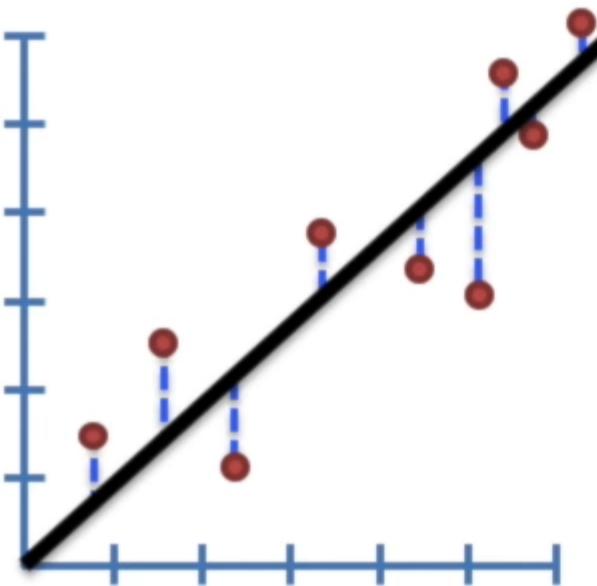


Linear regression

Supervised learning

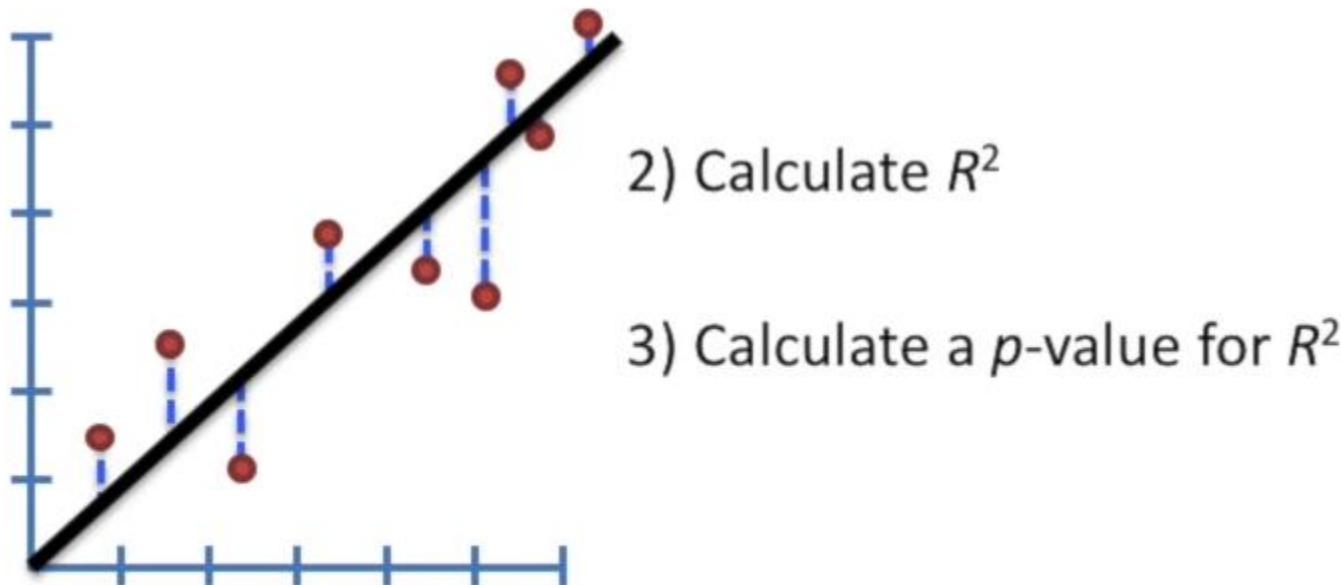
The Main Ideas!

- 1) Use least-squares to fit a line to the data.

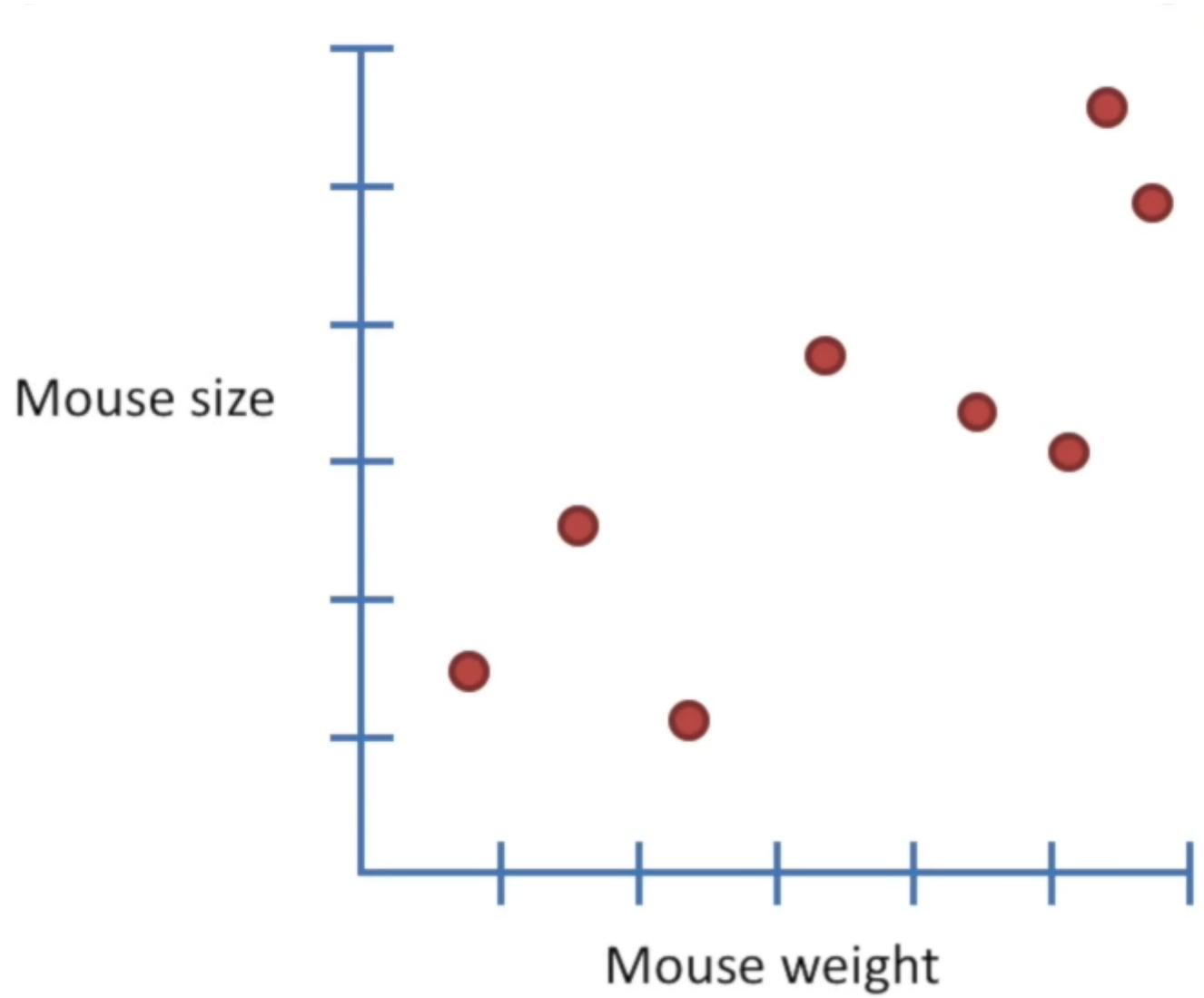


The Main Ideas!

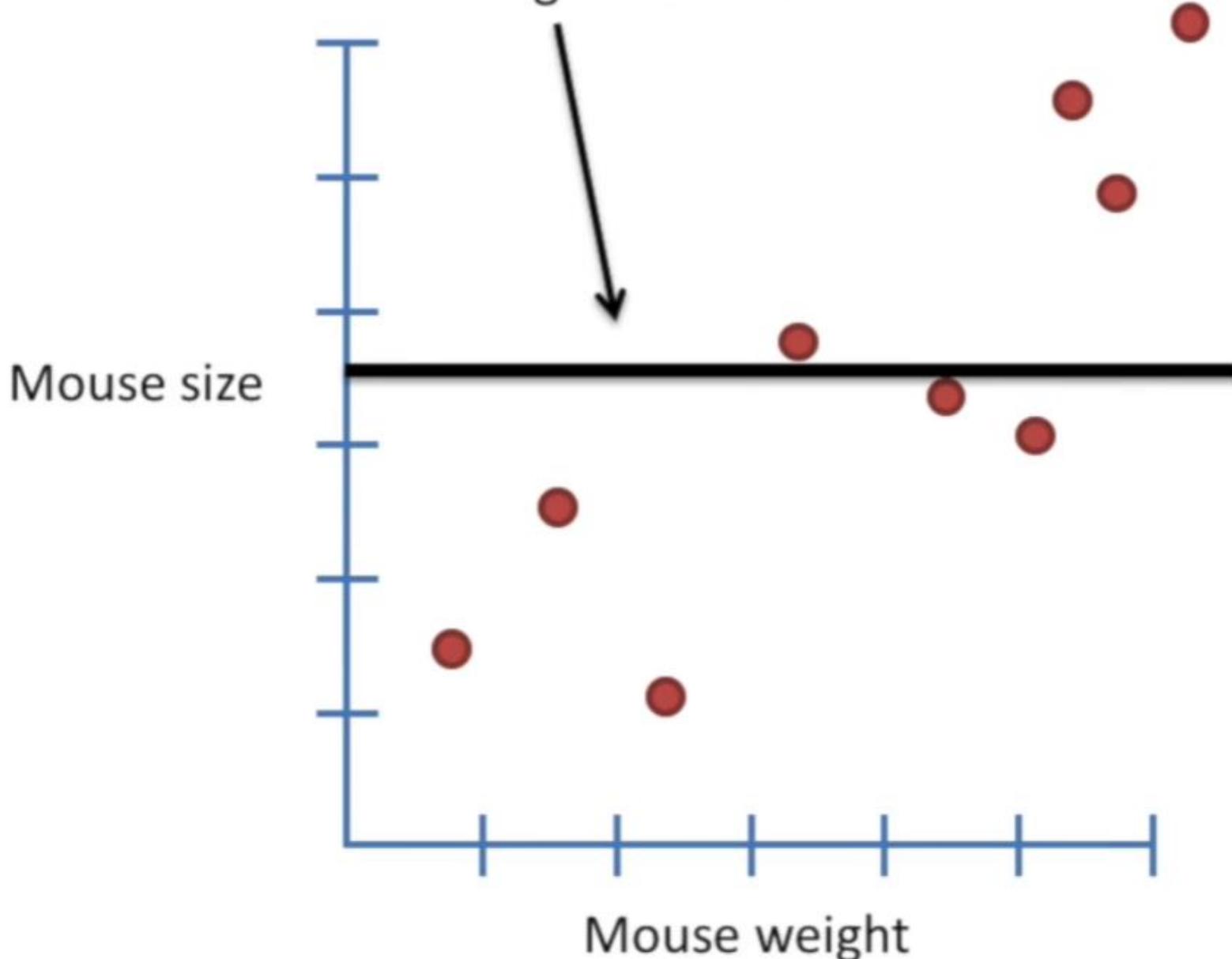
- 1) Use least-squares to fit a line to the data.

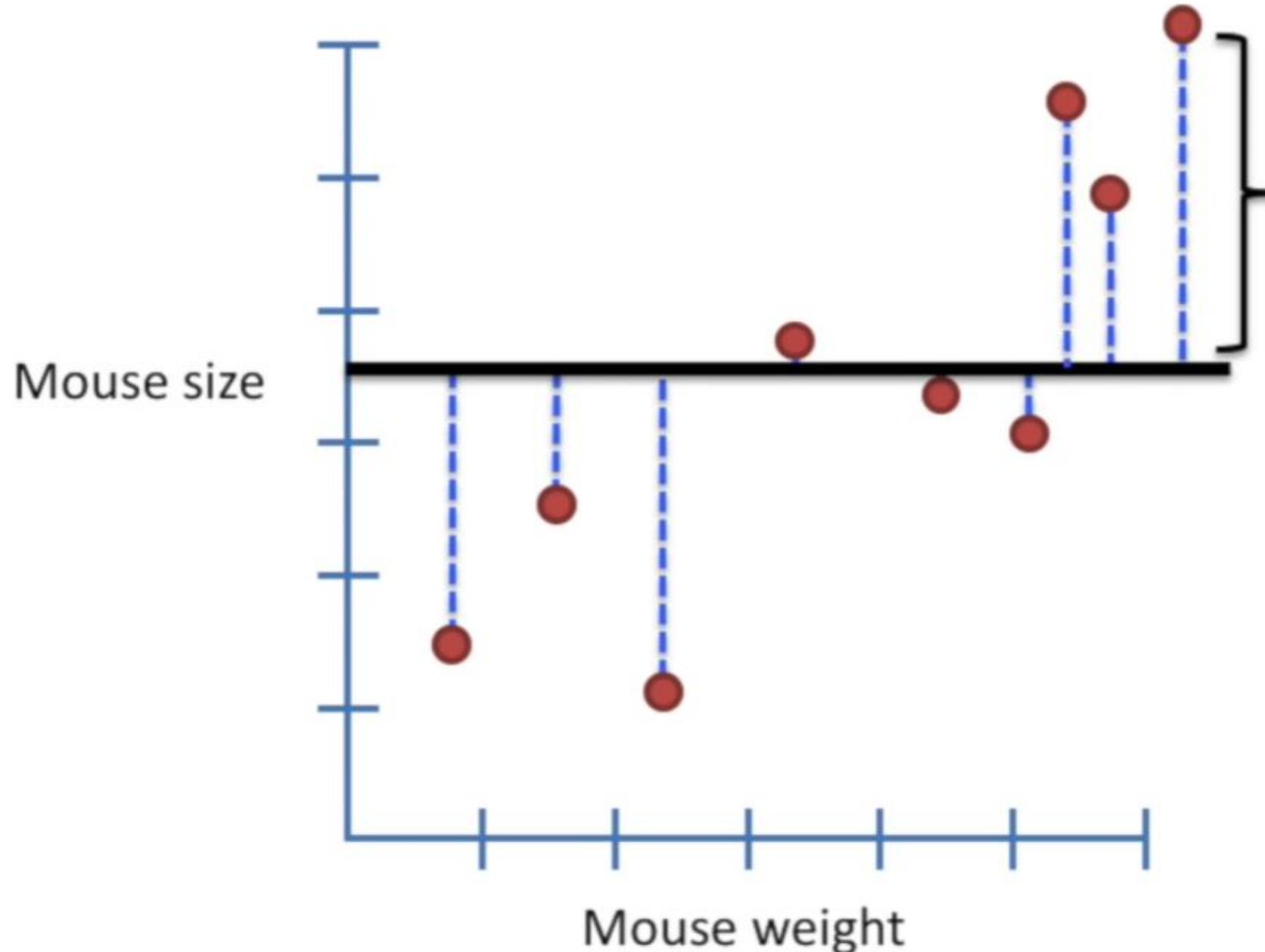


Let's do a quick review...



First, draw a line
through the data...

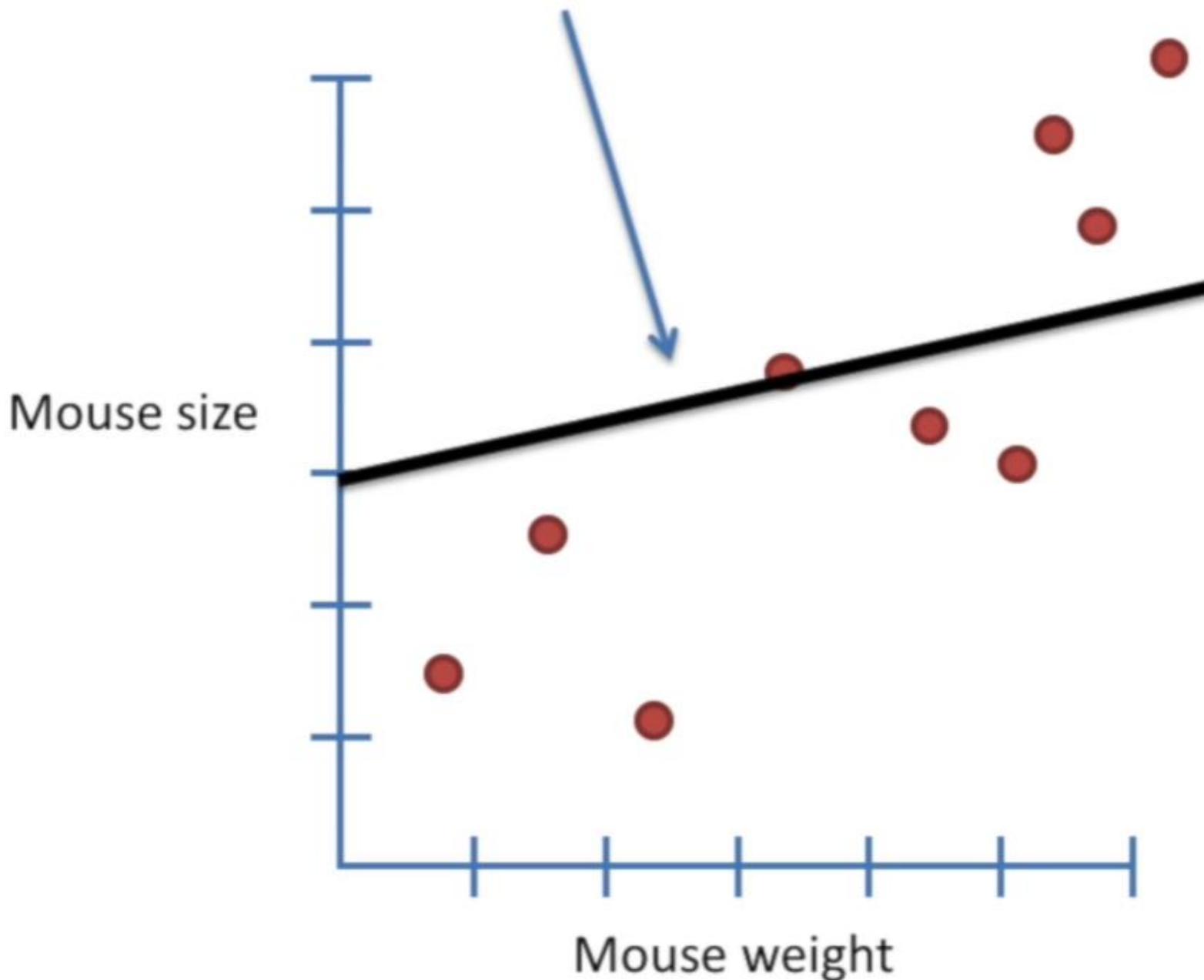


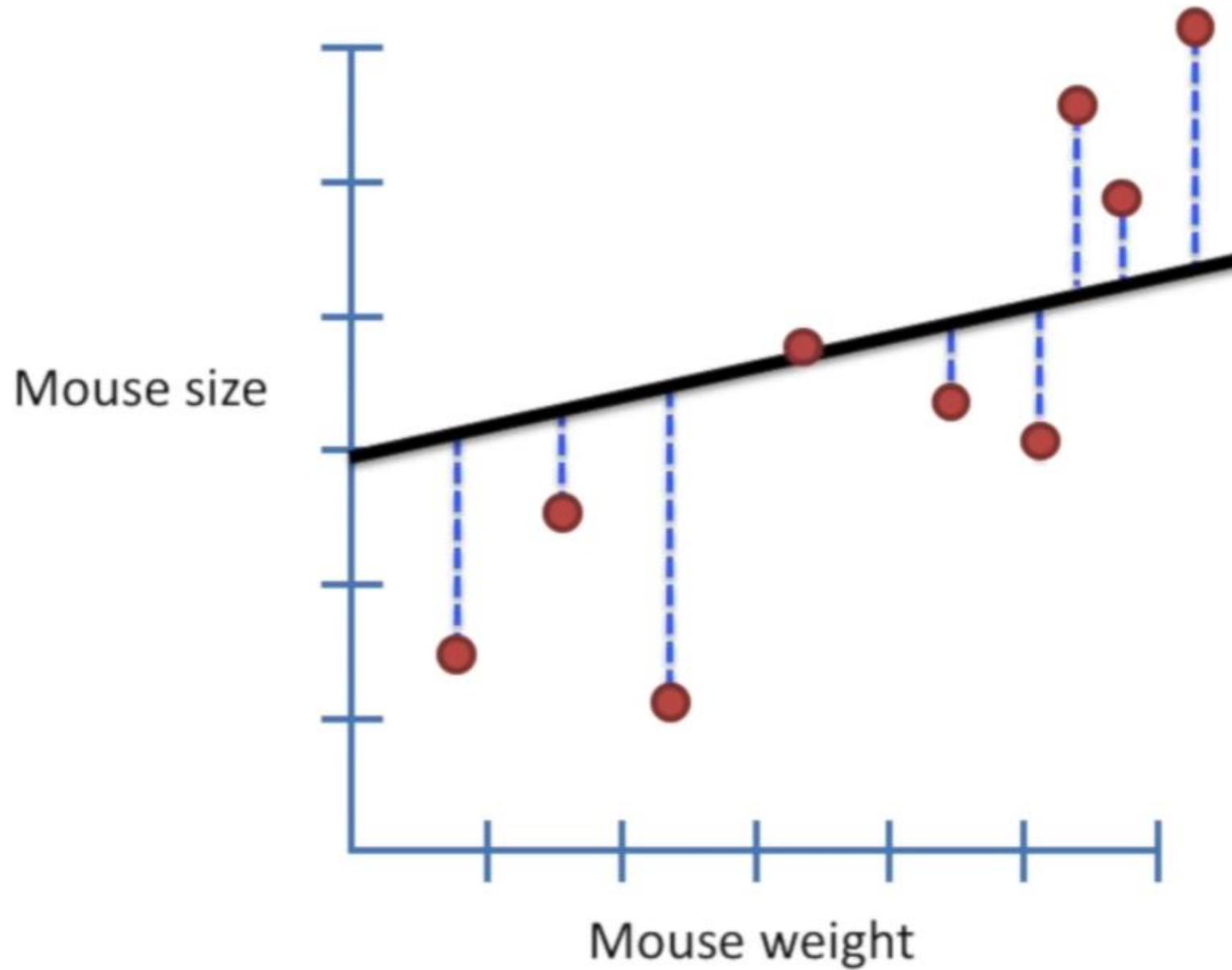


Second, measure the distance from the line to the data, square each distance, and then add them up.

Terminology alert!
The distance from a line to a data point is called a “**residual**”.

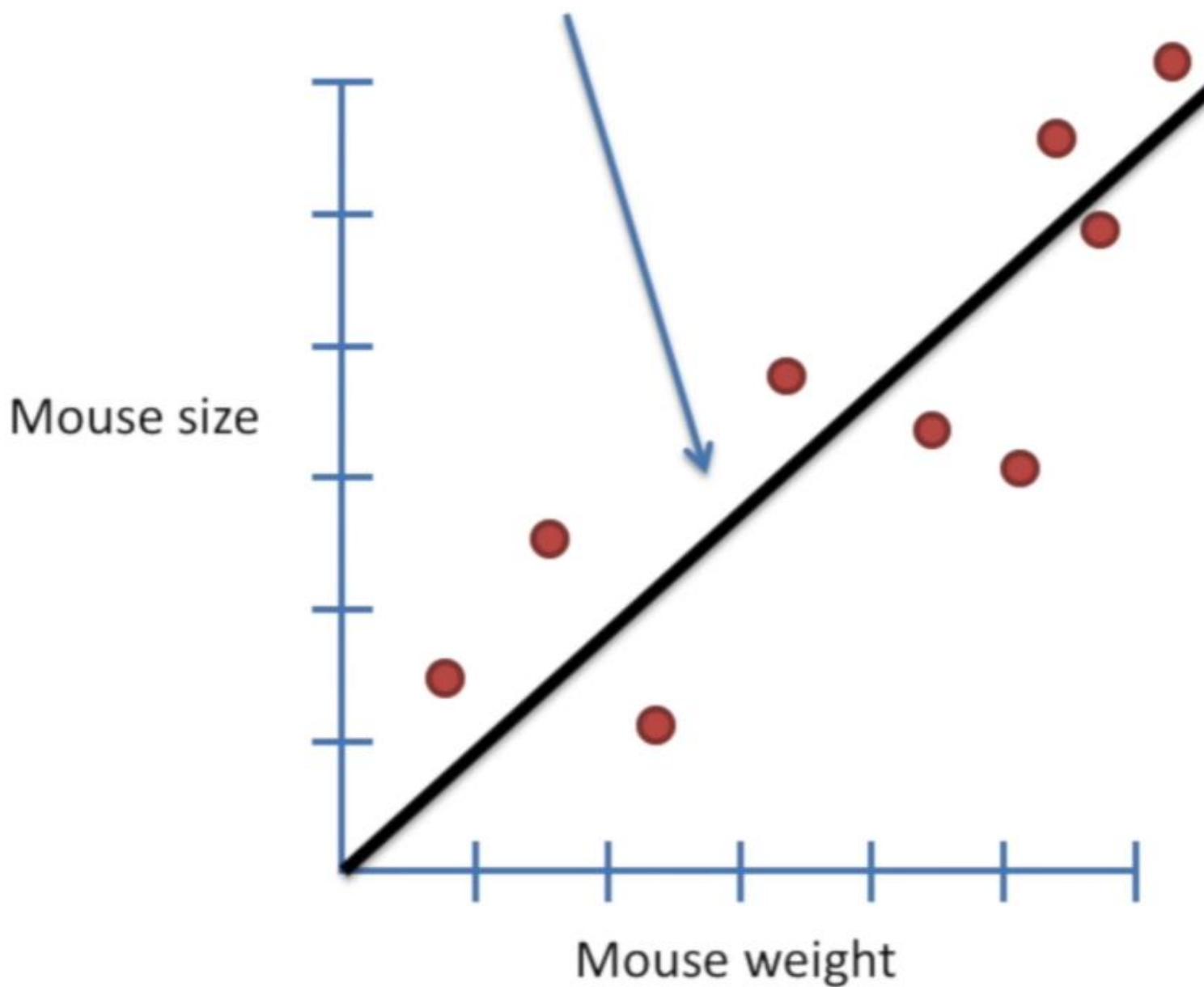
Third, rotate the line a little bit...

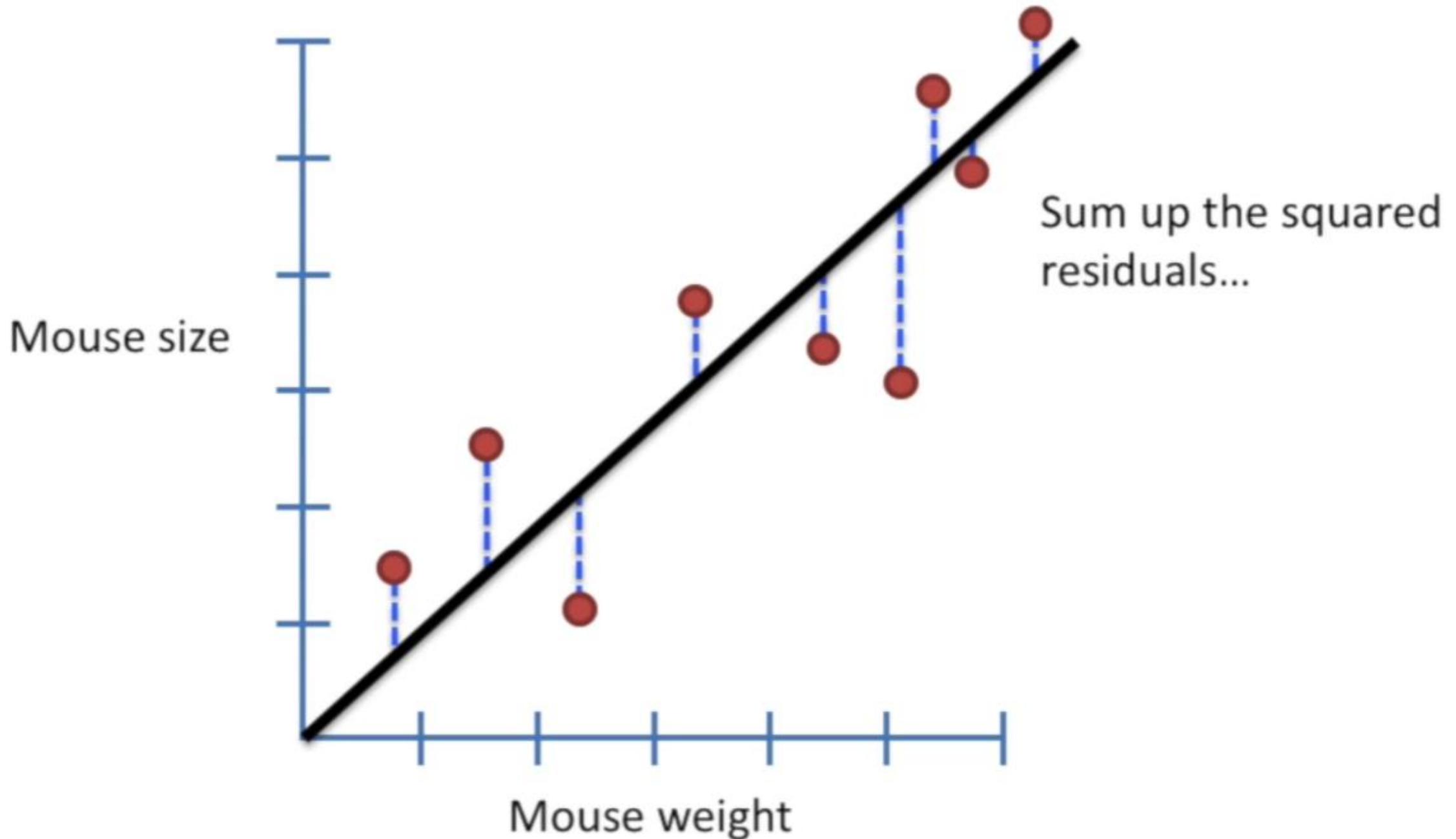


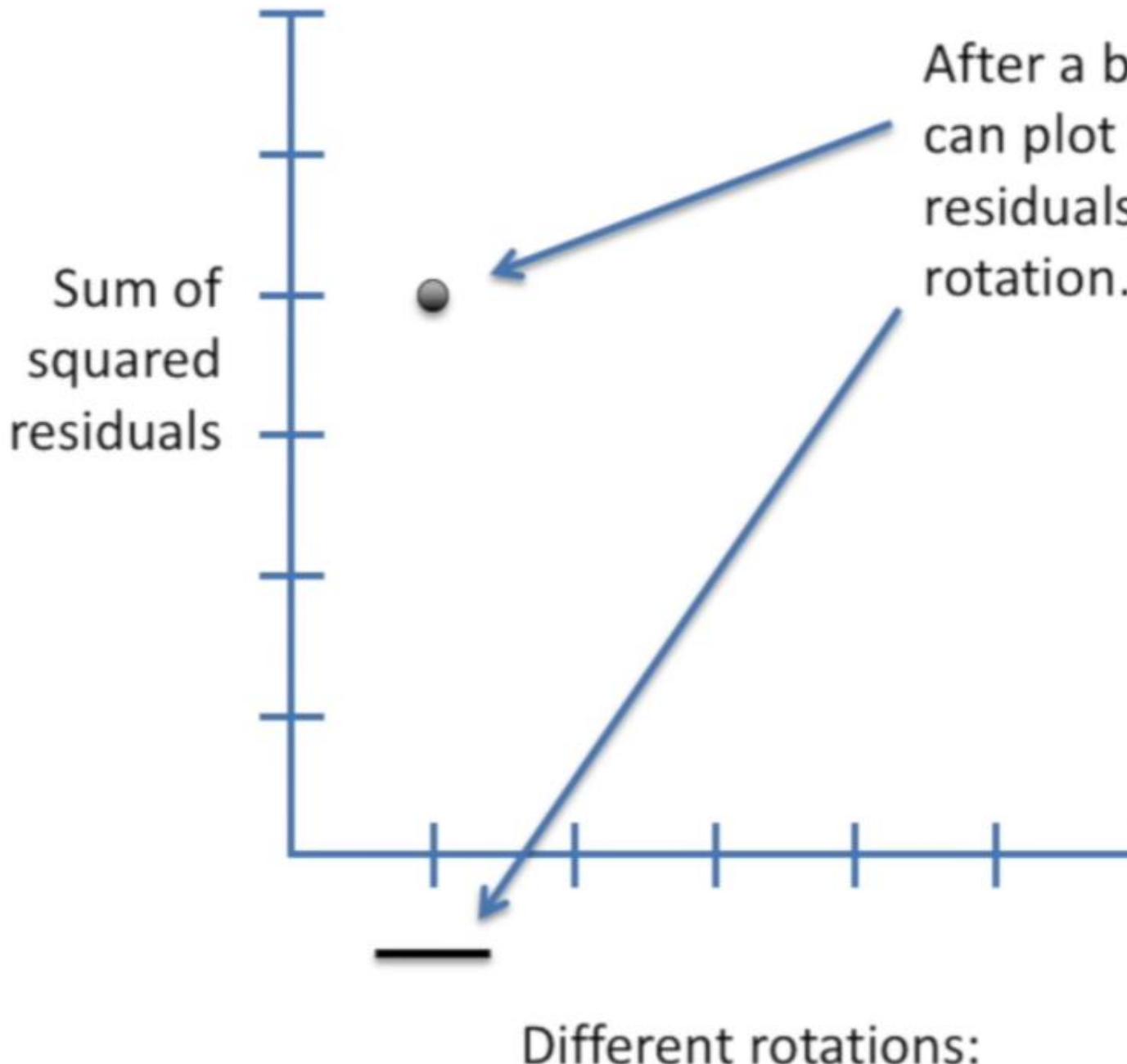


With the new line,
measure the
residuals, square
them, and then sum
up the squares.

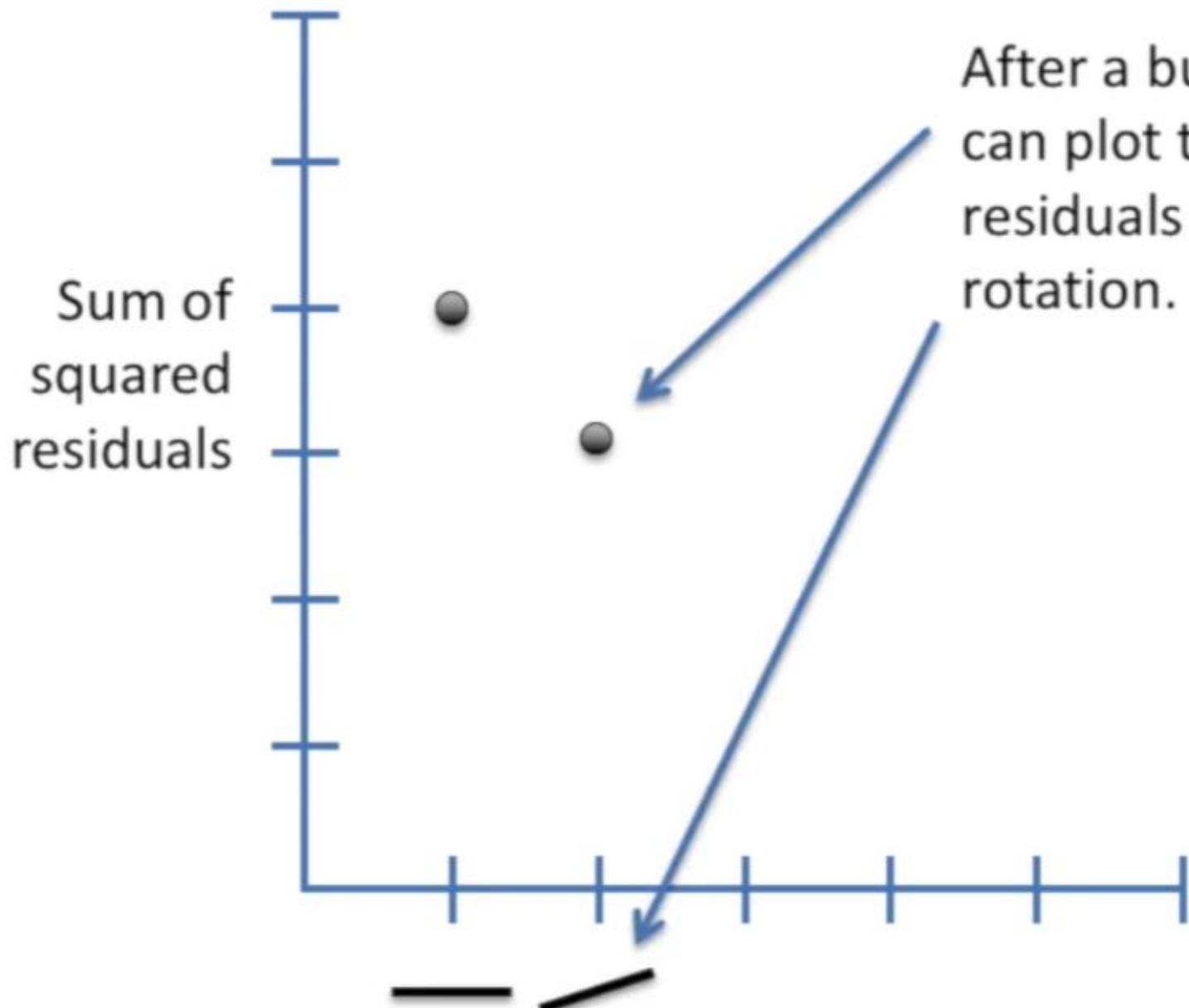
Rotate the line a little bit more...





