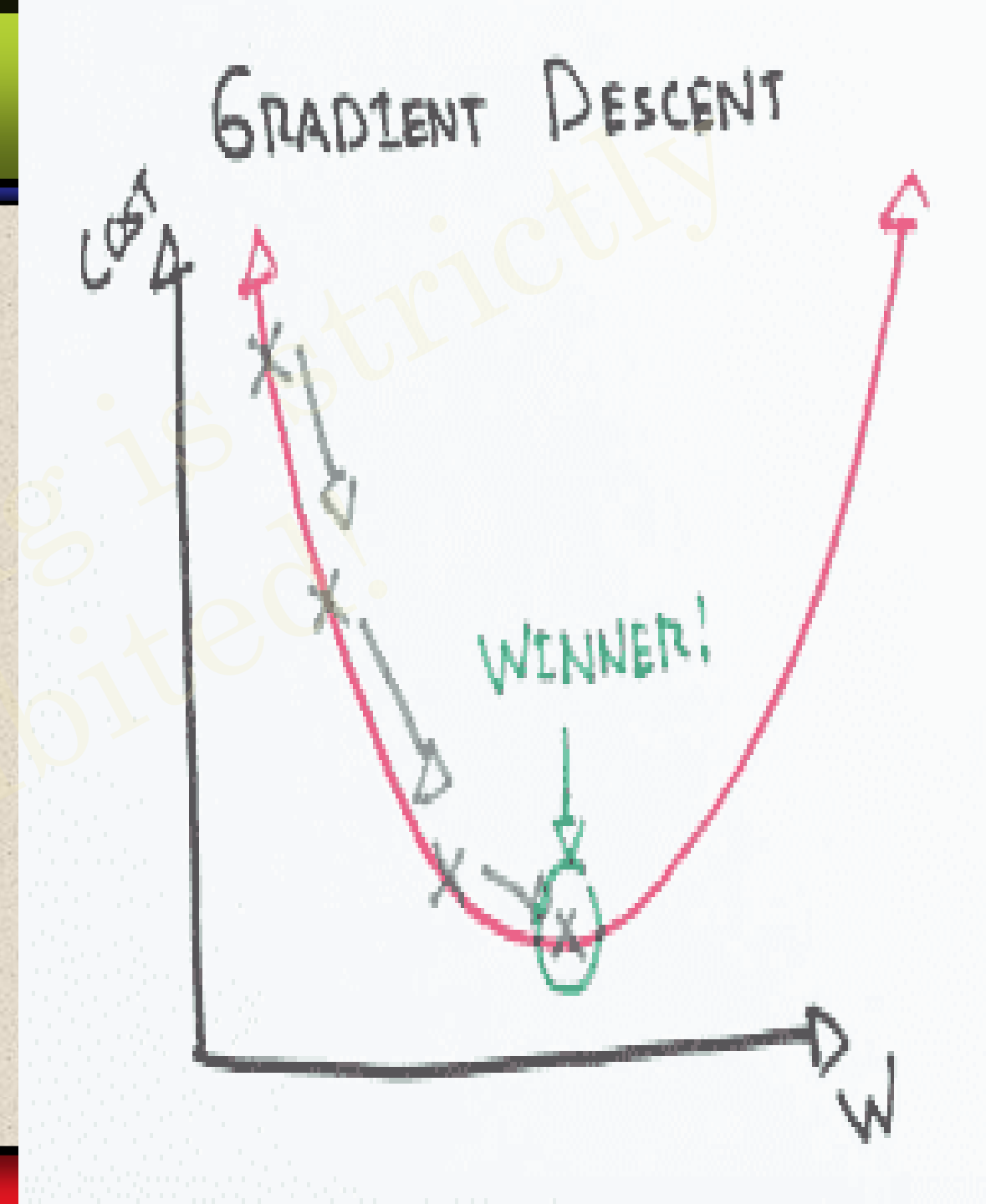Value of D is high
So take large steps

slope is less steep
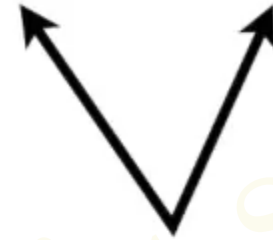Value of D is low
So take small steps

Goal

Mother of ML Algorithms

GRADIENT DESCENT

COST

WINNER!

W

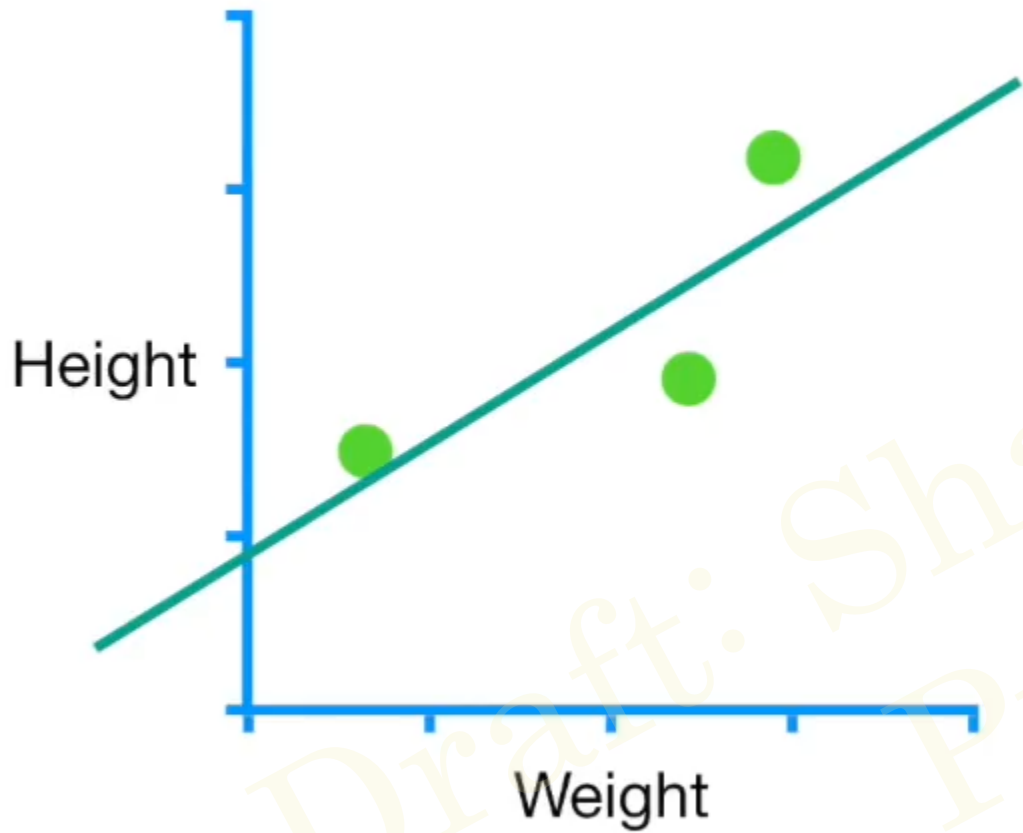# Gradient Descent

**Predicted Height** = intercept + slope × **Weight**

Height

Weight

So let's learn how **Gradient Descent** can fit a line to data by finding the optimal values for the **Intercept** and the **Slope**.

**Predicted Height** = intercept + slope × **Weight**

Actually, we'll start by using **Gradient Descent** to find the **Intercept**.

Height

Weight

**Predicted Height** = $\boxed{\text{intercept}}$ + $\boxed{\text{slope}}$ × **Weight**



Height

Weight

Then, once we understand how **Gradient Descent** works, we'll use it to solve for the **Intercept** *and* the **Slope**.

**Predicted Height** = intercept + slope × **Weight**



So for now, let's just plug in the **Least Squares** estimate for the **Slope**, **0.64**.

Height

Weight

**Predicted Height** = intercept + 0.64 × **Weight**



The first thing we do is pick a random value for the **Intercept**.

**Predicted Height** = intercept + 0.64 × **Weight**



Height

Weight

The first thing we do is pick a random value for the **Intercept**.

This is just an initial guess that gives **Gradient Descen**t something to improve upon.

**Predicted Height** = $\boxed{0}$ + 0.64 × **Weight**

In this case, we'll use **0**,
but any number will do.

Height

Weight

**Predicted Height =** $\boxed{0}$ **+ 0.64 × Weight**

Height

Weight

And that gives us the
equation for this line.

Height

(0.5, 1.4)

Weight

This datapoint represents a person with **Weight 0.5** and **Height 1.4**.

