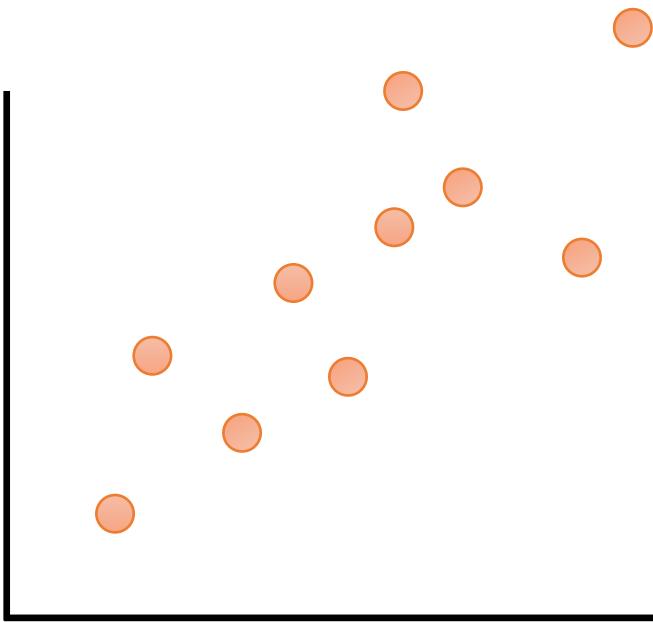


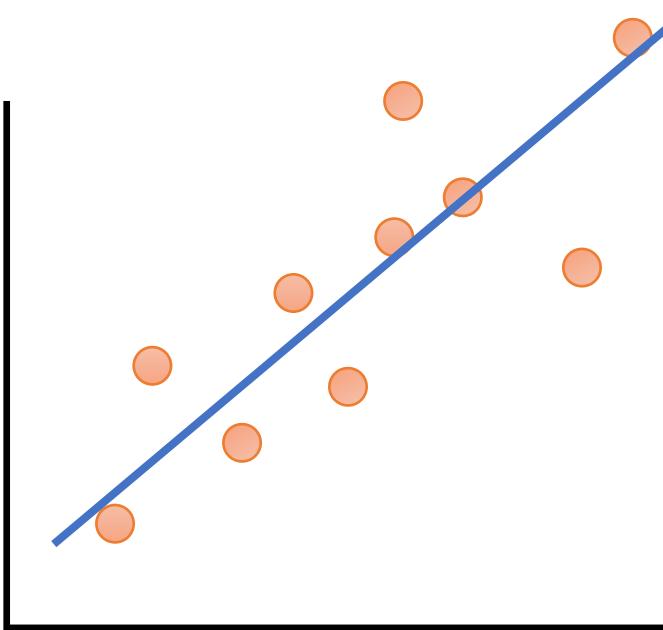
# Linear Models

Part II

Last time



# Last time



Call:

lm(formula = y ~ x, data = df)

Residuals:

1	2	3	4	5	6	7
5.4033	-4.9961	-0.1866	-0.8958	1.6027	-0.6310	-0.2965

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )		
(Intercept)	2.93686	6.12210	0.48	0.651679		
x	0.94070	0.09324	10.09	0.000164 ***		
---						
Signif. codes:	0 ‘***’	0.001 ‘**’	0.01 ‘*’	0.05 ‘.’	0.1 ‘ ’	1

Residual standard error: 3.407 on 5 degrees of freedom

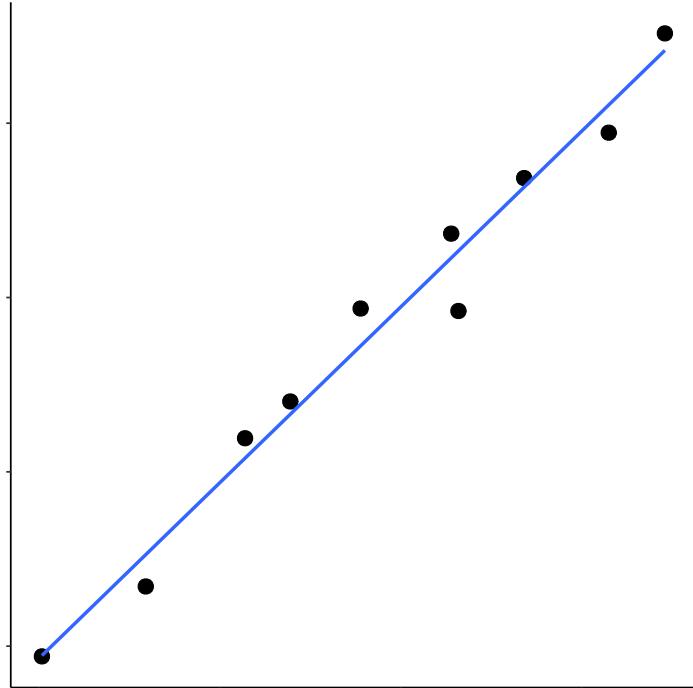
Multiple R-squared: 0.9532, Adjusted R-squared: 0.9438