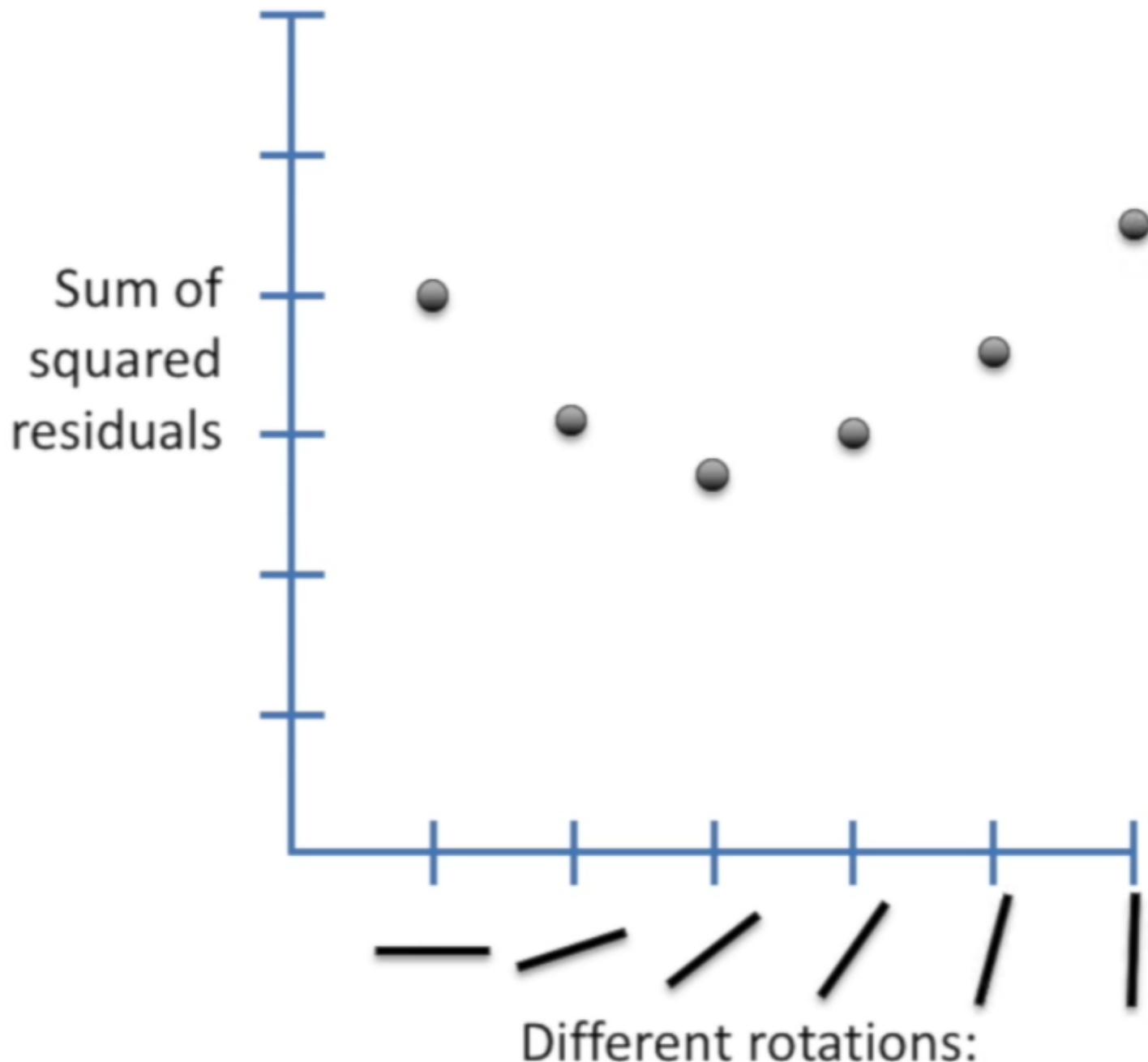


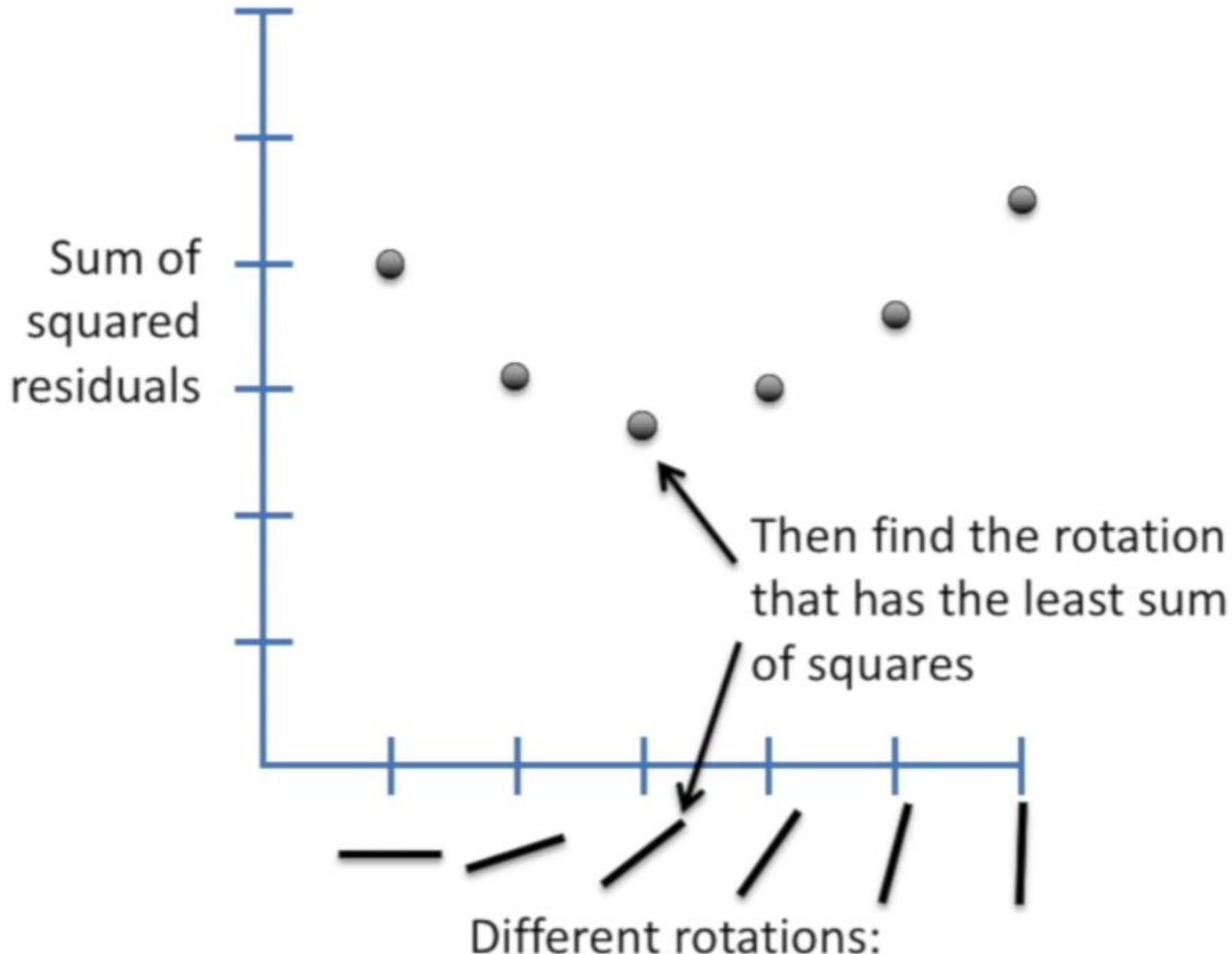
After a bunch of rotations, you can plot the sum of squared residuals and corresponding rotation.

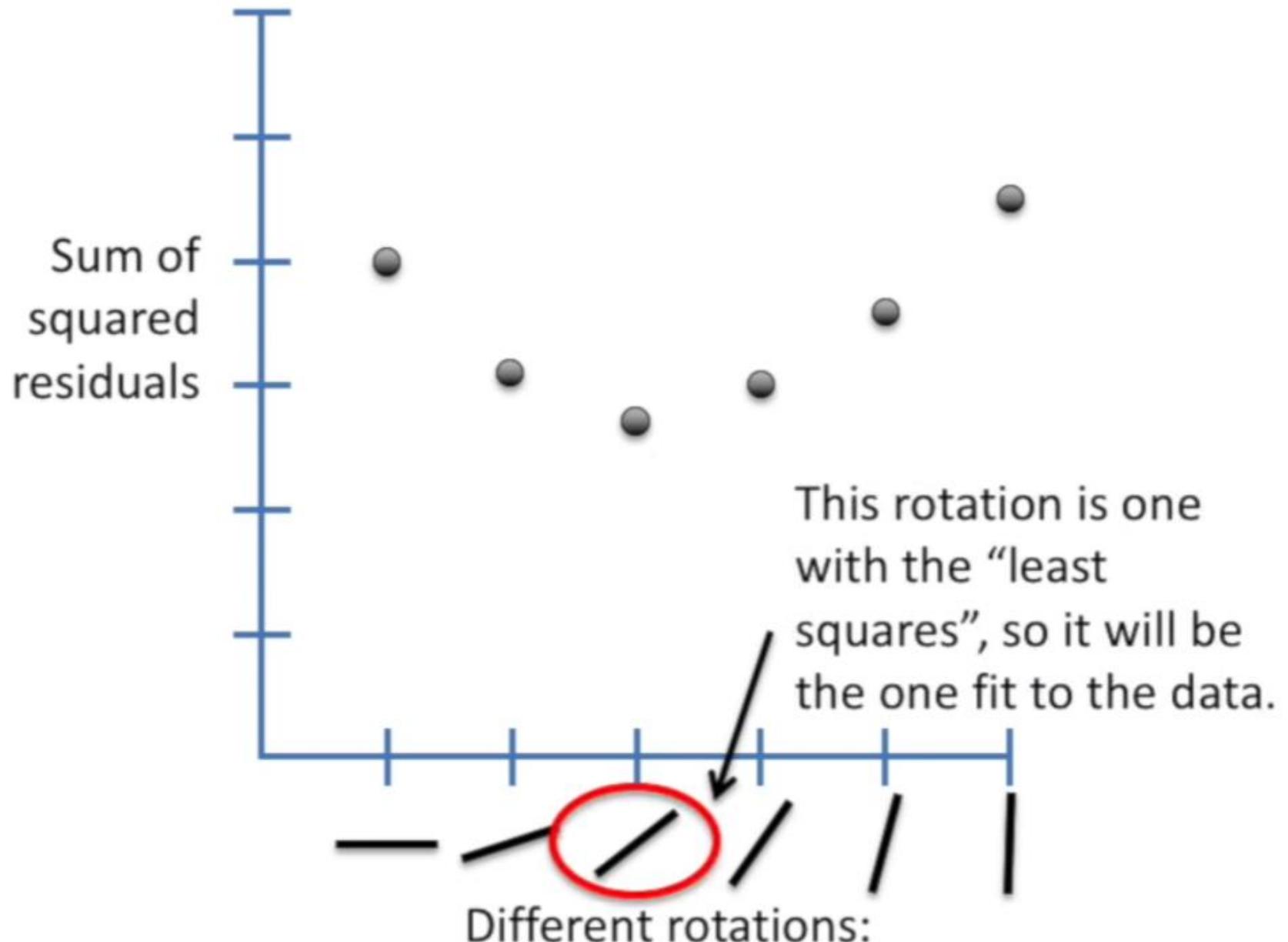


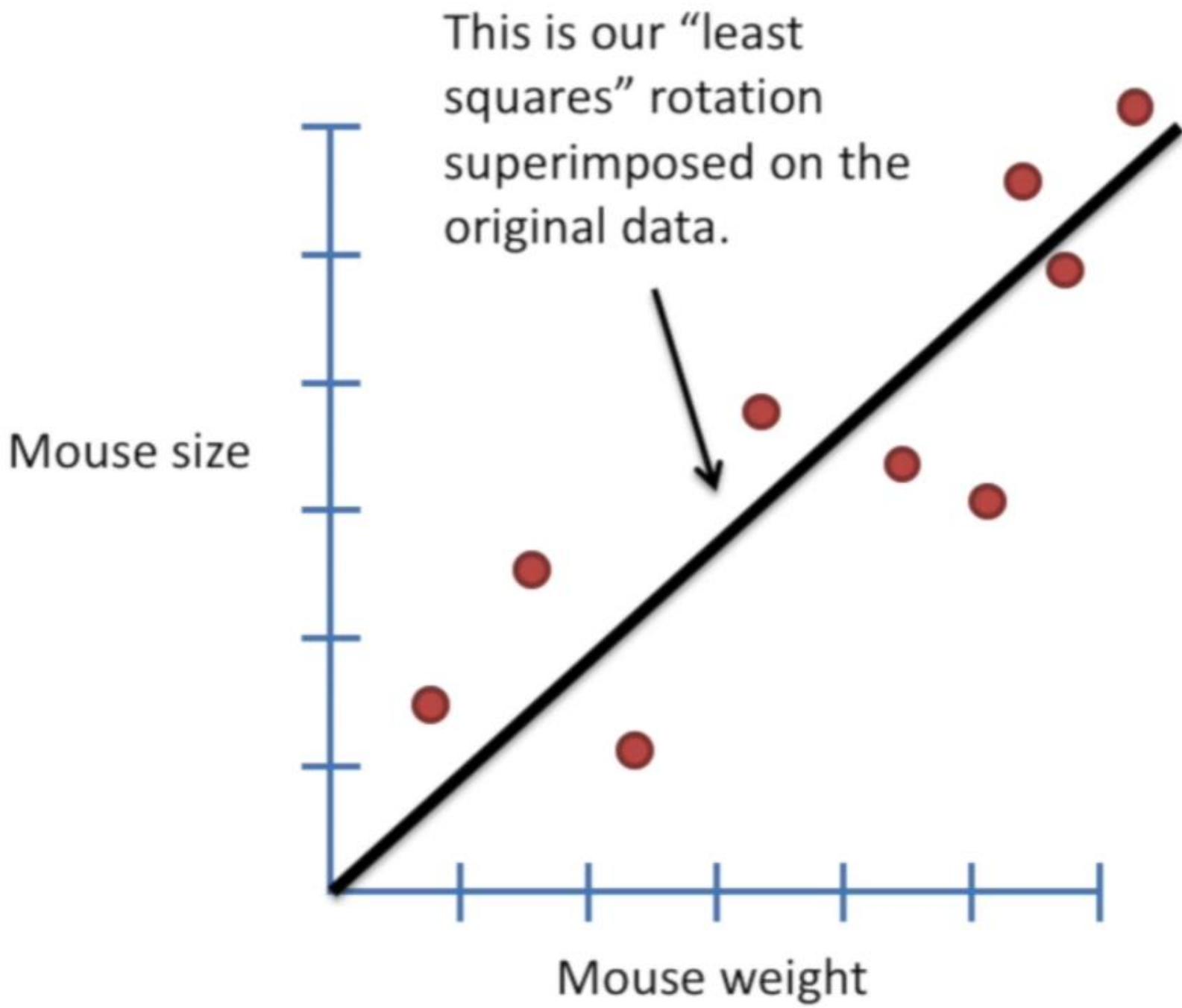
After a bunch of rotations, you can plot the sum of squared residuals and corresponding rotation.

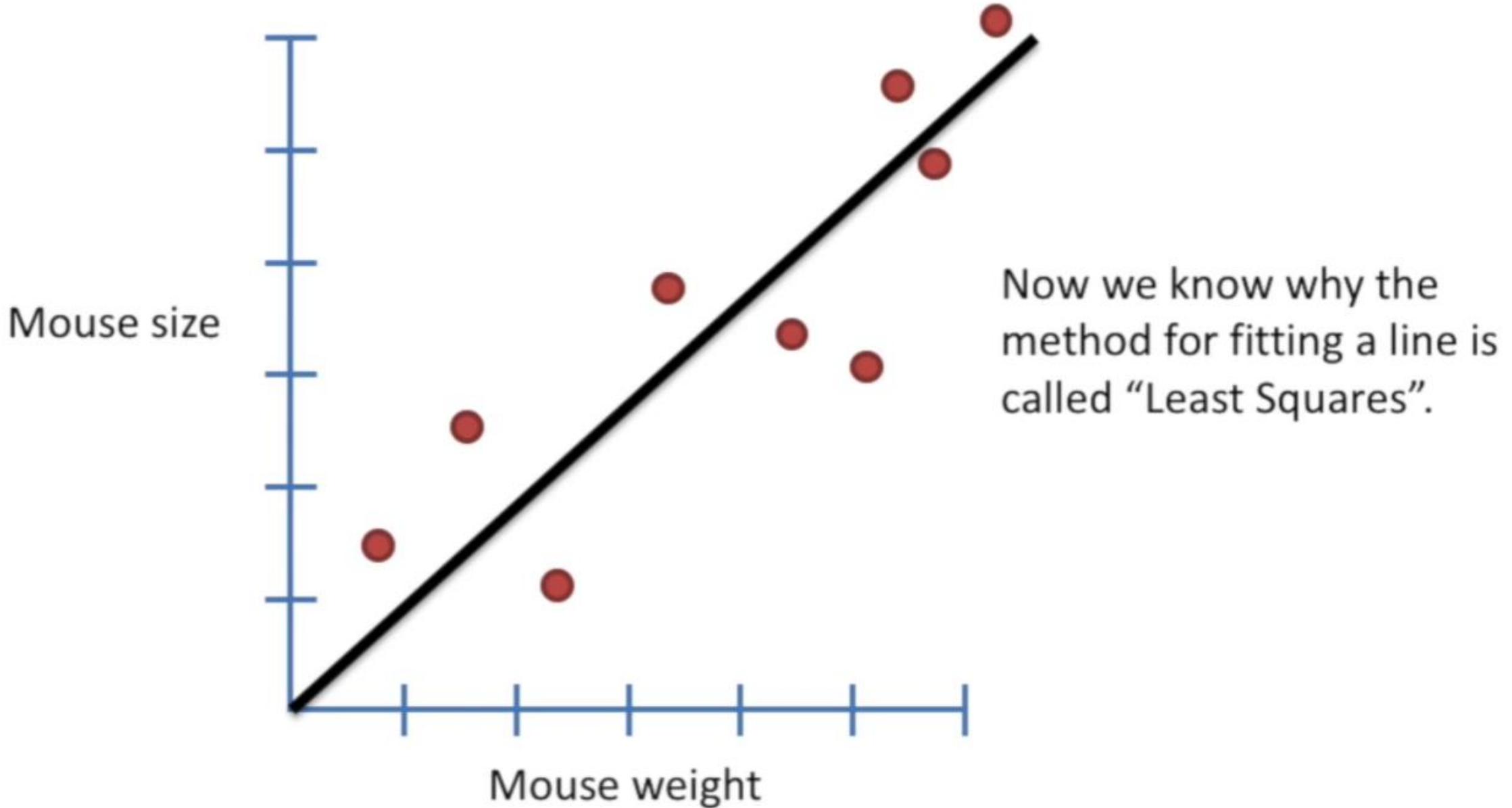
Different rotations:



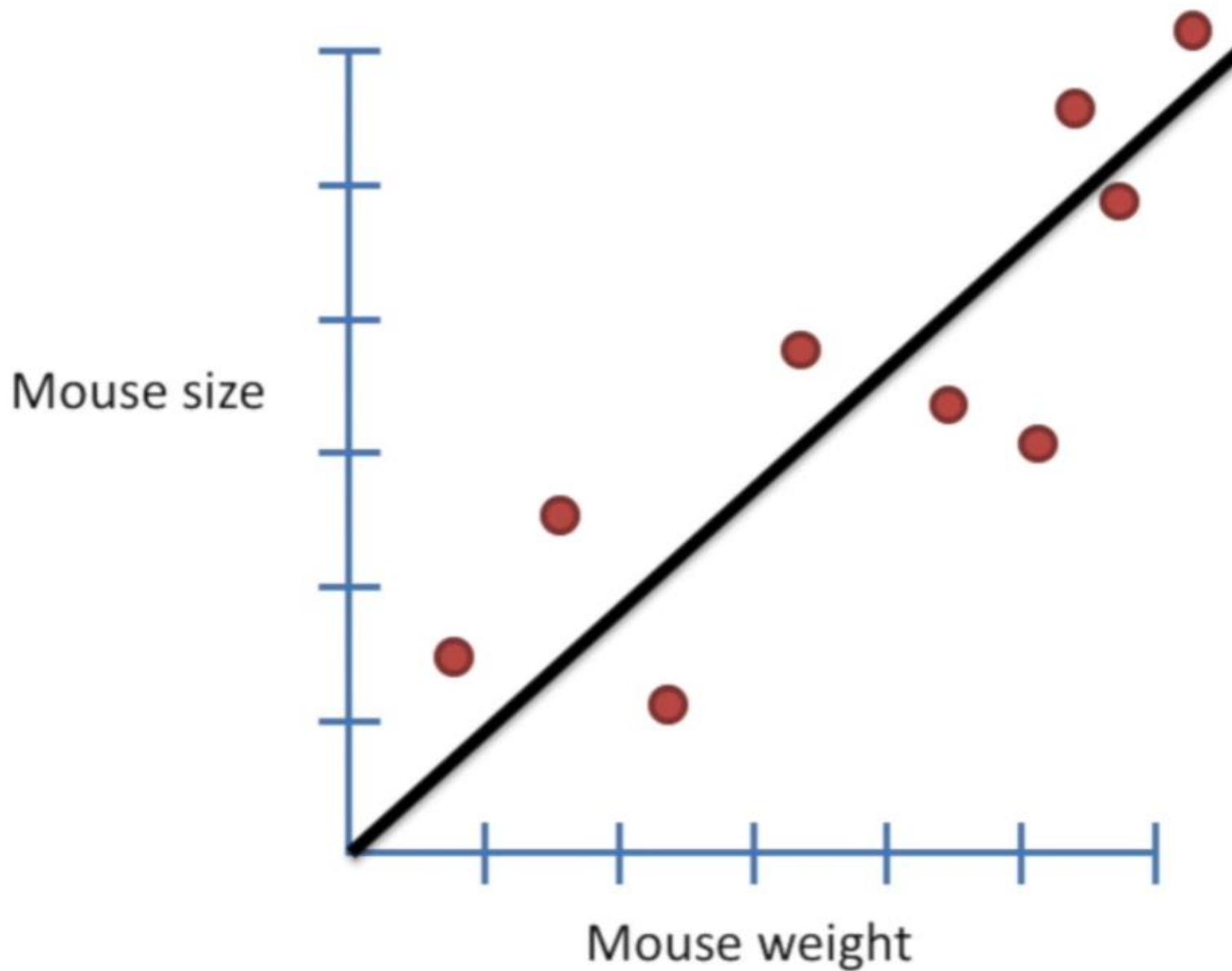




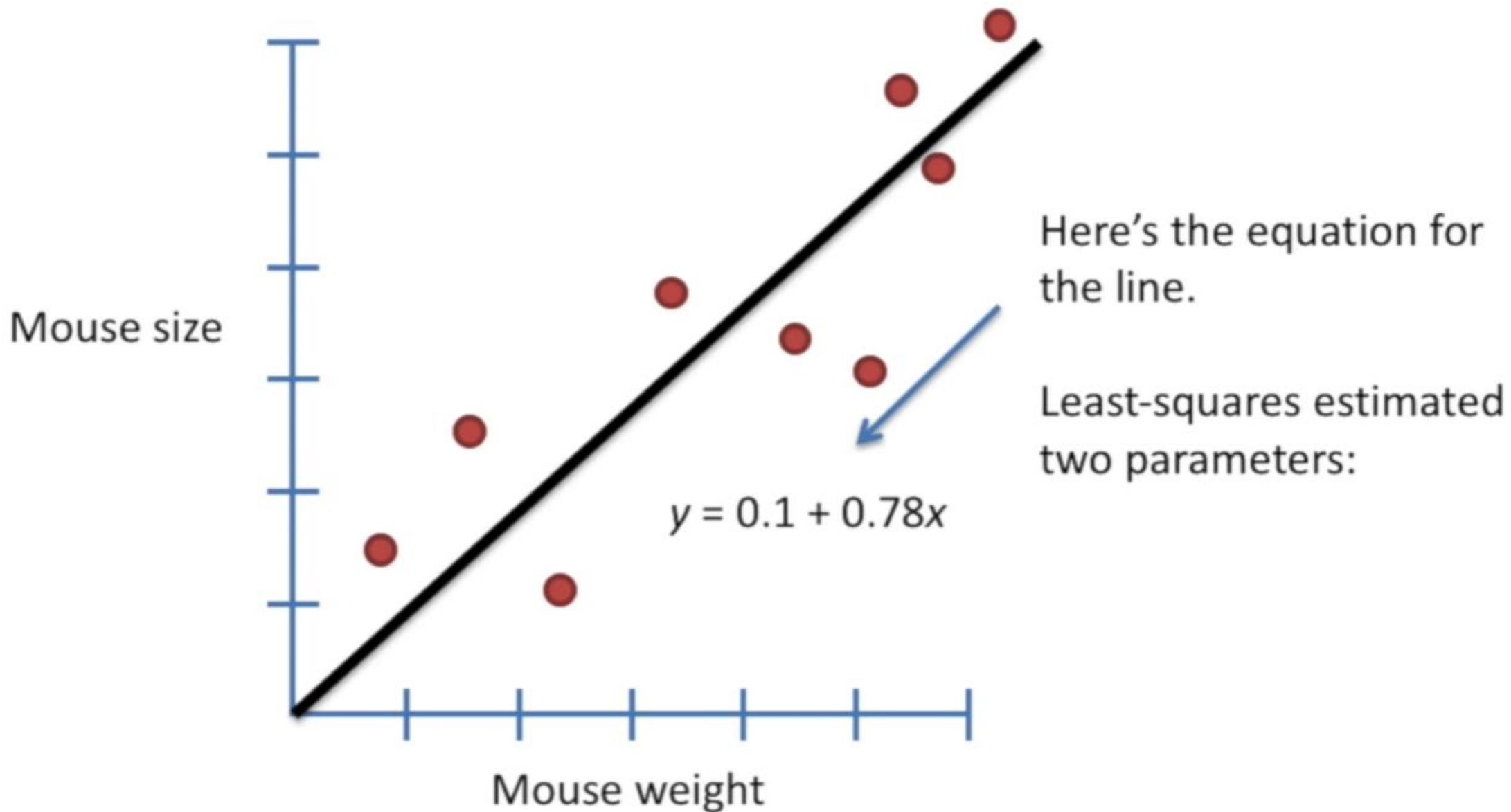




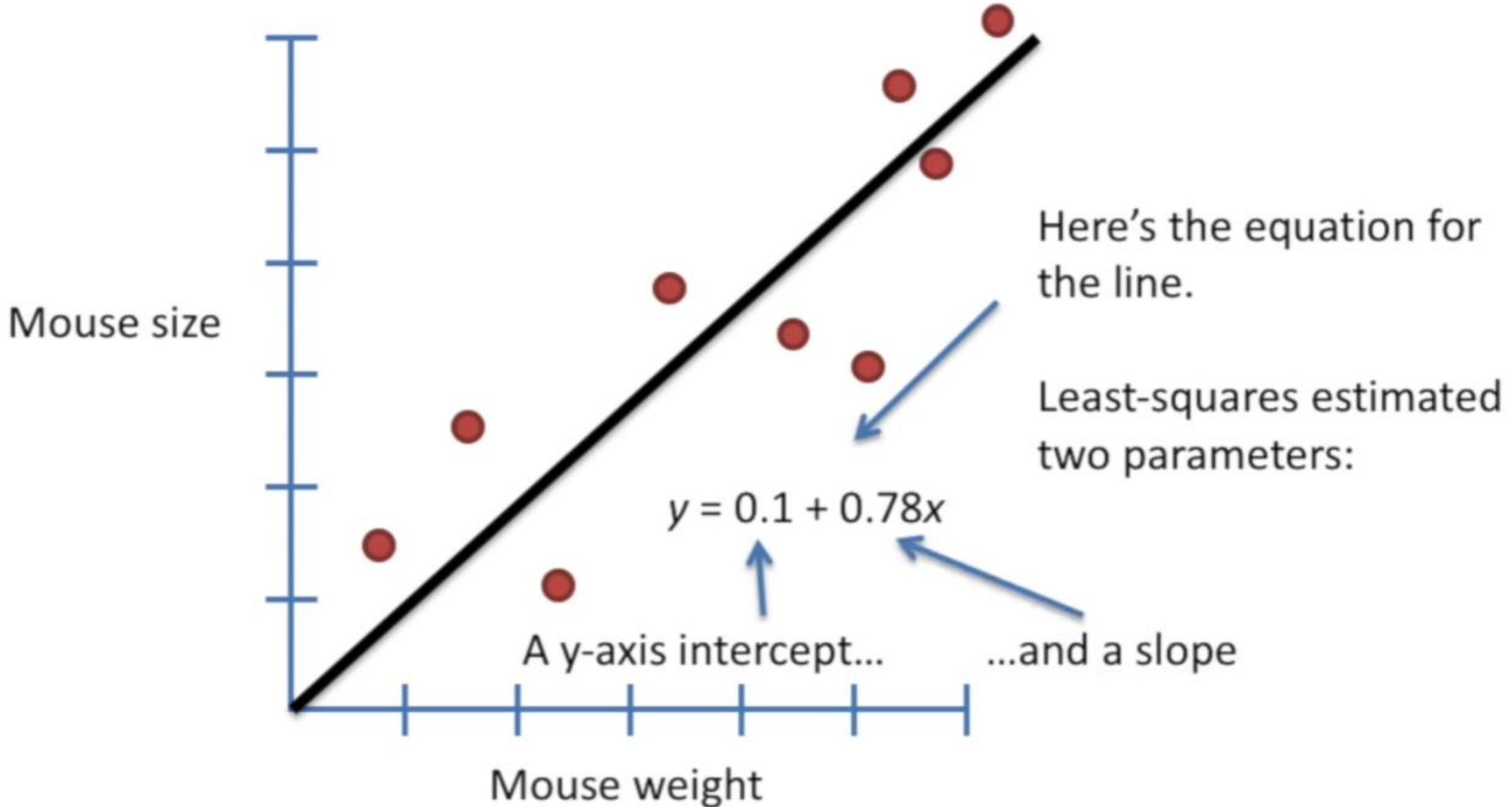
Now we have fit a line to the data! This is awesome!



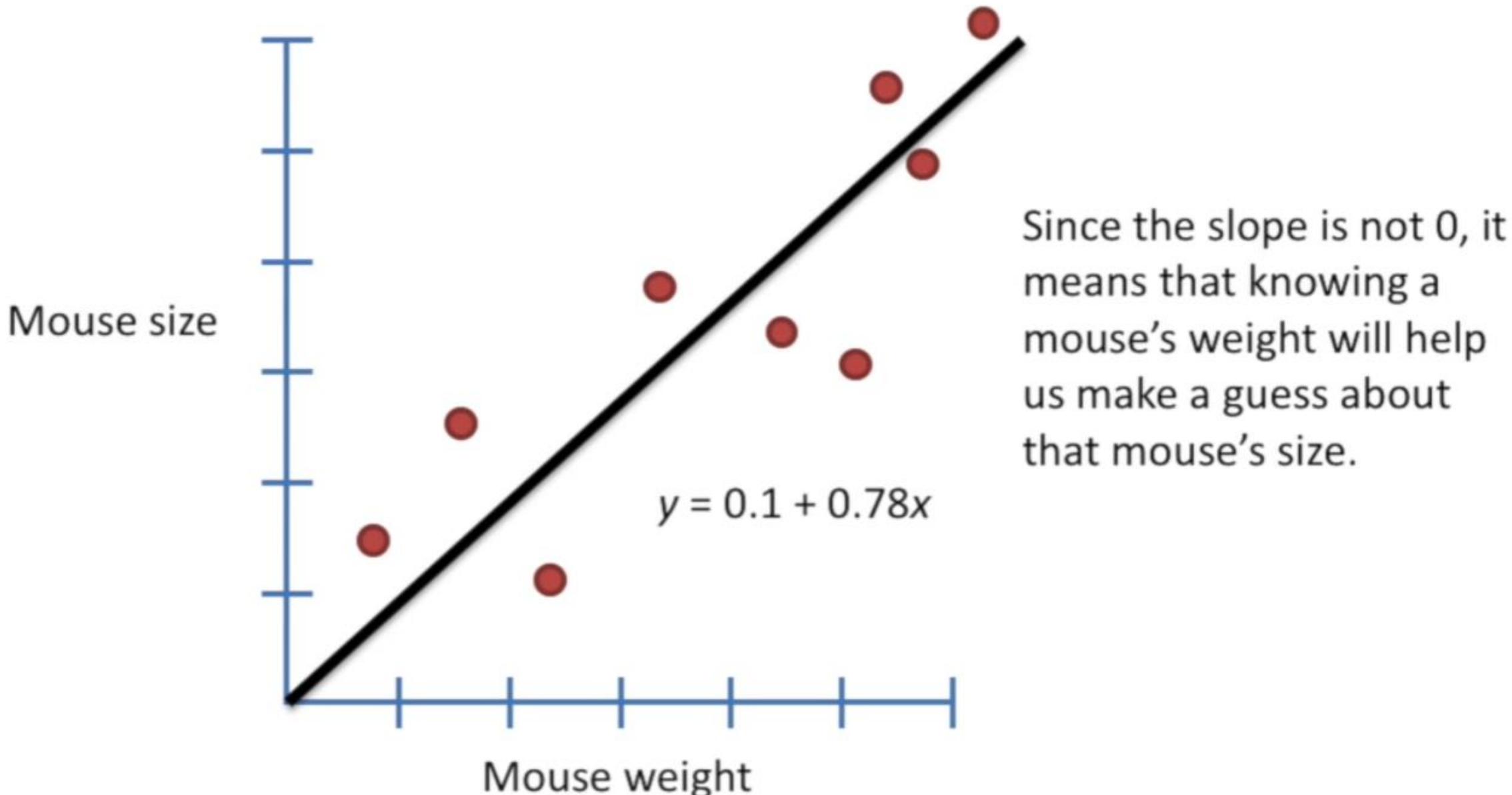
Now we have fit a line to the data! This is awesome!



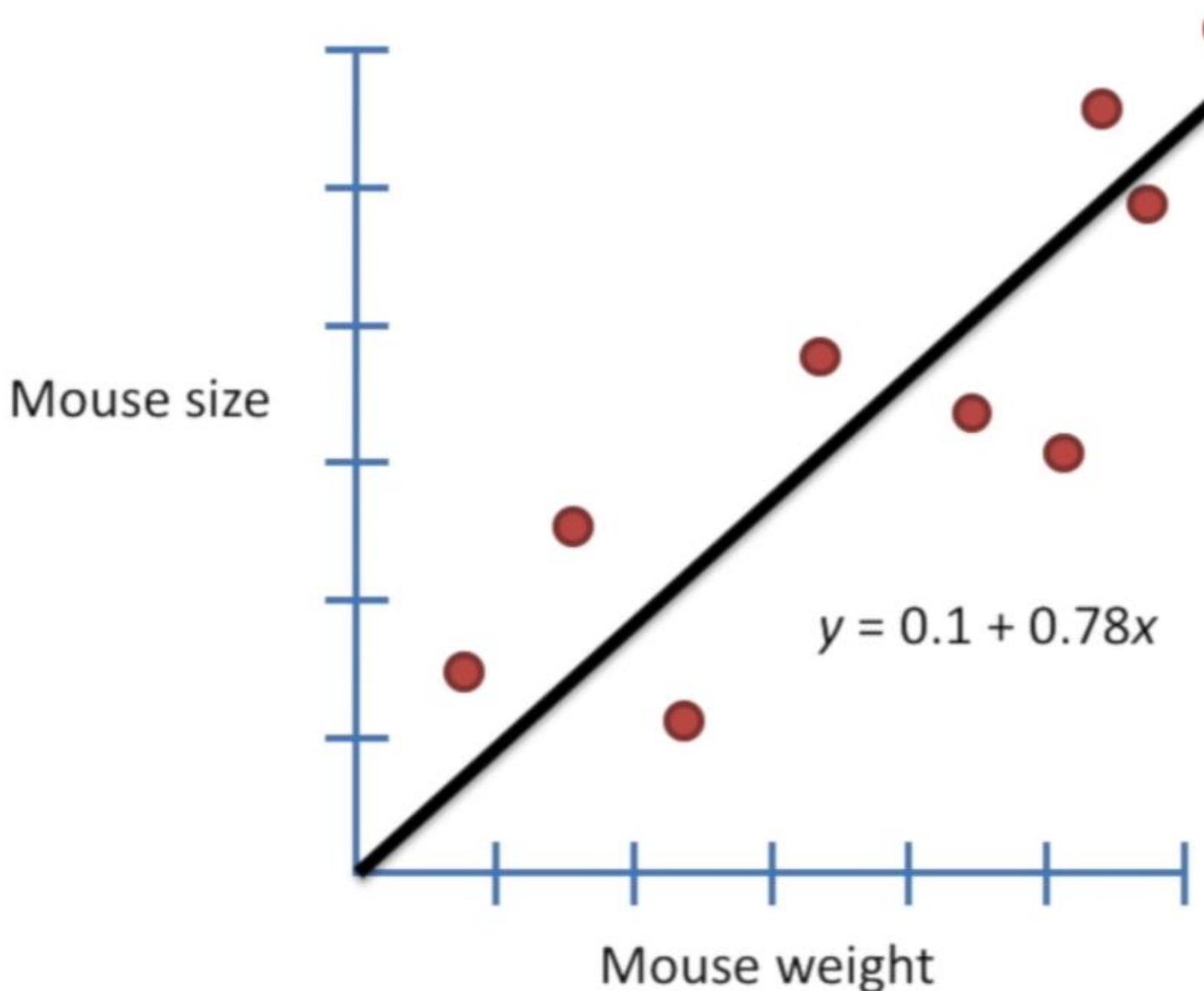
Now we have fit a line to the data! This is awesome!



Now we have fit a line to the data! This is awesome!



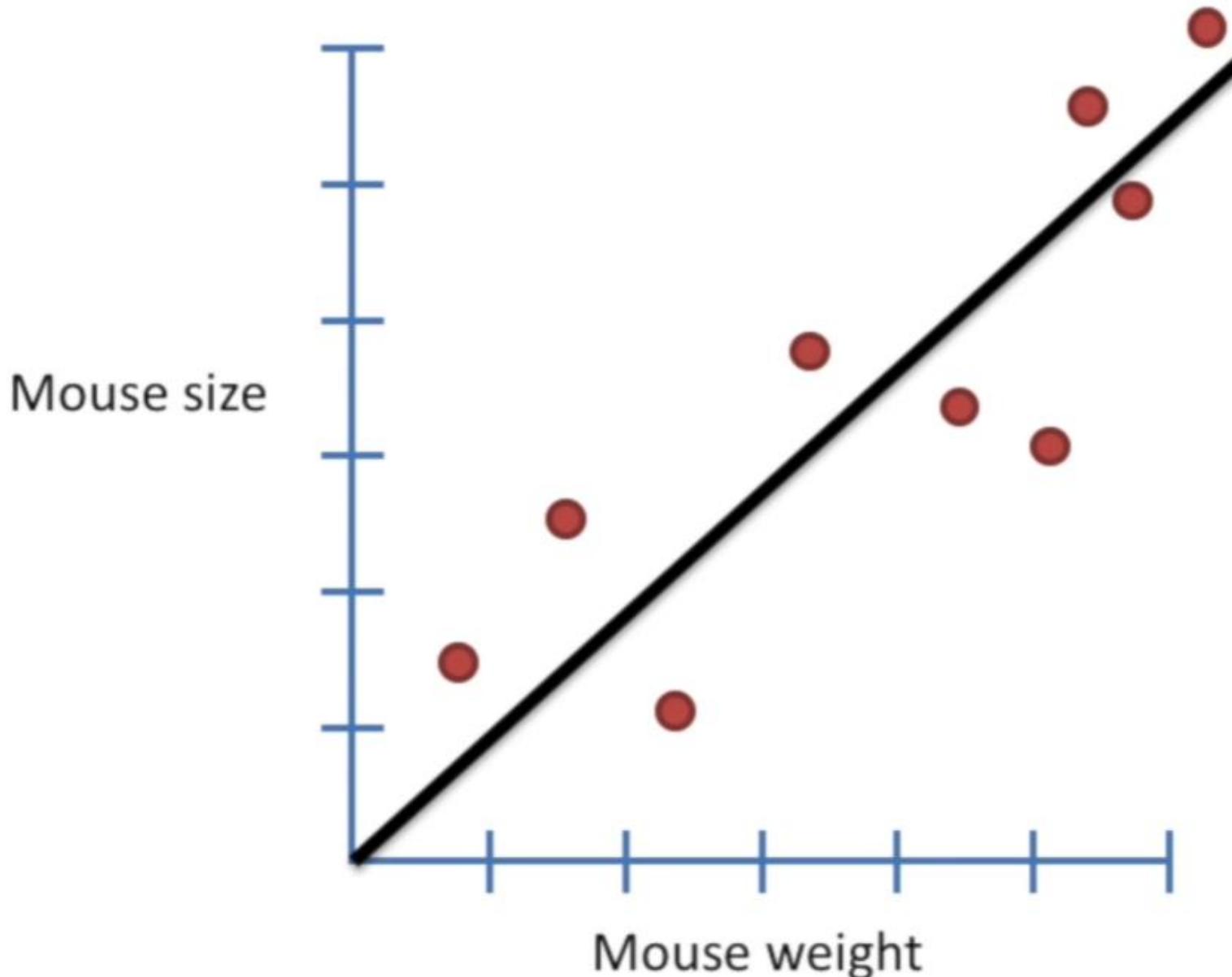
Now we have fit a line to the data! This is awesome!



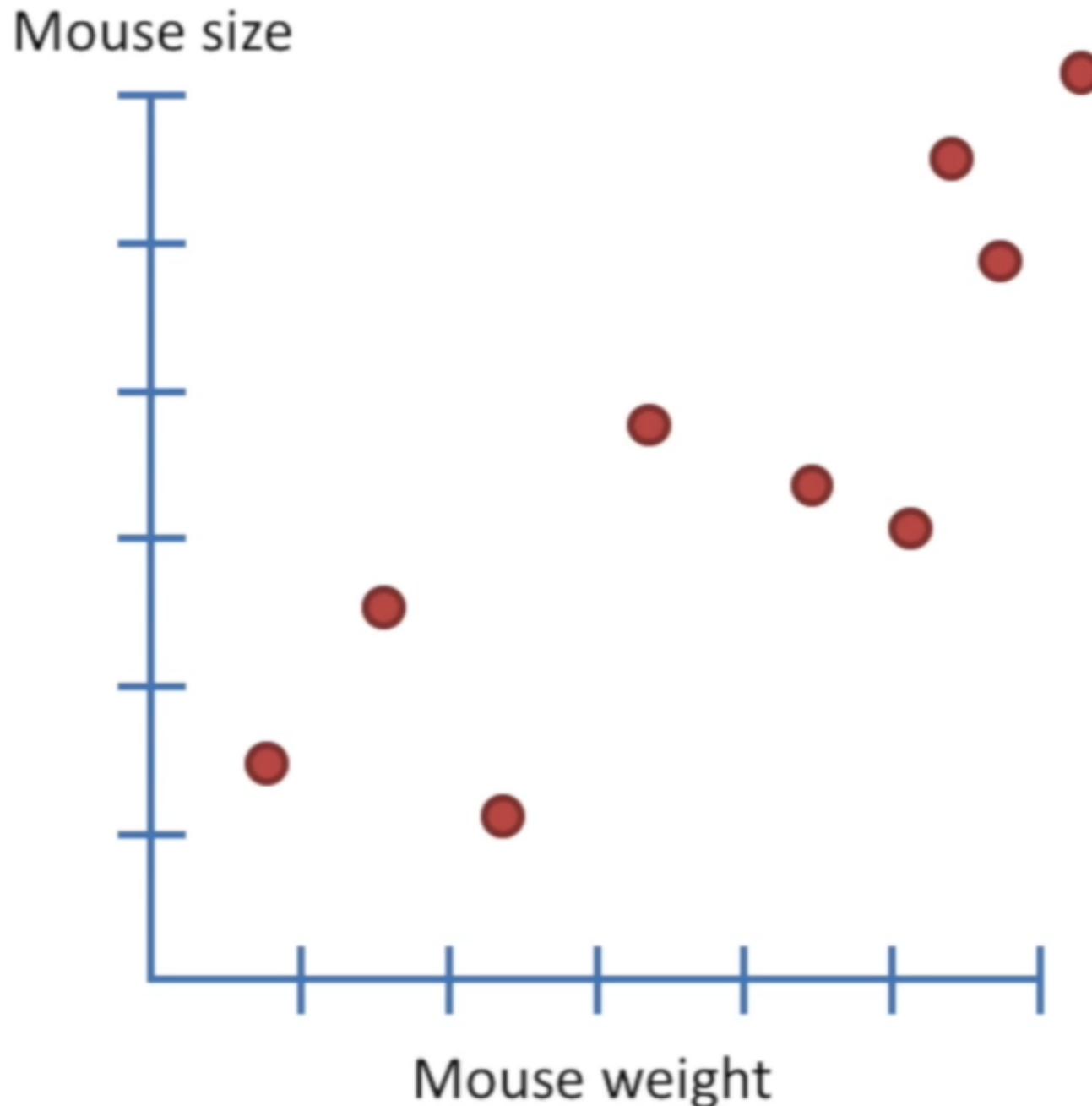
Since the slope is not 0, it means that knowing a mouse's weight will help us make a guess about that mouse's size.