We get the **Predicted Height**, the point on the line…

…by plugging **Weight = 0.5** into the equation for the line…

**Predicted Height** = 0 + 0.64 × **Weight**

We get the **Predicted Height**, the point on the line…

…by plugging **Weight = 0.5** into the equation for the line…

**Predicted Height** = 0 + 0.64 × **0.5**

...and the **Predicted Height** is **0.32**.

(0.5, 1.4)

Height

Weight

**Predicted Height** = 0 + 0.64 × **0.5** = **0.32**

The residual is the difference between the **Observed Height**, and the **Predicted Height**…

Height

(0.5, 1.4)

X

Weight

**Predicted Height** = 0 + 0.64 × **0.5** = **0.32**

Sum of squared residuals = $1.1^2 + 0.4^2$

Height

Weight

(2.3, 1.9)

...the second residual is **0.4**...

Sum of squared residuals = $1.1^2 + 0.4^2 + 1.3^2$



...and the third residual is **1.3**.

(2.9, 3.2)

Height

Weight

Sum of squared residuals = $1.1^2 + 0.4^2 + 1.3^2 =$ 3.1

In the end, **3.1** is the Sum of the Squared Residuals.

Height

Weight

Sum of squared residuals = $1.1^2 + 0.4^2 + 1.3^2 =$ 3.1

Now, just for fun, we can plot that value on a graph.

Height

Weight

Sum of Squared Residuals

0          1          2

Intercept

Sum of squared residuals = $1.1^2 + 0.4^2 + 1.3^2 = 3.1$



Height

Weight

Sum of Squared Residuals

This graph has the Sum of Squared Residuals on the y-axis...

0

1

2

Intercept

# Sum of squared residuals = $1.1^2 + 0.4^2 + 1.3^2 = 3.1$



Height

Weight

Sum of Squared Residuals

…and different values for the **Intercept** on the x-axis.

0    1    2

Intercept

Sum of squared residuals = $1.1^2 + 0.4^2 + 1.3^2 = 3.1$



This point represents the Sum of the Squared Residuals when the **Intercept = 0**.

However, if the
**Intercept = 0.25**…

Height

Weight

Sum of
Squared
Residuals

0　　　　　1　　　　　2

Intercept

And for increasing values for the **Intercept**, we get these points.

Height

Weight

Sum of Squared Residuals

0          1          2
Intercept

Of the points that we calculated for the graph, this one has the lowest Sum of Squared Residuals...



Sum of Squared Residuals

0                    1                    2

Intercept

What if the best value
for the **Intercept** is
somewhere between
these values?

Sum of
Squared
Residuals

0          1          2

Intercept

**Gradient Descent** identifies the optimal value by taking big steps when it is far away...

Sum of Squared Residuals

Intercept

0    1    2

So let's get back to using **Gradient Descent** to find the optimal value for the **Intercept**, starting from a random value. In this case, the random value was **0**.

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

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

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

Thus, we now have an equation for this curve…

Sum of Squared Residuals

0          1          2

Intercept

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

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

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

So let's take the derivative of the Sum of the Squared Residuals with respect to the **Intercept**.

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

$\quad$ -2(**1.4** - (intercept + 0.64 × **0.5**)

$\quad$ + -2(**1.9** - (intercept + 0.64 × **2.3**))

$\quad$ + -2(**3.2** - (intercept + 0.64 × **2.9**))

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



Sum of Squared Residuals

0     1     2

Intercept

$$\frac{d}{d\ intercept}$$

Sum of squared residuals =

    -2(**1.4** - (intercept + 0.64 × **0.5**)

    + -2(**1.9** - (intercept + 0.64 × **2.3**))

    + -2(**3.2** - (intercept + 0.64 × **2.9**))

**NOTE:** If we were using **Least Squares** to solve for the optimal value for the **Intercept**, we would simply find where the the slope of the curve = **0**.



Sum of Squared Residuals

0           1           2

Intercept

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

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

$$+ -2(\mathbf{1.9} - (intercept + 0.64 \times \mathbf{2.3}))$$

$$+ -2(\mathbf{3.2} - (intercept + 0.64 \times \mathbf{2.9}))$$

Sum of
Squared
Residuals

In contrast, **Gradient Descent** finds the minimum value by taking steps from an initial guess until it reaches the best value.

0          1          2

Intercept

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

   -2(**1.4** - (intercept + 0.64 × **0.5**)

 + -2(**1.9** - (intercept + 0.64 × **2.3**))

 + -2(**3.2** - (intercept + 0.64 × **2.9**))

This makes **Gradient Descent** very useful when it is not possible to solve for where the derivative = **0**, and this is why **Gradient Descent** can be used in so many different situations.

Sum of Squared Residuals

X   X XXX

0          1          2

Intercept

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

$\quad$ -2(**1.4** - (intercept + 0.64 × **0.5**)

$\quad$ + -2(**1.9** - (intercept + 0.64 × **2.3**))

$\quad$ + -2(**3.2** - (intercept + 0.64 × **2.9**))

Remember, we started by setting
the **Intercept** to a random number.
In this case, that was **0**.

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

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

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

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

So we plug **0** into
the derivative…



Sum of
Squared
Residuals

0          1          2

Intercept

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

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

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

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

$= -5.7$

...and we get **-5.7**.

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

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

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

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

$$= -5.7$$

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

Sum of
Squared
Residuals

0          1          2

Intercept

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

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

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

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

$$= -5.7$$

Sum of
Squared
Residuals

0          1          2

Intercept

**NOTE:** The closer we get to the optimal value for the **Intercept**, the closer the slope of the curve gets to **0**.

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

-2(**1.4** - (0 + 0.64 × **0.5**)

+ -2(**1.9** - (0 + 0.64 × **2.3**))

+ -2(**3.2** - (0 + 0.64 × **2.9**))

= -5.7

Sum of
Squared
Residuals

0          1          2

Intercept

This means that when
the slope of the curve is
close to **0**…

**X**

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

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

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

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

$$= -5.7$$



…then we should take baby steps, because we are close to the optimal value…

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

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

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

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

$= -5.7$

Sum of
Squared
Residuals

0    1    2

Intercept

...and when the slope is
far from **0**...

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

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

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

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

$$= -5.7$$

...then we should take big steps, because we are far from the optimal value.

Sum of Squared Residuals

0        1        2

Intercept

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

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

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

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

$$= -5.7$$

Sum of Squared Residuals

0        1        2

Intercept

However, if we take a super huge step...

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

-2(**1.4** - (0 + 0.64 × **0.5**)

+ -2(**1.9** - (0 + 0.64 × **2.3**))

+ -2(**3.2** - (0 + 0.64 × **2.9**))

= -5.7

Sum of
Squared
Residuals

0          1          2

Intercept

...then we would increase
the Sum of the Squared
Residuals!

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

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

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

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

$$= -5.7$$



So the size of the step should be related to the slope, since it tells us if we should take a baby step or a big step, but we need to make sure the big step is not too big.

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

-2(**1.4** - (0 + 0.64 × **0.5**)

+ -2(**1.9** - (0 + 0.64 × **2.3**))

+ -2(**3.2** - (0 + 0.64 × **2.9**))

= -5.7

**Step Size** = -5.7

**Gradient Descent** determines the
**Step Size** by multiplying the **slope**...



Sum of
Squared
Residuals

0          1          2
        Intercept

$$\frac{d}{d\ intercept}$$

Sum of squared residuals =
   -2(**1.4** - (0 + 0.64 × **0.5**)

   + -2(**1.9** - (0 + 0.64 × **2.3**))

   + -2(**3.2** - (0 + 0.64 × **2.9**))

   = -5.7

**Step Size** = -5.7 × 0.1

...by a small number called
**The Learning Rate**.



Sum of
Squared
Residuals

0          1          2
          Intercept

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

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

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

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

$= -5.7$

**Step Size** = -5.7 × 0.1 = -0.57

When the **Intercept** = **0**, the **Step Size** = **-0.57**.