F-statistic: 101.8 on 1 and 5 DF, p-value: 0.0001638

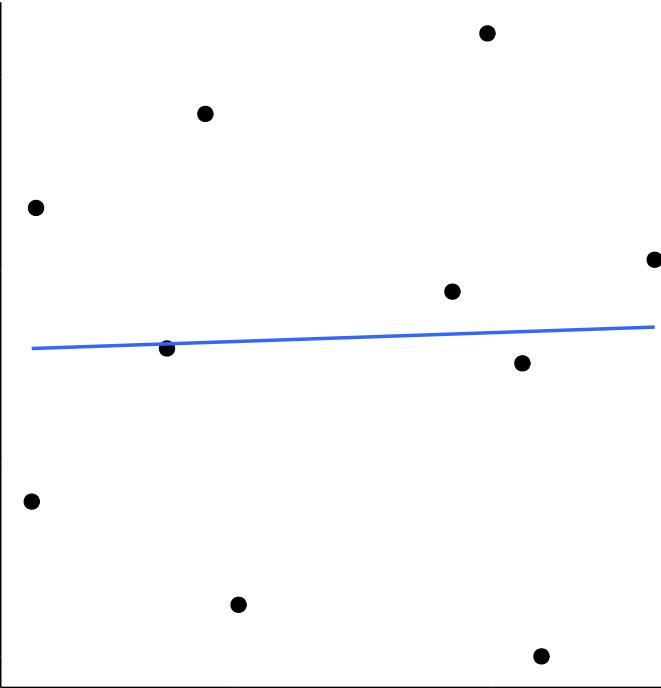
```
Call:  
lm(formula = y ~ x, data = df)  
  
Residuals:  
    1      2      3      4      5      6      7  
 5.4033 -4.9961 -0.1866 -0.8958  1.6027 -0.6310 -0.2965  
  
Coefficients:  
              Estimate Std. Error t value Pr(>|t|)  
(Intercept) 2.93686   6.12210   0.48  0.651679  
x            0.94070   0.09324  10.09 0.000164 ***  
---  
Signif. codes:  0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1  
  
Residual standard error: 3.407 on 5 degrees of freedom  
Multiple R-squared:  0.9532,    Adjusted R-squared:  0.9438  
F-statistic: 101.8 on 1 and 5 DF,  p-value: 0.0001638
```