How good is that guess?

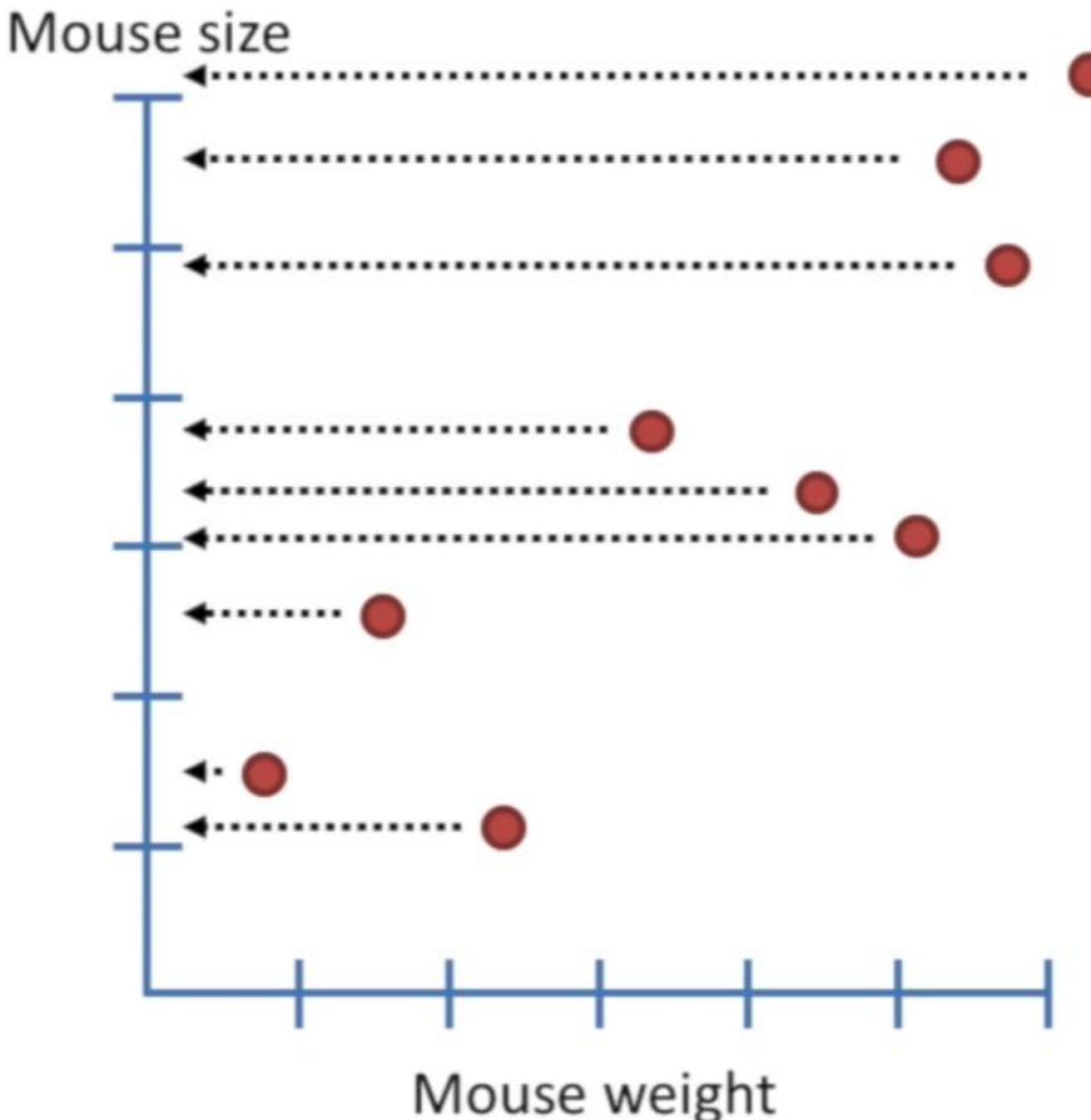
Calculating R^2 is the first step in determining how good that guess will be.



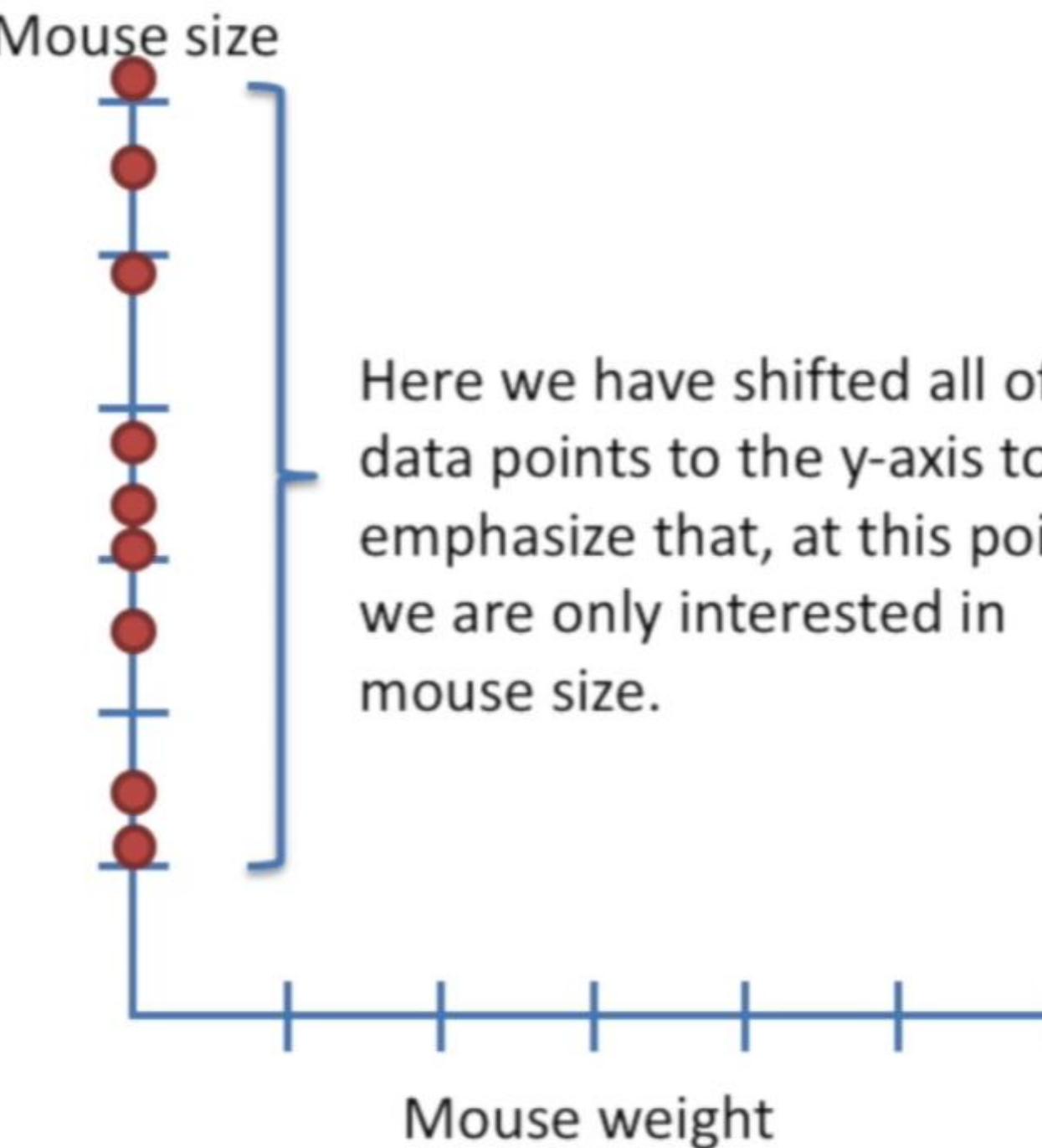
First, calculate the average mouse size.



First, calculate the average mouse size.

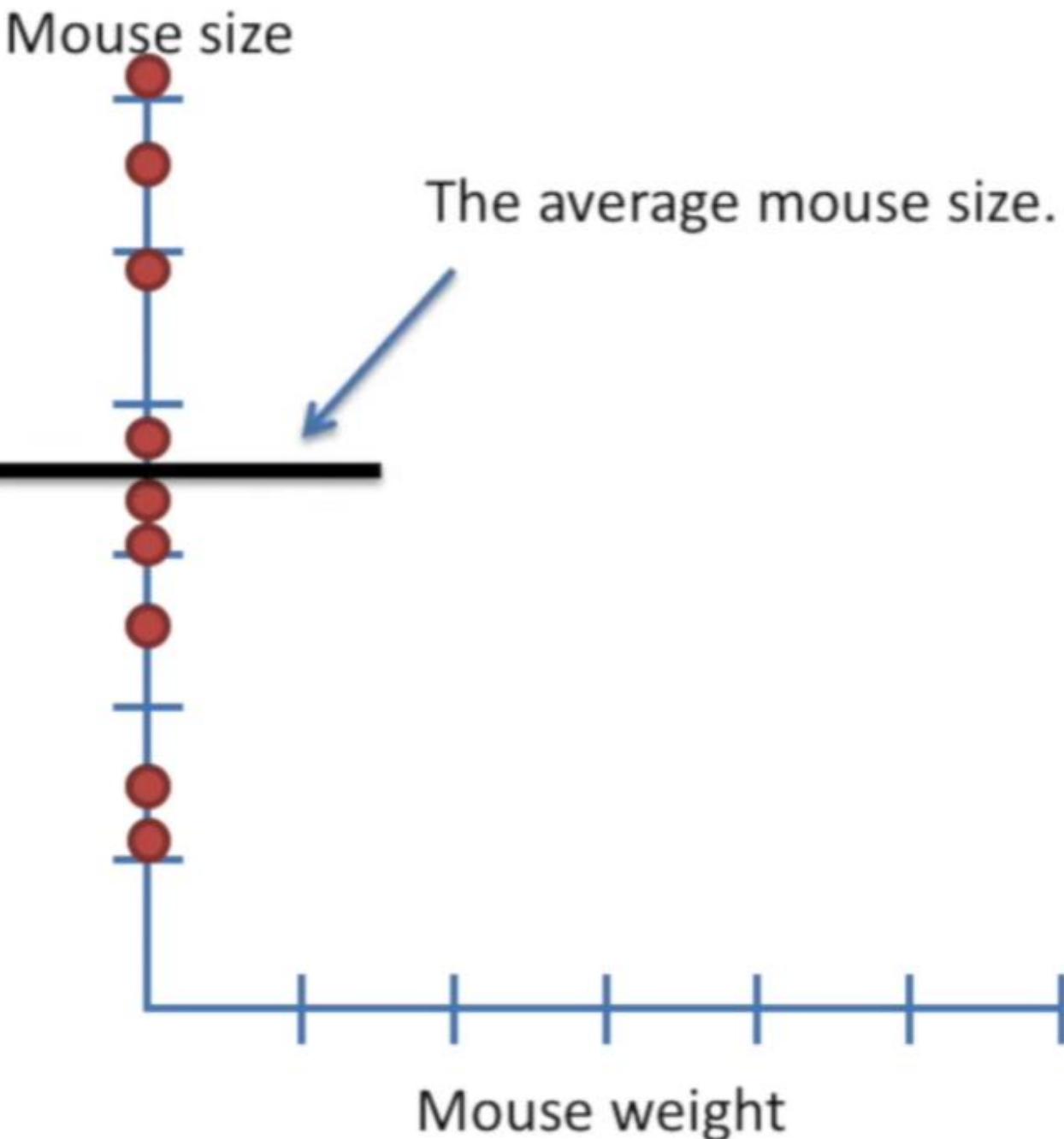


First, calculate the average mouse size.

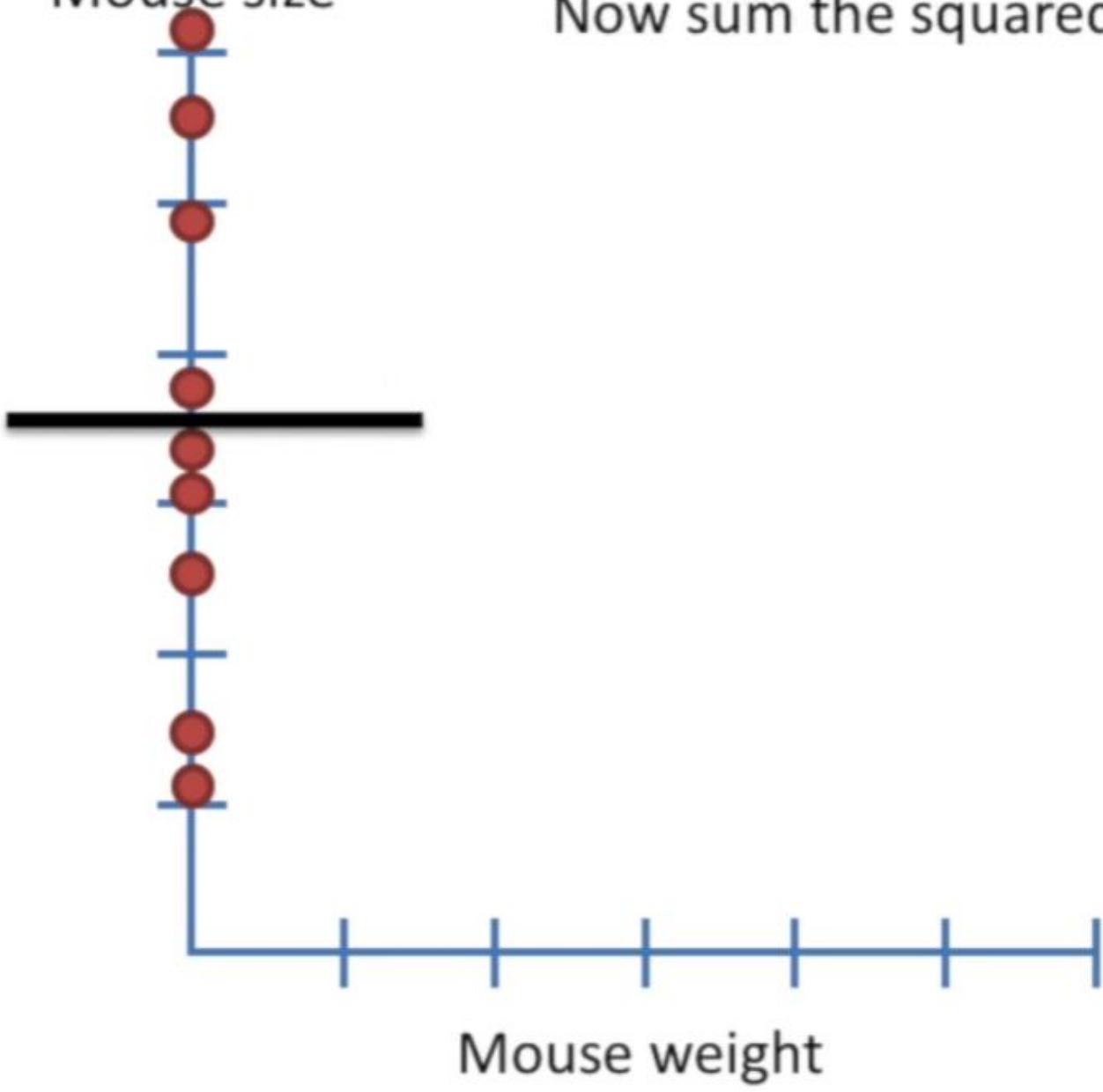


Here we have shifted all of the data points to the y-axis to emphasize that, at this point, we are only interested in mouse size.

First, calculate the average mouse size.

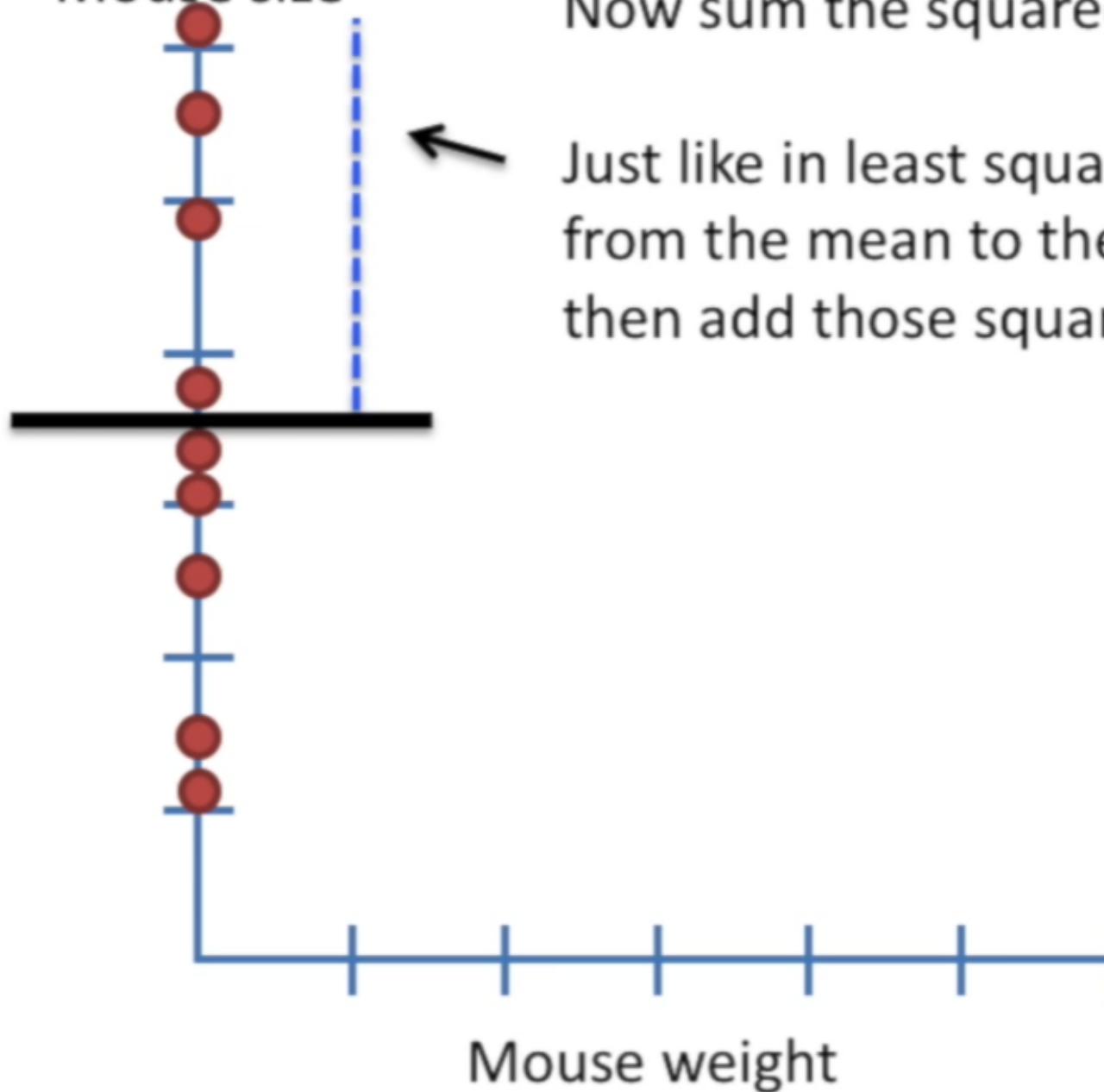


Mouse size



Now sum the squared residuals...

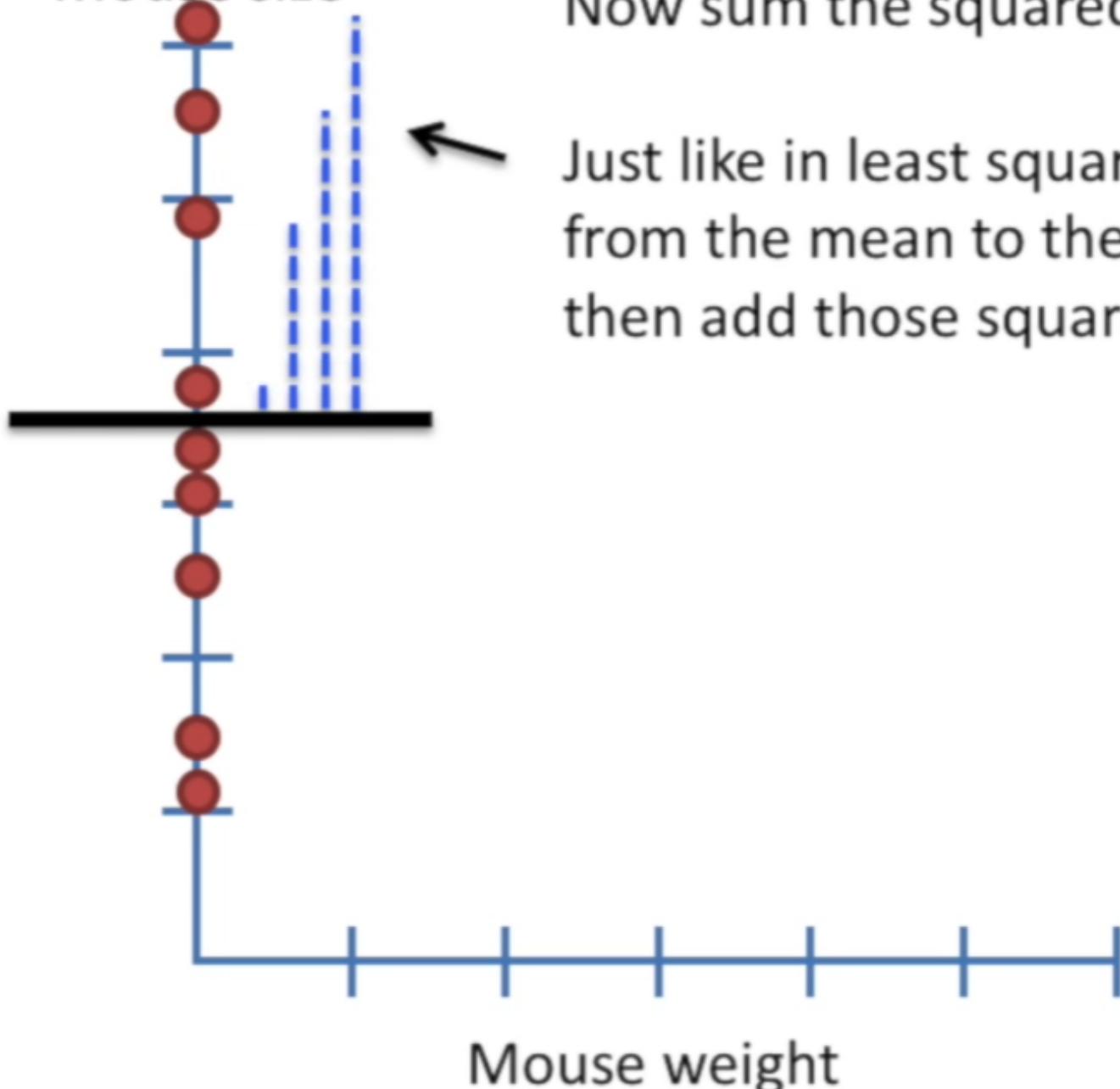
Mouse size



Now sum the squared residuals...

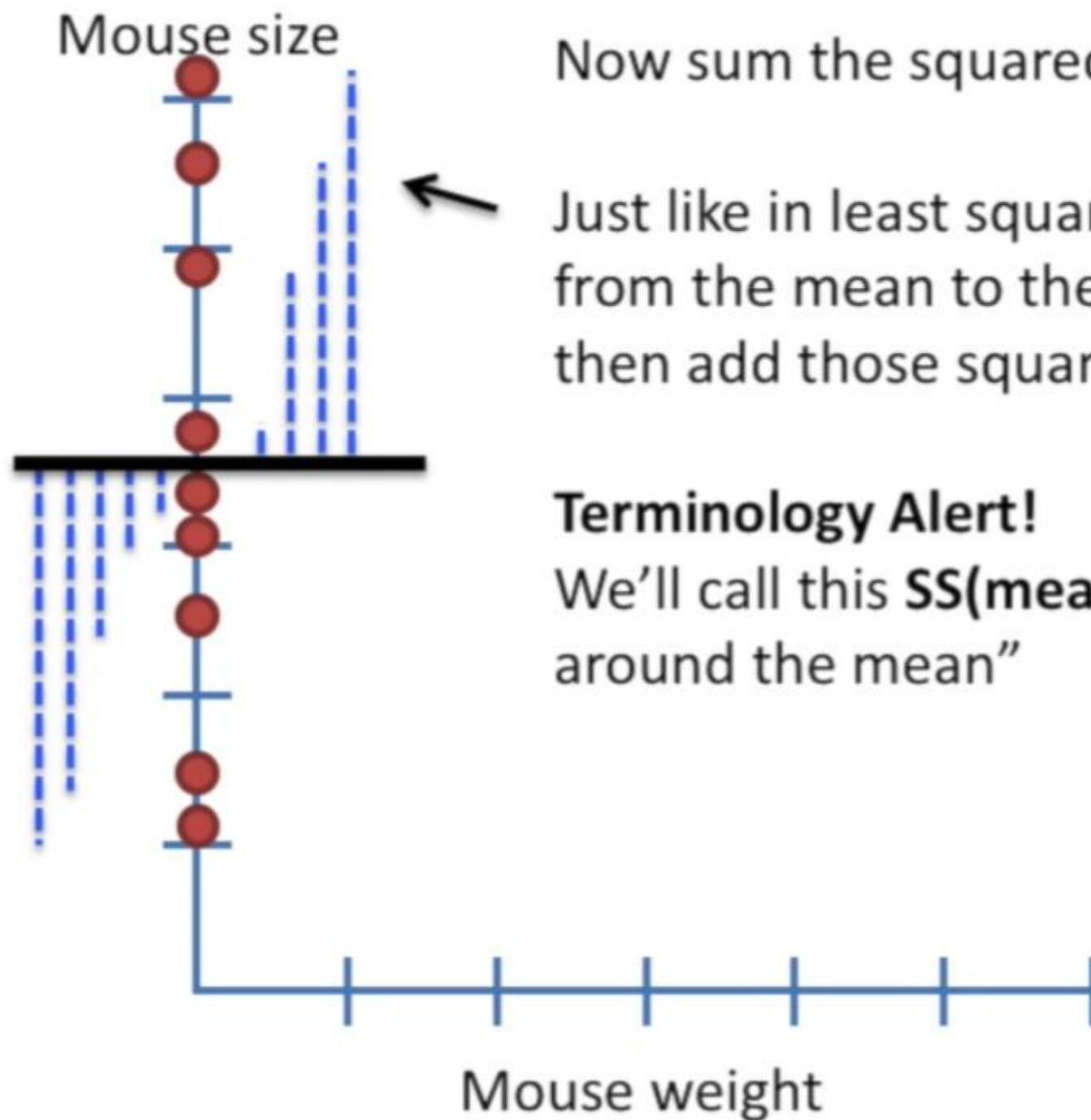
Just like in least squares, we measure the distance from the mean to the data point and square it, then add those squares together.

Mouse size



Now sum the squared residuals...

Just like in least squares, we measure the distance from the mean to the data point and square it, then add those squares together.



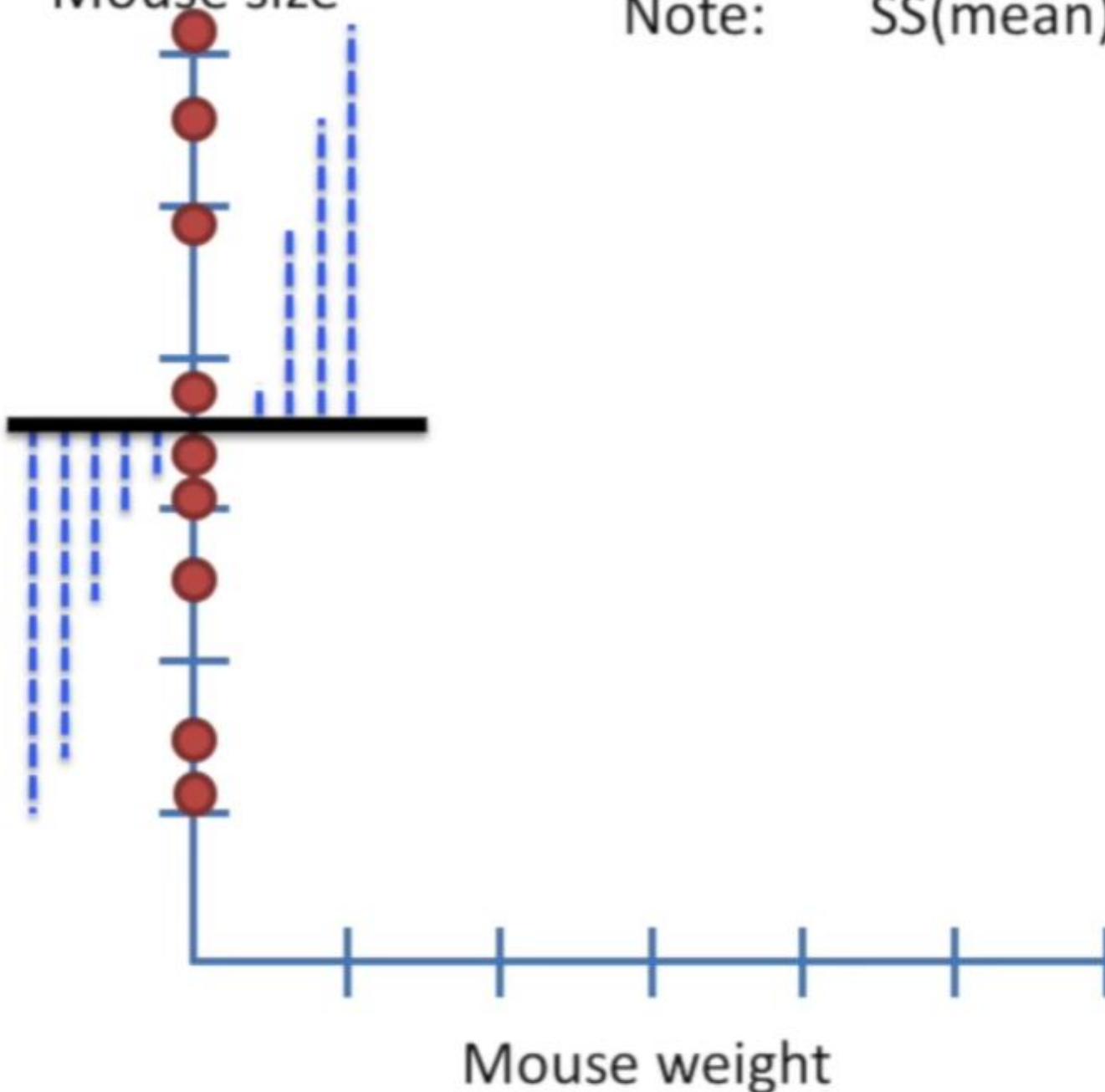
Now sum the squared residuals...

Just like in least squares, we measure the distance from the mean to the data point and square it, then add those squares together.

Terminology Alert!

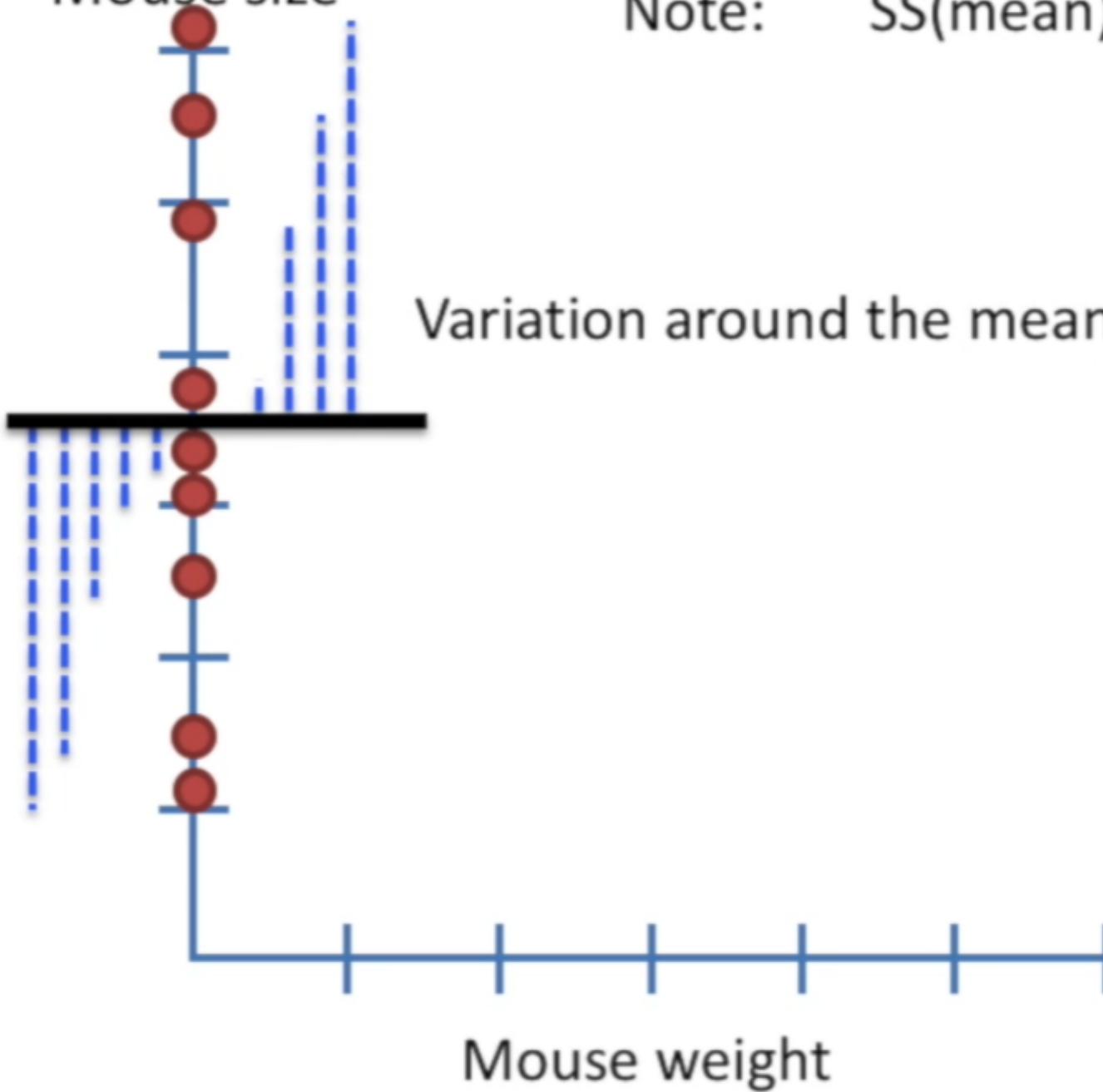
We'll call this **SS(mean)**, for "sum of squares around the mean"

Mouse size



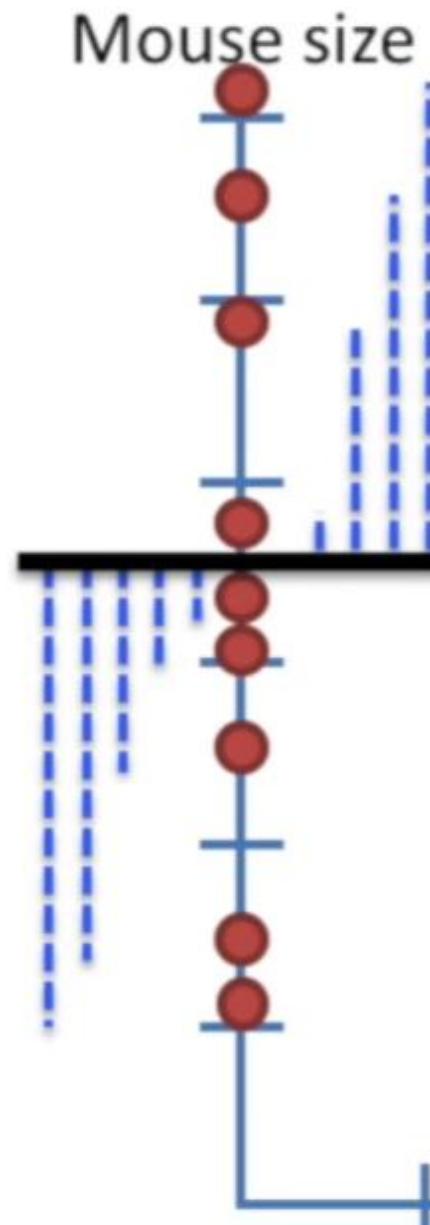
Note: $SS(\text{mean}) = (\text{data} - \text{mean})^2$

Mouse size



Note: $SS(\text{mean}) = (\text{data} - \text{mean})^2$

$$\frac{(\text{data} - \text{mean})^2}{n}$$



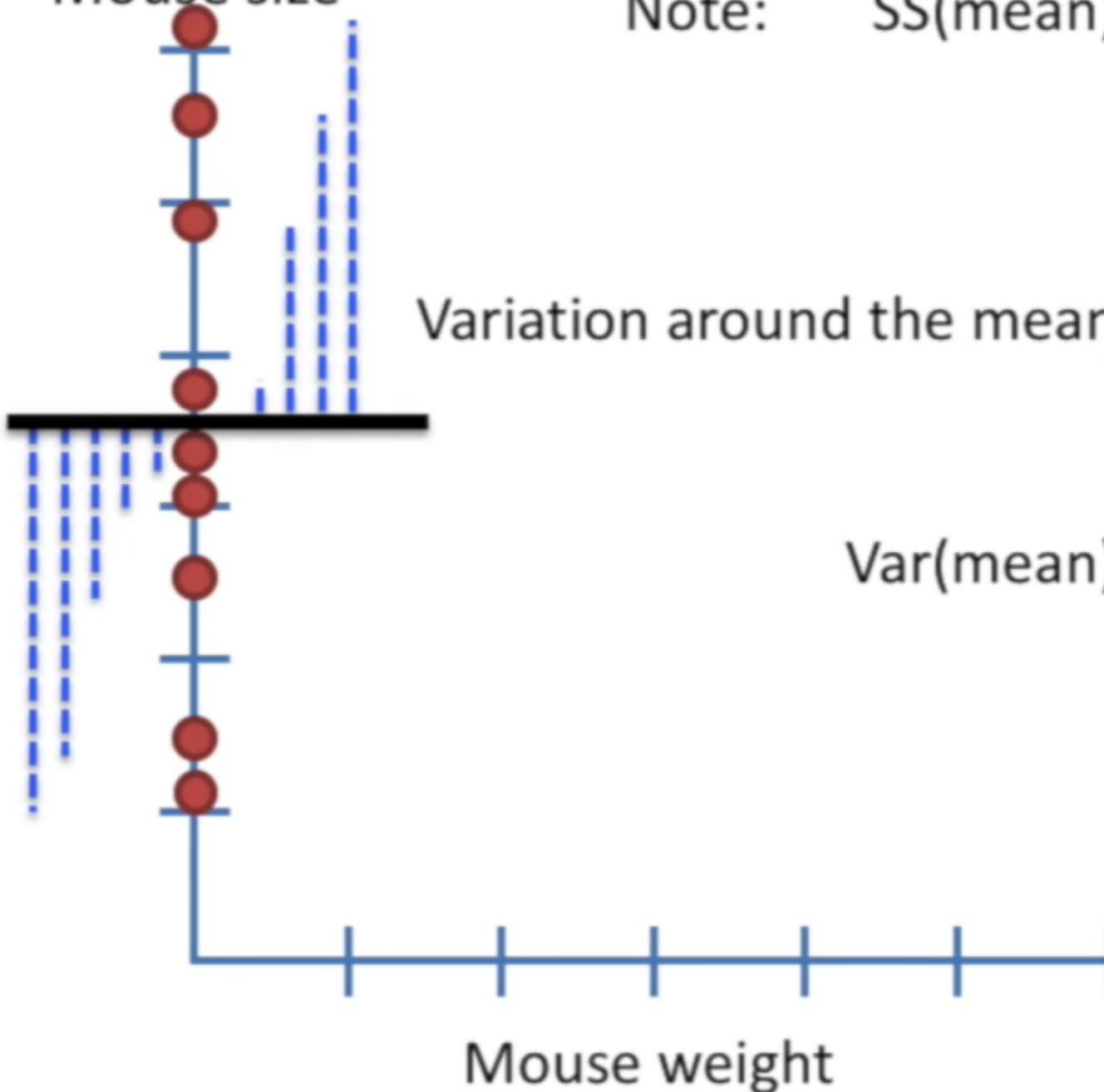
Note: $SS(\text{mean}) = (\text{data} - \text{mean})^2$