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

-2(**1.4** - (0 + 0.64 × **0.5**)

+ -2(**1.9** - (0 + 0.64 × **2.3**))

+ -2(**3.2** - (0 + 0.64 × **2.9**))

= -5.7

**Step Size** = -5.7 × 0.1 = -0.57

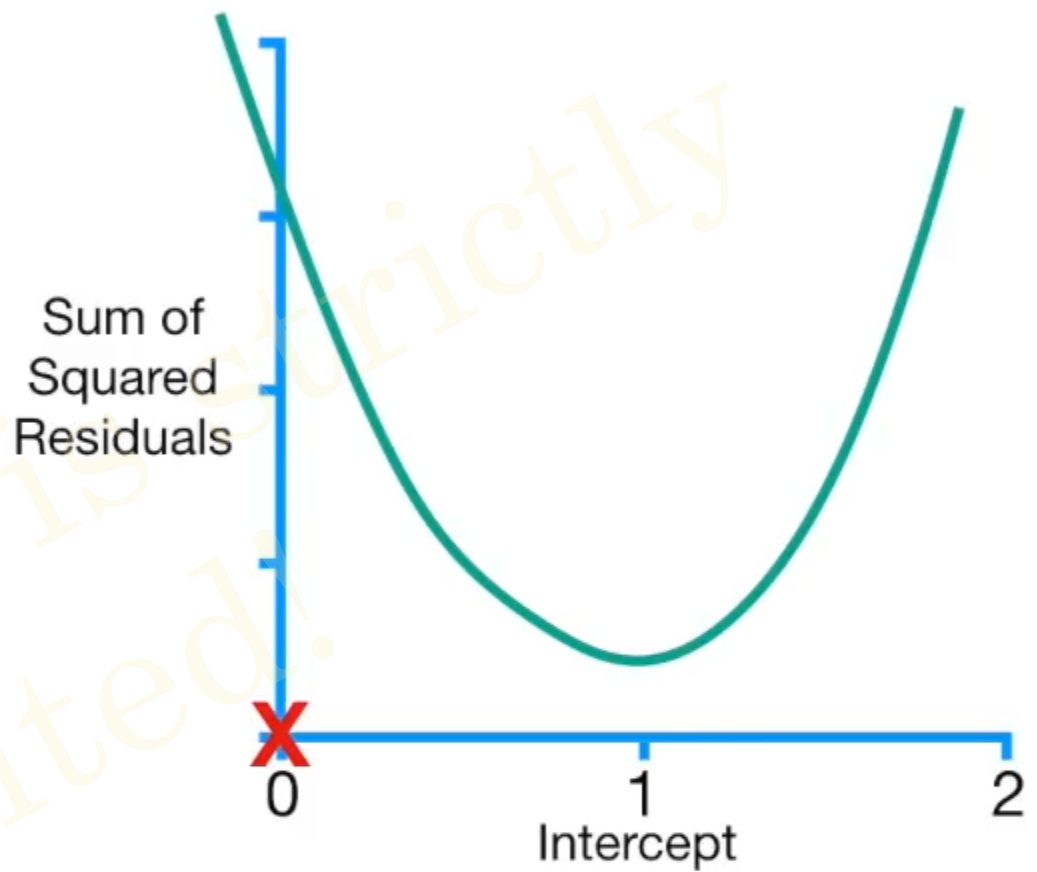**New Intercept** = ⟵

With the **Step Size**, we can calculate a **New Intercept**.

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

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

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

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

$= -5.7$

**Step Size** = -5.7 × 0.1 = -0.57

**New Intercept = Old Intercept - Step Size**

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

Sum of
Squared
Residuals

0          1          2

Intercept

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

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

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

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

$$= -5.7$$

**Step Size** = -5.7 × 0.1 = -0.57

**New Intercept** = 0 - (-0.57) = **0.57**

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

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

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

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

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

$$= -5.7$$

**Step Size** = -5.7 × 0.1 = -0.57

**New Intercept** = 0 - (-0.57) = **0.57**

Sum of
Squared
Residuals

0          1          2

Intercept

In one big step, we moved
much closer to the optimal
value for the **Intercept**.

Going back to the original data and the original line, with the **Intercept** = **0**...

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

Now let's take another step
closer to the optimal value
for the **Intercept**.

$$\frac{d}{d\ intercept}$$ Sum of squared residuals =
$$-2(\mathbf{1.4} - (0.57 + 0.64 \times \mathbf{0.5})$$
$$+\ -2(\mathbf{1.9} - (0.57 + 0.64 \times \mathbf{2.3}))$$
$$+\ -2(\mathbf{3.2} - (0.57 + 0.64 \times \mathbf{2.9}))$$

To take another step, we go back to the derivative and plug in the **New Intercept** (**0.57**)...



Sum of Squared Residuals

0          1          2

Intercept

$$\frac{d}{d \text{ intercept}}$$

Sum of squared residuals =

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

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

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

$= \boxed{-2.3}$

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

Sum of Squared Residuals

0          1          2

Intercept

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

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

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

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

$$= -2.3$$

**Step Size = Slope × Learning Rate**

Now let's calculate the
**Step Size**…



Sum of
Squared
Residuals

Intercept

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

-2(**1.4** - (0.57 + 0.64 × **0.5**)

+ -2(**1.9** - (0.57 + 0.64 × **2.3**))

+ -2(**3.2** - (0.57 + 0.64 × **2.9**))

= -2.3

**Step Size** = -2.3 × 0.1 = **-0.23**

Ultimately, the **Step Size** is **-0.23**...



Sum of
Squared
Residuals

0          1          2

Intercept

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

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

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

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

$$=\ -2.3$$

**Step Size** = -2.3 × 0.1 = **-0.23**



Sum of
Squared
Residuals

Intercept

**New Intercept** = 0.57 - (-0.23) = **0.8**

…and the **New Intercep**t = **0.8**

Height

Weight

Sum of
Squared
Residuals

0          1          2

Intercept

Now we can compare the
residuals when the
**Intercept** = **0.57**…

Height

Weight

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

Sum of
Squared
Residuals

0    1    2

Intercept

Height

Weight

Sum of
Squared
Residuals

0          1          2

Intercept

Overall, the Sum of the
Squared Residuals is getting
smaller.

$$\frac{d}{d\ intercept}$$

Sum of squared residuals =

-2(**1.4** - (0.8 + 0.64 × **0.5**)

+ -2(**1.9** - (0.8 + 0.64 × **2.3**))

+ -2(**3.2** - (0.8 + 0.64 × **2.9**))

Now let's calculate the derivative at the **New Intercept** (**0.8**)...



Sum of Squared Residuals

0          1          2

Intercept

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

$-2(\textbf{1.4} - (0.8 + 0.64 \times \textbf{0.5})$

$+ -2(\textbf{1.9} - (0.8 + 0.64 \times \textbf{2.3}))$

$+ -2(\textbf{3.2} - (0.8 + 0.64 \times \textbf{2.9}))$

$= \boxed{\textbf{-0.9}}$

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

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

$$-2(\mathbf{1.4} - (0.8 + 0.64 \times \mathbf{0.5})$$

$$+ -2(\mathbf{1.9} - (0.8 + 0.64 \times \mathbf{2.3}))$$

$$+ -2(\mathbf{3.2} - (0.8 + 0.64 \times \mathbf{2.9}))$$

$$= \mathbf{-0.9}$$

**Step Size** = -0.9 × 0.1 = **-0.09**

The **Step Size** = **-0.09**…

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

$$-2(\mathbf{1.4} - (0.8 + 0.64 \times \mathbf{0.5})$$

$$+ \ -2(\mathbf{1.9} - (0.8 + 0.64 \times \mathbf{2.3}))$$

$$+ \ -2(\mathbf{3.2} - (0.8 + 0.64 \times \mathbf{2.9}))$$

$$= \mathbf{-0.9}$$

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

**New Intercept** = 0.8 - (-0.09) = **0.89**

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



Sum of
Squared
Residuals

Intercept

0  1  2

Sum of Squared Residuals

0    1    2

Intercept

Notice how each step gets smaller and smaller the closer we get to the bottom of the curve.

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

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

**Step Size** = **Slope** × **Learning Rate**



Sum of Squared Residuals

Intercept

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**.

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**.

**Step Size** = **Slope** × **Learning Rate**

plug in

**0.009** for the **Slope** and **0.1**
for the **Learning Rate**..

**Step Size** = **0.009** × **0.1**



Sum of
Squared
Residuals

0          1          2

Intercept

...and get **0.0009**, which is smaller than **0.001**, so **Gradient Descent** would stop.

**Step Size** = **0.009** × **0.1** = 0.0009

That said, **Gradient Descent** also includes a limit on the number of steps it will take before giving up.

That said, **Gradient Descent** also includes a limit on the number of steps it will take before giving up.

In practice, the **Maximum Number of Steps** = **1,000** or greater.