# Correlation-Coefficient “R”

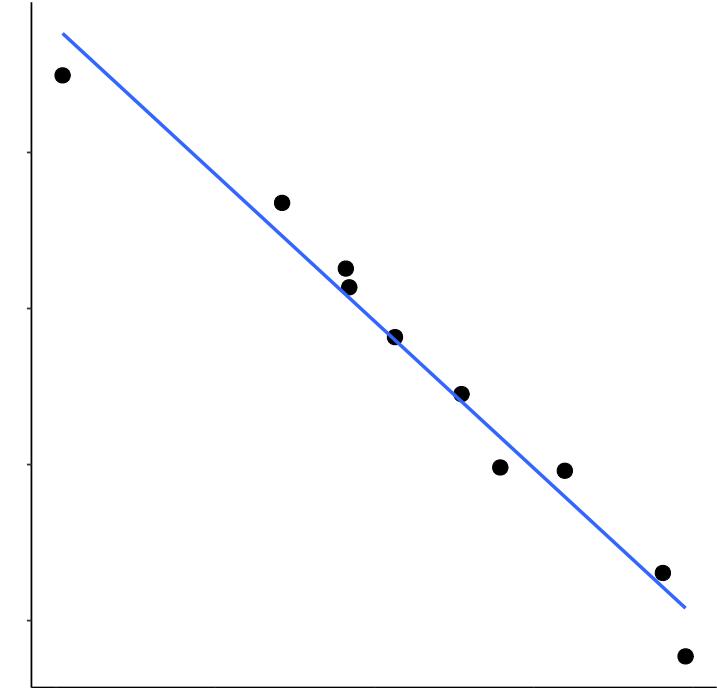
$R = 1$



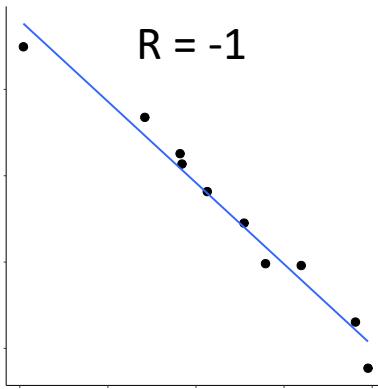
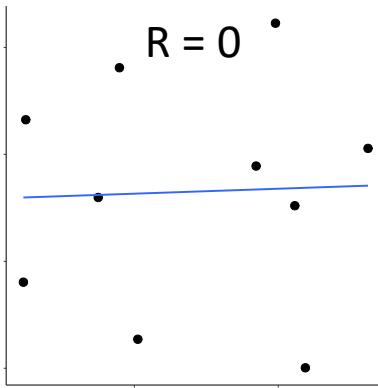
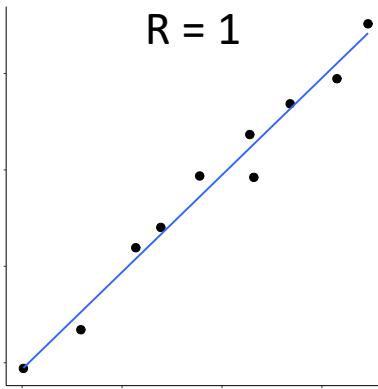
$R = 0$