Variation around the mean =

$$\frac{(\text{data} - \text{mean})^2}{n}$$

'n' is the sample size (in this case, n= 9)

Mouse size

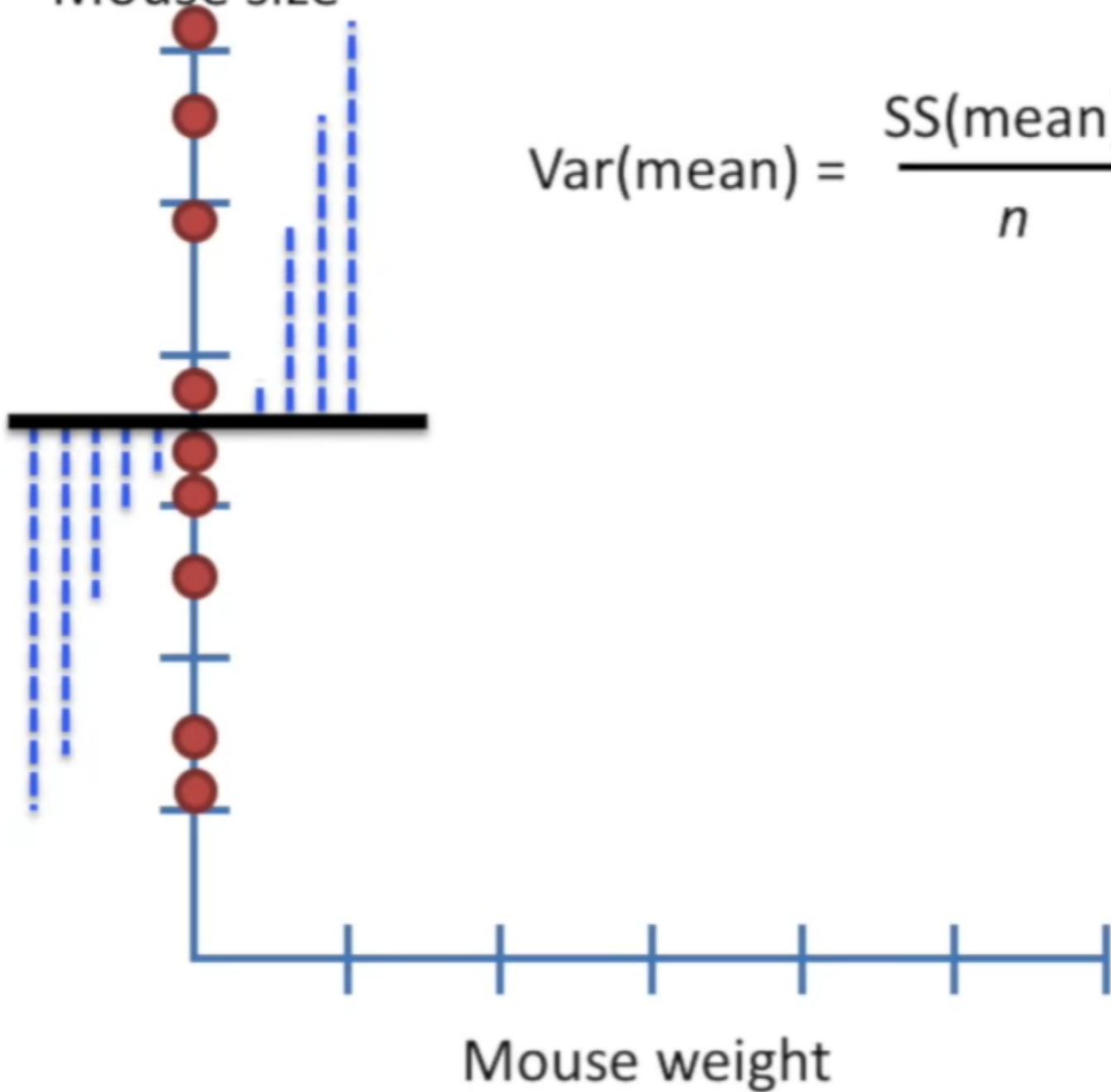


Note: $SS(\text{mean}) = (\text{data} - \text{mean})^2$

Variation around the mean = $\frac{(\text{data} - \text{mean})^2}{n}$

$\text{Var}(\text{mean}) = \frac{SS(\text{mean})}{n}$

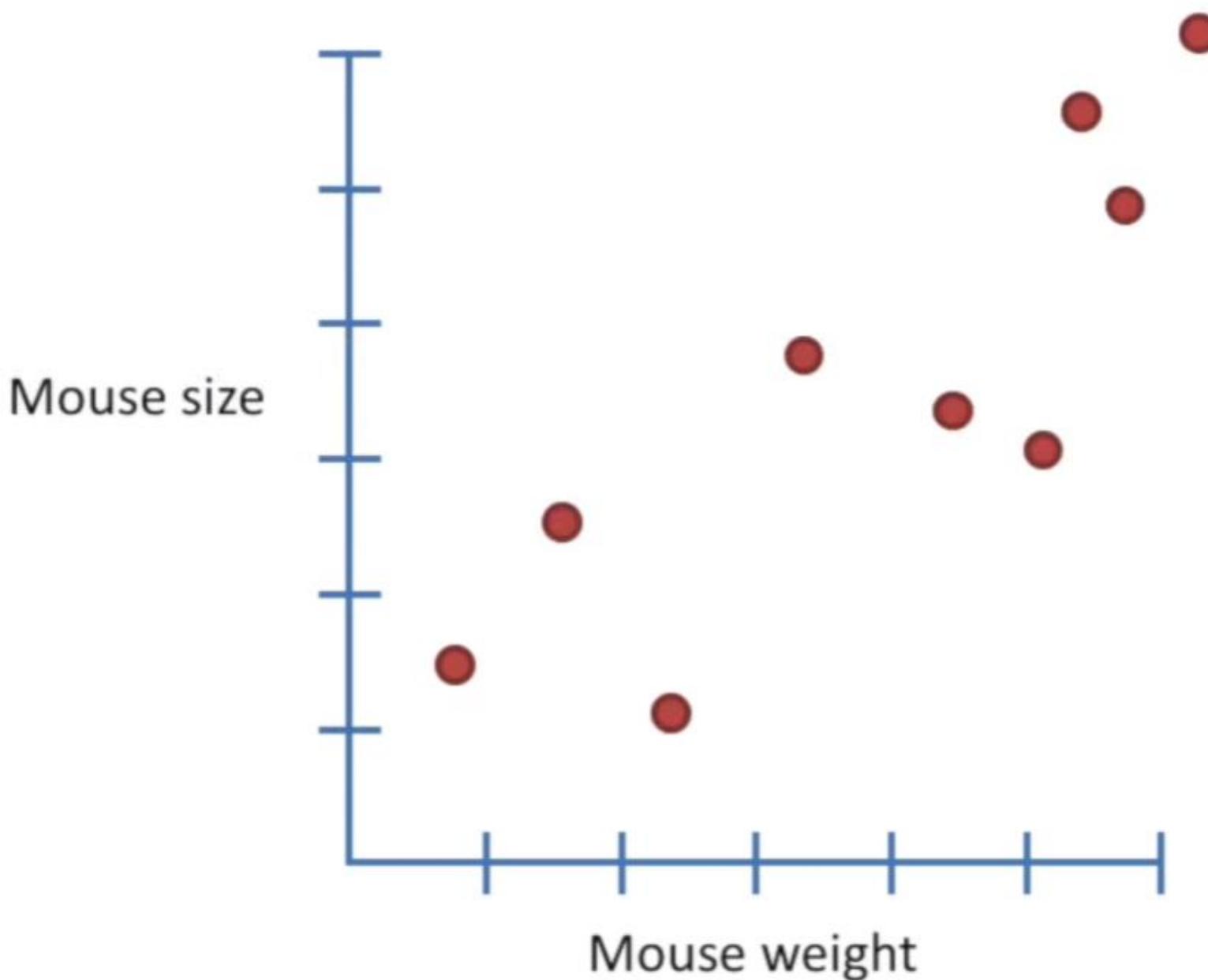
Mouse size



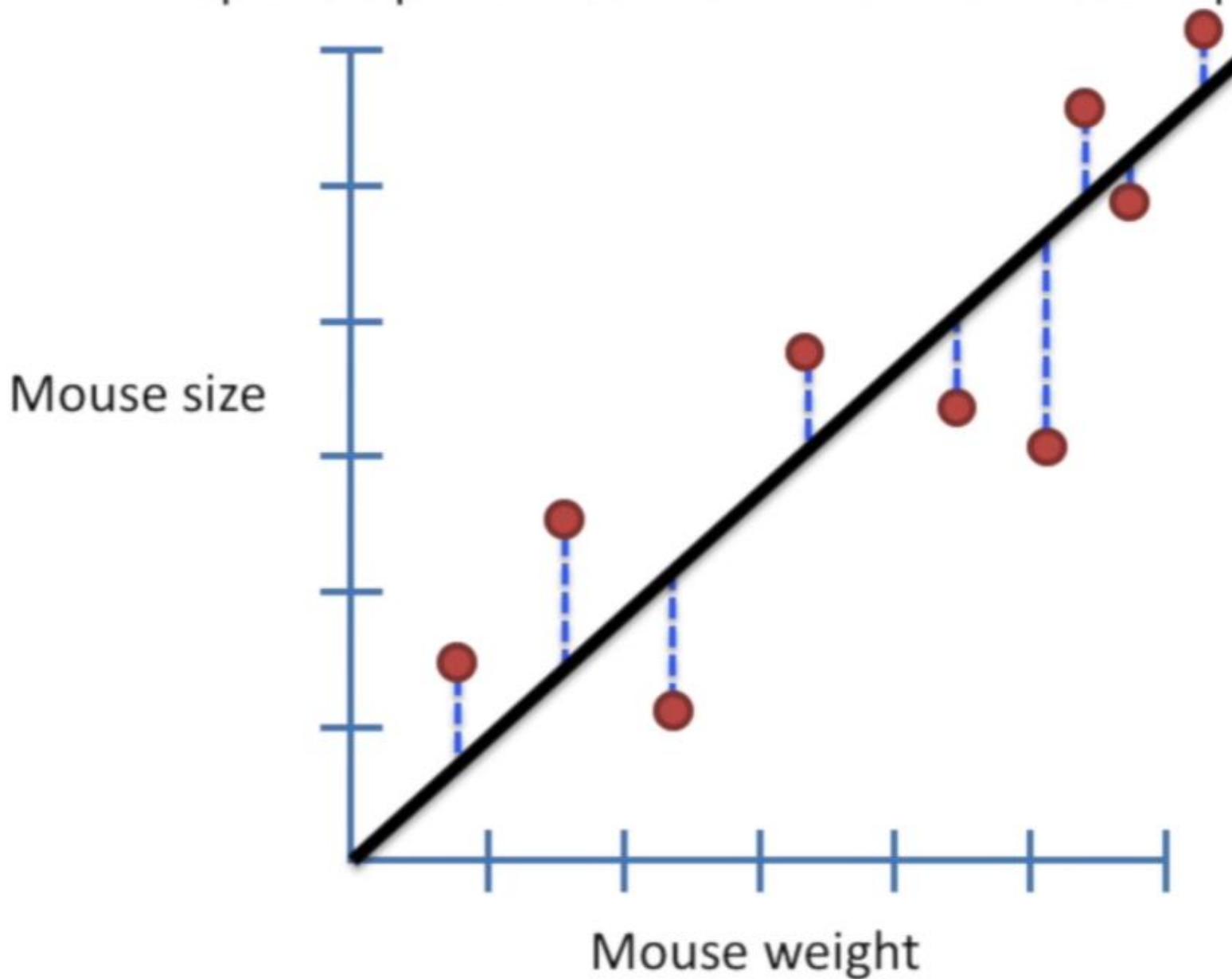
$$\text{Var}(\text{mean}) = \frac{\text{SS}(\text{mean})}{n}$$

Another way to think about variance is as the average sum of squares per mouse.

Now go back to the original plot.

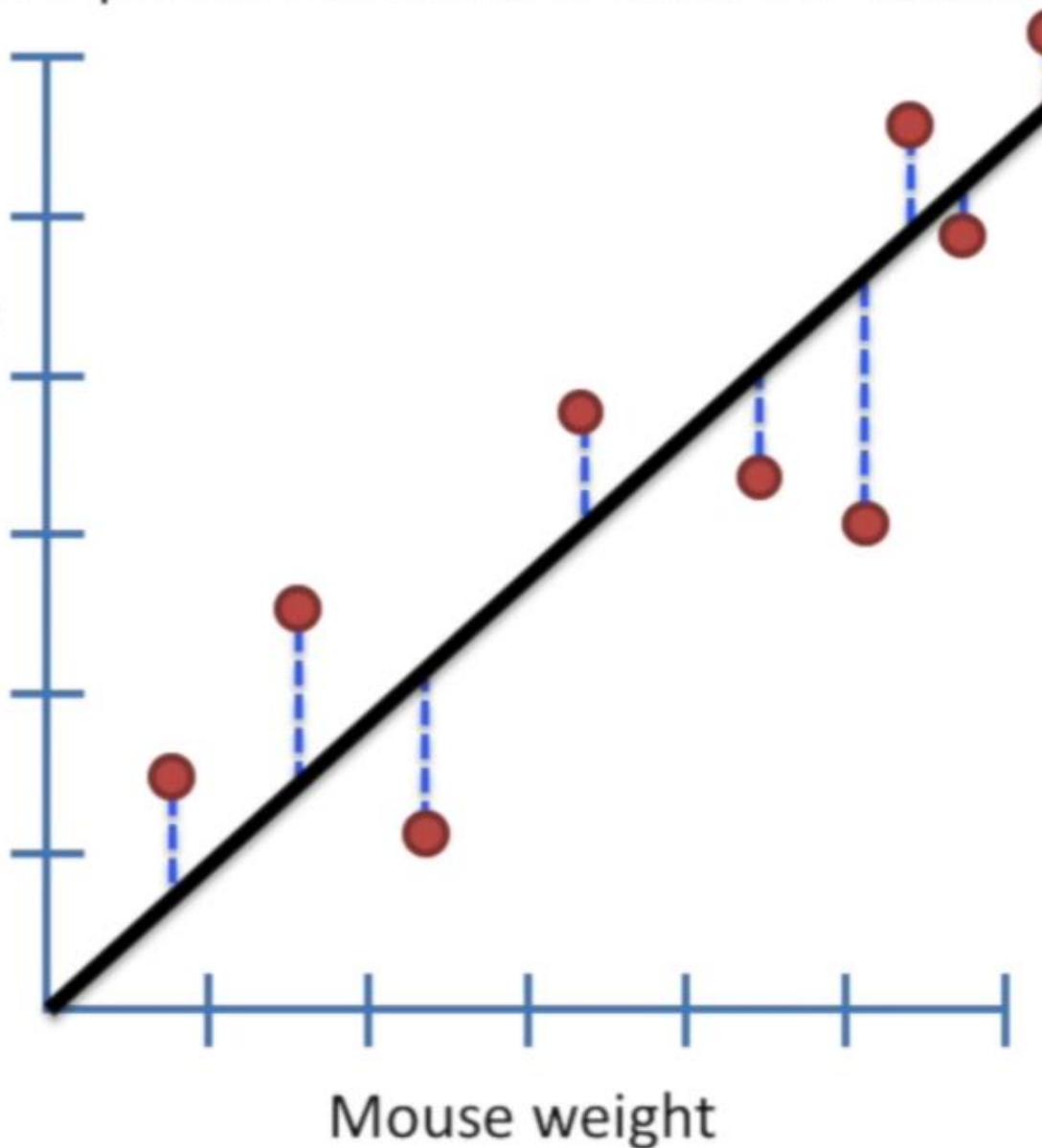


Now go back to the original plot.
Sum up the squared residuals around our least-squares fit.



Now go back to the original plot.
Sum up the squared residuals around our least-squares fit.

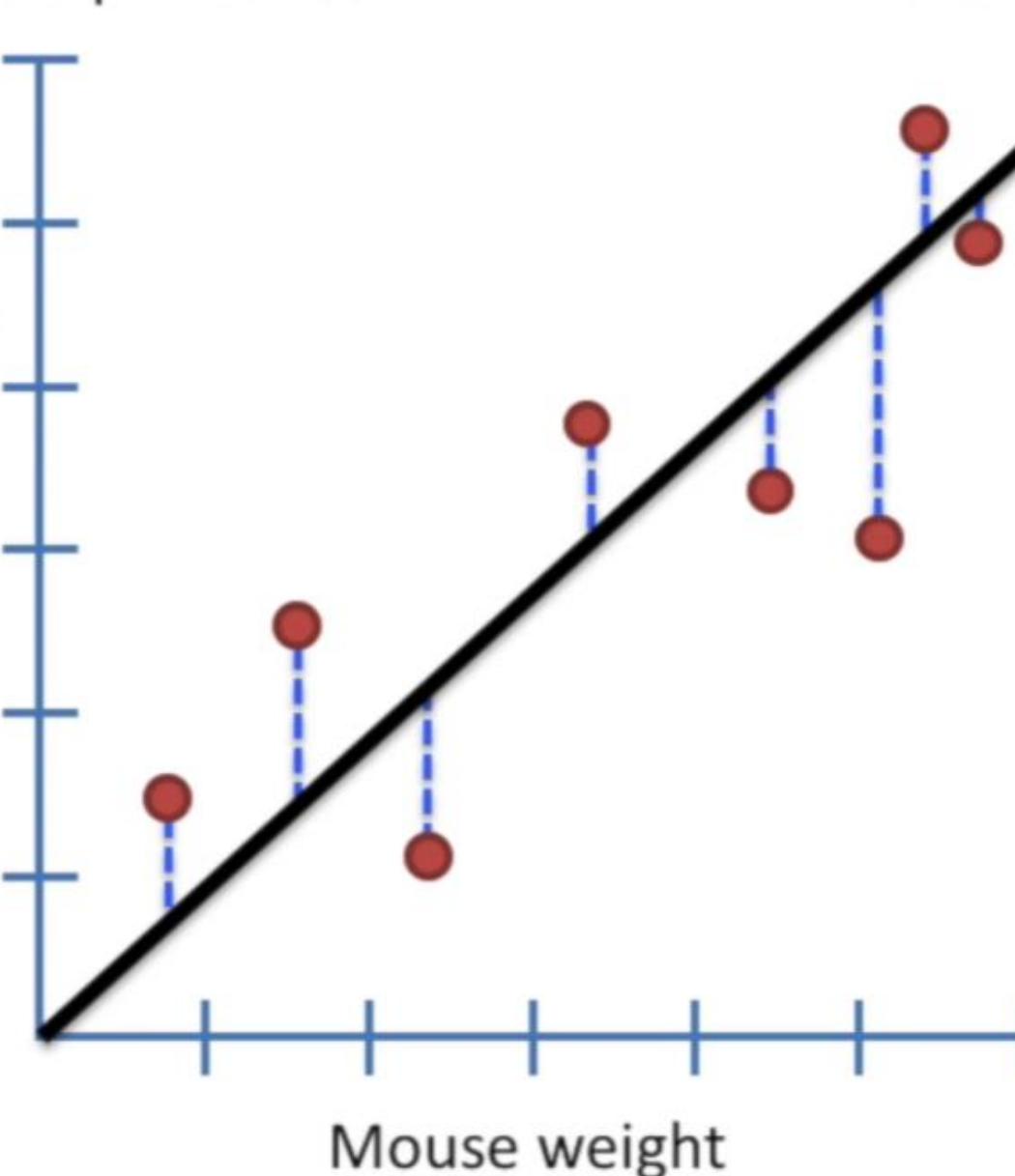
We'll call this **SS(fit)**, for the sum of squares around the least-squares fit.



Now go back to the original plot.
Sum up the squared residuals around our least-squares fit.

We'll call this **SS(fit)**, for the sum of squares around the least-squares fit.

$$SS(\text{fit}) = (\text{data} - \text{line})^2$$

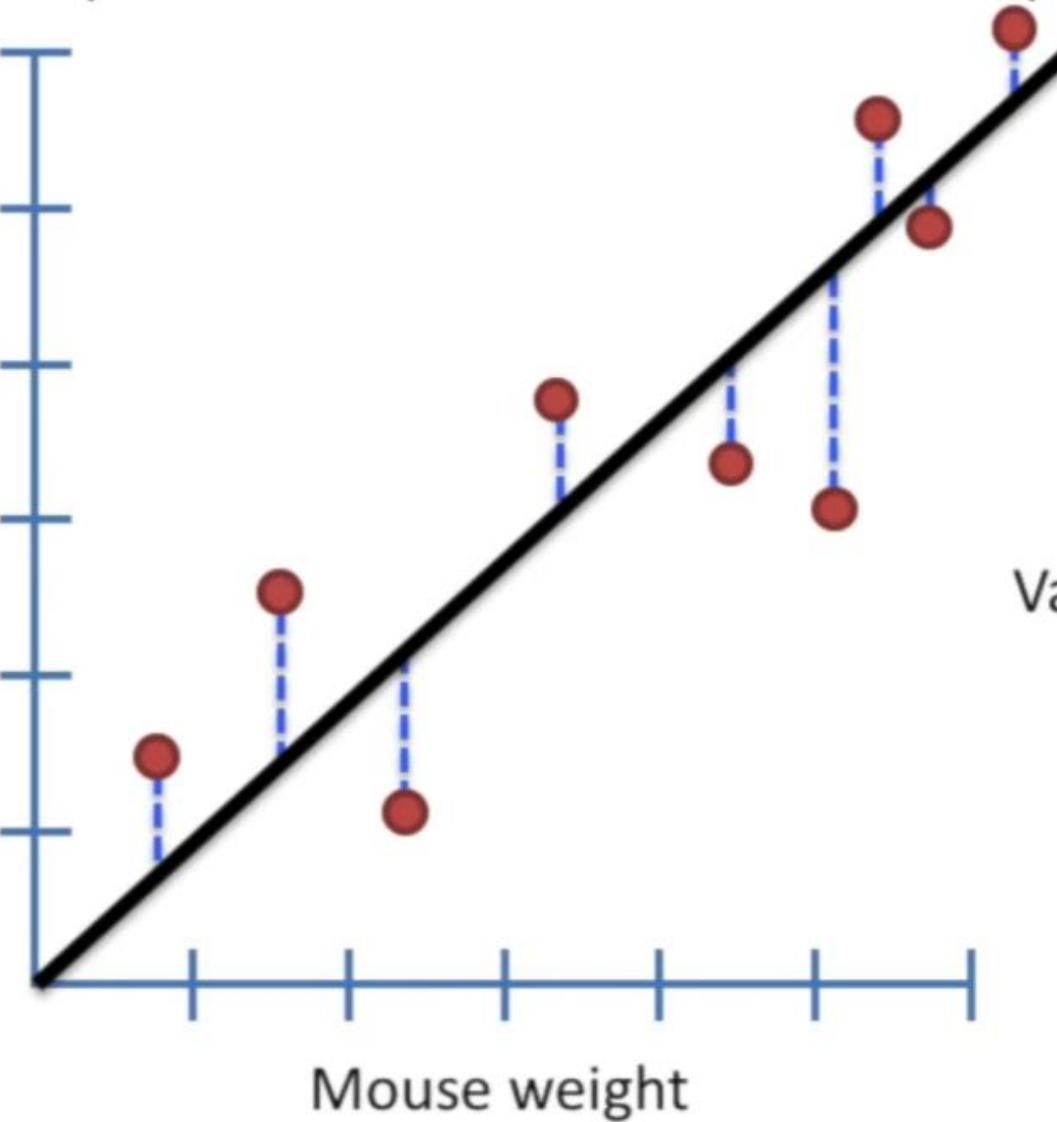


Now go back to the original plot.

Sum up the squared residuals around our least-squares fit.

We'll call this **SS(fit)**, for the sum of squares around the least-squares fit.

$$SS(\text{fit}) = (\text{data} - \text{line})^2$$



Just like with the mean, the variance around the fit...

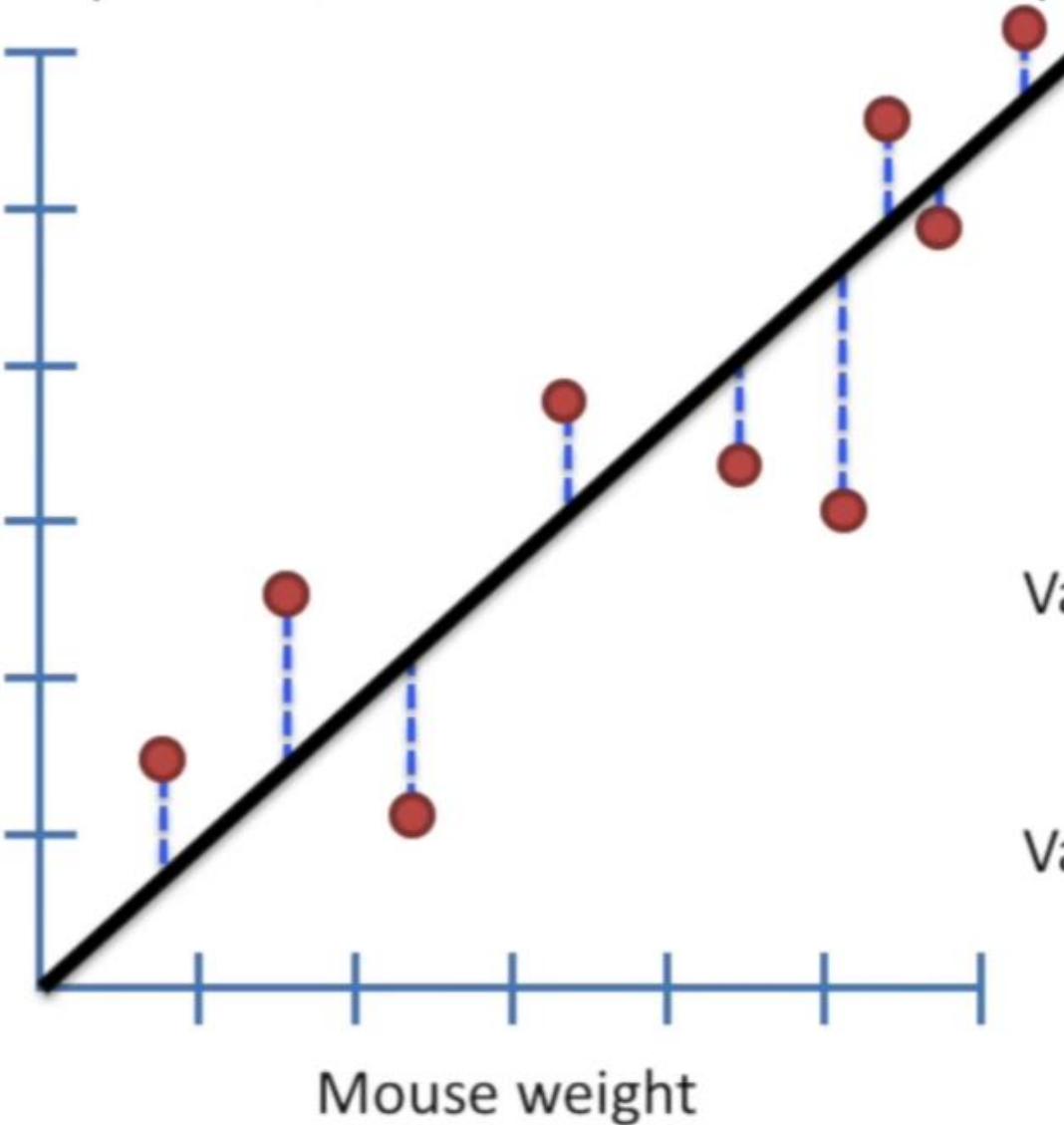
$$\text{Var}(\text{fit}) = \frac{(\text{data} - \text{line})^2}{n}$$

Now go back to the original plot.

Sum up the squared residuals around our least-squares fit.

We'll call this **SS(fit)**, for the sum of squares around the least-squares fit.

$$SS(\text{fit}) = (\text{data} - \text{line})^2$$



Just like with the mean, the variance around the fit...

$$\text{Var}(\text{fit}) = \frac{(\text{data} - \text{line})^2}{n}$$

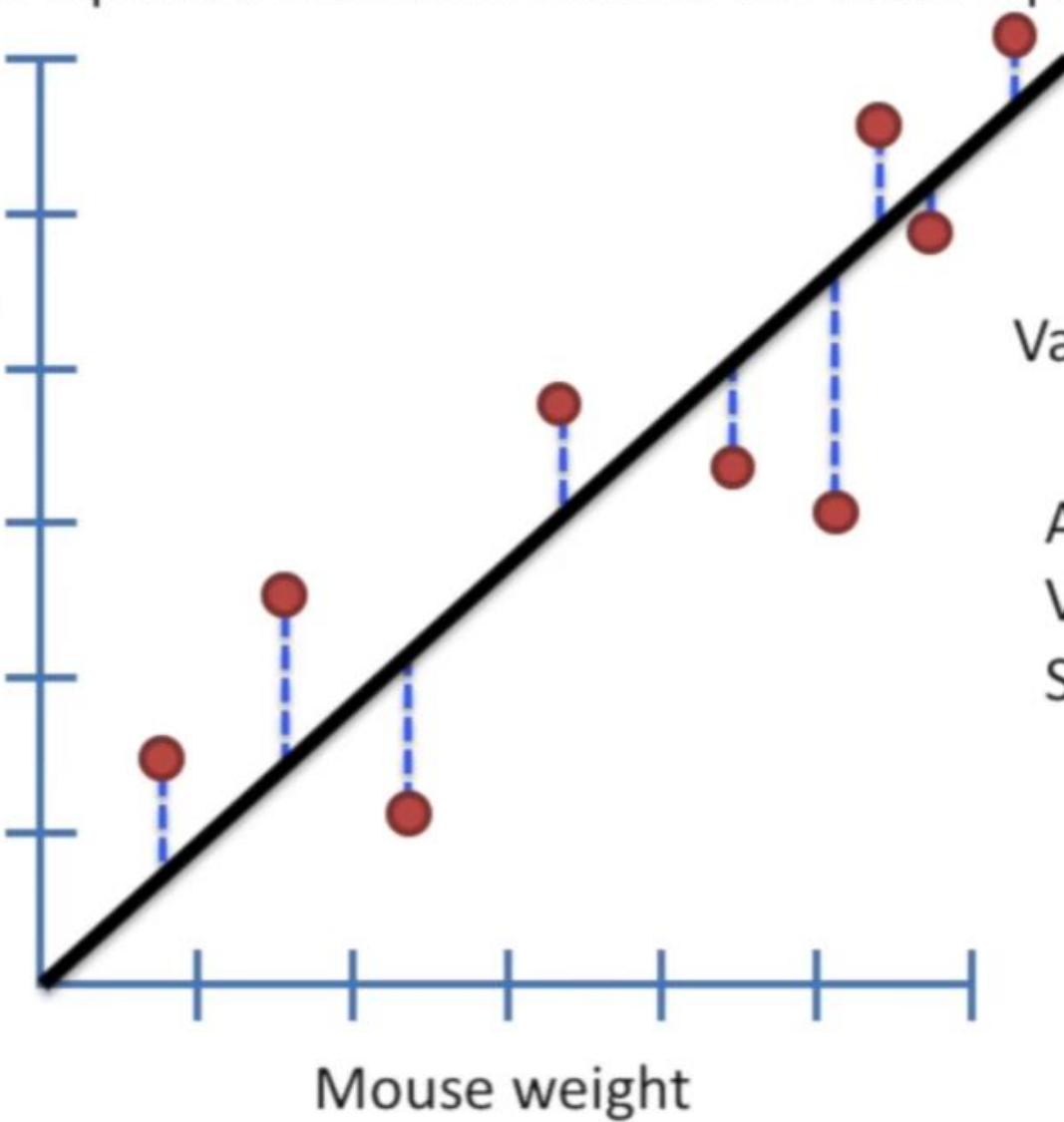
$$\text{Var}(\text{fit}) = \frac{SS(\text{fit})}{n}$$

Now go back to the original plot.

Sum up the squared residuals around our least-squares fit.

We'll call this **SS(fit)**, for the sum of squares around the least-squares fit.

$$SS(\text{fit}) = (\text{data} - \text{line})^2$$



$$\text{Var}(\text{fit}) = \frac{SS(\text{fit})}{n}$$

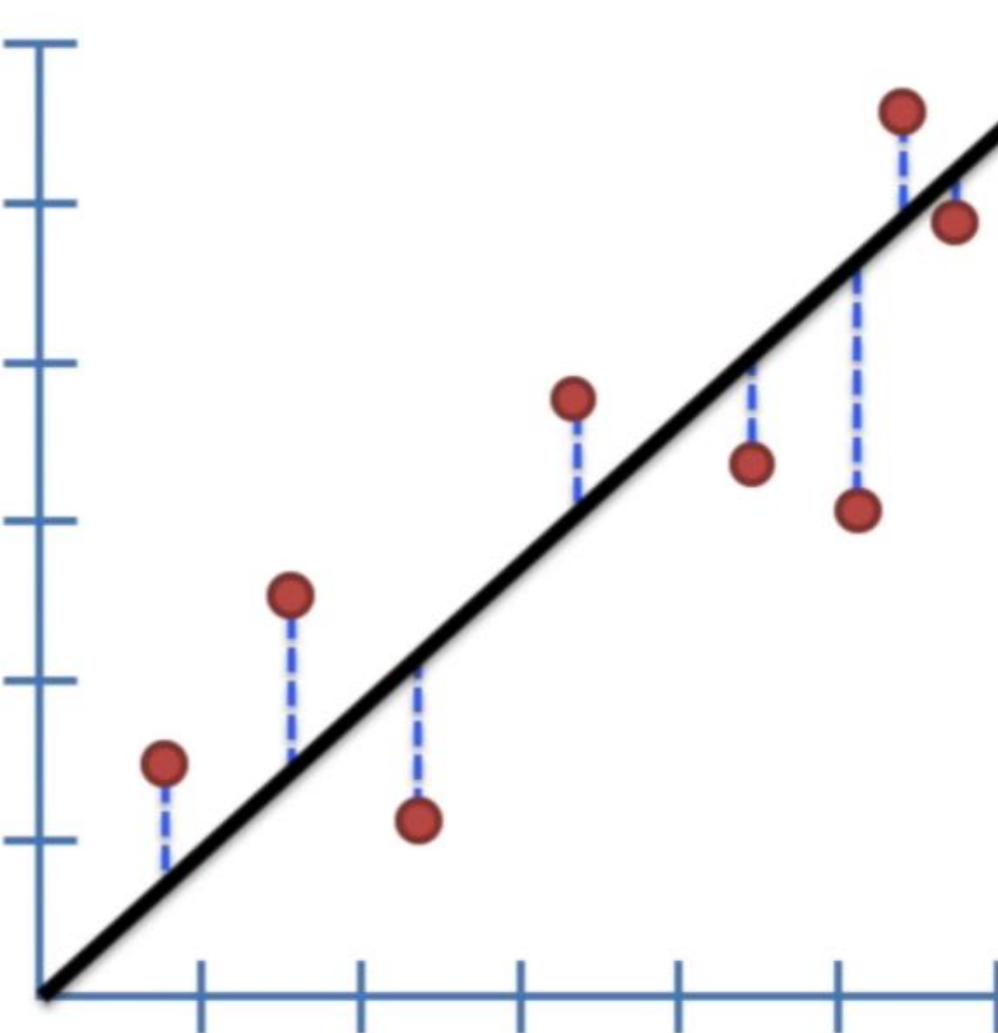
Again, we can think of $\text{Var}(\text{fit})$ as the average $SS(\text{fit})$ for each mouse.

Now go back to the original plot.

Sum up the squared residuals around our least-squares fit.

We'll call this **SS(fit)**, for the sum of squares around the least-squares fit.

$$SS(\text{fit}) = (\text{data} - \text{line})^2$$



$$\text{Var}(\text{fit}) = \frac{SS(\text{fit})}{n}$$

Again, we can think of $\text{Var}(\text{fit})$ as the average $SS(\text{fit})$ for each mouse.

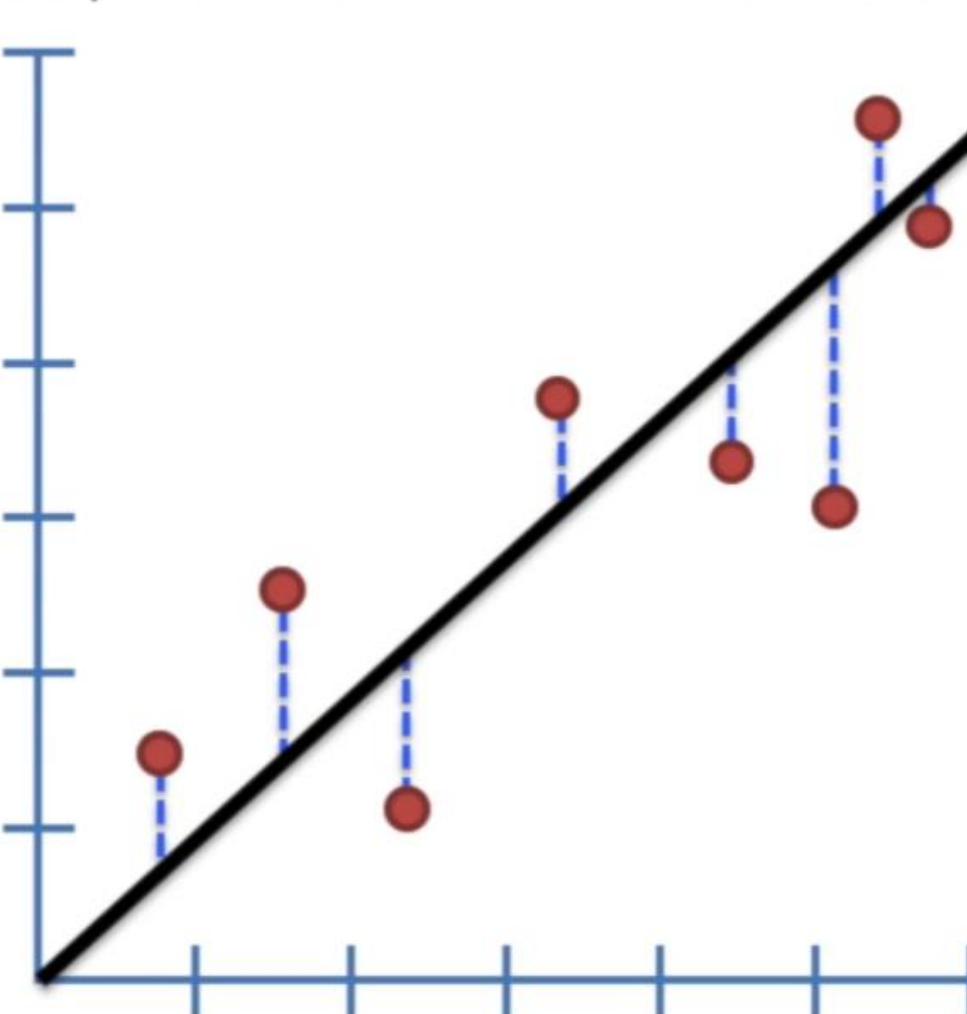
In general: $\text{Variance}(\text{something}) = \frac{\text{Sums of squares}}{\text{The number of those things}}$

Now go back to the original plot.

Sum up the squared residuals around our least-squares fit.