Sum of Squared Residuals

0          1          2

Intercept

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

OK, let's review what we've learned so far…

The first thing we did is decide to use the Sum of the Squared Residuals as the **Loss Function** to evaluate how well a line fits the data…

Height

Weight

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

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

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

...then we took the derivative of the Sum of the Squared Residuals. In other words, we took the derivative of the **Loss Function**...

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

-2(**1.4** - (intercept + 0.64 × **0.5**)

+ -2(**1.9** - (intercept + 0.64 × **2.3**))

+ -2(**3.2** - (intercept + 0.64 × **2.9**))

Sum of
Squared
Residuals

0                    1                    2

Intercept

…then we picked a random value for the **Intercept**, in this case we set the **Intercept** = 0…

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

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

$$+ -2(\mathbf{1.9} - (\text{intercept} + 0.64 \times \mathbf{2.3}))$$

$$+ -2(\mathbf{3.2} - (\text{intercept} + 0.64 \times \mathbf{2.9}))$$

Sum of Squared Residuals

0       1       2

Intercept

...then we calculated the derivative
when the **Intercept** = 0...



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

-2(**1.4** - (intercept + 0.64 × **0.5**)

+ -2(**1.9** - (intercept + 0.64 × **2.3**))

+ -2(**3.2** - (intercept + 0.64 × **2.9**))

Sum of
Squared
Residuals

0          1          2

Intercept

...plugged that slope into the **Step Size** calculation...

**Step Size = Slope × Learning Rate**

Sum of Squared Residuals

0   1   2

Intercept

...then calculated the **New Intercept**, the difference between the **Old Intercept** and the **Step Size**.

**Step Size = Slope × Learning Rate**

**New Intercept = Old Intercept - Step Size**

Sum of Squared Residuals

0          1          2

Intercept

Lastly, we plugged the **New Intercept** into the derivative and repeated everything until **Step Size** was close to **0**.

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

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

$$+ -2(\mathbf{1.9} - (\text{intercept} + 0.64 \times \mathbf{2.3}))$$

$$+ -2(\mathbf{3.2} - (\text{intercept} + 0.64 \times \mathbf{2.9}))$$

Sum of Squared Residuals

0    1    2

Intercept

**Step Size = Slope × Learning Rate**

**New Intercept = Old Intercept - Step Size**

Lastly, we plugged the **New Intercept**
into the derivative and repeated
everything until **Step Size** was close to **0**.



Sum of
Squared
Residuals

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

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

$$+ -2(\textbf{1.9} - (\text{intercept} + 0.64 \times \textbf{2.3}))$$

$$+ -2(\textbf{3.2} - (\text{intercept} + 0.64 \times \textbf{2.9}))$$

0               1               2

Intercept

**Step Size = Slope × Learning Rate**

**New Intercept = Old Intercept - Step Size**

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

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

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



Just like before, we will use the Sum of the Squared Residuals as the **Loss Function**

Height

Weight

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

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

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

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

Here's the derivative of the Sum of the Squared Residuals with respect to the **Intercept**…

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

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

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

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

Here's the derivative of the Sum of the Squared Residuals with respect to the **Intercept**…

…and here's the derivative with respect to the **Slope**.

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

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

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

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

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

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

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

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

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

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

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

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

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

$\dfrac{d}{d\ intercept}$ Sum of squared residuals =

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

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

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

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…

$\dfrac{d}{d\ slope}$ Sum of squared residuals =

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

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

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

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

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

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

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

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!**

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

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

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

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

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

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

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

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

Just like before, we will start by picking a random number for the **Intercept**. In this case we'll set the **Intercept = 0**…

…and we'll pick a random number for the **Slope**. In this case we'll set the **Slope = 1**.

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

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

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

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

$\dfrac{d}{d\ intercept}$ Sum of squared residuals =

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

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

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

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

$\dfrac{d}{d\ slope}$ Sum of squared residuals =

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

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

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



Height

Weight

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

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

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

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

Now let's plug in **0** for the **Intercept** and **1** for the **Slope**…

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

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

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

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

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

$$-2(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5})$$

$$+ -2(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3}))$$

$$+ -2(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9})) \boxed{= \mathbf{-1.6}}$$

...and that gives us two **Slopes**...

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

$$-2 \times \mathbf{0.5}(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5}))$$

$$+ -2 \times \mathbf{2.9}(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9}))$$

$$+ -2 \times \mathbf{2.3}(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3})) \boxed{= \mathbf{-0.8}}$$

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

$$-2(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5})$$

$$+ -2(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3}))$$

$$+ -2(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9})) \boxed{= \mathbf{-1.6}}$$

**Step Size**$_{\text{Intercept}}$ = **Slope** × **Learning Rate**

…now we plug the **Slopes** into the **Step Size** formulas…

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

$$-2 \times \mathbf{0.5}(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5}))$$

$$+ -2 \times \mathbf{2.9}(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9}))$$

$$+ -2 \times \mathbf{2.3}(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3})) \boxed{= \mathbf{-0.8}}$$

**Step Size**$_{\text{Slope}}$ = **Slope** × **Learning Rate**

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

-2(**1.4** - (0 + 1 × **0.5**)

+ -2(**1.9** - (0 + 1 × **2.3**))

+ -2(**3.2** - (0 + 1 × **2.9**)) $\boxed{= \textbf{-1.6}}$

**Step Size**Intercept = -1.6 × **Learning Rate**

...now we plug the **Slopes** into the **Step Size** formulas...

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

-2 × **0.5**(**1.4** - (0 + 1 × **0.5**))

+ -2 × **2.9**(**3.2** - (0 + 1 × **2.9**))

+ -2 × **2.3**(**1.9** - (0 + 1 × **2.3**)) $\boxed{= \textbf{-0.8}}$

**Step Size**Slope = -0.8 × **Learning Rate**

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

$$-2(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5})$$

$$+ -2(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3}))$$

$$+ -2(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9})) = \mathbf{-1.6}$$

**Step Size**$_{\text{Intercept}}$ = -1.6 × **Learning Rate**

…and multiply by the **Learning Rate**, which this time we set to **0.01**…

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

$$-2 \times \mathbf{0.5}(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5}))$$

$$+ -2 \times \mathbf{2.9}(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9}))$$

$$+ -2 \times \mathbf{2.3}(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3})) = \mathbf{-0.8}$$

**Step Size**$_{\text{Slope}}$ = -0.8 × **Learning Rate**

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

$-2(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5})$

$+ -2(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3}))$

$+ -2(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9})) = \mathbf{-1.6}$

**Step Size**<sub>Intercept</sub> = -1.6 × 0.01

**NOTE:** The larger **Learning Rate** that we used in the first example doesn't work this time. Even after a bunch of steps, **Gradient Descent** doesn't arrive at the correct answer.

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

$-2 \times \mathbf{0.5}(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5}))$

$+ -2 \times \mathbf{2.9}(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9}))$

$+ -2 \times \mathbf{2.3}(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3})) = \mathbf{-0.8}$

**Step Size**<sub>Slope</sub> = -0.8 × 0.01

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

$$-2(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5})$$

$$+ -2(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3}))$$

$$+ -2(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9})) = \mathbf{-1.6}$$

**Step Size**$_{\text{Intercept}}$ = -1.6 × 0.01

This means that **Gradient Descent** can be very sensitive to the **Learning Rate**.

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

$$-2 \times \mathbf{0.5}(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5}))$$

$$+ -2 \times \mathbf{2.9}(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9}))$$

$$+ -2 \times \mathbf{2.3}(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3})) = \mathbf{-0.8}$$

**Step Size**$_{\text{Slope}}$ = -0.8 × 0.01

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

$$-2(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5}))$$

$$+ -2(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3}))$$

$$+ -2(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9})) = \mathbf{-1.6}$$

**Step Size**$_{\text{Intercept}}$ = -1.6 × 0.01

The good news is that in practice, a reasonable **Learning Rate** can be determined automatically by starting large and getting smaller with each step.

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

$$-2 \times \mathbf{0.5}(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5}))$$

$$+ -2 \times \mathbf{2.9}(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9}))$$

$$+ -2 \times \mathbf{2.3}(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3})) = \mathbf{-0.8}$$

**Step Size**$_{\text{Slope}}$ = -0.8 × 0.01

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

$-2(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5})$

$+ \ -2(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3}))$

$+ \ -2(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9})) \ = \mathbf{-1.6}$

**Step Size**$_{\text{Intercept}}$ = -1.6 × 0.01

So, in theory, you shouldn't have to worry too much about the **Learning Rate**.