$R = -1$



# R and R<sup>2</sup>

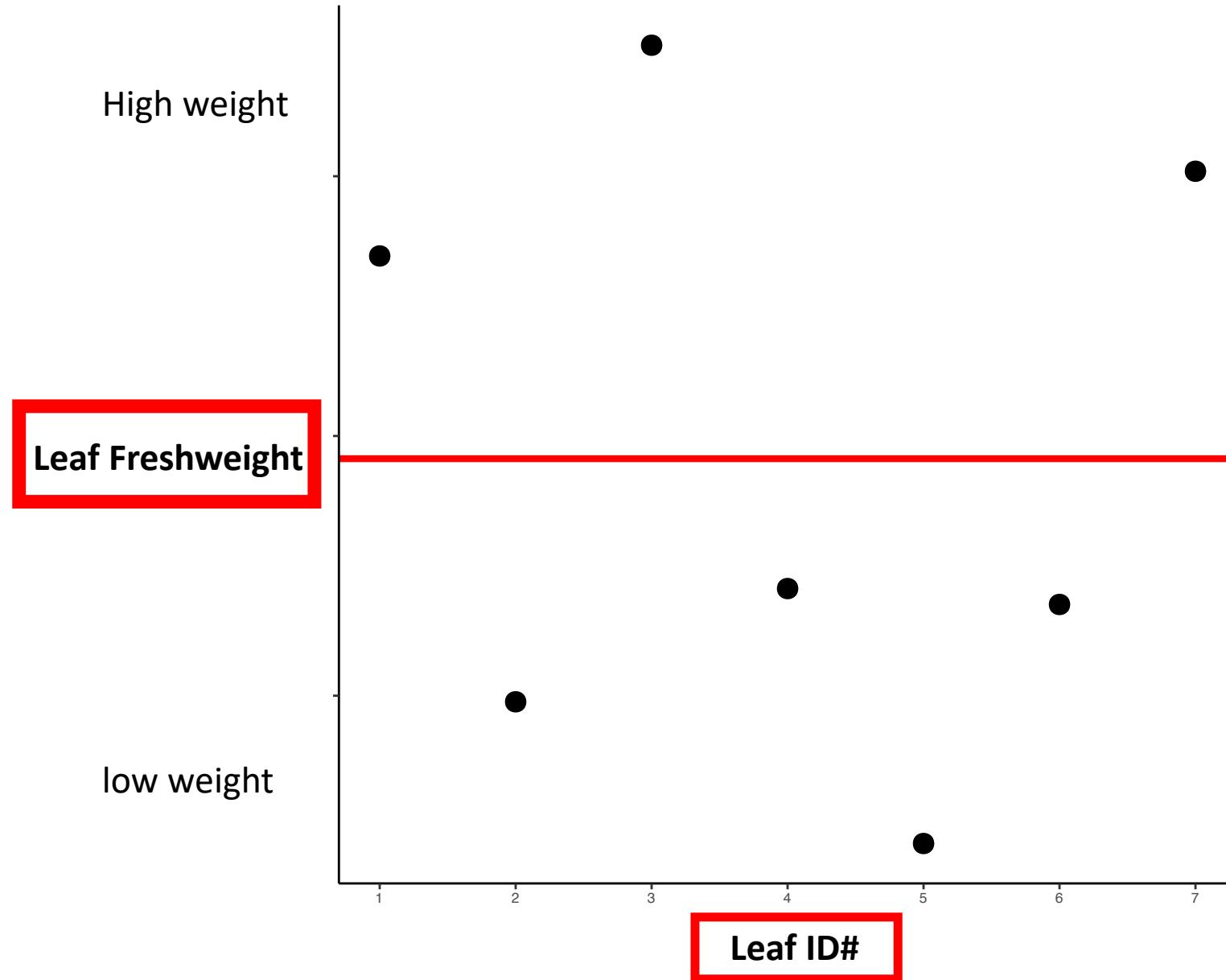


**You probably know about correlation (regular “R”)**

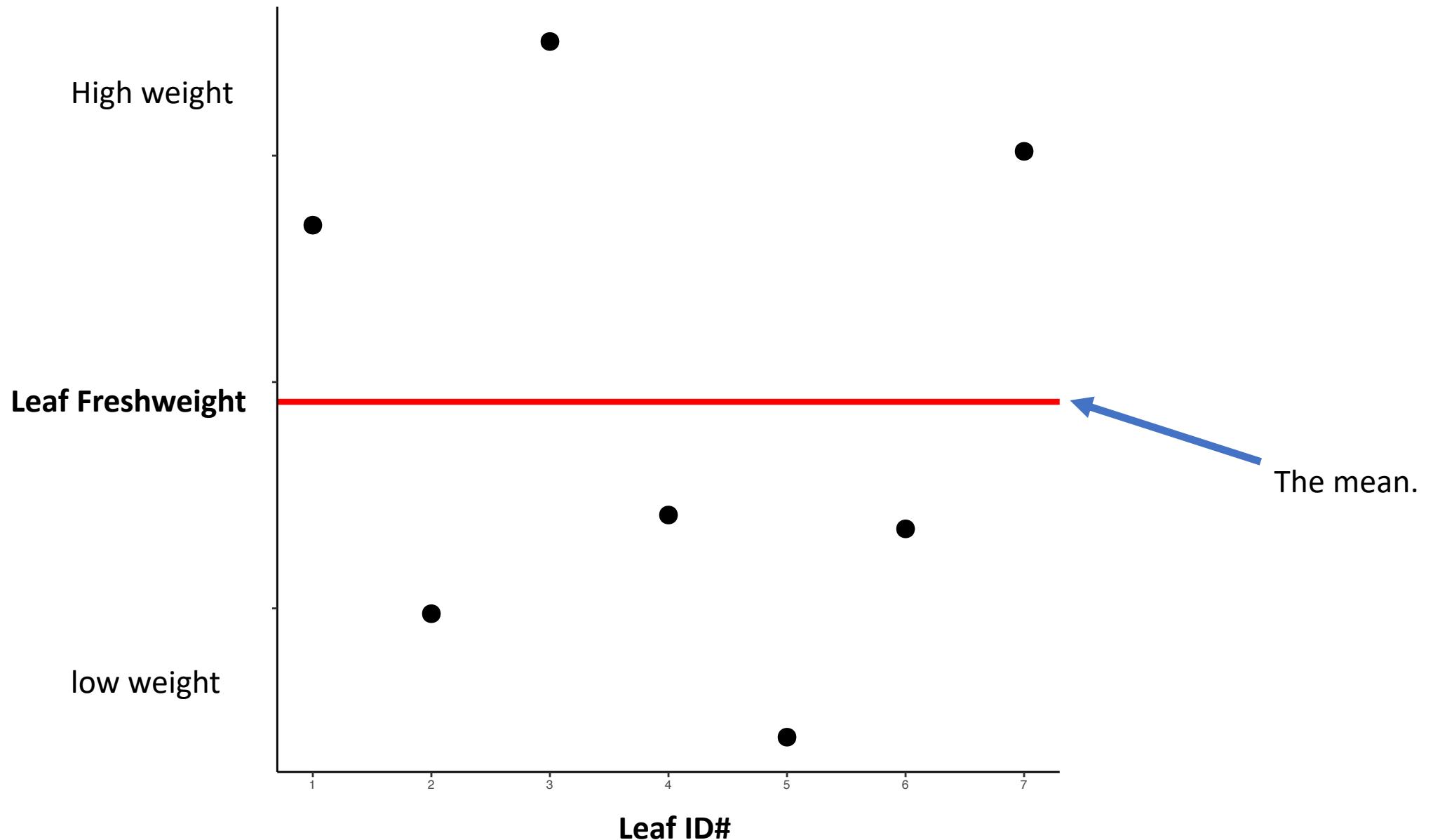
- Correlation values close to 1 or -1 are good and tell you two quantitative variables are strongly related.
- Correlation values close to 0 are lame