We'll call this **SS(fit)**, for the sum of squares around the least-squares fit.

$$SS(\text{fit}) = (\text{data} - \text{line})^2$$



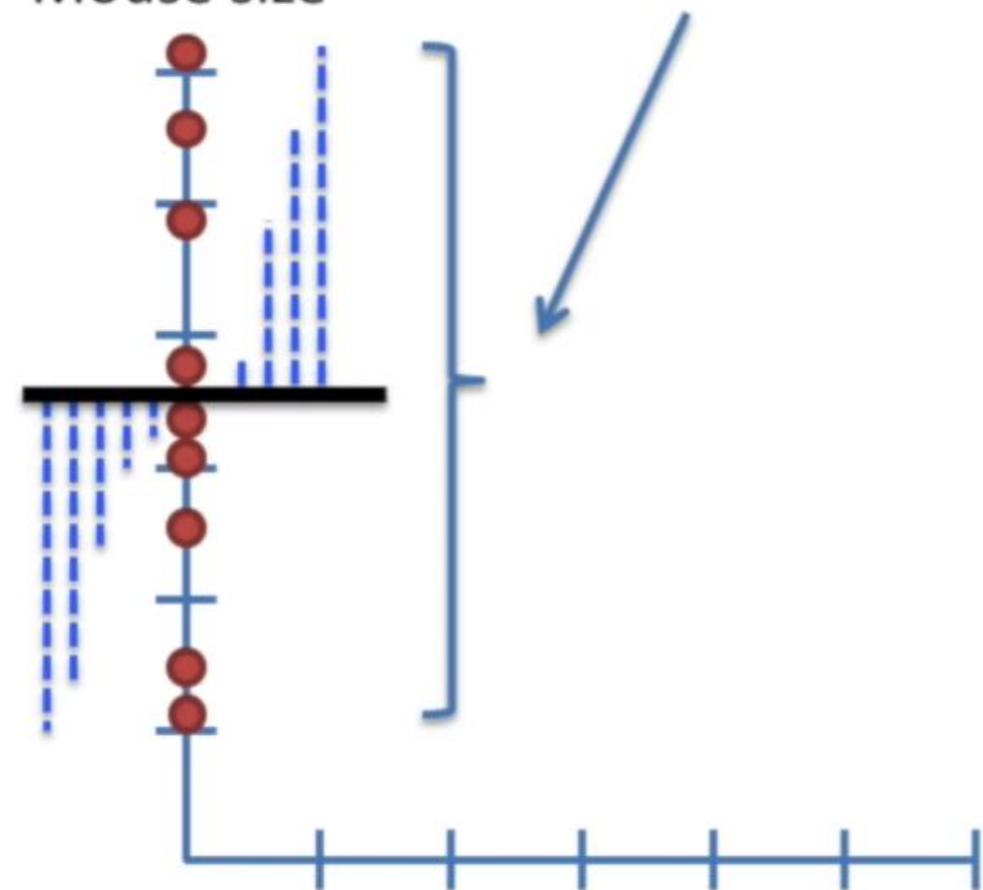
$$\text{Var}(\text{fit}) = \frac{SS(\text{fit})}{n}$$

Again, we can think of $\text{Var}(\text{fit})$ as the average $SS(\text{fit})$ for each mouse.

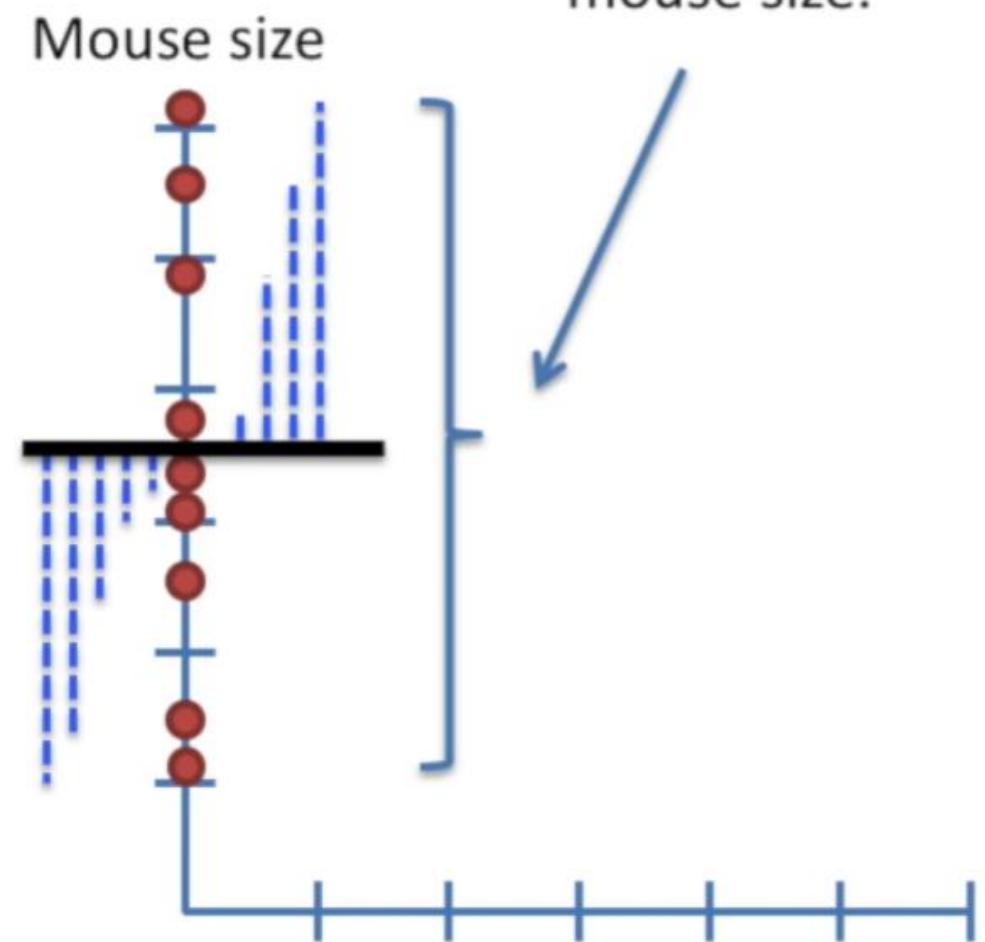
In general: $\text{Variance}(\text{something}) = \frac{\text{Sums of squares}}{\text{The number of those things}} = \text{Average sum of squares.}$

This is the raw variation in mouse size.

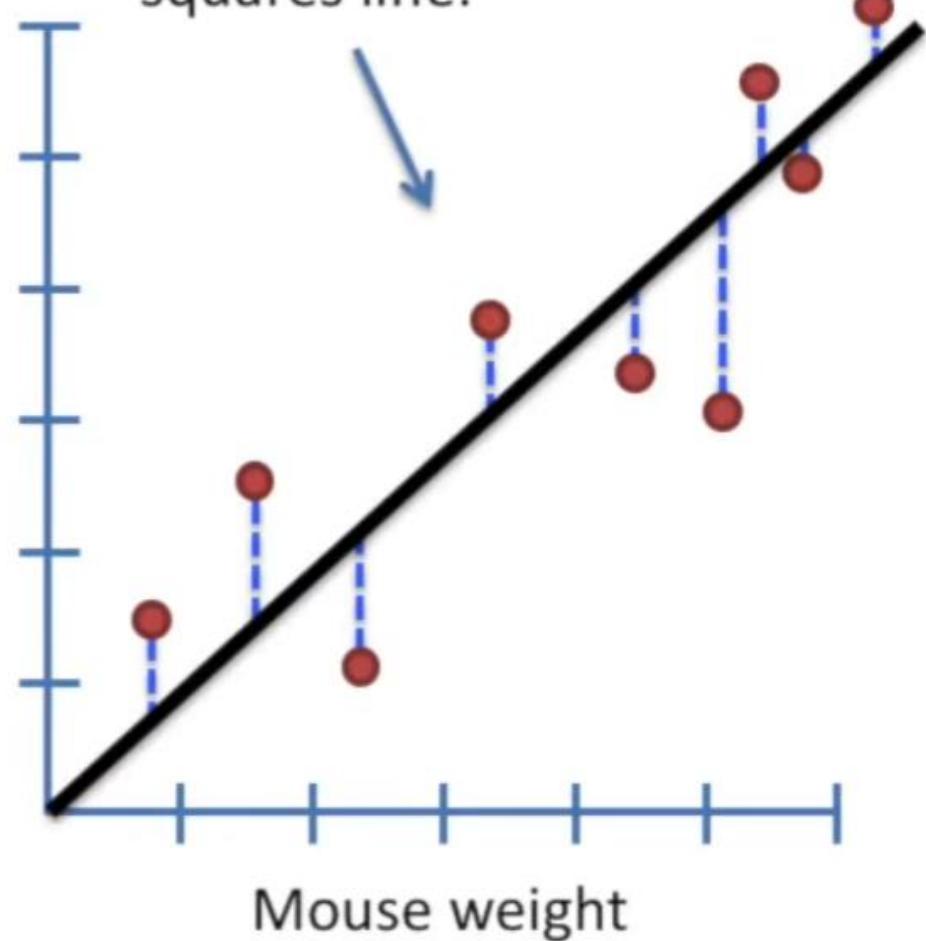
Mouse size



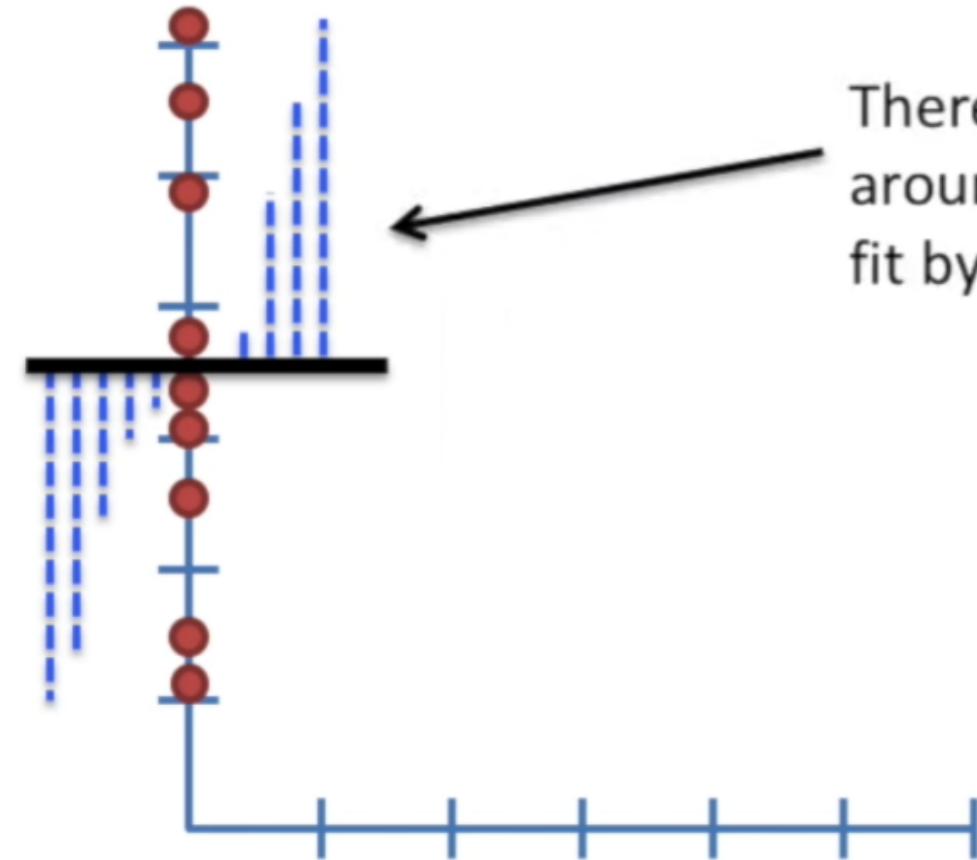
This is the raw variation in mouse size.



This is the variation around the least squares line.

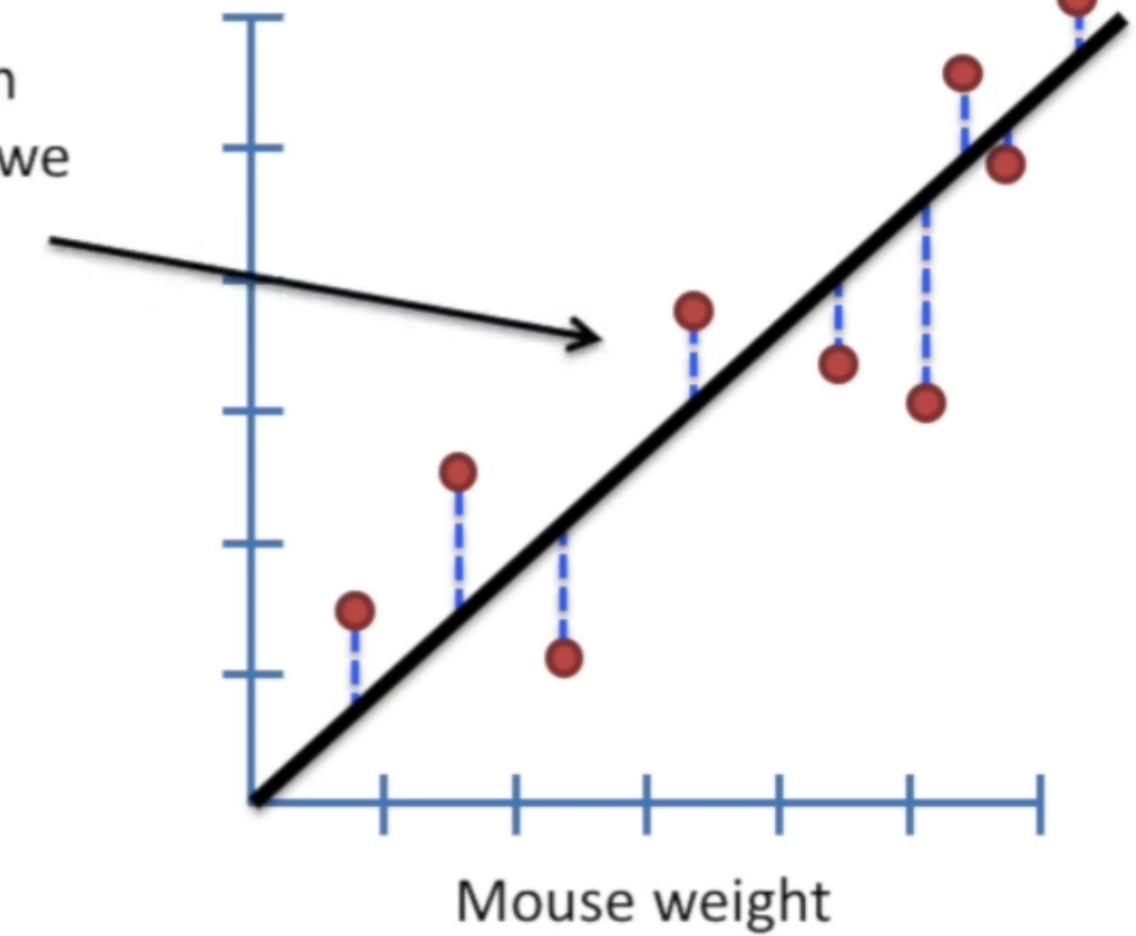


Mouse size

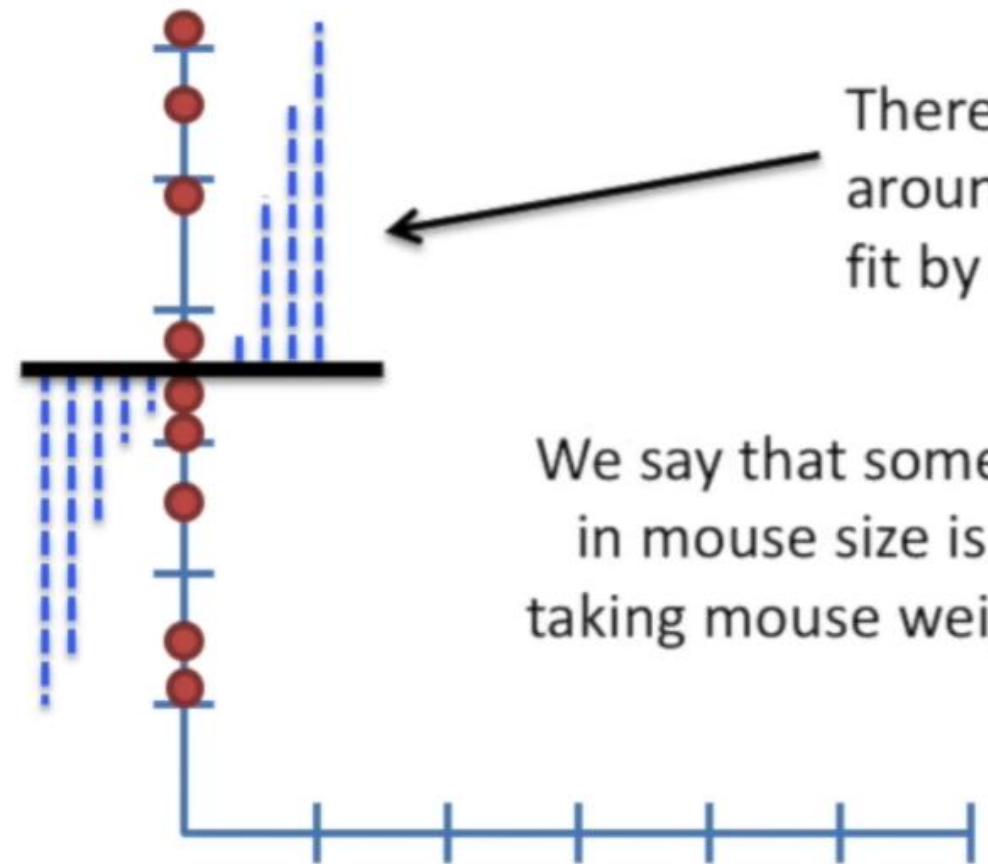


There is less variation around the line that we fit by least-squares.

Mouse size



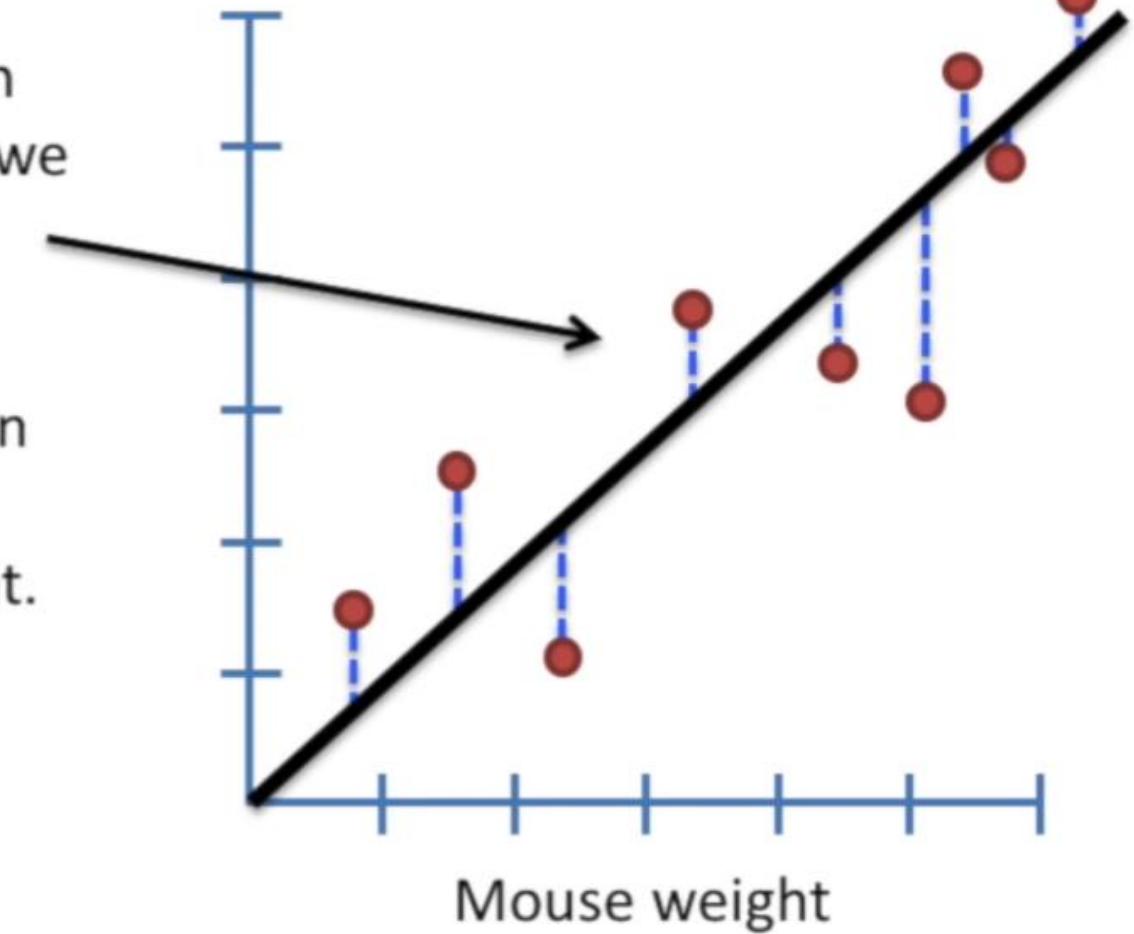
Mouse size

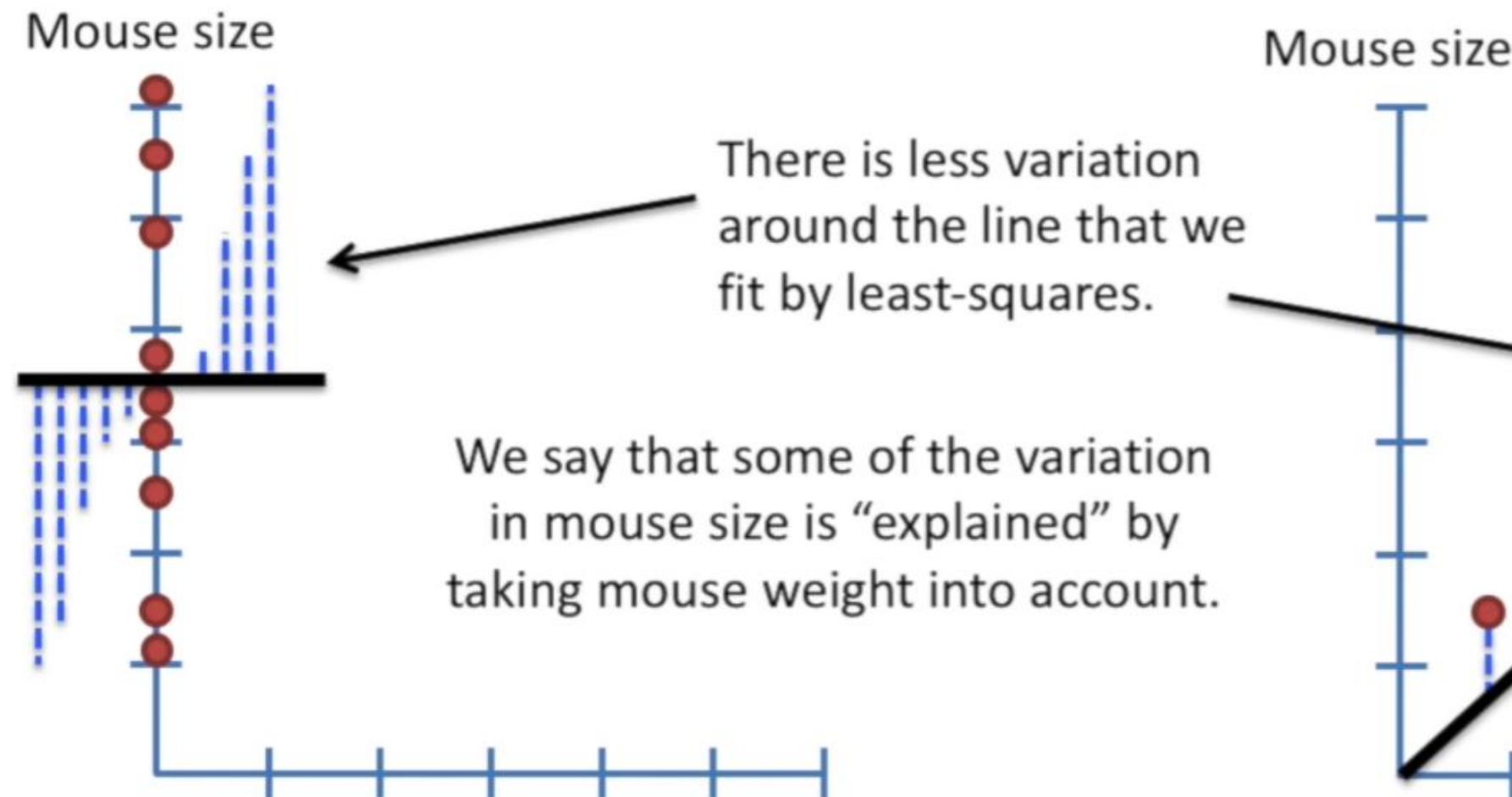


There is less variation around the line that we fit by least-squares.

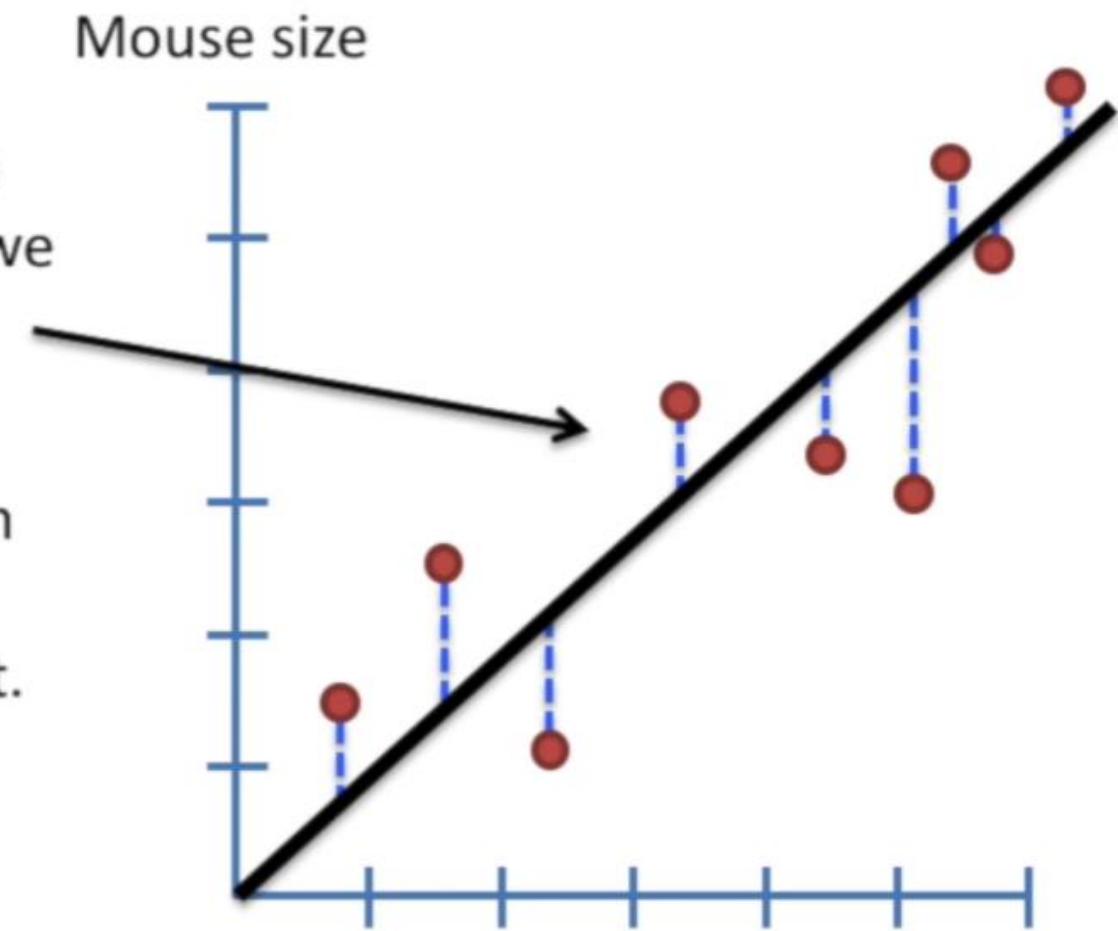
We say that some of the variation in mouse size is “explained” by taking mouse weight into account.

Mouse size

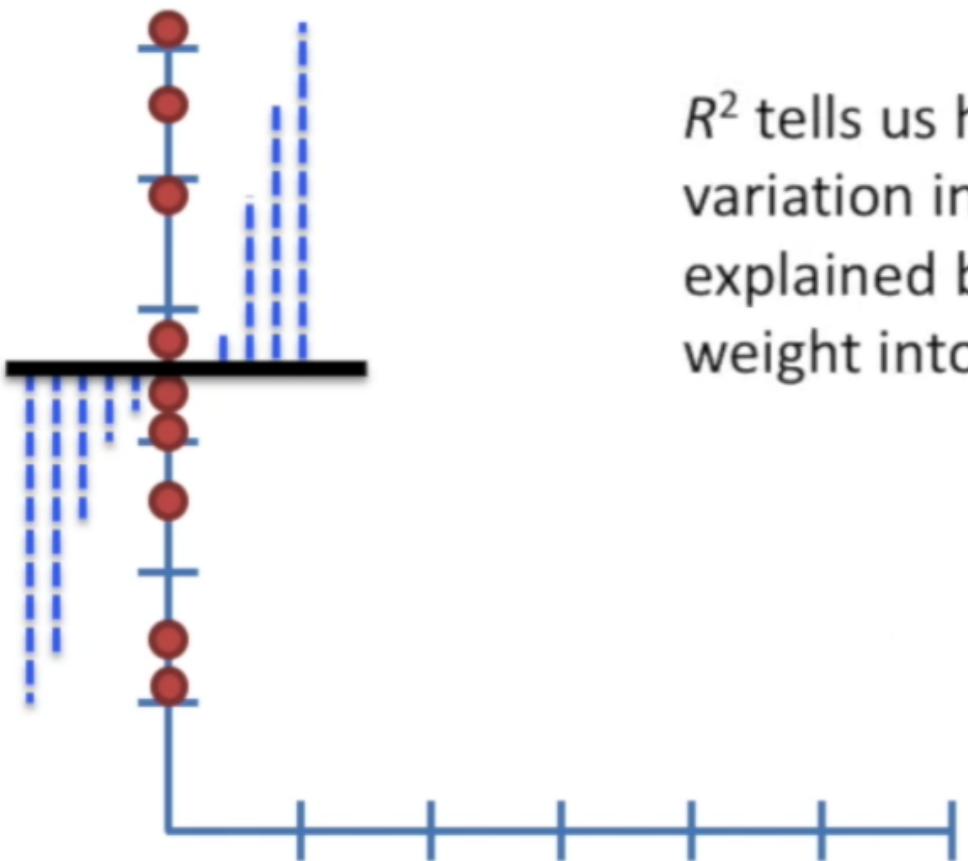




Heavier mice are bigger.
Lighter mice are smaller.

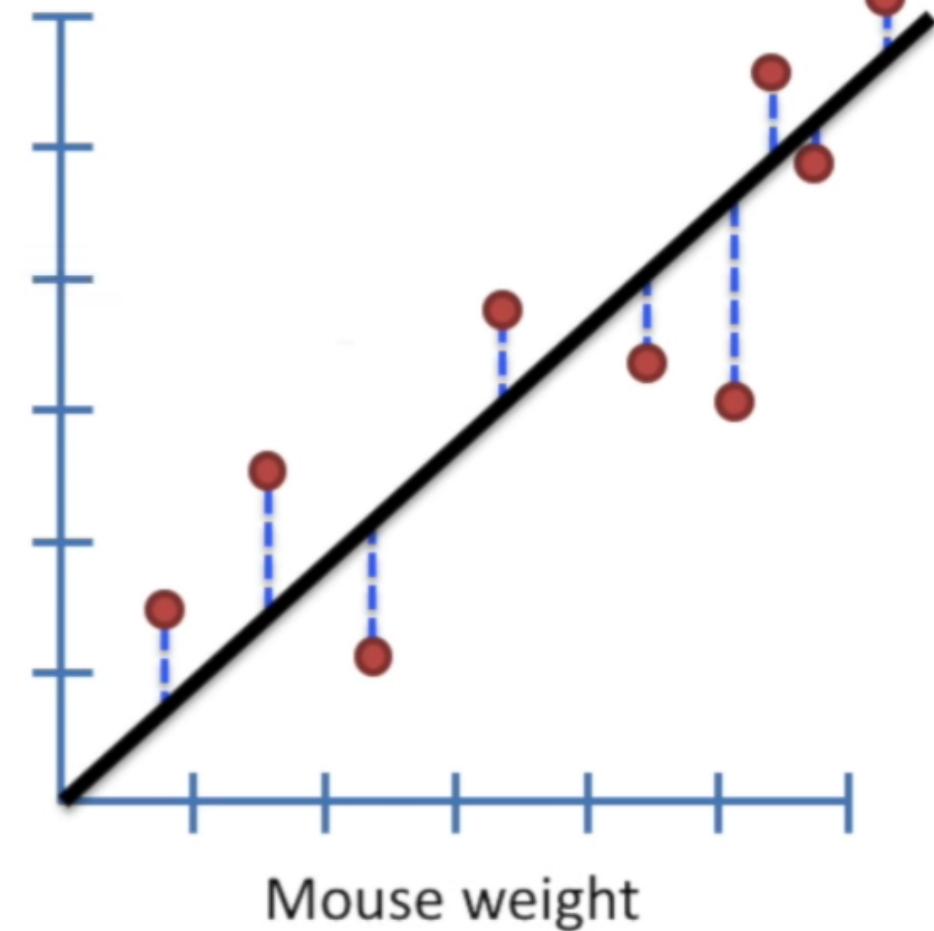


Mouse size

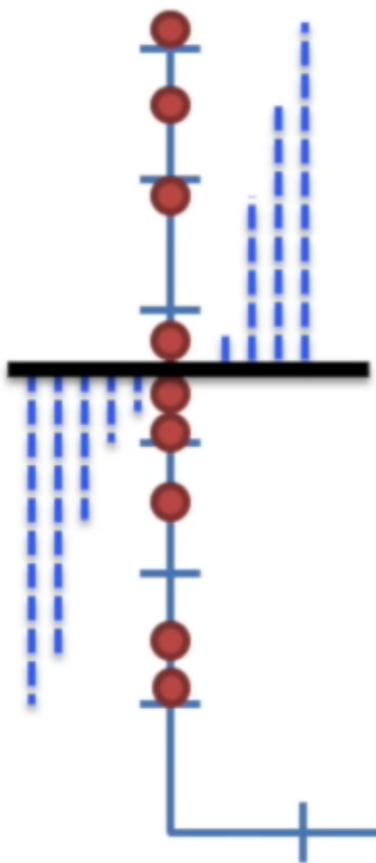


R^2 tells us how much of the variation in mouse size can be explained by taking mouse weight into account.

Mouse size



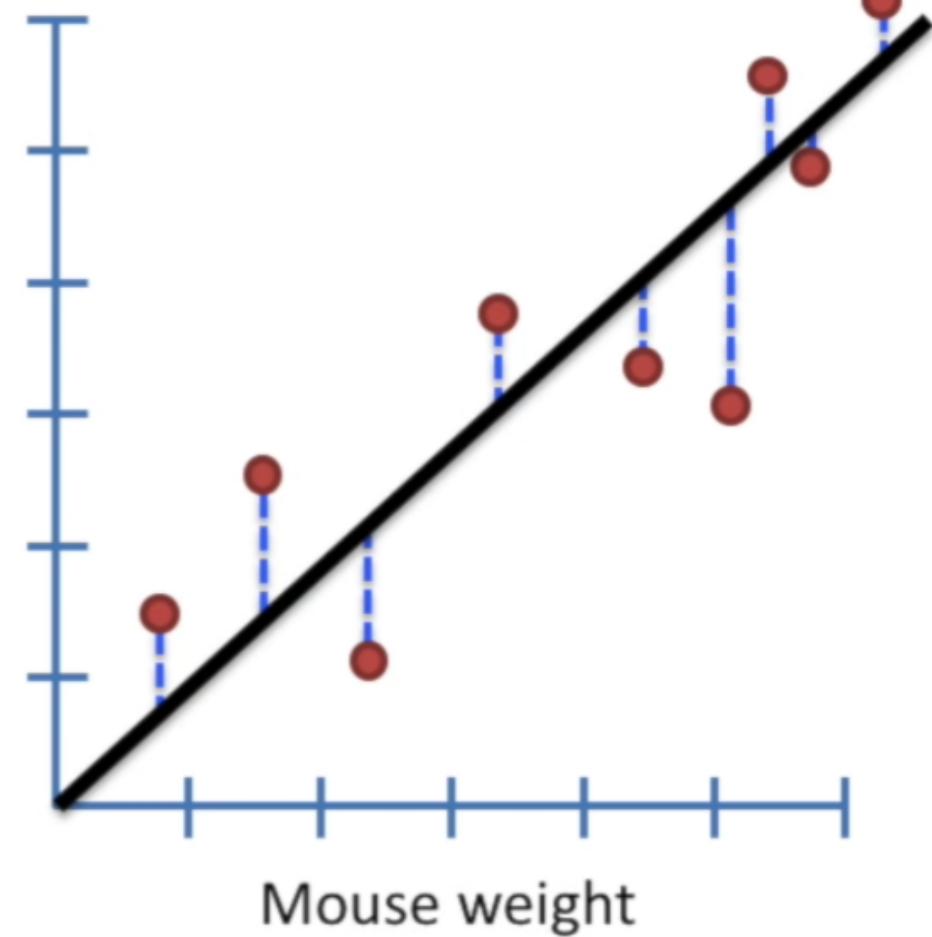
Mouse size



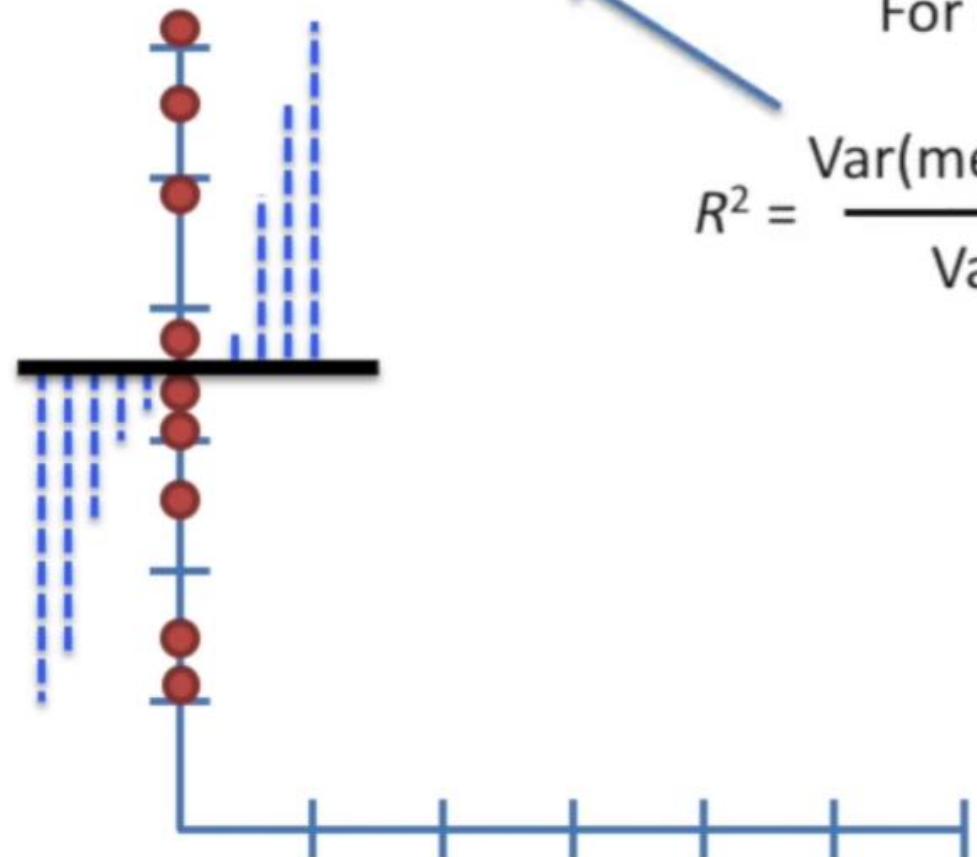
R^2 tells us how much of the variation in mouse size can be explained by taking mouse weight into account.