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

$-2 \times \mathbf{0.5}(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5}))$

$+ \ -2 \times \mathbf{2.9}(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9}))$

$+ \ -2 \times \mathbf{2.3}(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3})) \ = \mathbf{-0.8}$

**Step Size**$_{\text{Slope}}$ = -0.8 × 0.01

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

$$-2(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5})$$

$$+ -2(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3}))$$

$$+ -2(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9})) = \mathbf{-1.6}$$

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

Anyway, we do the math and get two **Step Sizes**.

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

$$-2 \times \mathbf{0.5}(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5}))$$

$$+ -2 \times \mathbf{2.9}(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9}))$$

$$+ -2 \times \mathbf{2.3}(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3})) = \mathbf{-0.8}$$

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

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

$-2(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5})$

$+ -2(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3}))$

$+ -2(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9})) = \mathbf{-1.6}$

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

**New Intercept = Old Intercept - Step Size**

Now we calculate the **New Intercept** and **New Slope** by plugging in the **Old Intercept** and the **Old Slope**…

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

$-2 \times \mathbf{0.5}(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5}))$

$+ -2 \times \mathbf{2.9}(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9}))$

$+ -2 \times \mathbf{2.3}(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3})) = \mathbf{-0.8}$

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

**New Slope = Old Slope - Step Size**

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

$$-2(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5})$$

$$+ -2(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3}))$$

$$+ -2(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9})) = \mathbf{-1.6}$$

**Step Size**<sub>Intercept</sub> = -1.6 × 0.01 = **-0.016**

**New Intercept** = 0 - (-0.016) ←

…and the
**Step Sizes**…

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

$$-2 \times \mathbf{0.5}(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5}))$$

$$+ -2 \times \mathbf{2.9}(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9}))$$

$$+ -2 \times \mathbf{2.3}(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3})) = \mathbf{-0.8}$$

**Step Size**<sub>Slope</sub> = -0.8 × 0.01 = **-0.008**

**New Slope** = 1 - (-0.008) ←

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

$$-2(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5})$$

$$+ -2(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3}))$$

$$+ -2(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9})) = \mathbf{-1.6}$$

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

**New Intercept** = 0 - (-0.016) = 0.016

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

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

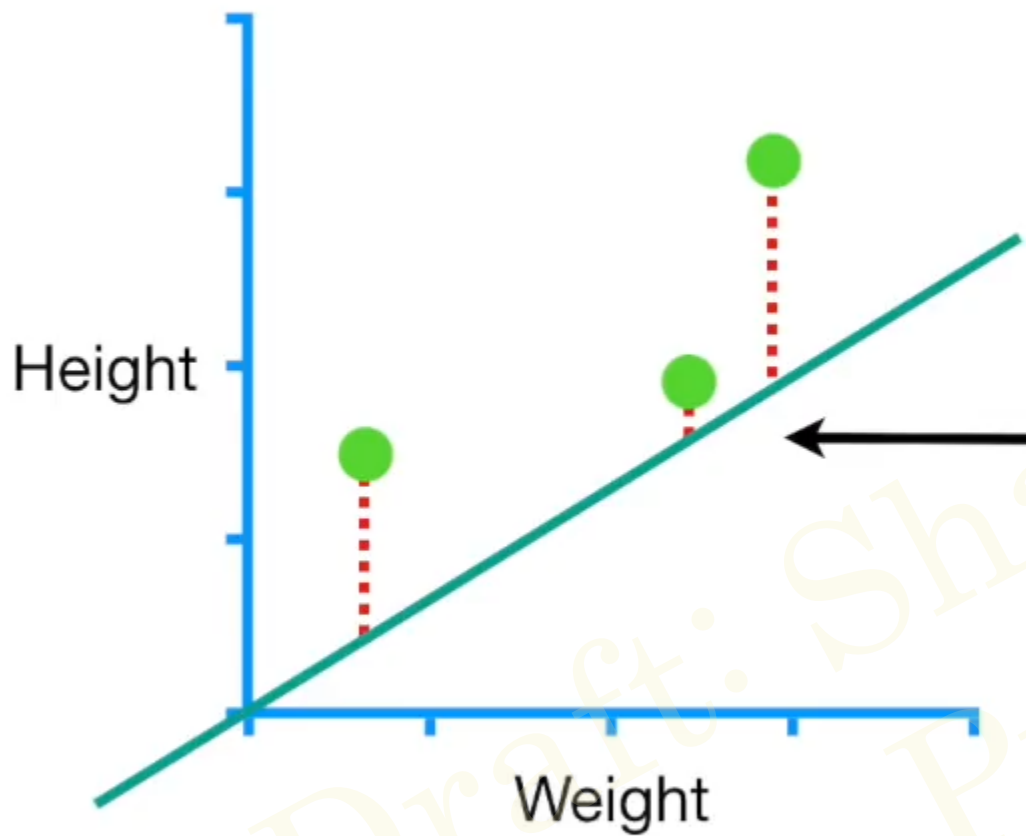$$-2 \times \mathbf{0.5}(\mathbf{1.4} - (0 + 1 \times \mathbf{0.5}))$$

$$+ -2 \times \mathbf{2.9}(\mathbf{3.2} - (0 + 1 \times \mathbf{2.9}))$$

$$+ -2 \times \mathbf{2.3}(\mathbf{1.9} - (0 + 1 \times \mathbf{2.3})) = \mathbf{-0.8}$$
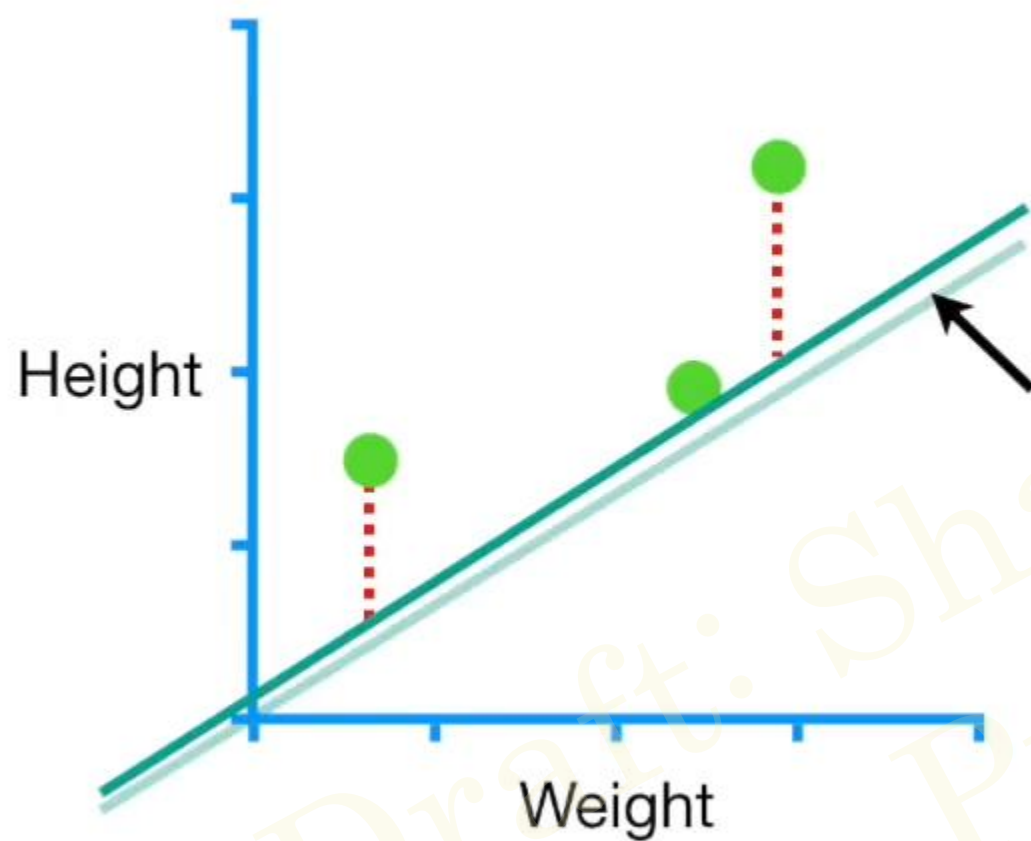
**Step Size**$_{\text{Slope}}$ = -0.8 × 0.01 = **0.008**