**Why should we care about R<sup>2</sup>? (and what even is that?)**

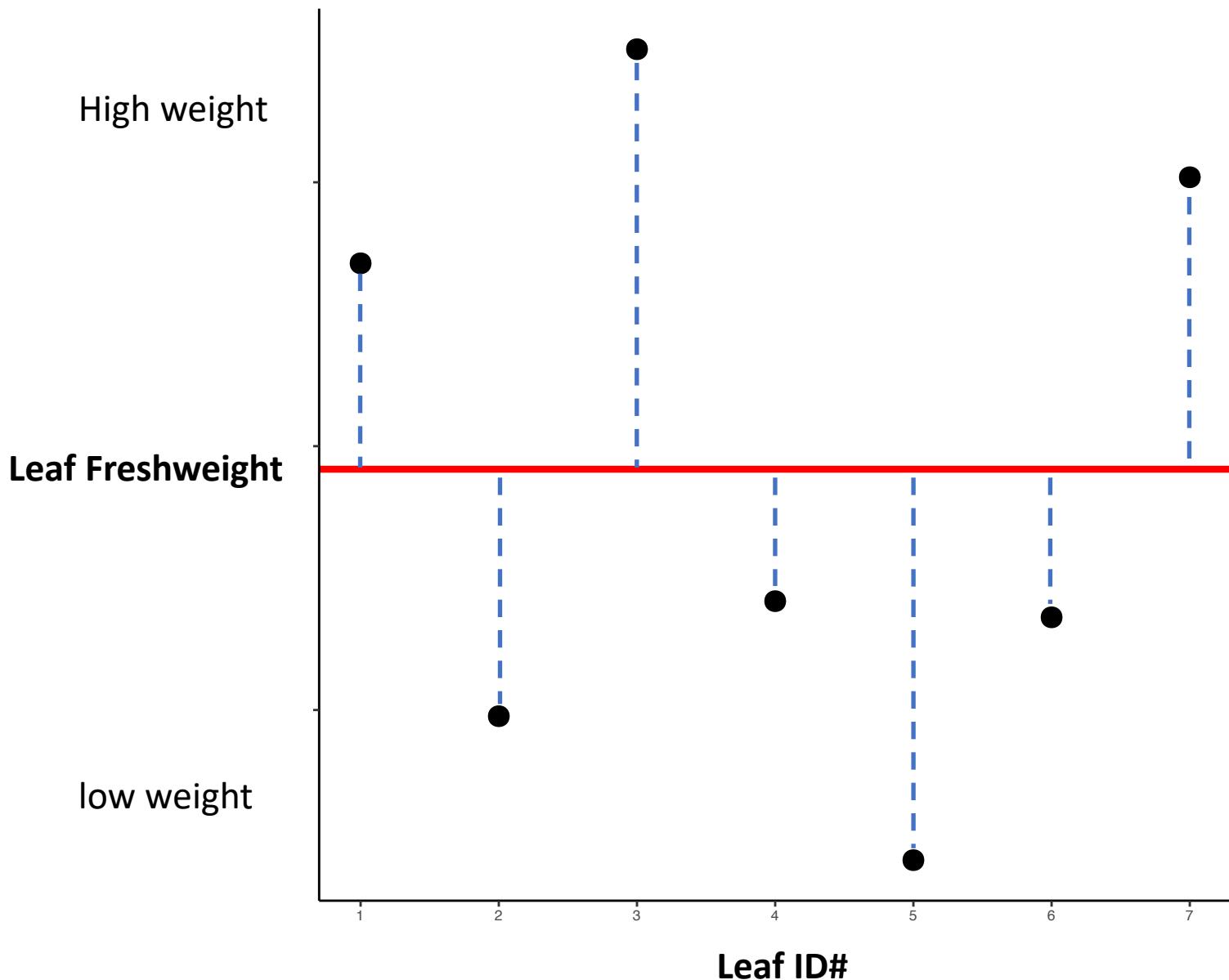
- $R^2$  is very similar to its hipper cousin, R, but..
- Interpretation is easier:
  - It's not obvious that (regular)  $R = 0.7$  is twice as good as  $R = 0.5$
  - However,  $R^2=0.7$  is what it looks like, 1.4 times as good as  $R^2=0.5$