$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{fit})}{\text{Var}(\text{mean})}$$

Mouse size



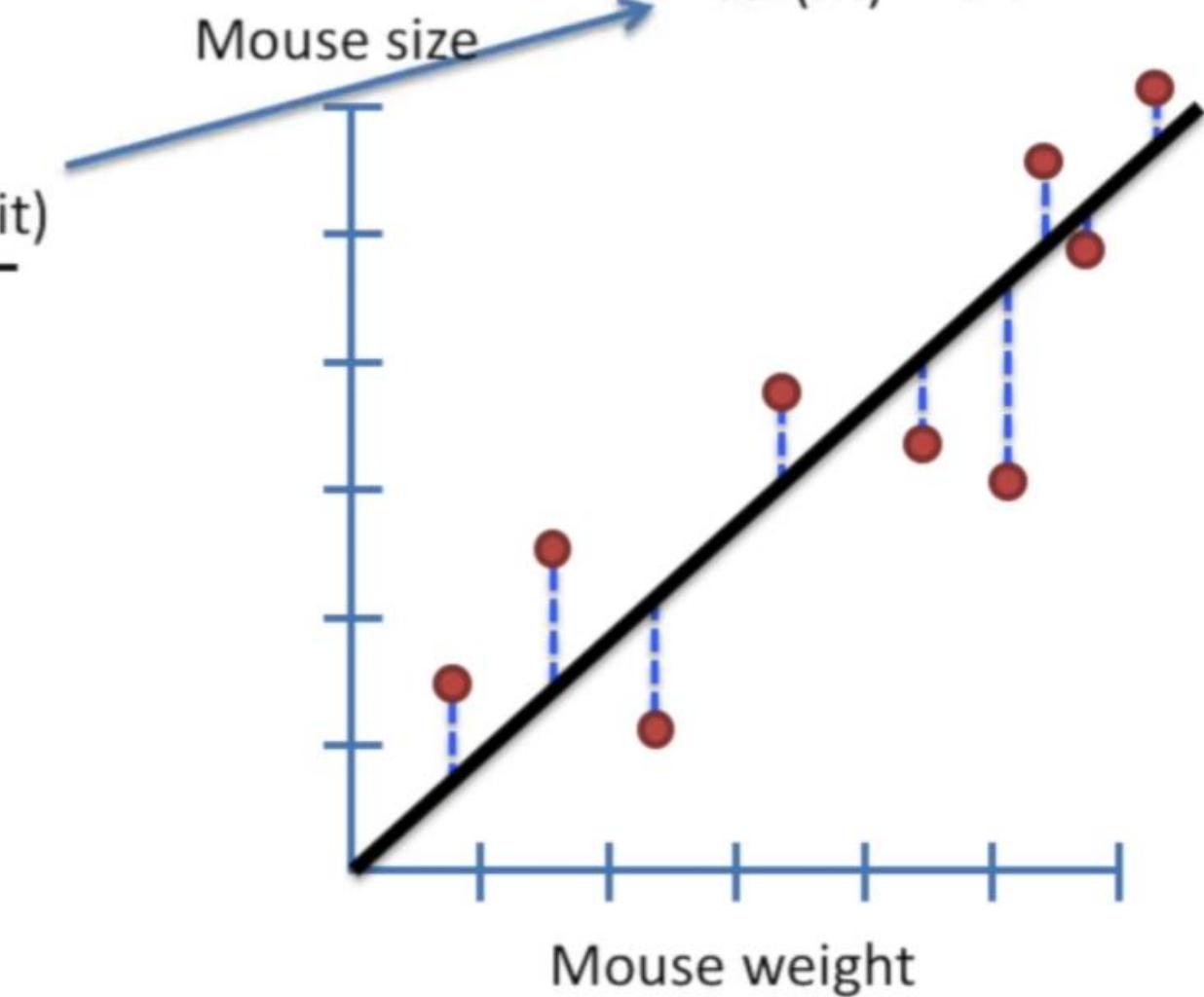
$\text{Var}(\text{mean}) = 11.1$



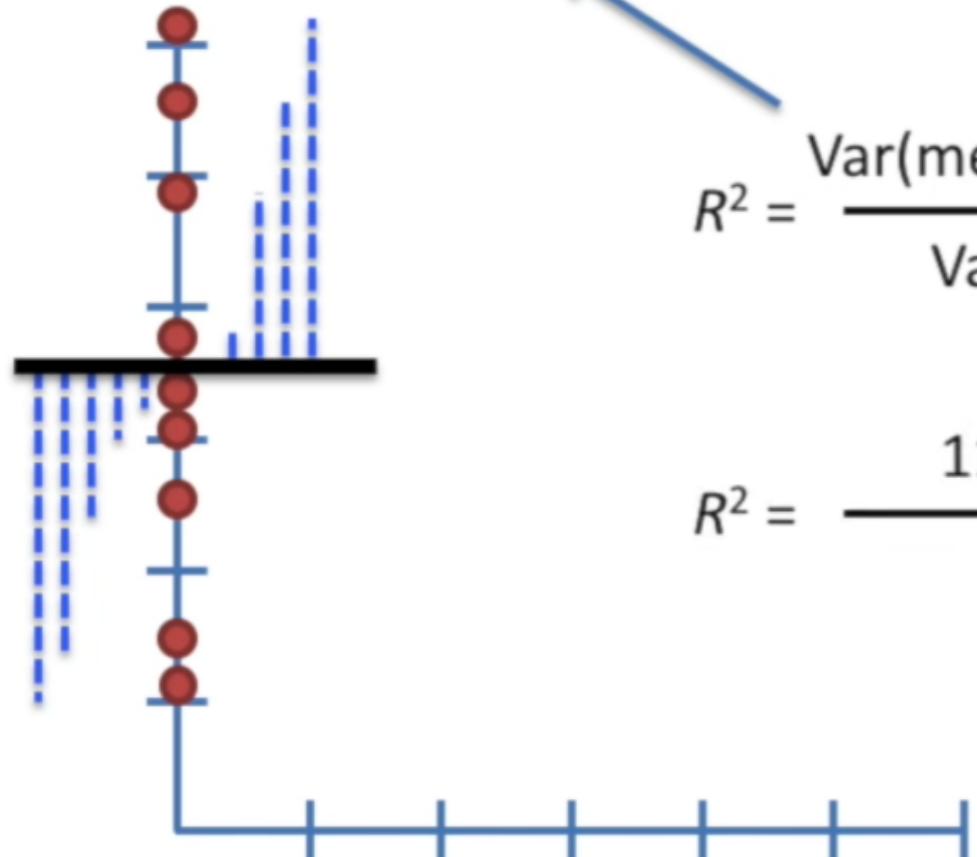
For example...

$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{fit})}{\text{Var}(\text{mean})}$$

$\text{Var}(\text{fit}) = 4.4$

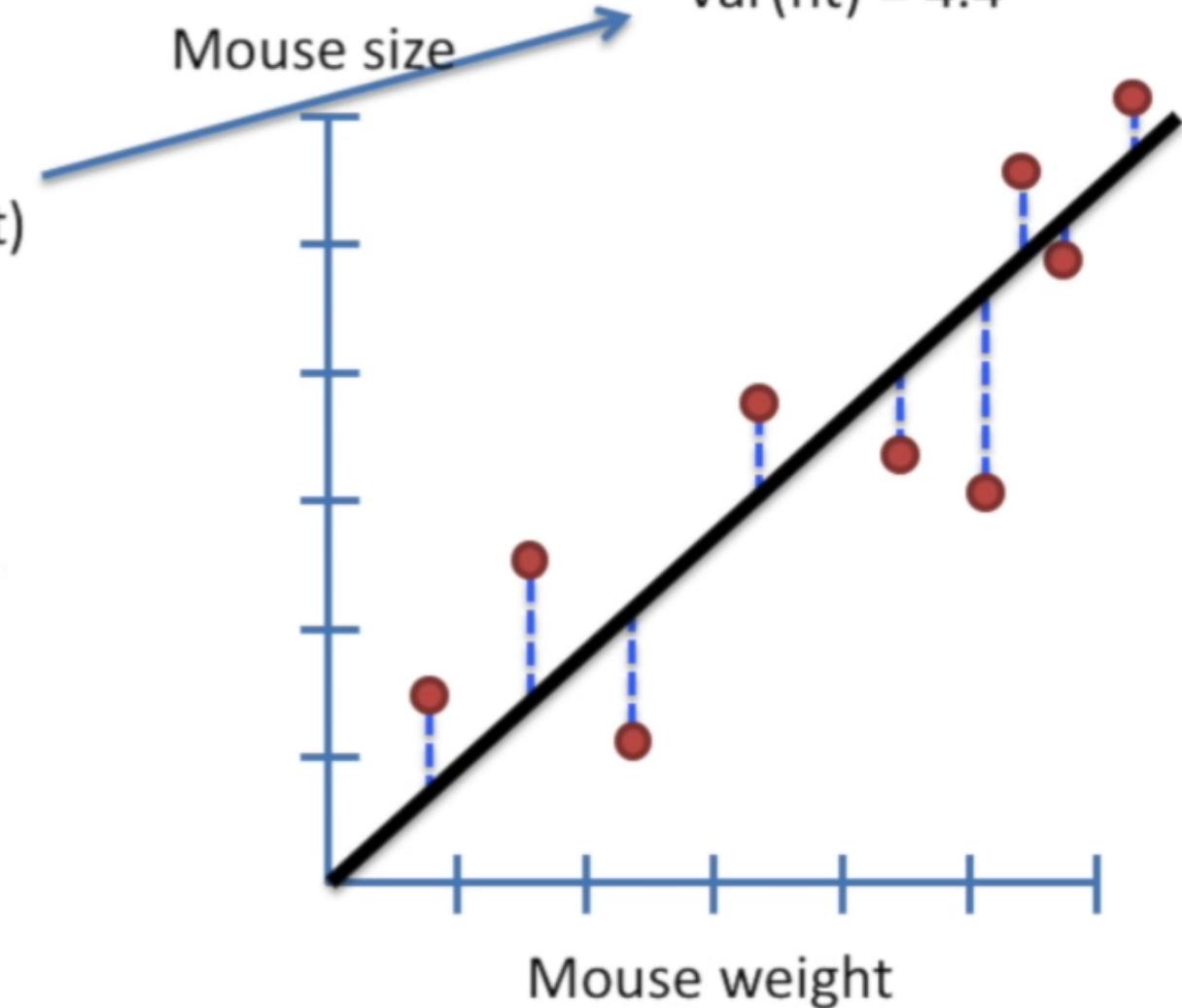


$\text{Var}(\text{mean}) = 11.1$

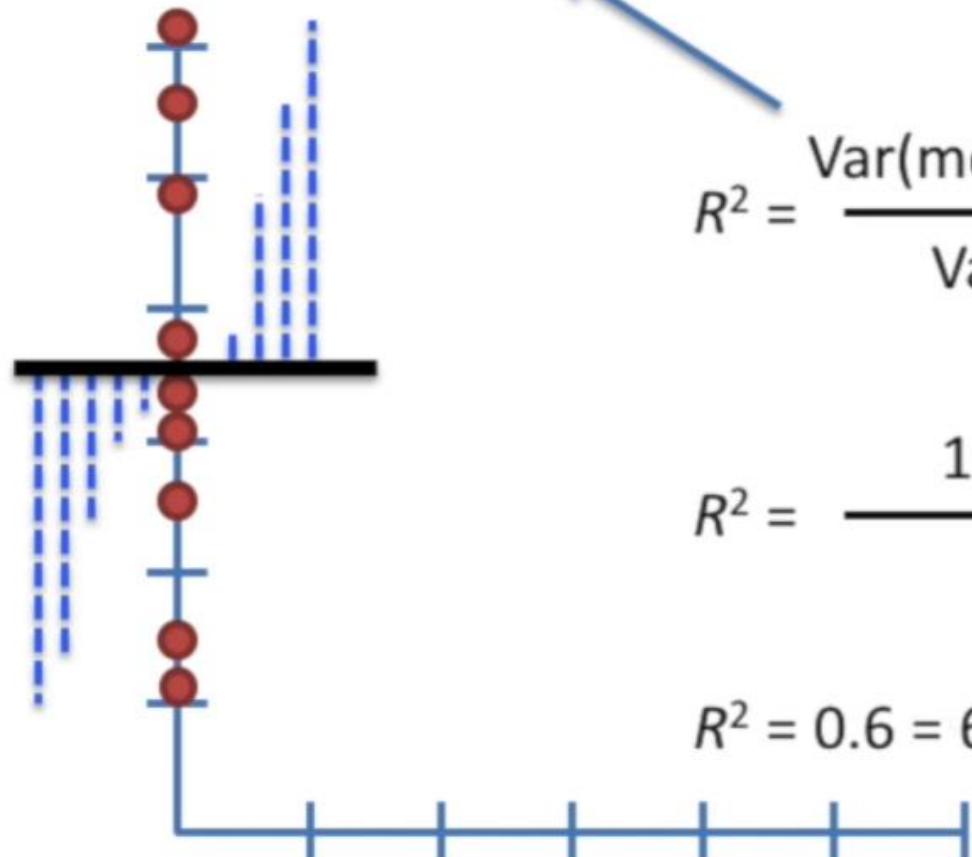


$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{fit})}{\text{Var}(\text{mean})}$$
$$R^2 = \frac{11.1 - 4.4}{11.1}$$

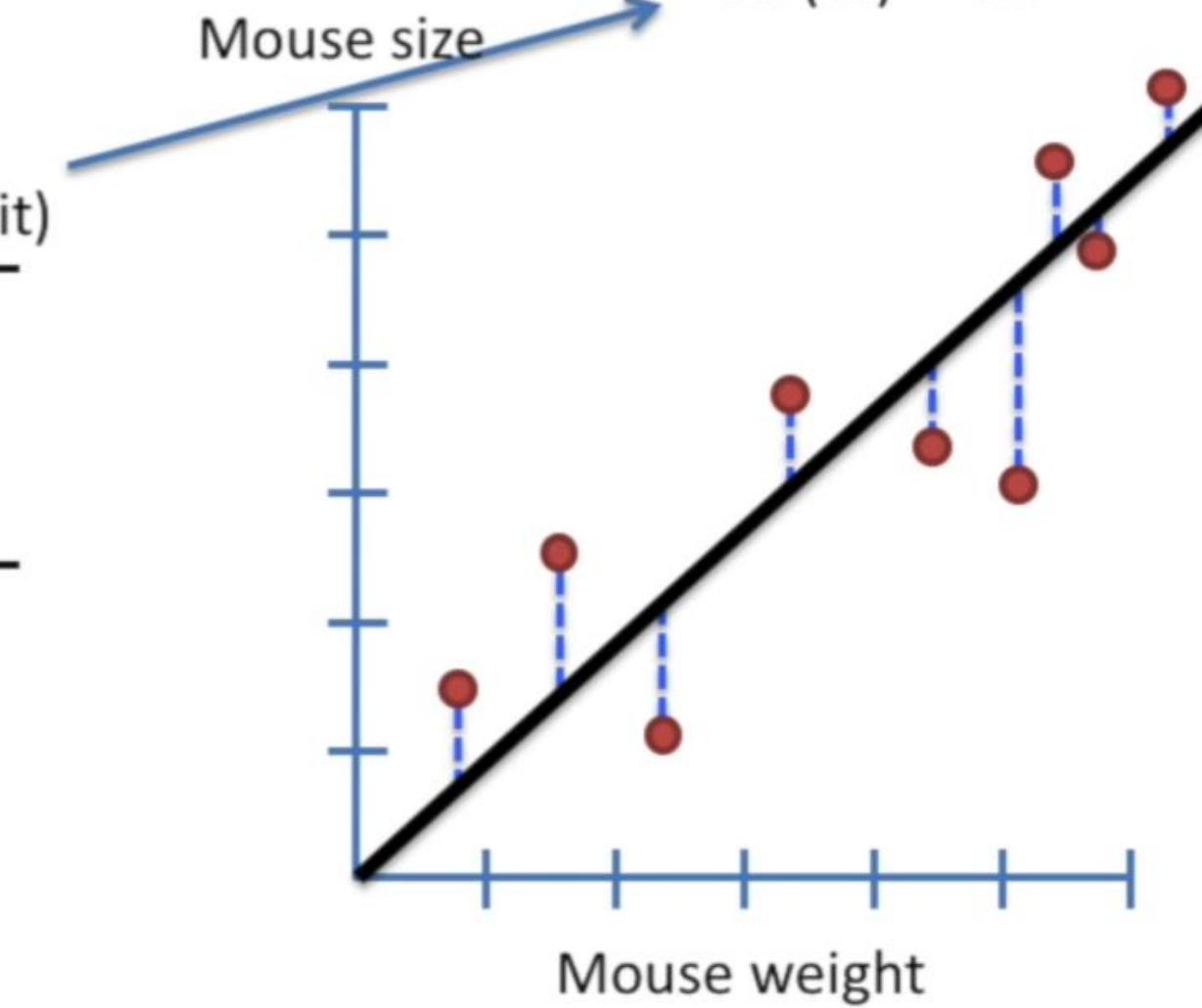
$\text{Var}(\text{fit}) = 4.4$



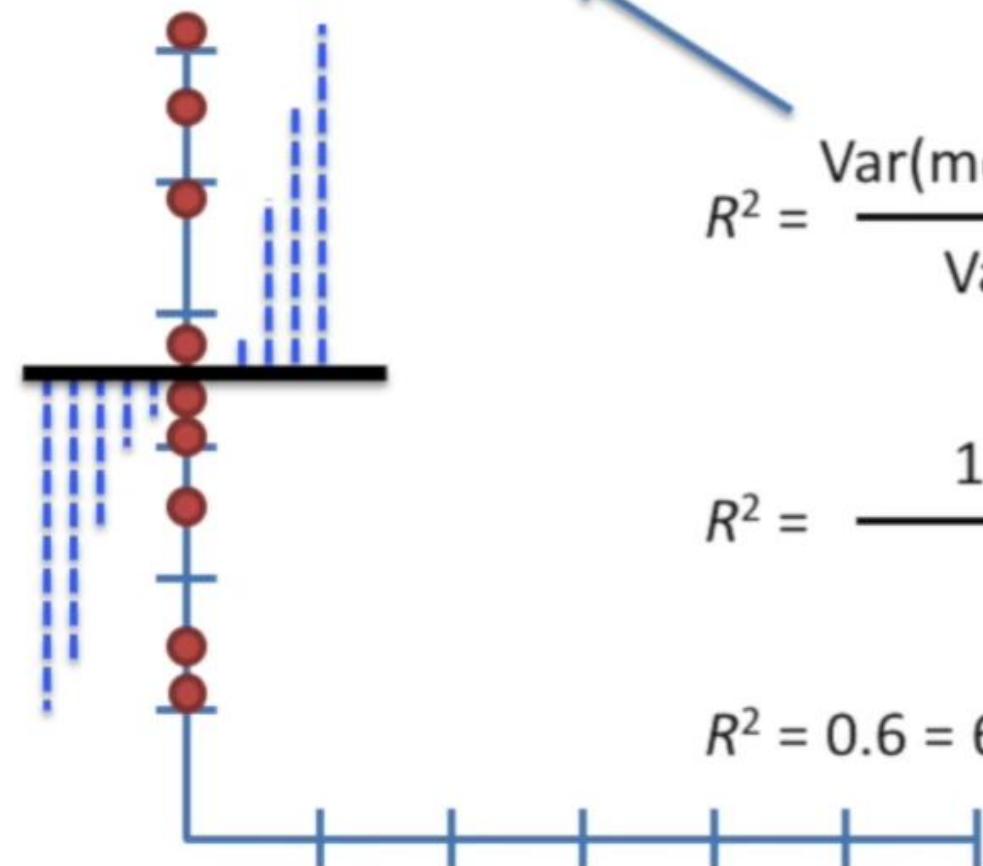
$\text{Var}(\text{mean}) = 11.1$



$\text{Var}(\text{fit}) = 4.4$

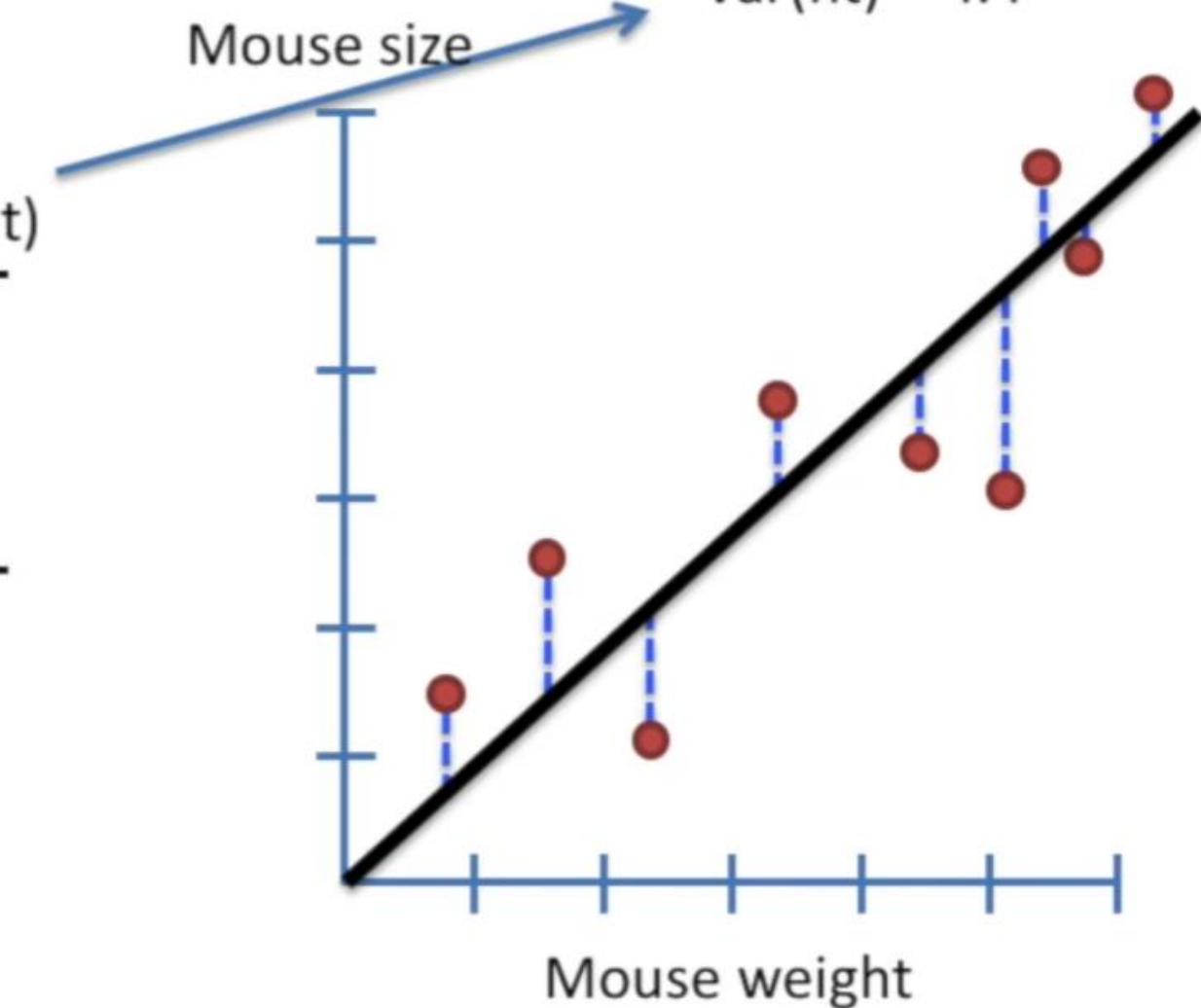


$\text{Var}(\text{mean}) = 11.1$

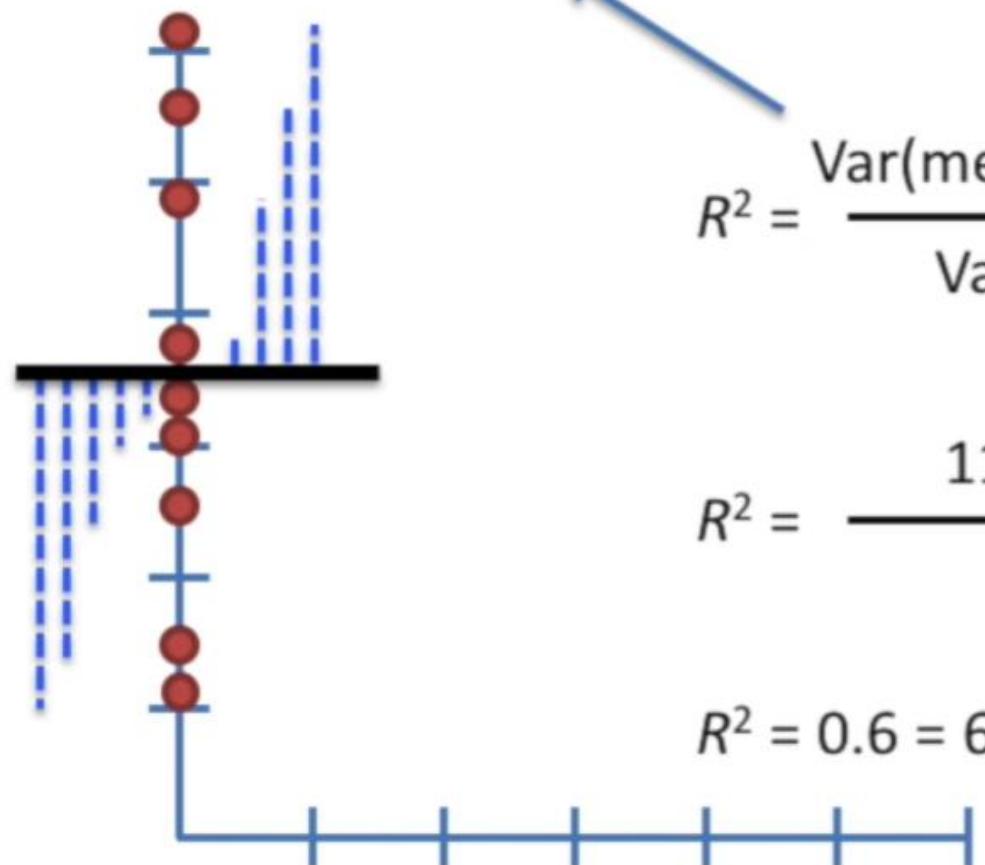


There is a 60% reduction in variance when we take the mouse weight into account.

$\text{Var}(\text{fit}) = 4.4$



$\text{Var}(\text{mean}) = 11.1$



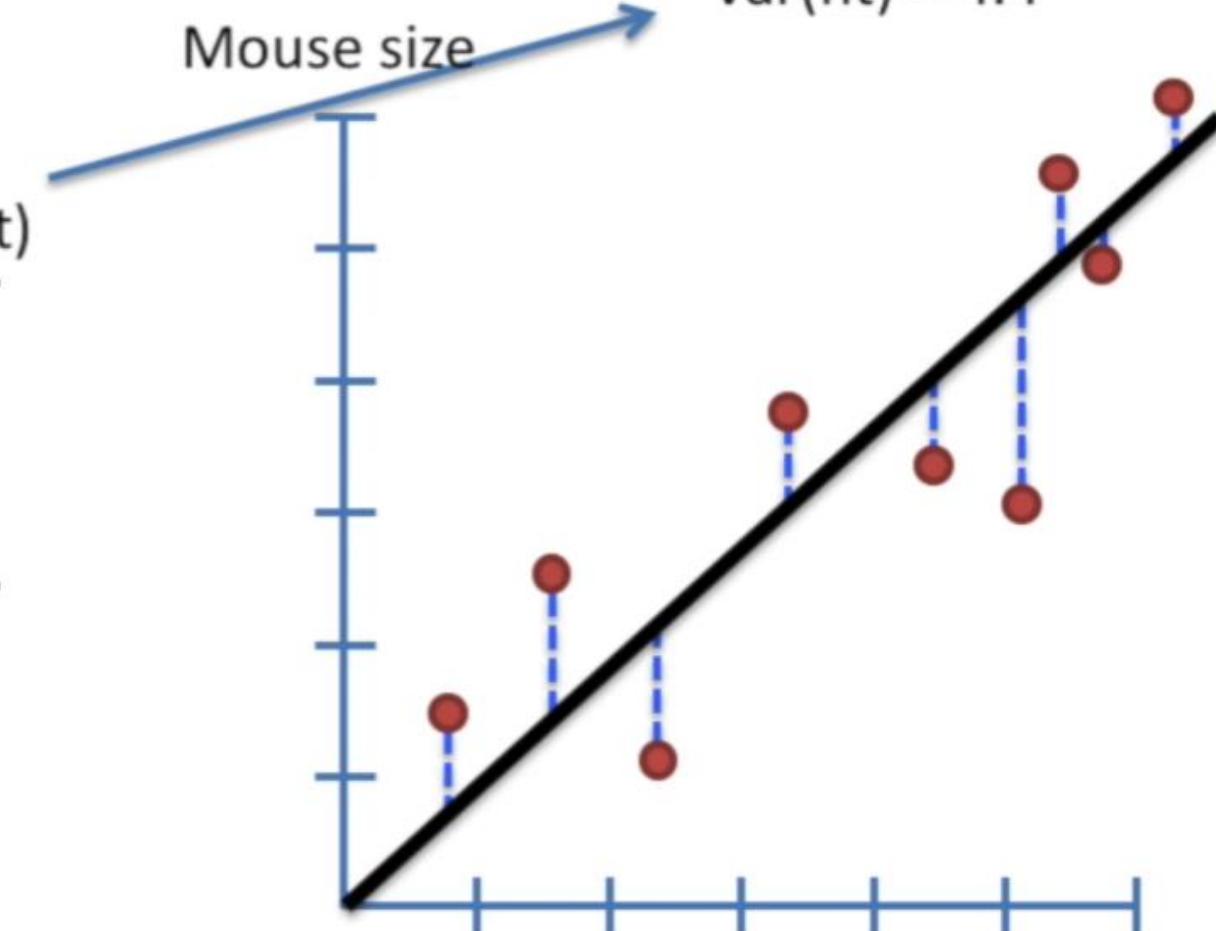
$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{fit})}{\text{Var}(\text{mean})}$$

$$R^2 = \frac{11.1 - 4.4}{11.1}$$

$$R^2 = 0.6 = 60\%$$

There is a 60% reduction in variance when we take the mouse weight into account.

$\text{Var}(\text{fit}) = 4.4$



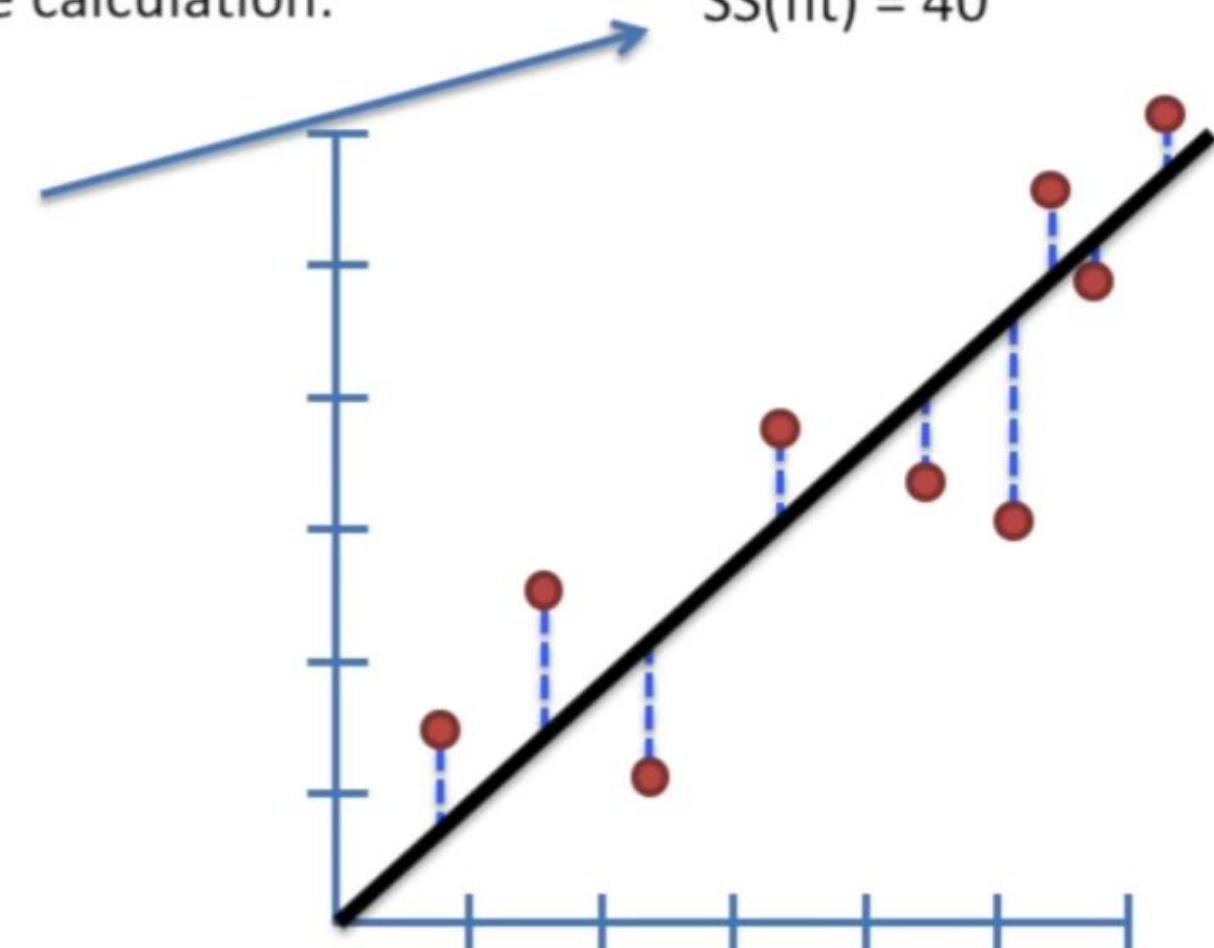
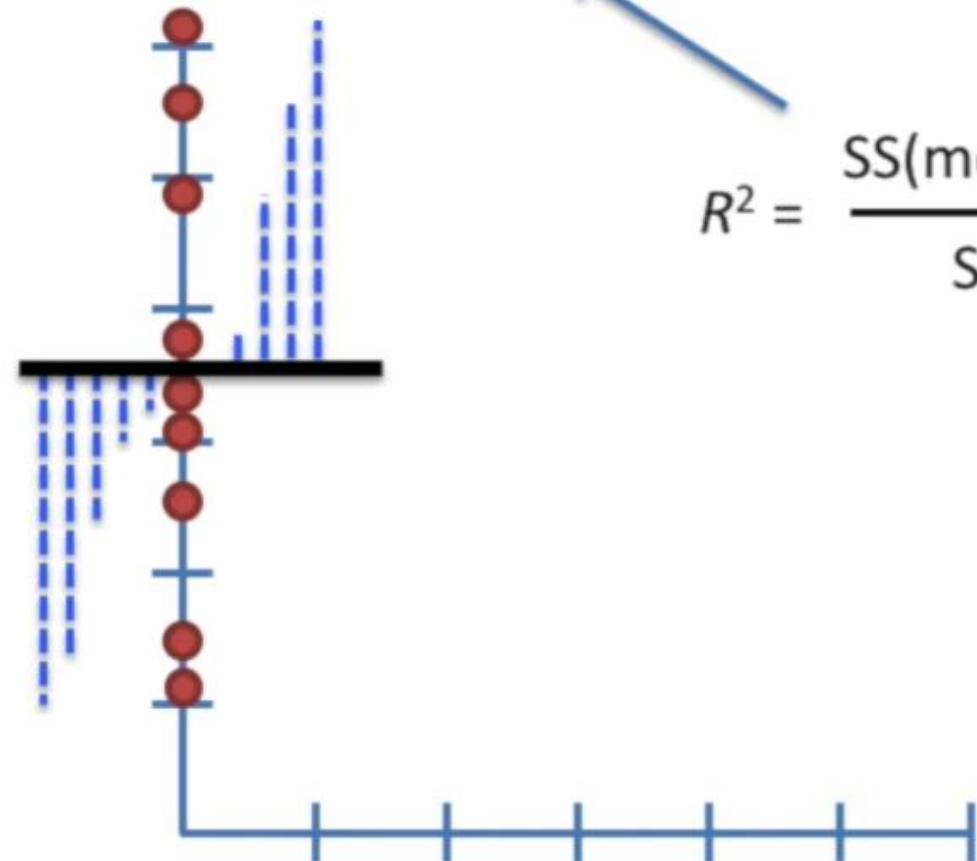
Alternatively, we can say that mouse weight “explains” 60% of the variation in mouse size.

We can also use the sums of squares to make the same calculation.

$$SS(\text{mean}) = 100$$

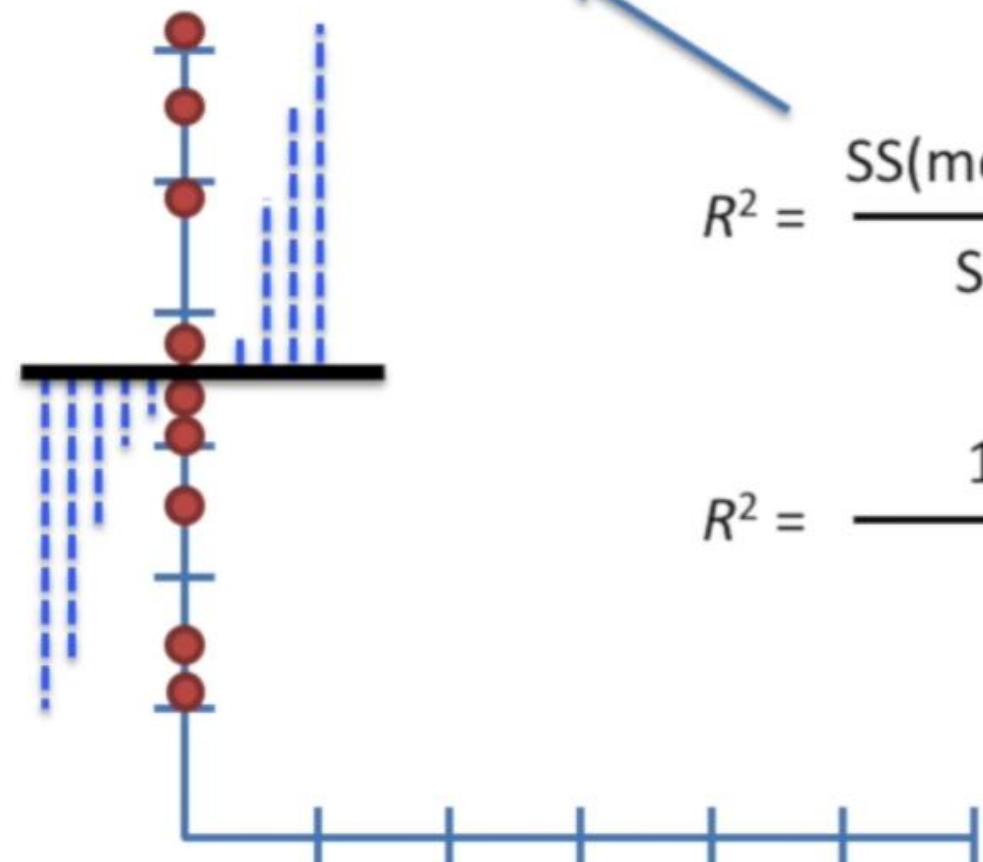
$$R^2 = \frac{SS(\text{mean}) - SS(\text{fit})}{SS(\text{mean})}$$

$$SS(\text{fit}) = 40$$



We can also use the sums of squares to make the same calculation.