**New Slope** = 1 - (-0.008) = 1.008

**New Intercept** = 0 - (-0.016) = 0.016

Height

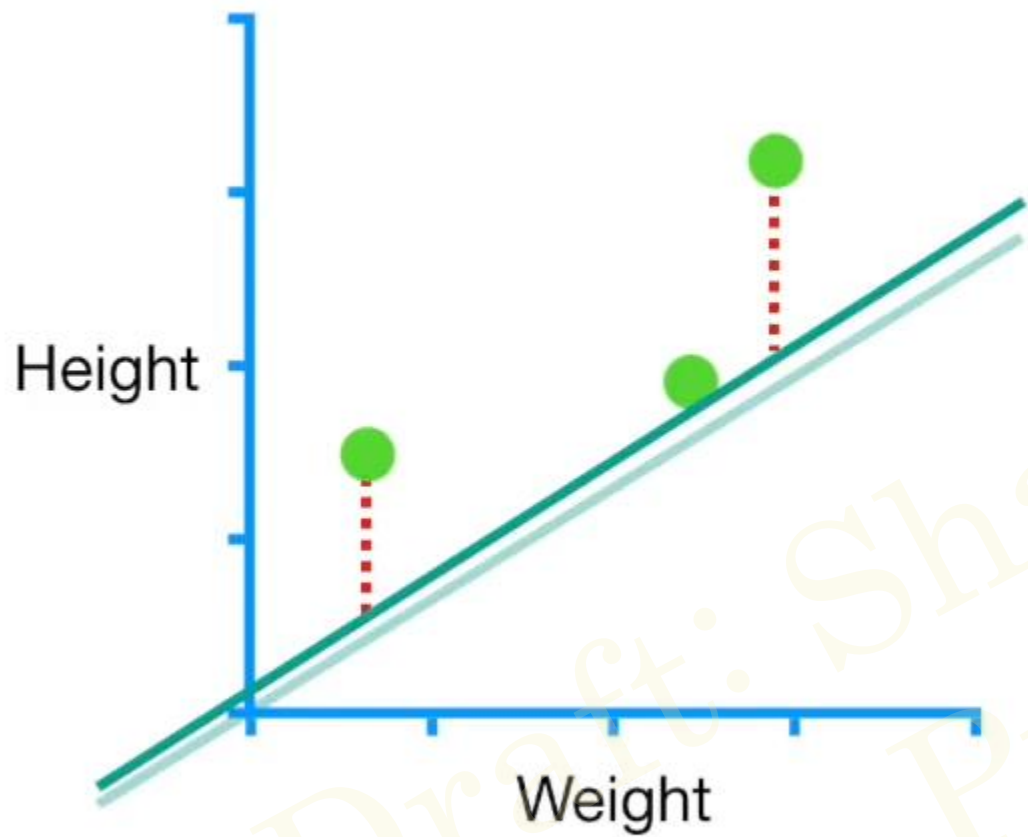This is the line we started with… (**Slope = 1** and **Intercept = 0**)

Weight

**New Slope** = 1 - (-0.008) = 1.008

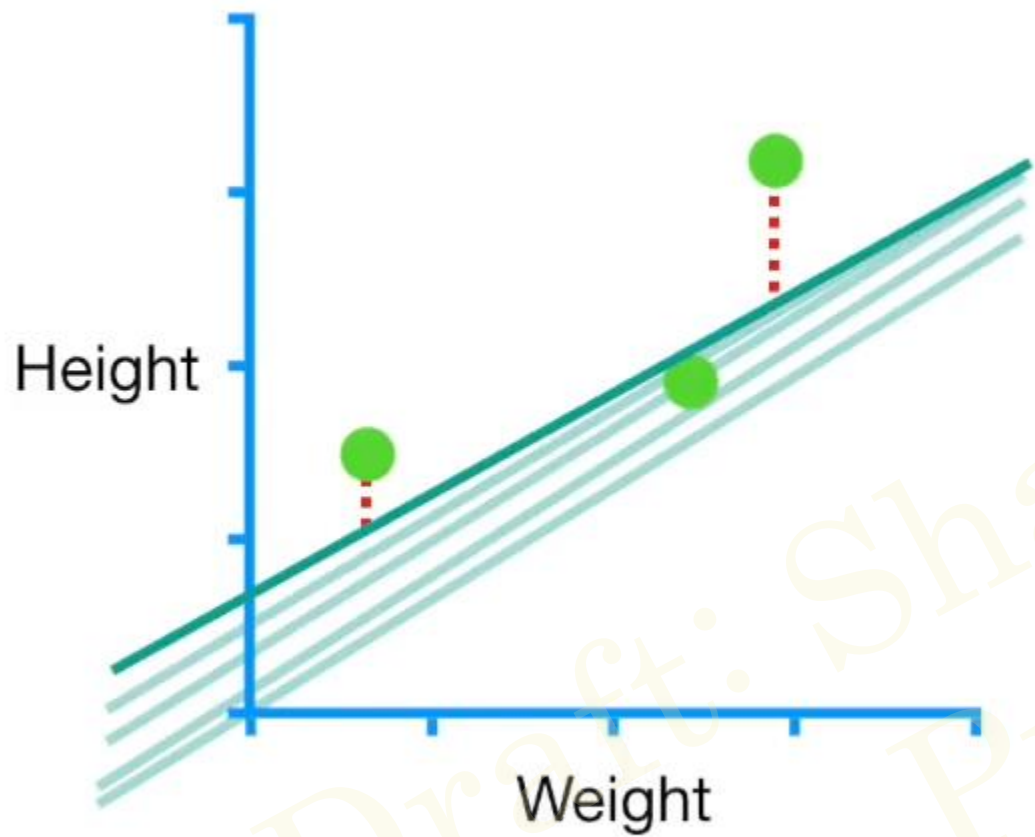**New Intercept** = 0 - (-0.016) = 0.016

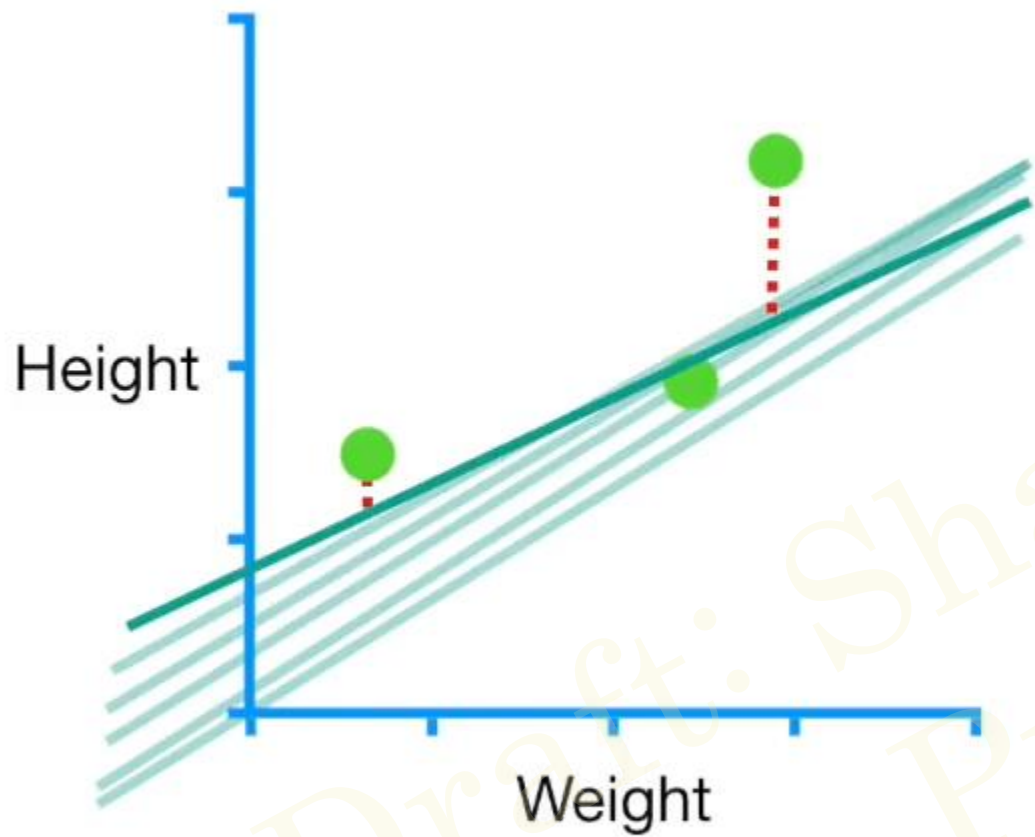...and this is the new line (with **Slope = 1.008** and **Intercept = 0.016**) after the first step.