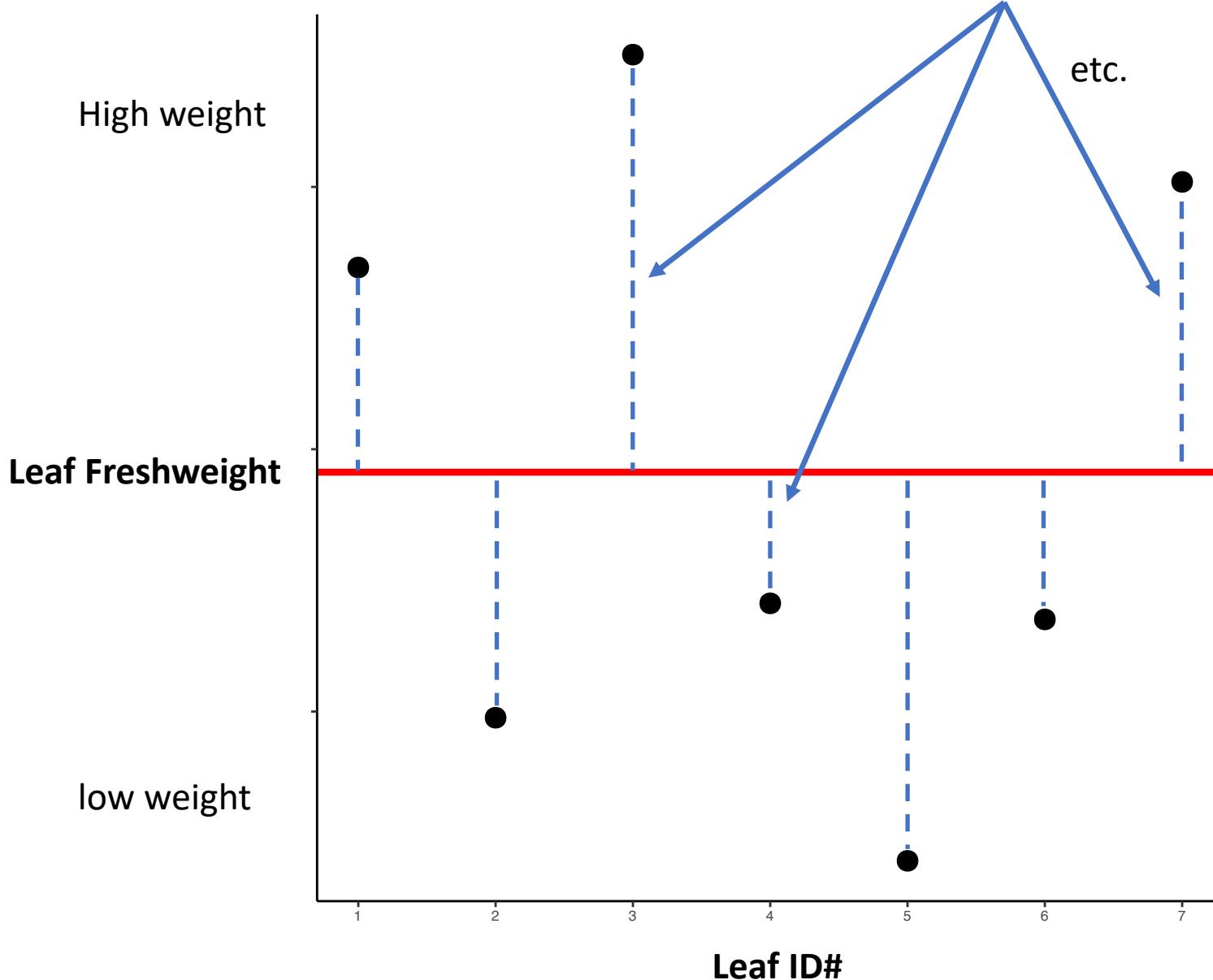
## The Mean (or Average) of leaf weights



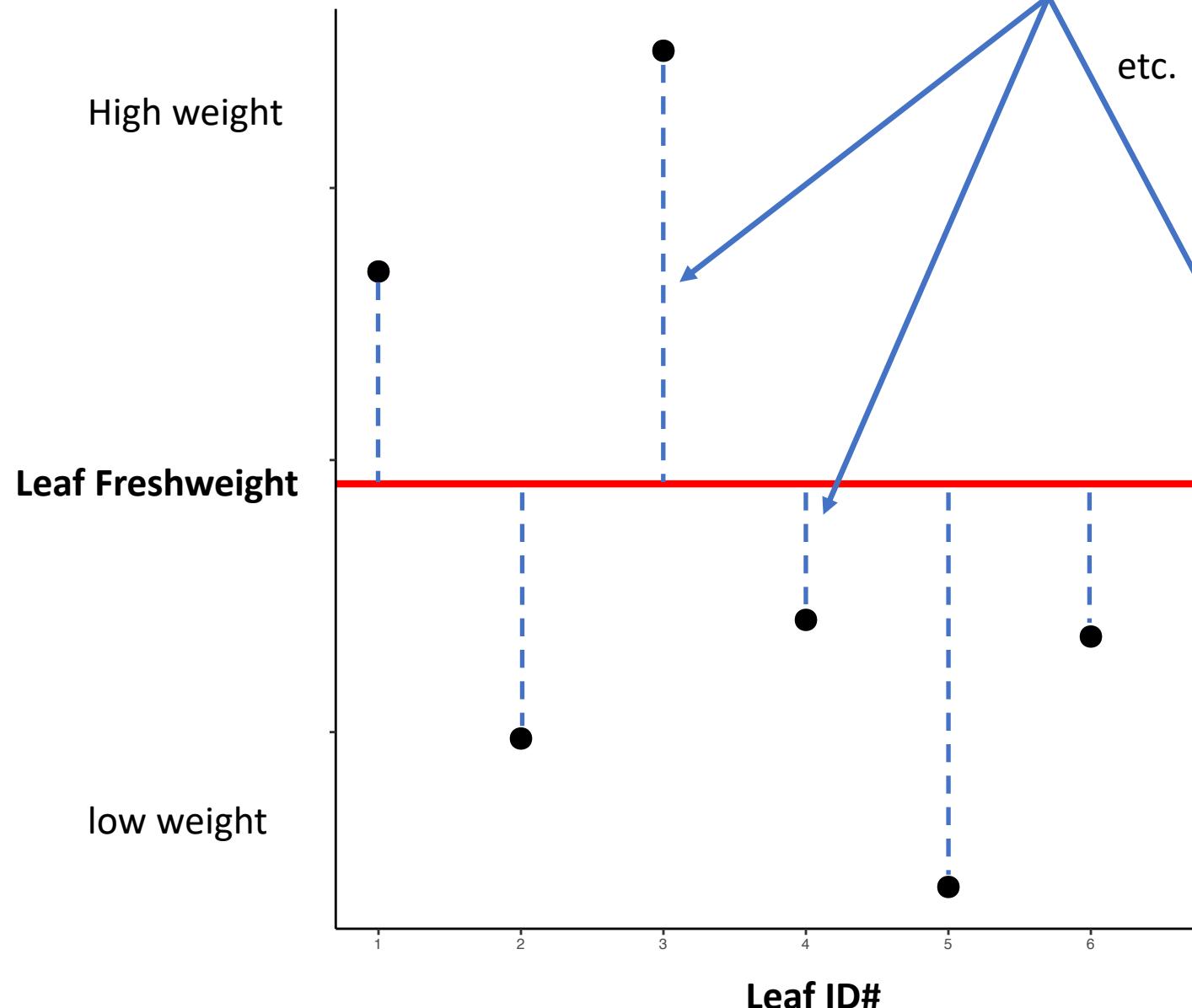
The Variation of the data = Sum(weight – mean)<sup>2</sup>