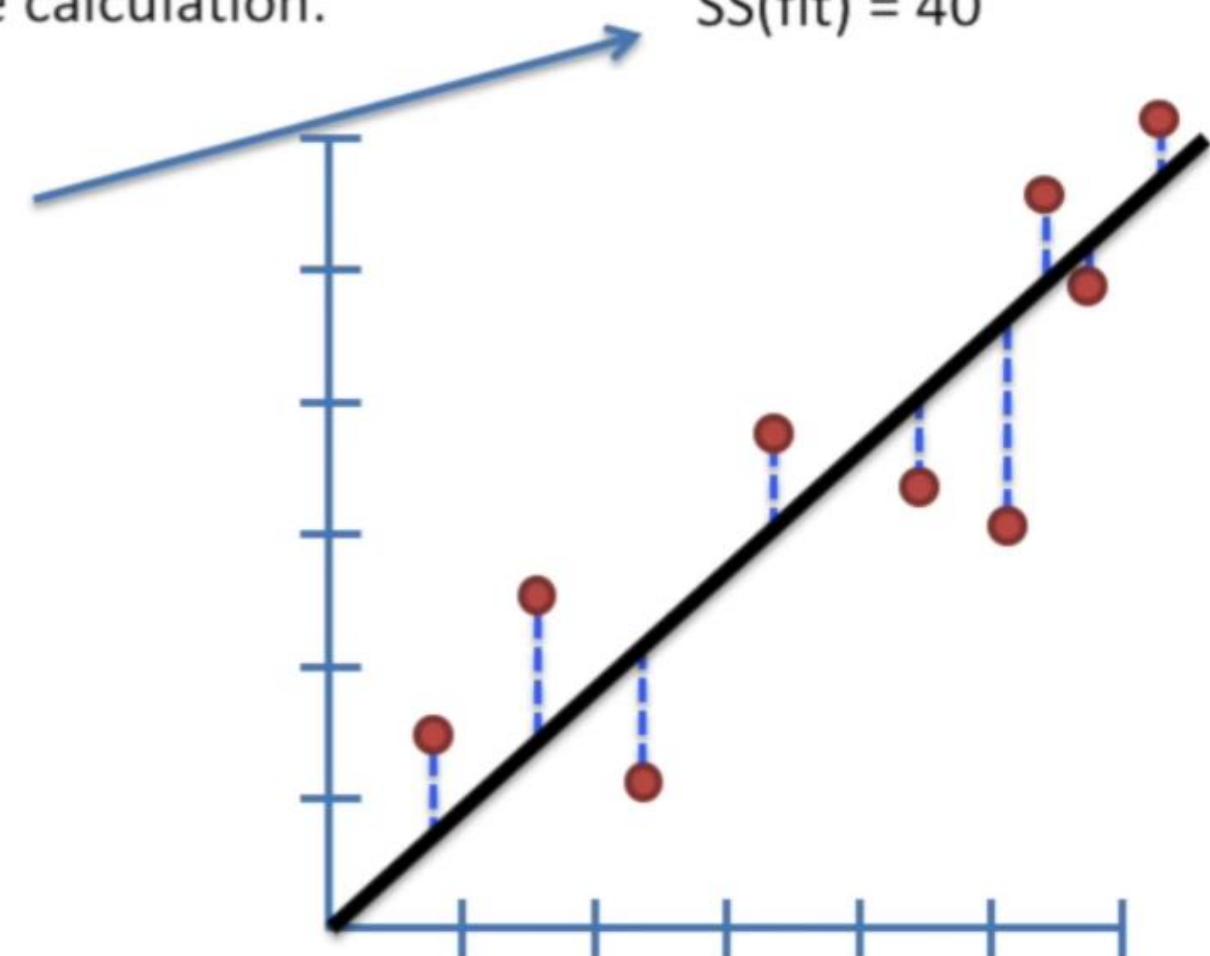
$$SS(\text{mean}) = 100$$



$$R^2 = \frac{SS(\text{mean}) - SS(\text{fit})}{SS(\text{mean})}$$

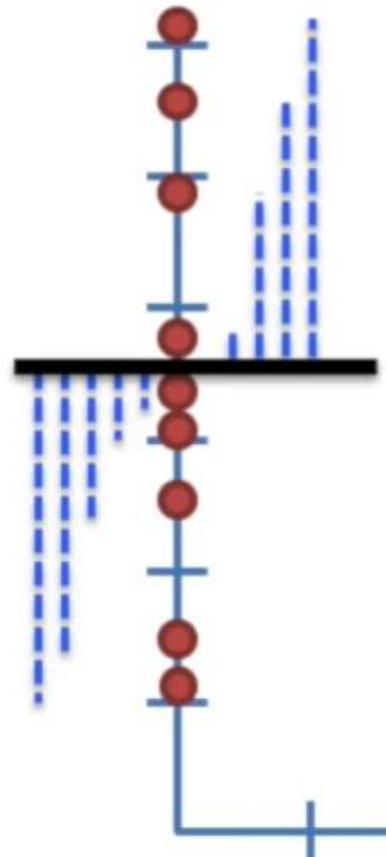
$$R^2 = \frac{100 - 40}{100}$$

$$SS(\text{fit}) = 40$$



We can also use the sums of squares to make the same calculation.

$$SS(\text{mean}) = 100$$

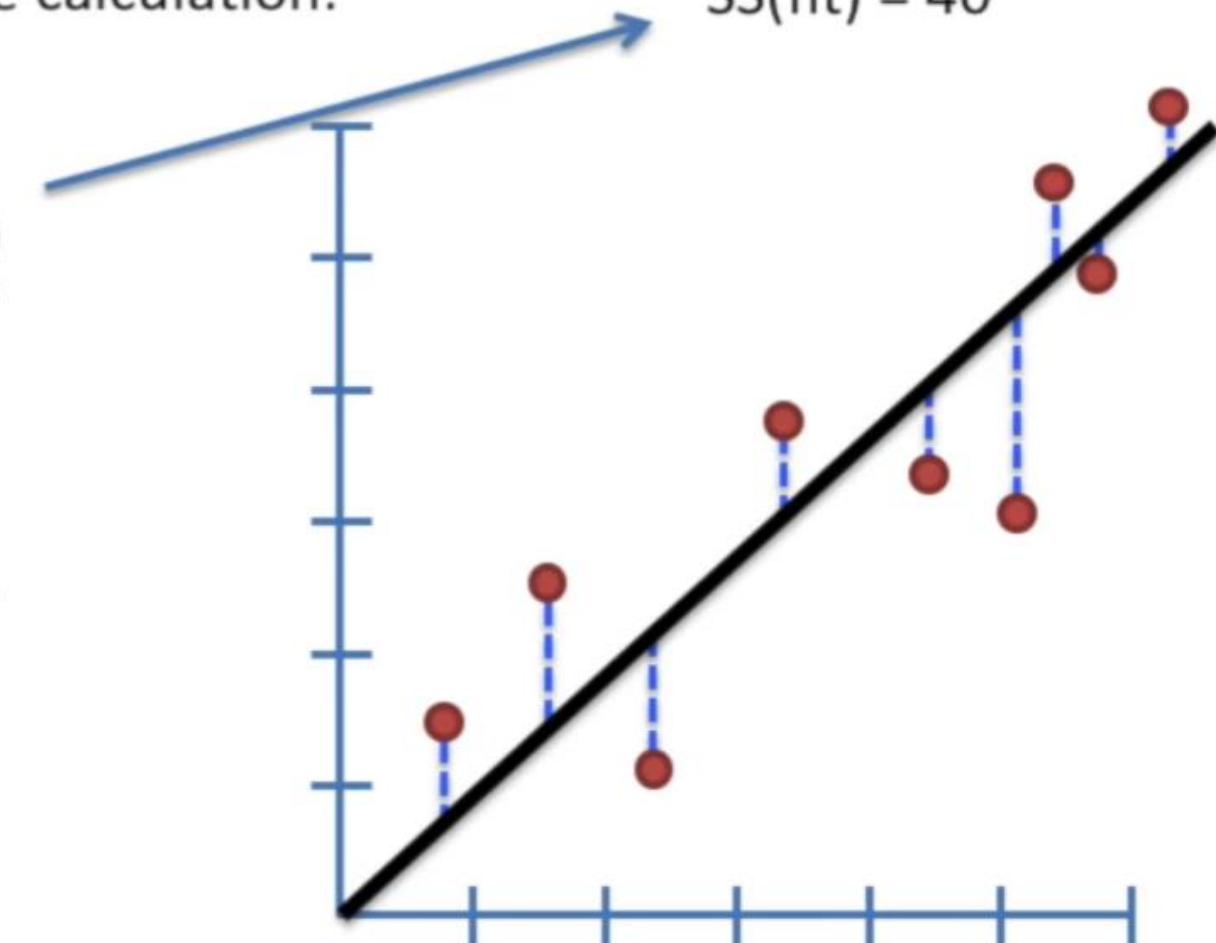


$$R^2 = \frac{SS(\text{mean}) - SS(\text{fit})}{SS(\text{mean})}$$

$$R^2 = \frac{100 - 40}{100}$$

$$R^2 = 0.6 = 60\%$$

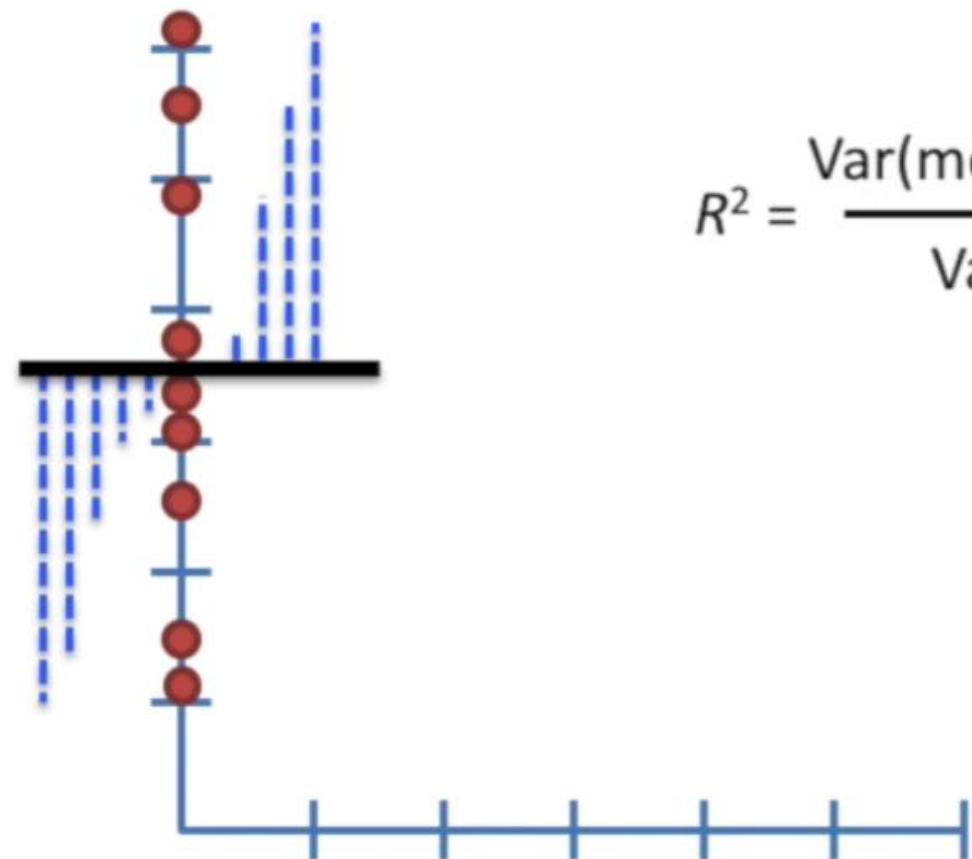
$$SS(\text{fit}) = 40$$



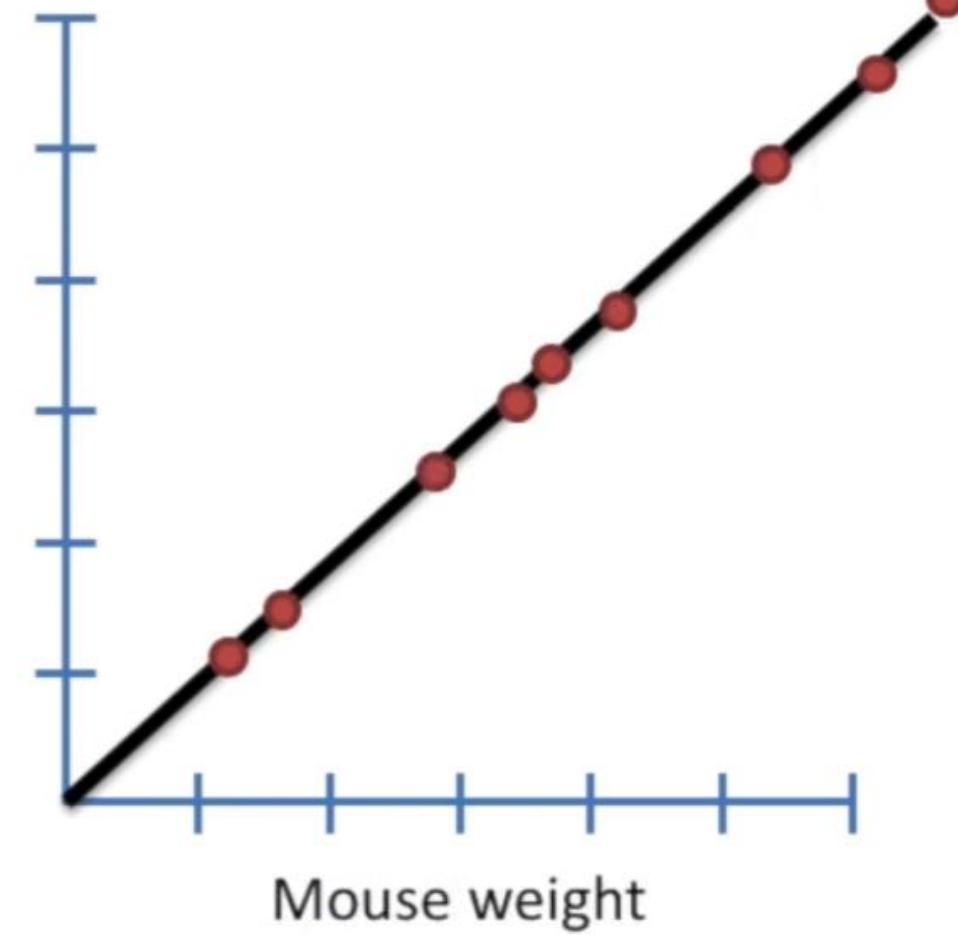
60% of the sums of squares of the mouse size can be explained by mouse weight..

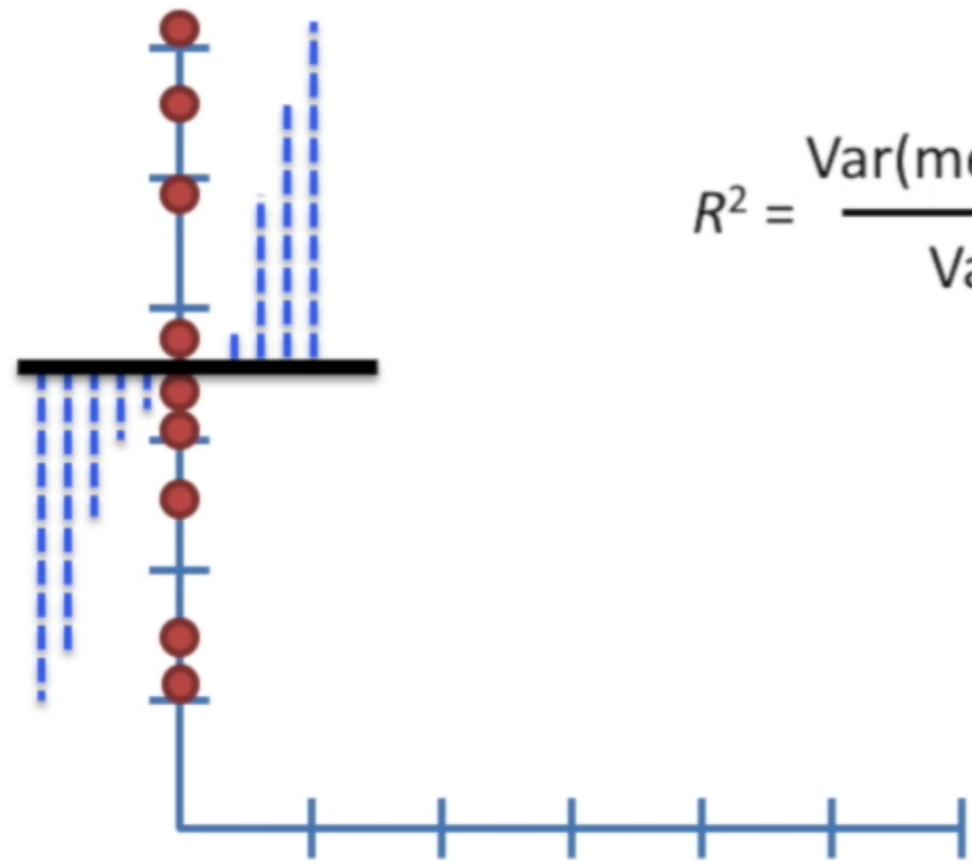
Another Example...

$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{fit})}{\text{Var}(\text{mean})}$$

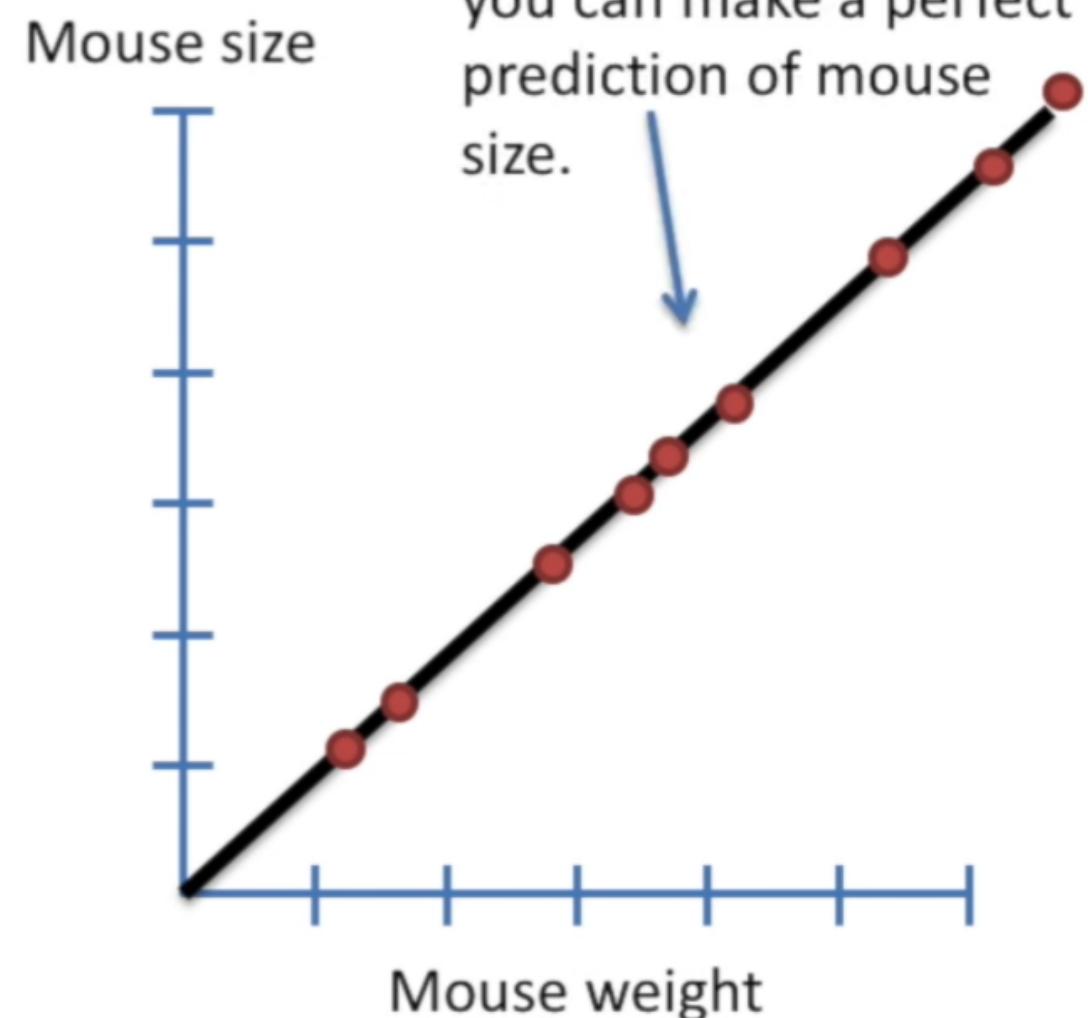


Mouse size

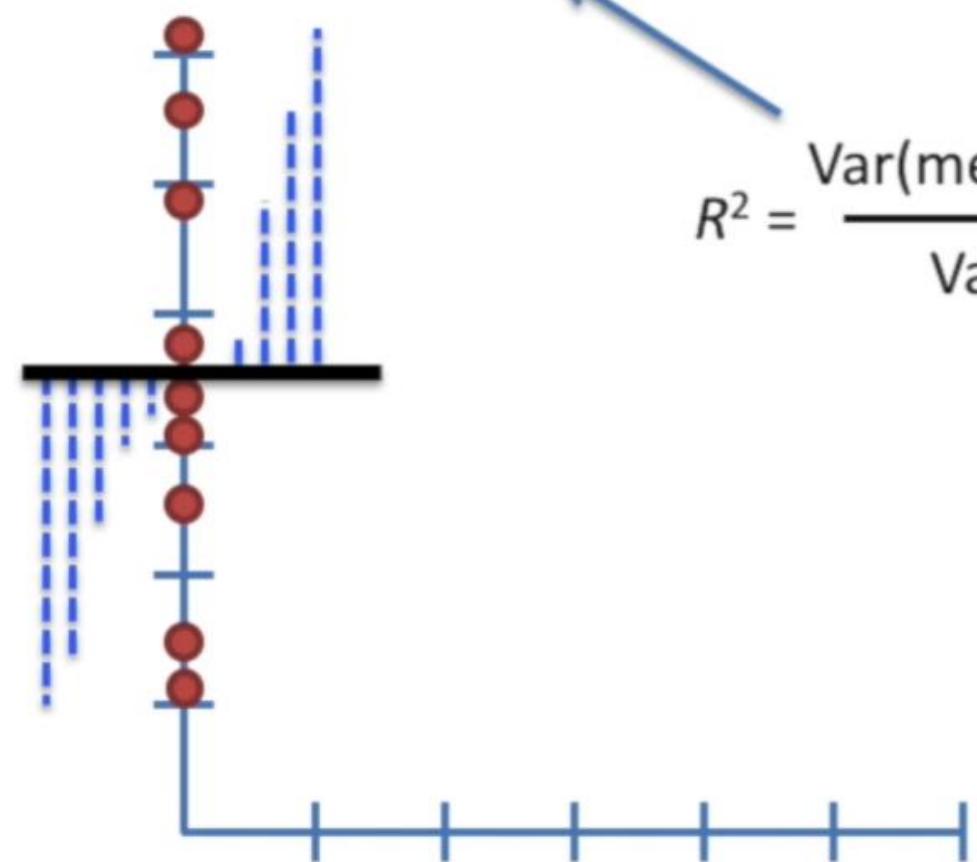




$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{fit})}{\text{Var}(\text{mean})}$$

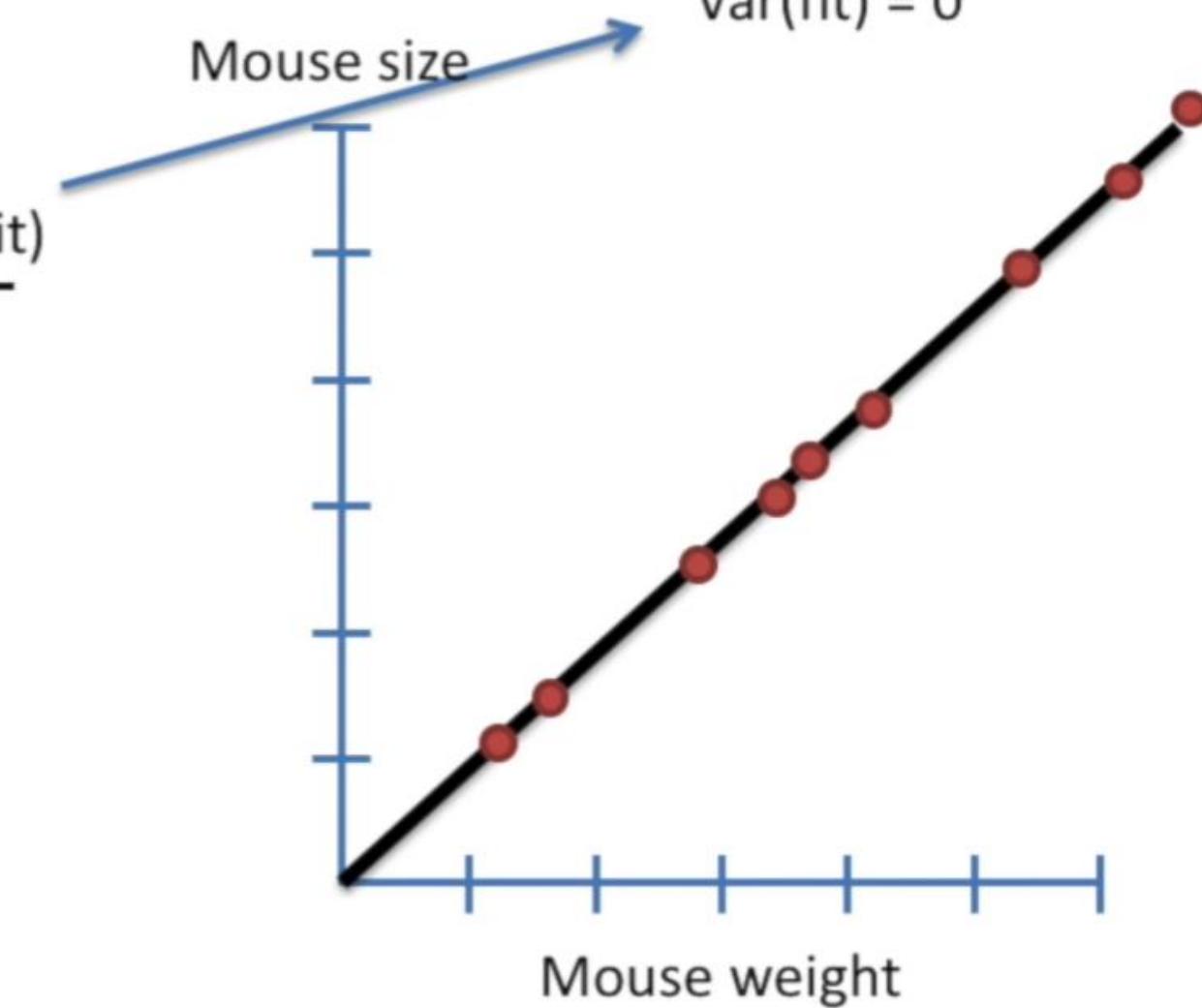


$\text{Var}(\text{mean}) = 11.1$

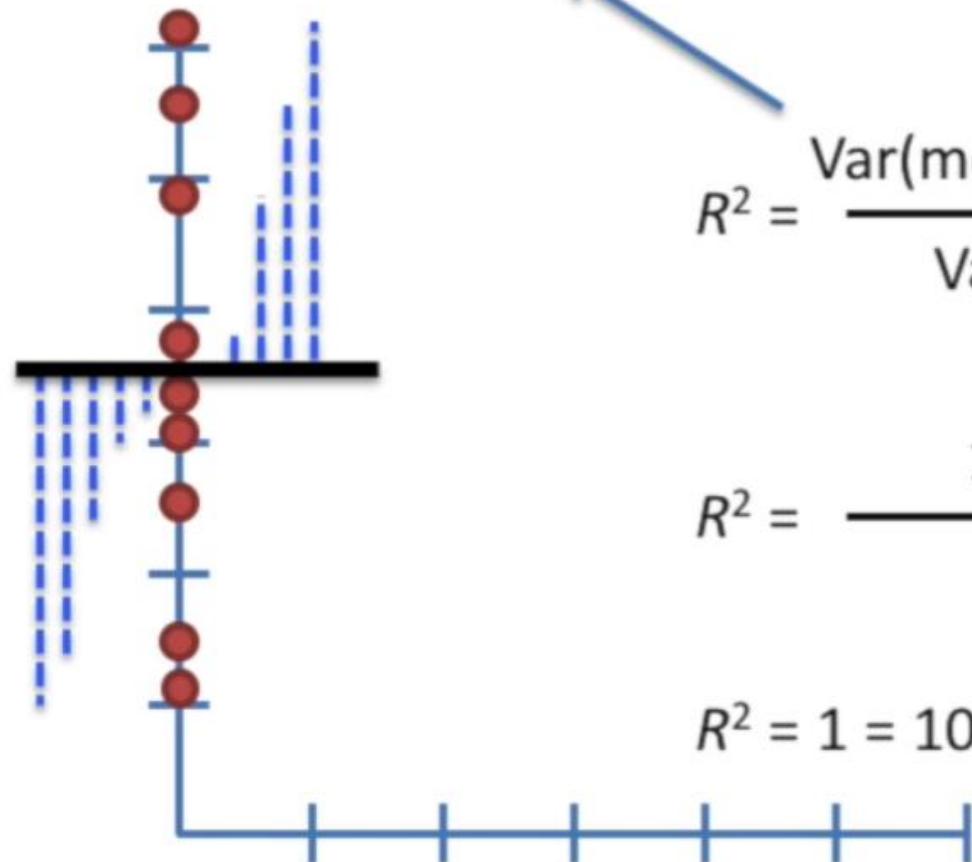


$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{fit})}{\text{Var}(\text{mean})}$$

$\text{Var}(\text{fit}) = 0$



$\text{Var}(\text{mean}) = 11.1$

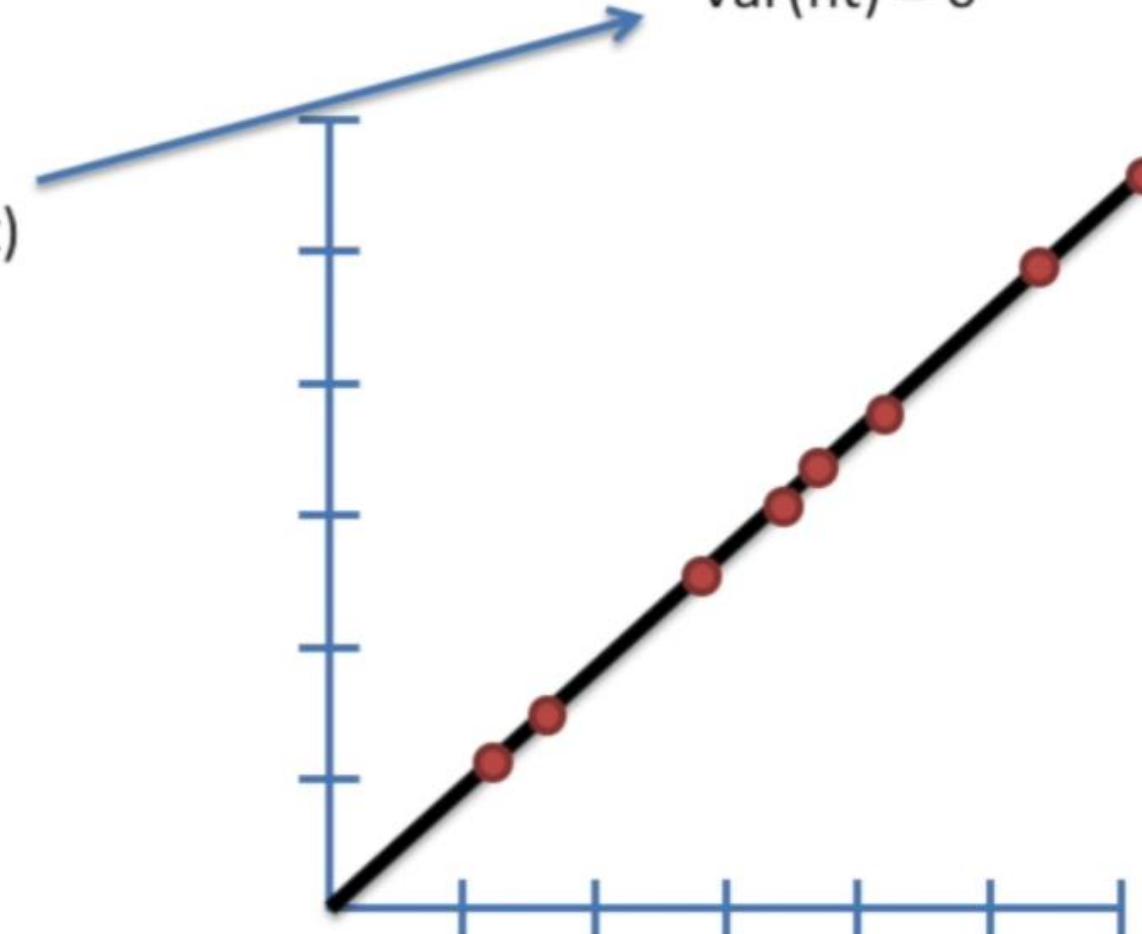


$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{fit})}{\text{Var}(\text{mean})}$$

$$R^2 = \frac{11.1 - 0}{11.1}$$

$$R^2 = 1 = 100\%$$

$\text{Var}(\text{fit}) = 0$

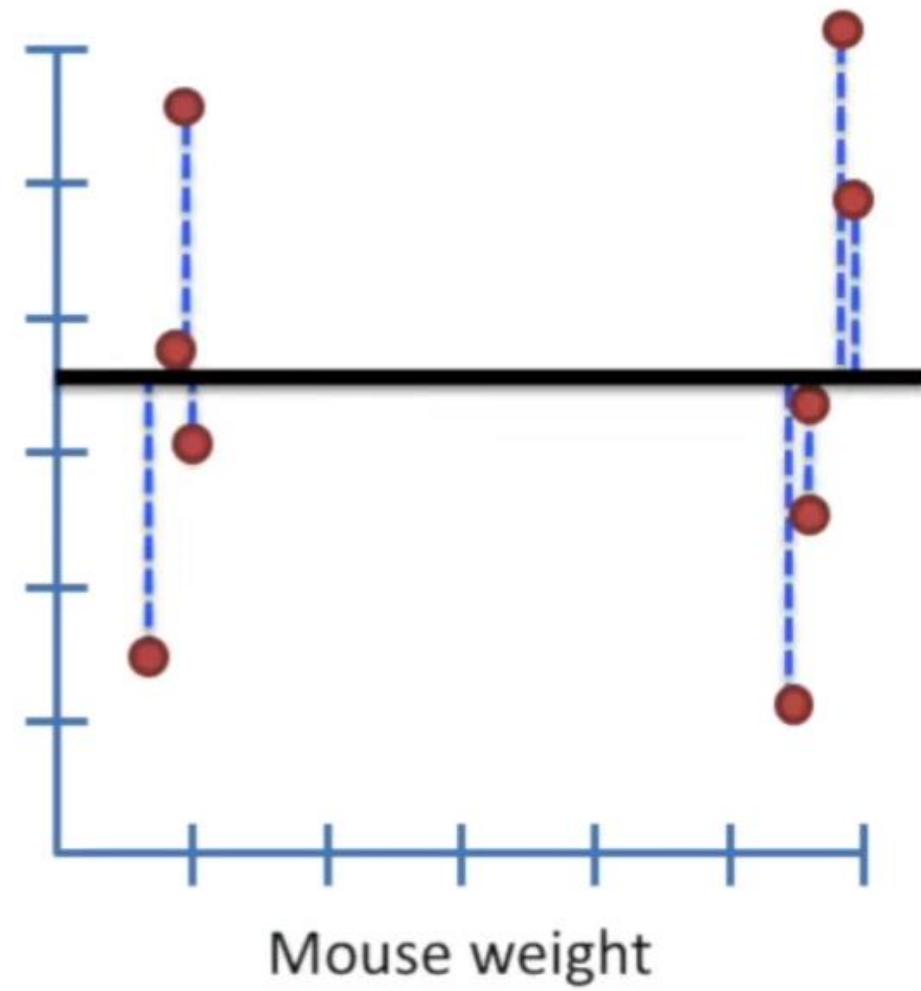
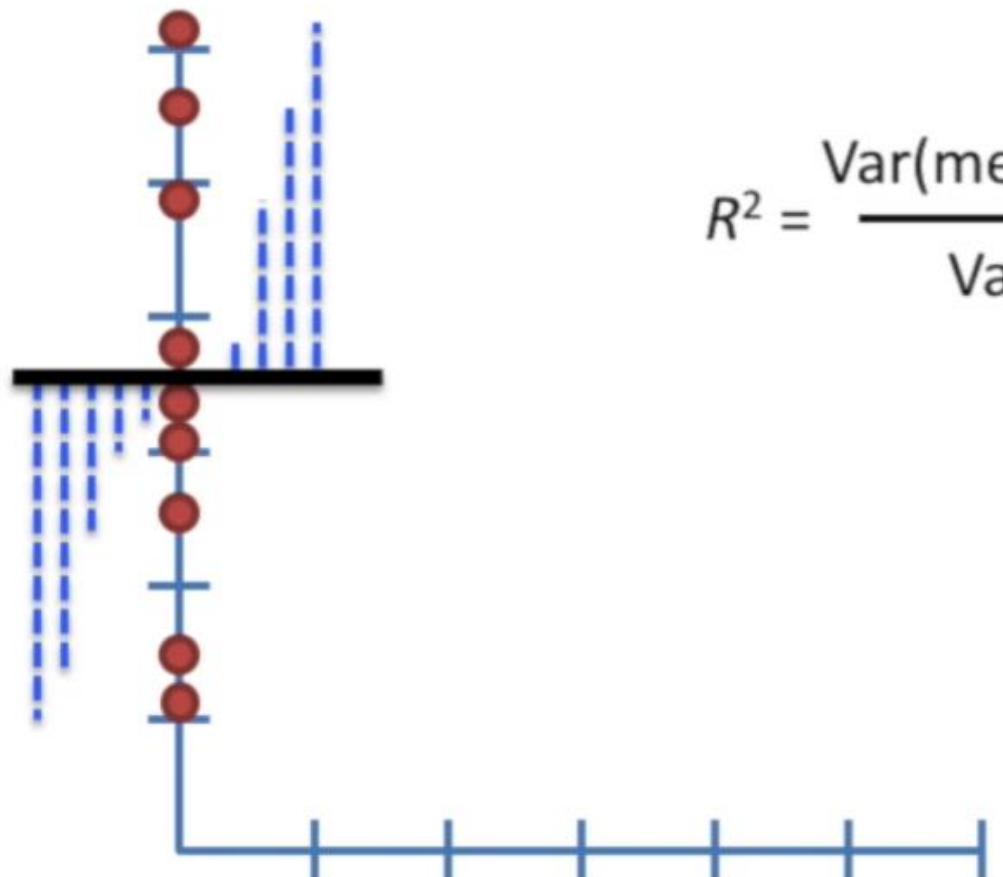


In this case, mouse weight “explains” 100% of the variation in mouse size.

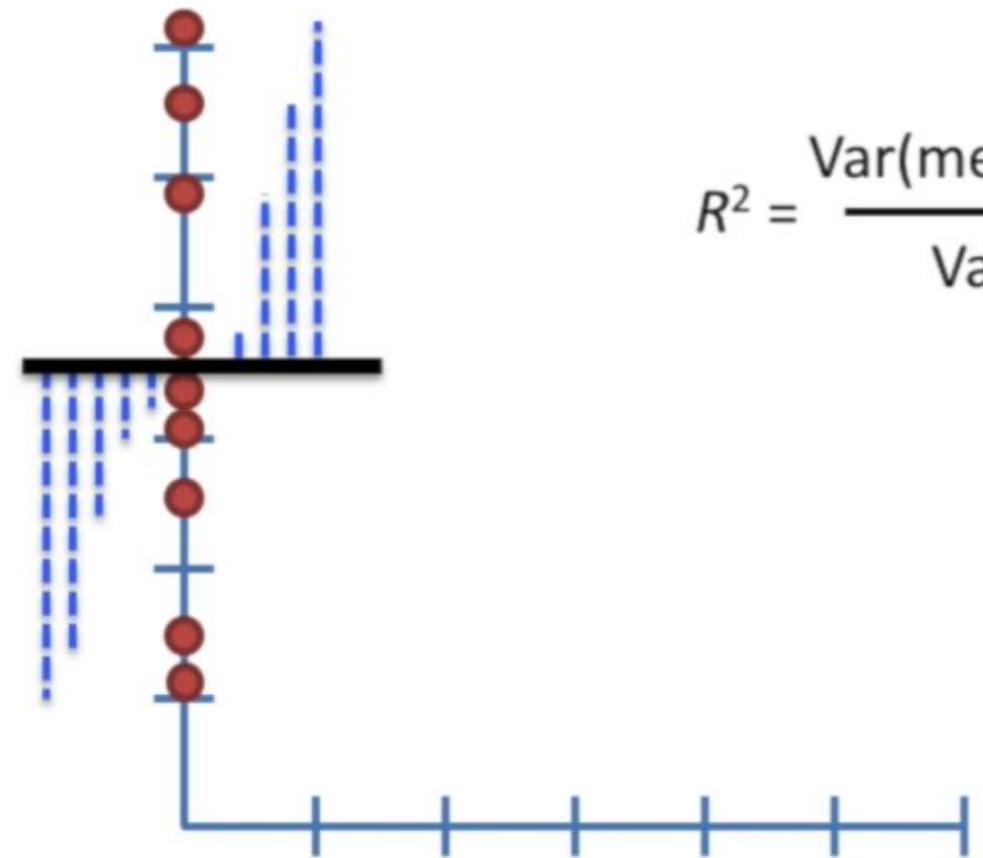
$\text{Var}(\text{mean}) = 11.1$

One last example...