**New Slope** = 1 - (-0.008) = 1.008
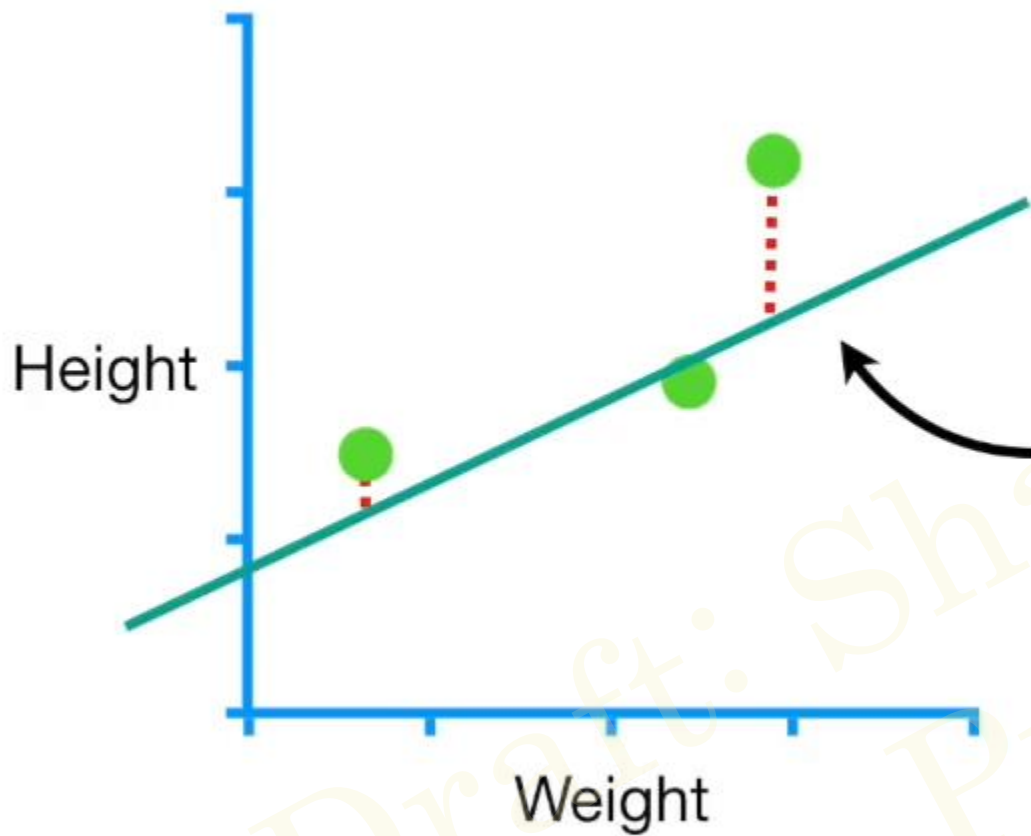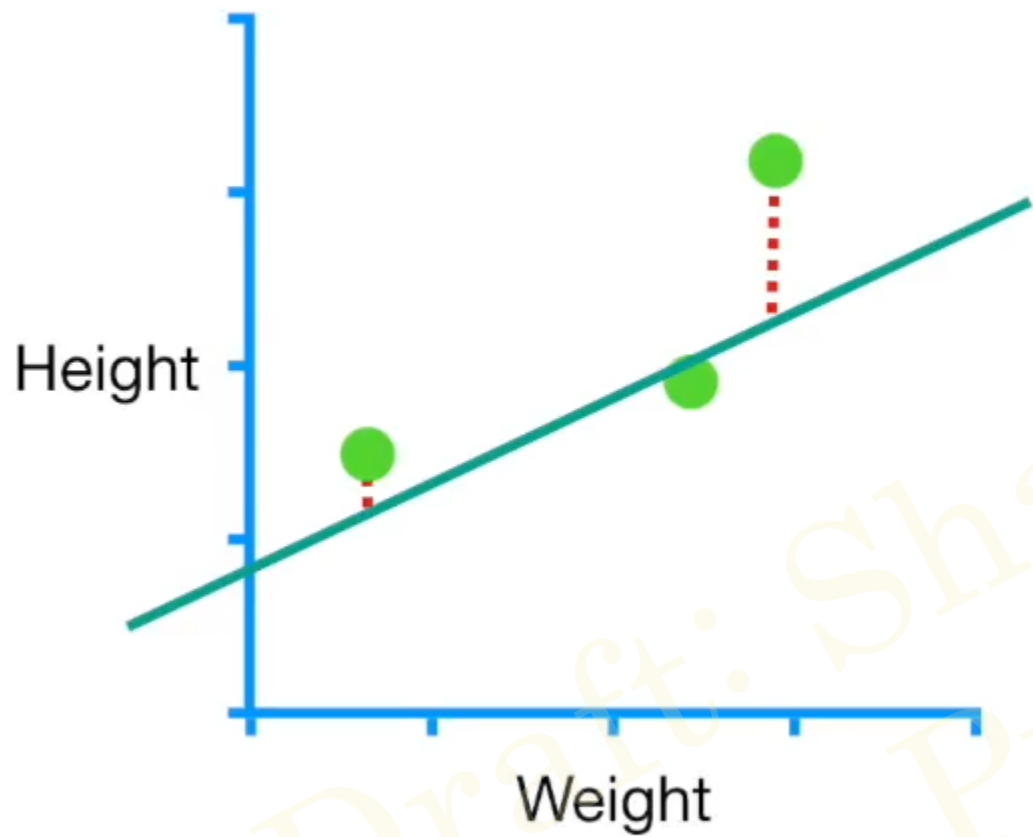
Height

Weight

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

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

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

This is the best fitting line, with **Intercept** = **0.95** and **Slope** = **0.64**, the same values we get from **Least Squares**.

We now know how **Gradient Descent** optimizes two parameters, the **Slope** and **Intercept**.

**Step 1:** Take the derivative of the **Loss Function** for each parameter in it. In fancy Machine Learning Lingo, take the **Gradient** of the **Loss Function**.

**Step 1:** Take the derivative of the **Loss Function** for each parameter in it. In fancy Machine Learning Lingo, take the **Gradient** of the **Loss Function**.

**Step 2:** Pick random values for the parameters.

**Step 1:** Take the derivative of the **Loss Function** for each parameter in it. In fancy Machine Learning Lingo, take the **Gradient** of the **Loss Function**.

**Step 2:** Pick random values for the parameters.

**Step 3:** Plug the parameter values into the derivatives (ahem, the **Gradient**).

**Step 1:** Take the derivative of the **Loss Function** for each parameter in it. In fancy Machine Learning Lingo, take the **Gradient** of the **Loss Function**.

**Step 2:** Pick random values for the parameters.

**Step 3:** Plug the parameter values into the derivatives (ahem, the **Gradient**).

**Step 4:** Calculate the Step Sizes:  **Step Size = Slope × Learning Rate**

**Step 1:** Take the derivative of the **Loss Function** for each parameter in it. In fancy Machine Learning Lingo, take the **Gradient** of the **Loss Function**.