$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{fit})}{\text{Var}(\text{mean})}$$

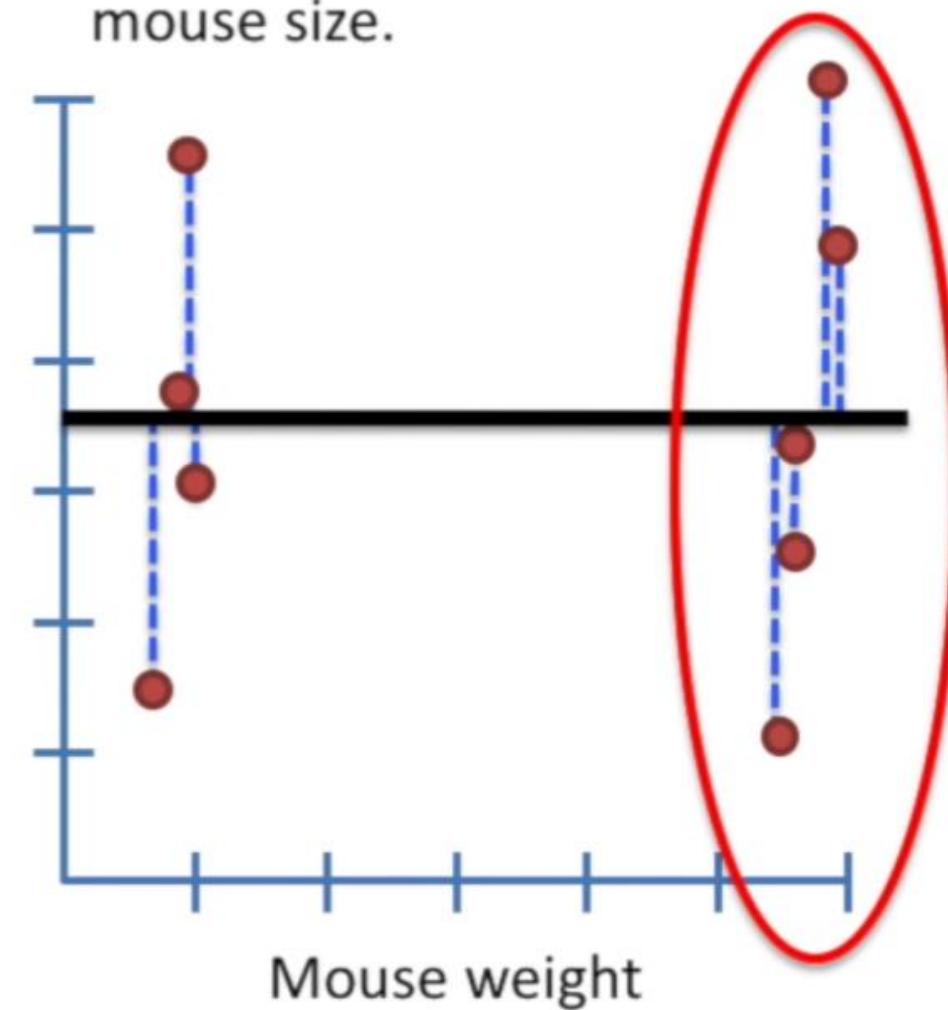


$\text{Var}(\text{mean}) = 11.1$

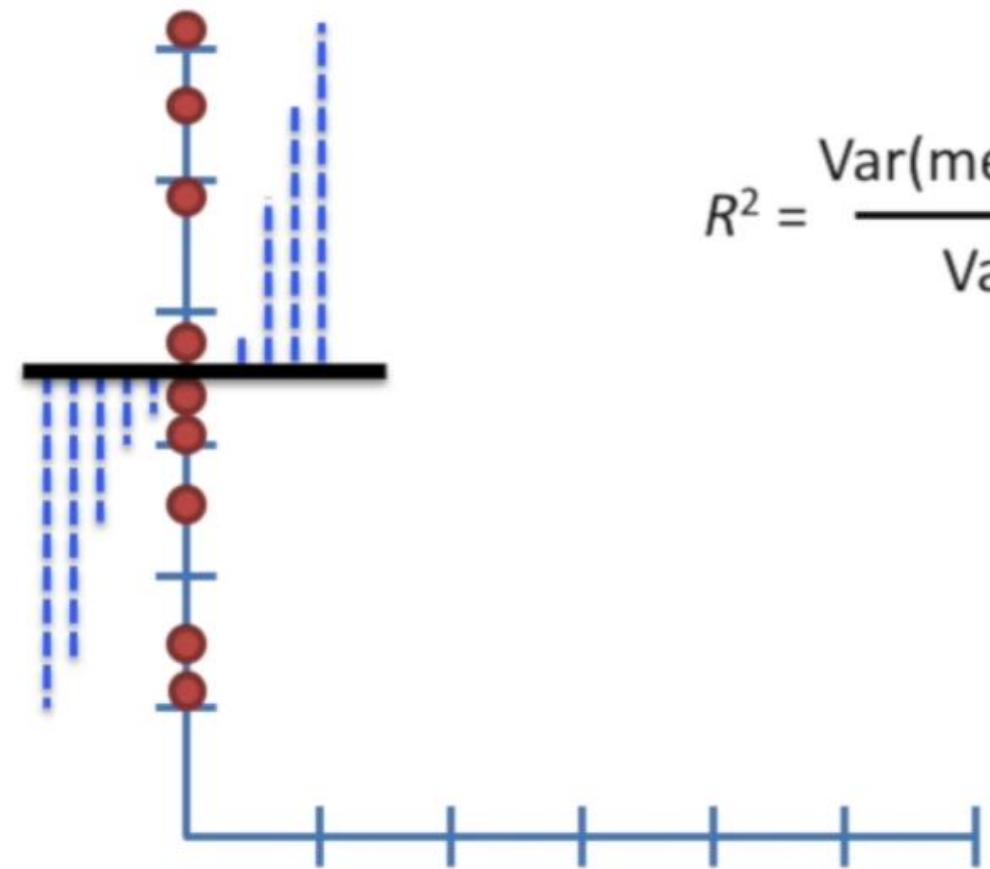


$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{fit})}{\text{Var}(\text{mean})}$$

In this case, knowing mouse weight doesn't help us predict mouse size.

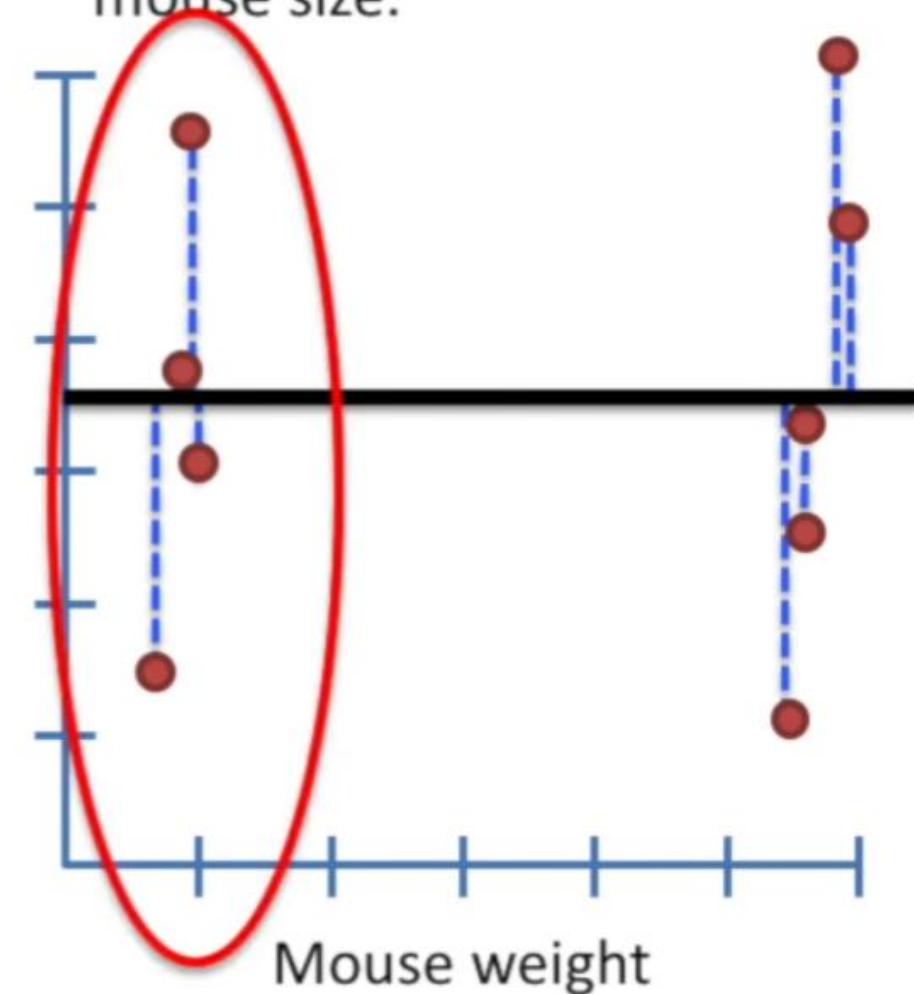


$\text{Var}(\text{mean}) = 11.1$

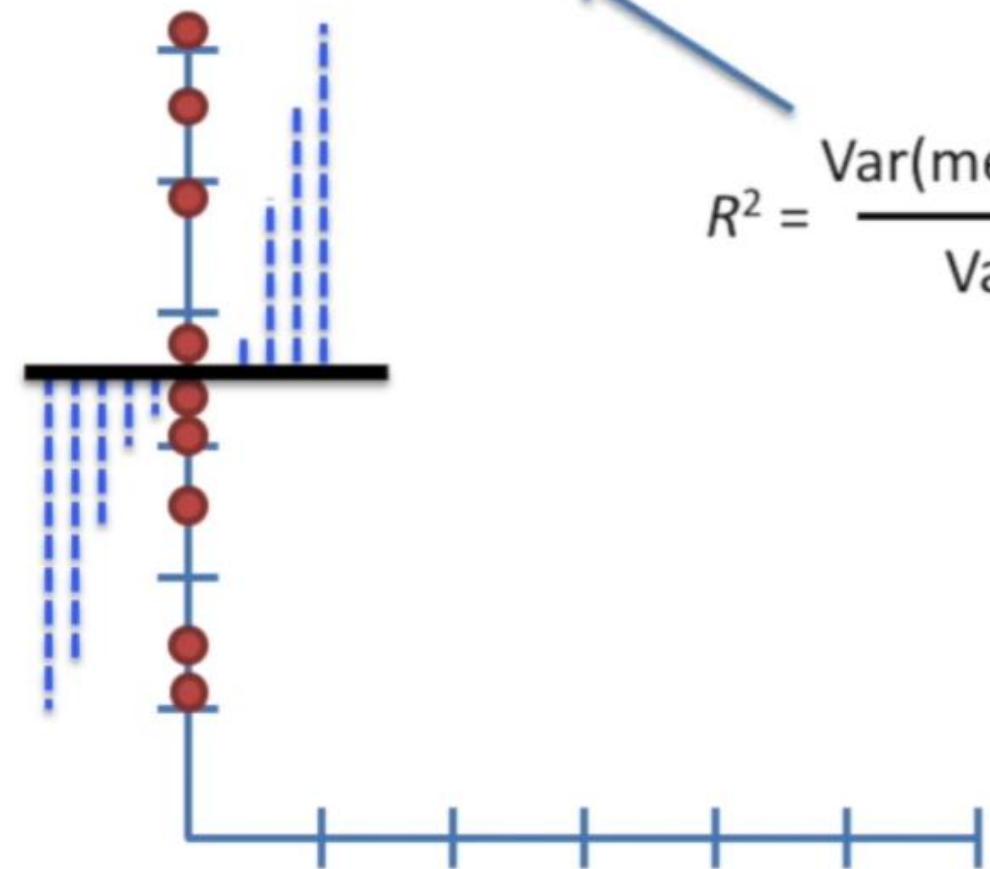


$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{fit})}{\text{Var}(\text{mean})}$$

In this case, knowing mouse weight doesn't help us predict mouse size.

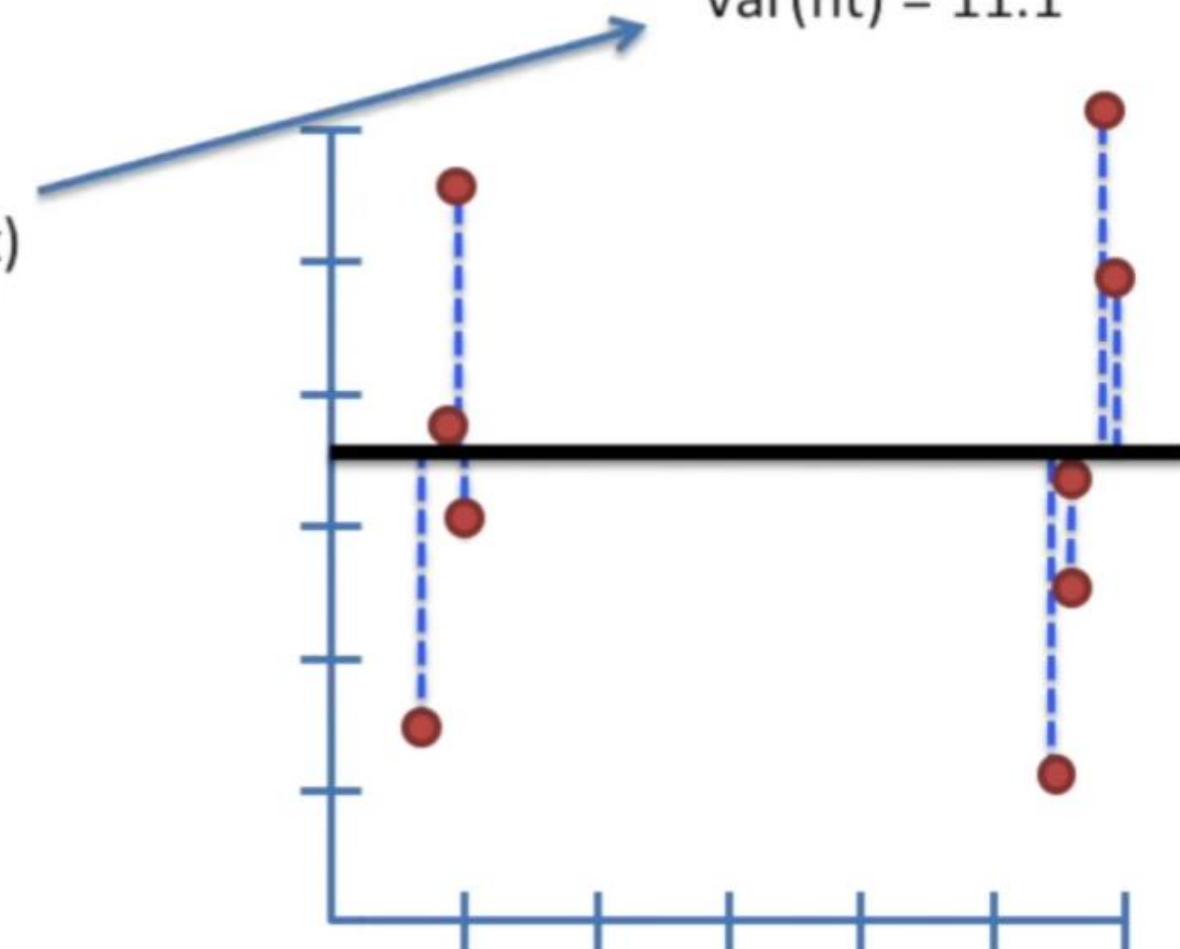


$\text{Var}(\text{mean}) = 11.1$

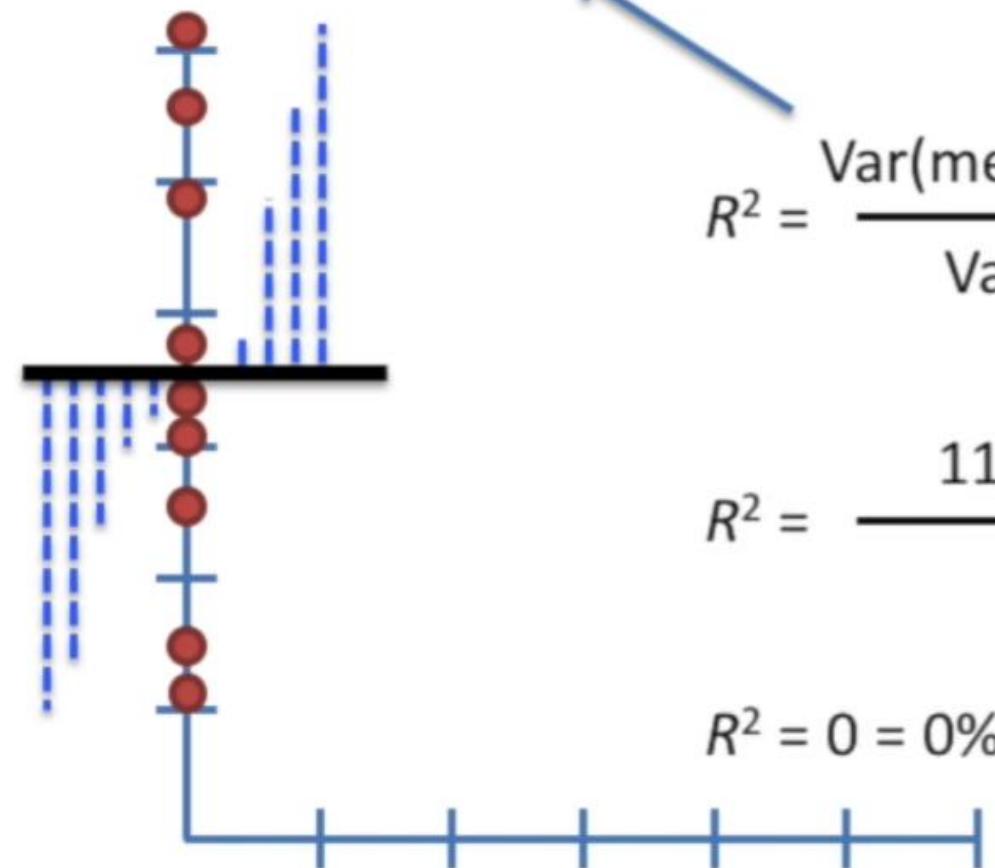


$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{fit})}{\text{Var}(\text{mean})}$$

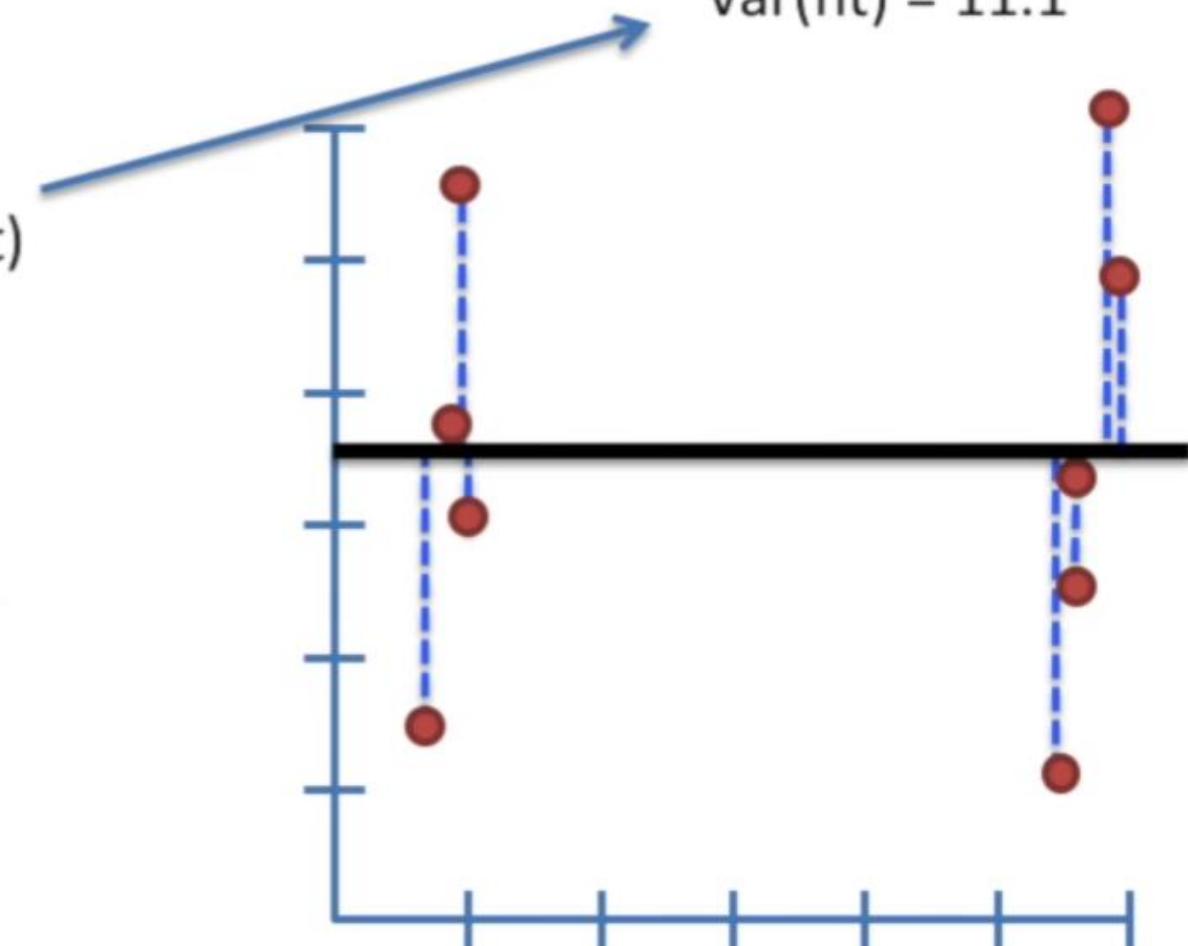
$\text{Var}(\text{fit}) = 11.1$



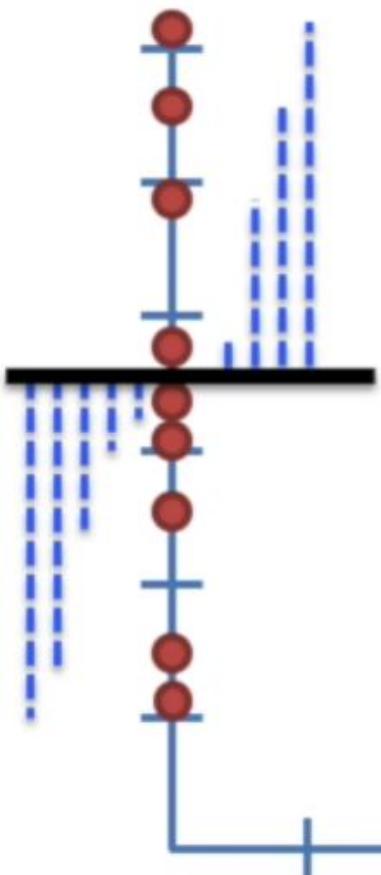
$\text{Var}(\text{mean}) = 11.1$



$\text{Var}(\text{fit}) = 11.1$



Var(mean) = 11.1

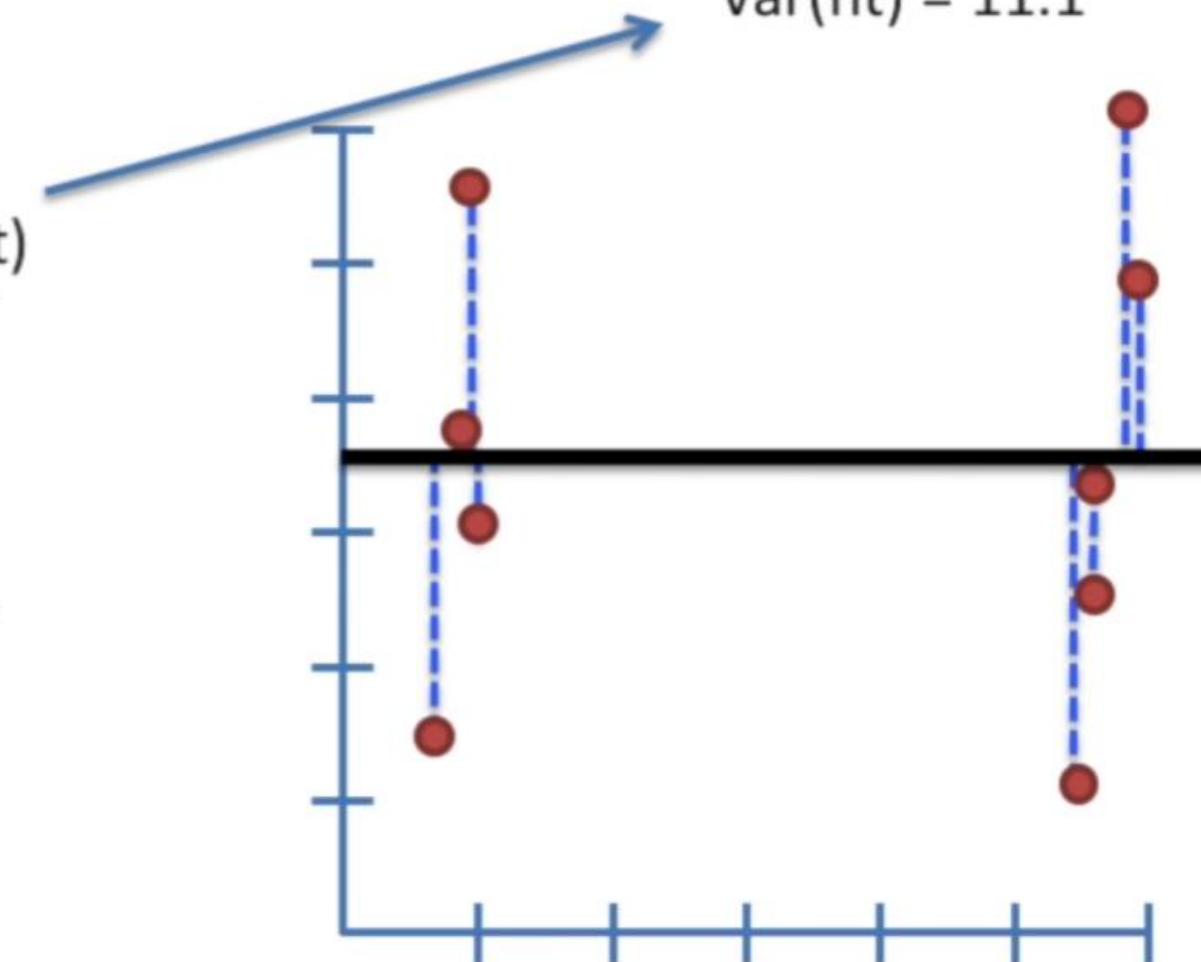


$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{fit})}{\text{Var}(\text{mean})}$$

$$R^2 = \frac{11.1 - 11.1}{11.1}$$

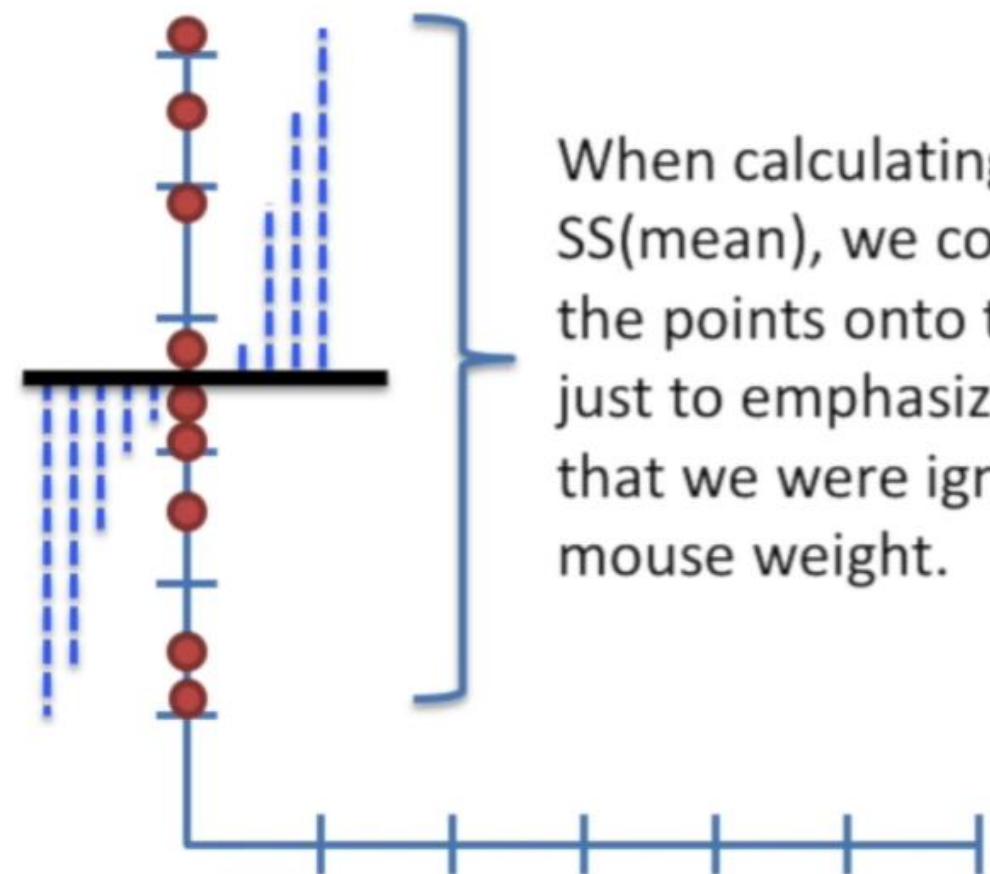
$$R^2 = 0 = 0\%$$

Var(fit) = 11.1



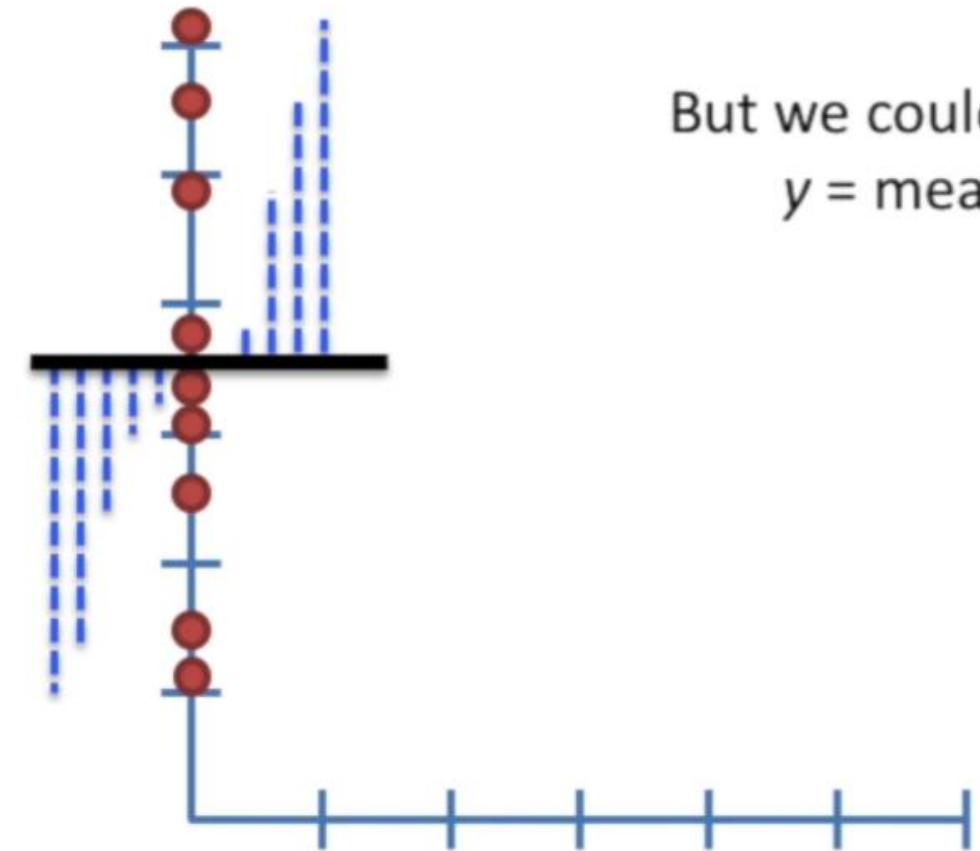
In this case, mouse weight doesn't "explain" any of the variation around the mean.

$SS(\text{mean})$



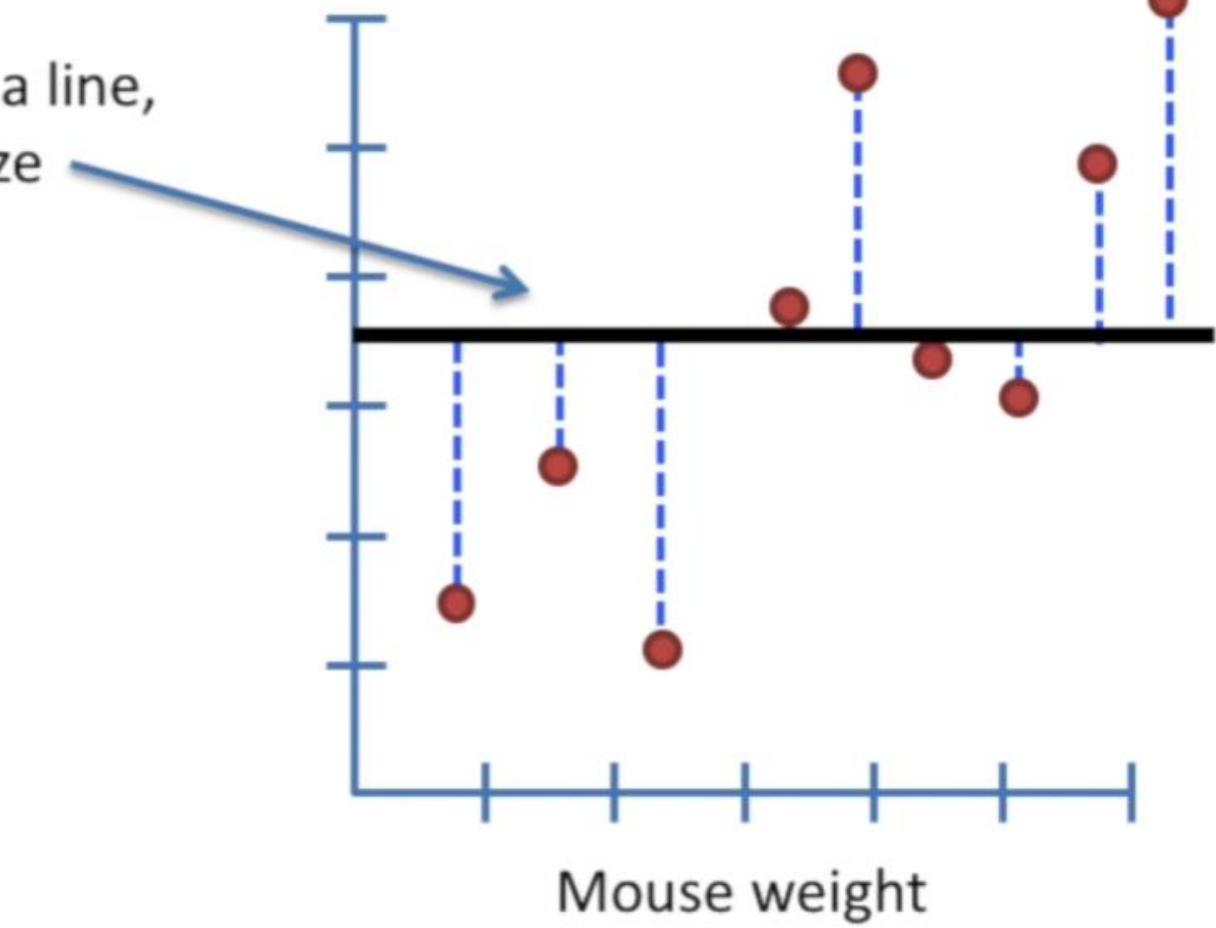
When calculating $SS(\text{mean})$, we collapsed the points onto the y-axis just to emphasize the fact that we were ignoring mouse weight.

$SS(\text{mean})$

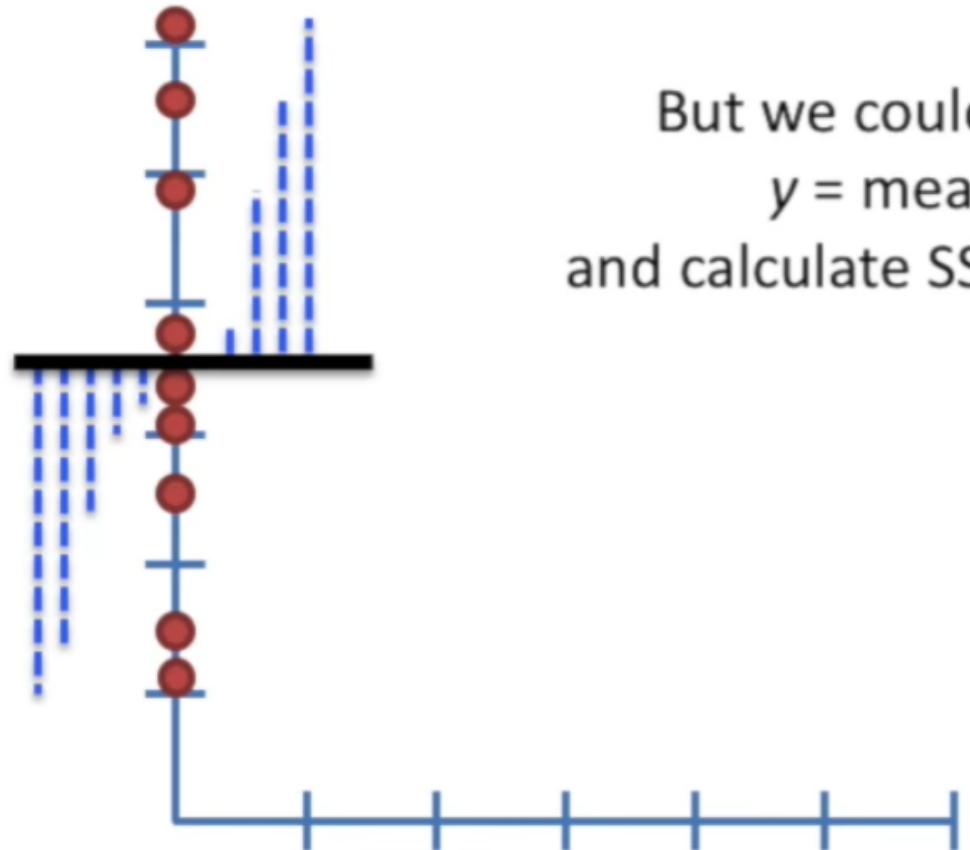


But we could just draw a line,
 $y = \text{mean mouse size}$

Mouse size

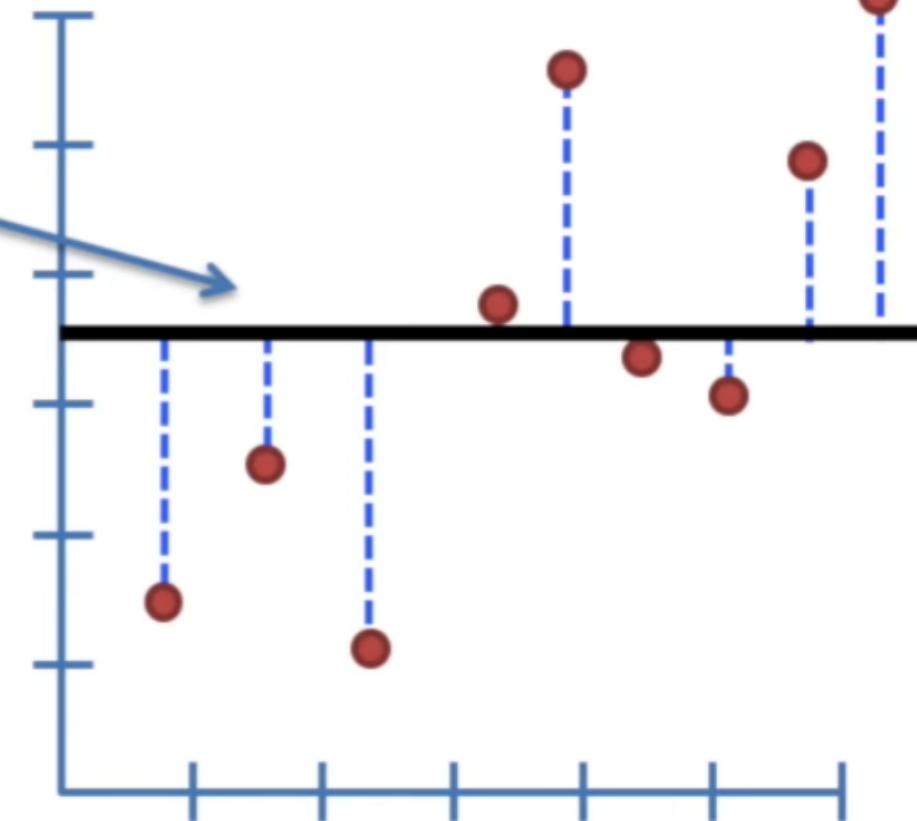


$SS(\text{mean})$



But we could just draw a line,
 $y = \text{mean mouse size}$
and calculate $SS(\text{mean})$ around that.

Mouse size



Mouse weight