**Step 2:** Pick random values for the parameters.

**Step 3:** Plug the parameter values into the derivatives (ahem, the **Gradient**).

**Step 4:** Calculate the Step Sizes:  **Step Size = Slope × Learning Rate**

**Step 5:** Calculate the New Parameters:

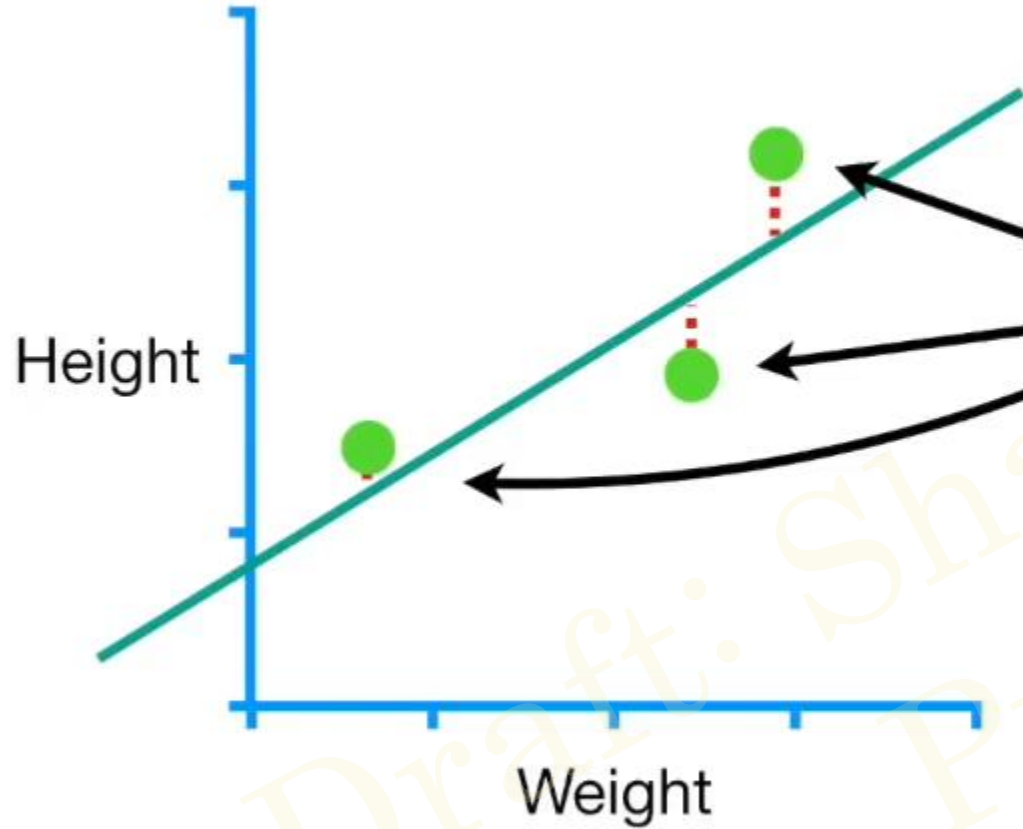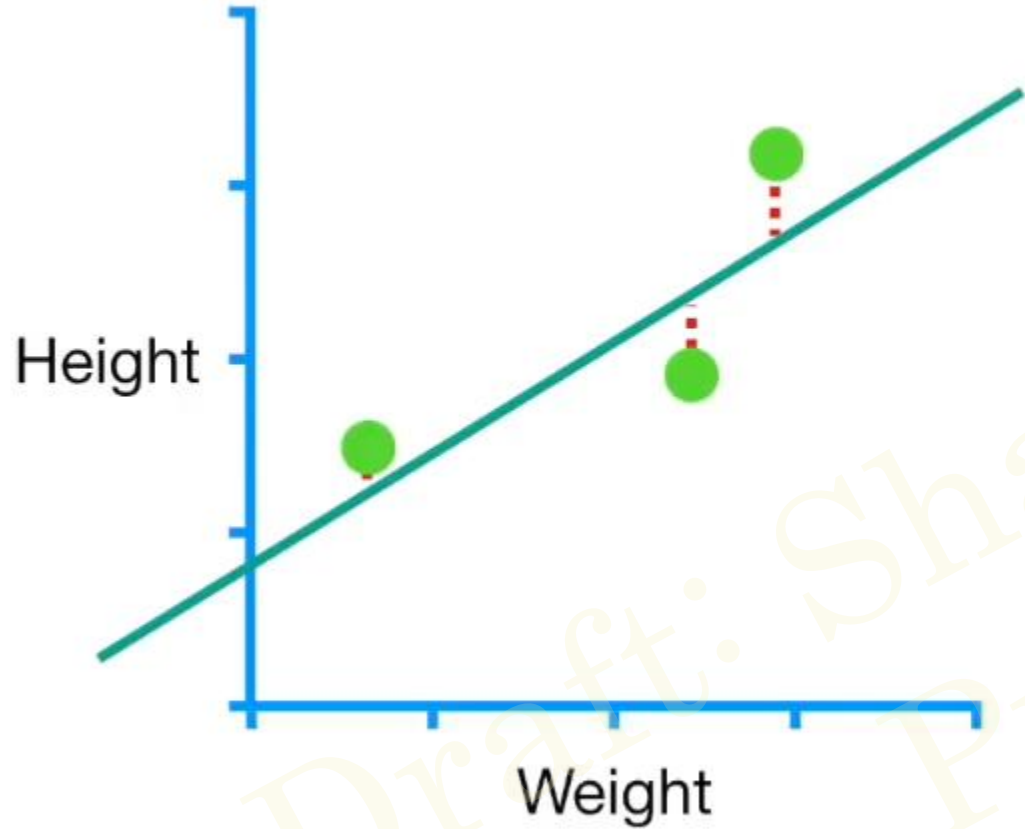**New Parameter = Old Parameter - Step Size**

Now go back to **Step 3** and repeat until
**Step Size** is very small, or you reach
the **Maximum Number of Steps**.

**Step 3:** Plug the parameter values into the derivatives (ahem, the **Gradient**).

**Step 4:** Calculate the Step Sizes:  **Step Size** = **Slope** × **Learning Rate**

**Step 5:** Calculate the New Parameters:

**New Parameter** = **Old Parameter** - **Step Size**

In our example, we only had three data points, so the math didn't take very long...
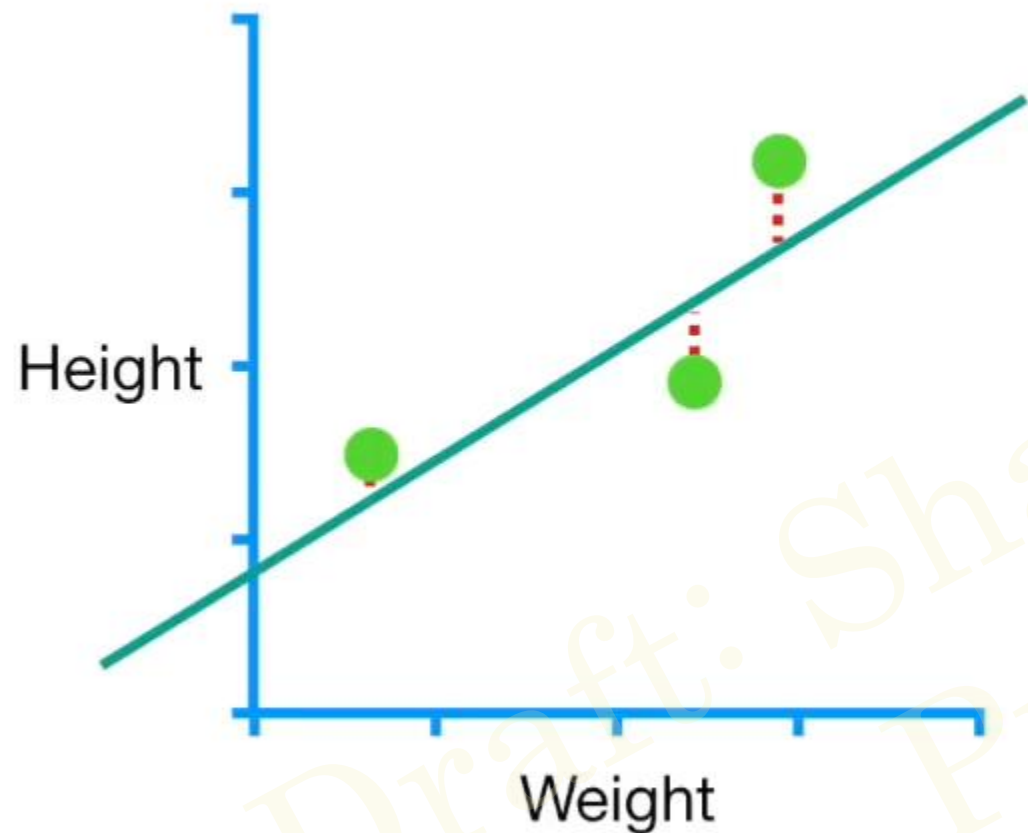
So there is a thing called **Stochastic Gradient Descent** that uses a randomly selected subset of the data at every step rather than the full dataset.

So there is a thing called **Stochastic Gradient Descent** that uses a randomly selected subset of the data at every step rather than the full dataset.

This reduces the time spent calculating the derivatives of the **Loss Function**.

So there is a thing called **Stochastic Gradient Descent** that uses a randomly selected subset of the data at every step rather than the full dataset.

This reduces the time spent calculating the derivatives of the **Loss Function**.

That's all.

**Stochastic Gradient Descent** sounds fancy, but it's no big deal.

# THANK YOU!