The Variation of the data = Sum(weight – mean)<sup>2</sup>

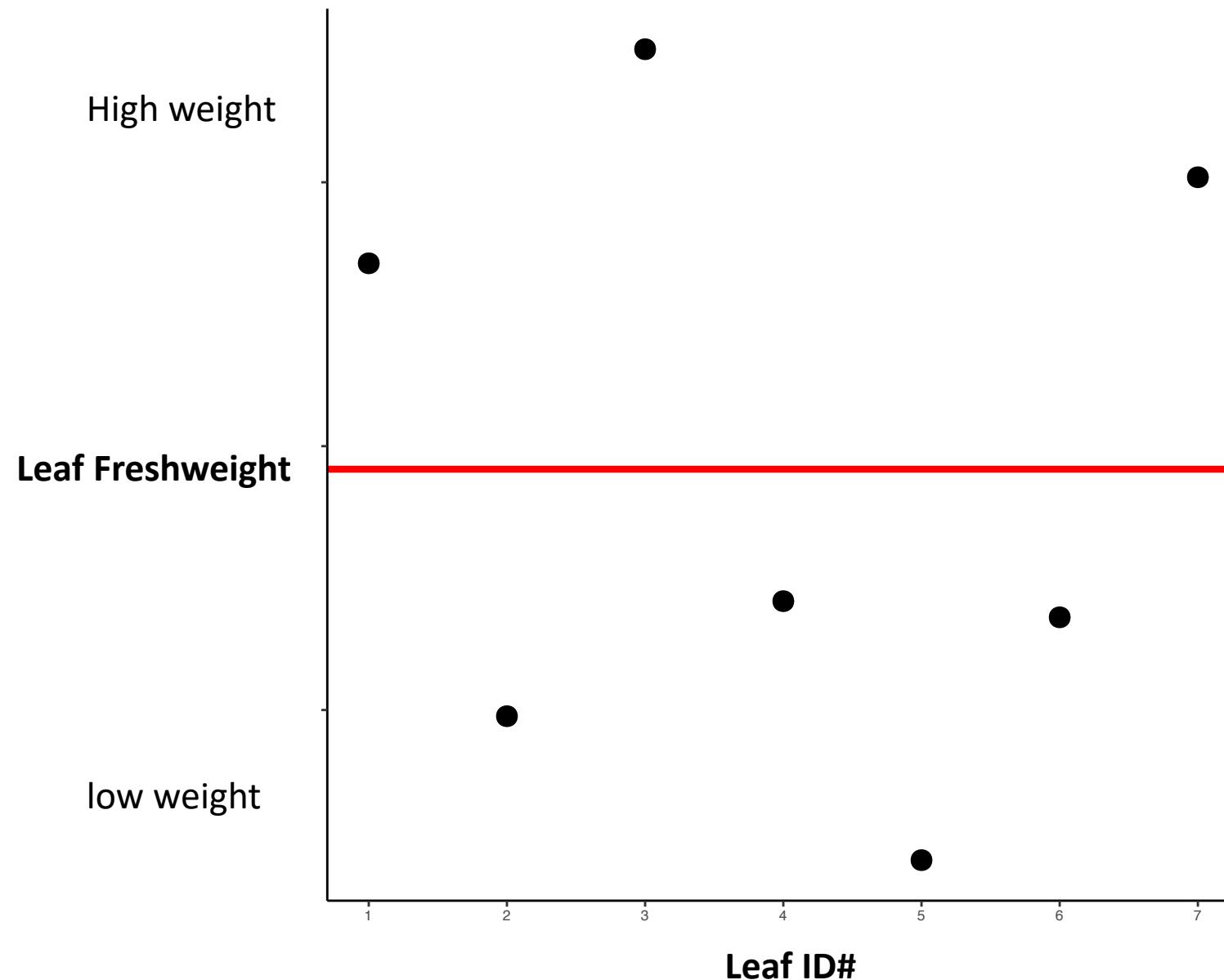


The Variation of the data = Sum(weight – mean)<sup>2</sup>

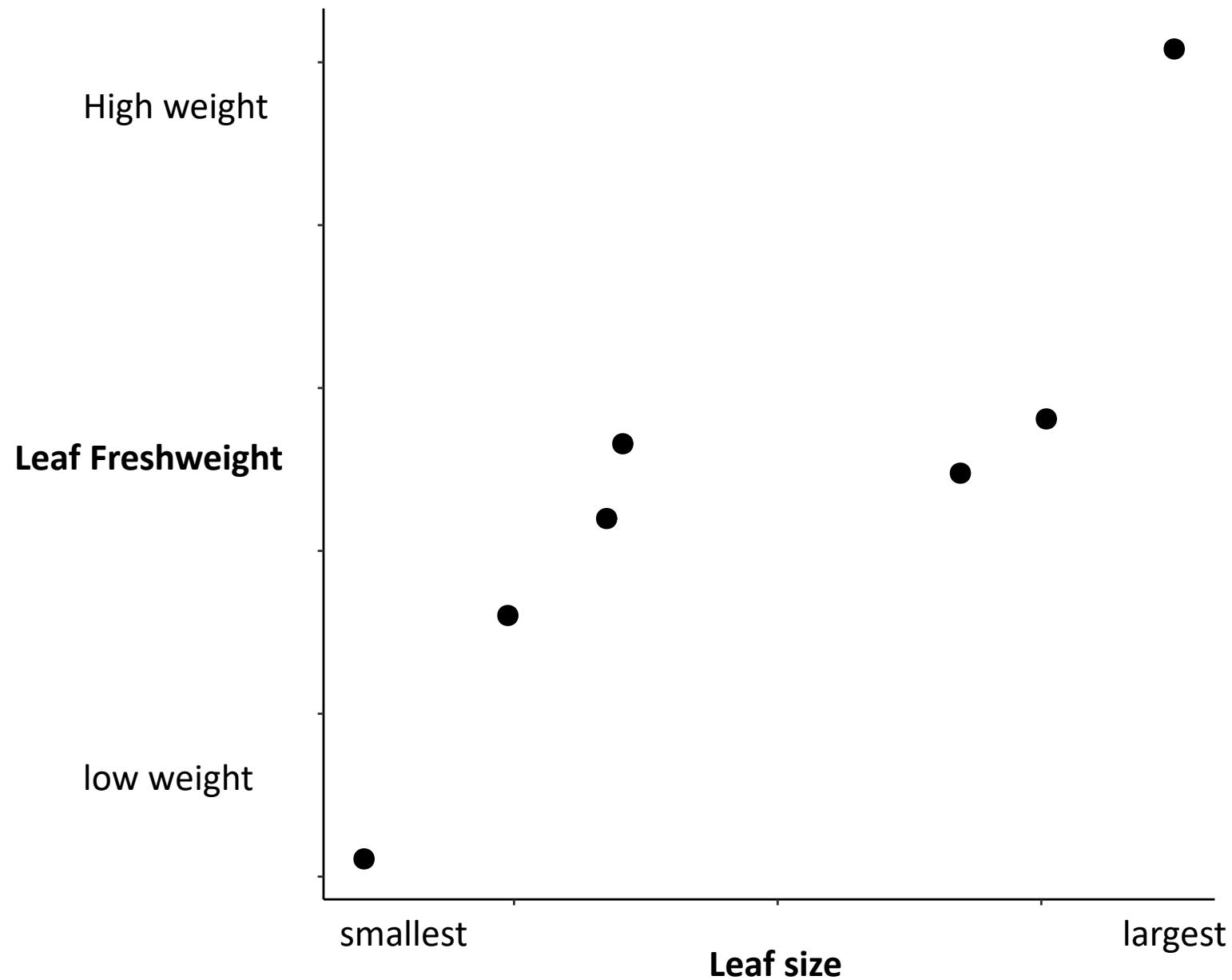


The difference between each data point is squared so that points below the mean don't cancel out points above the mean.

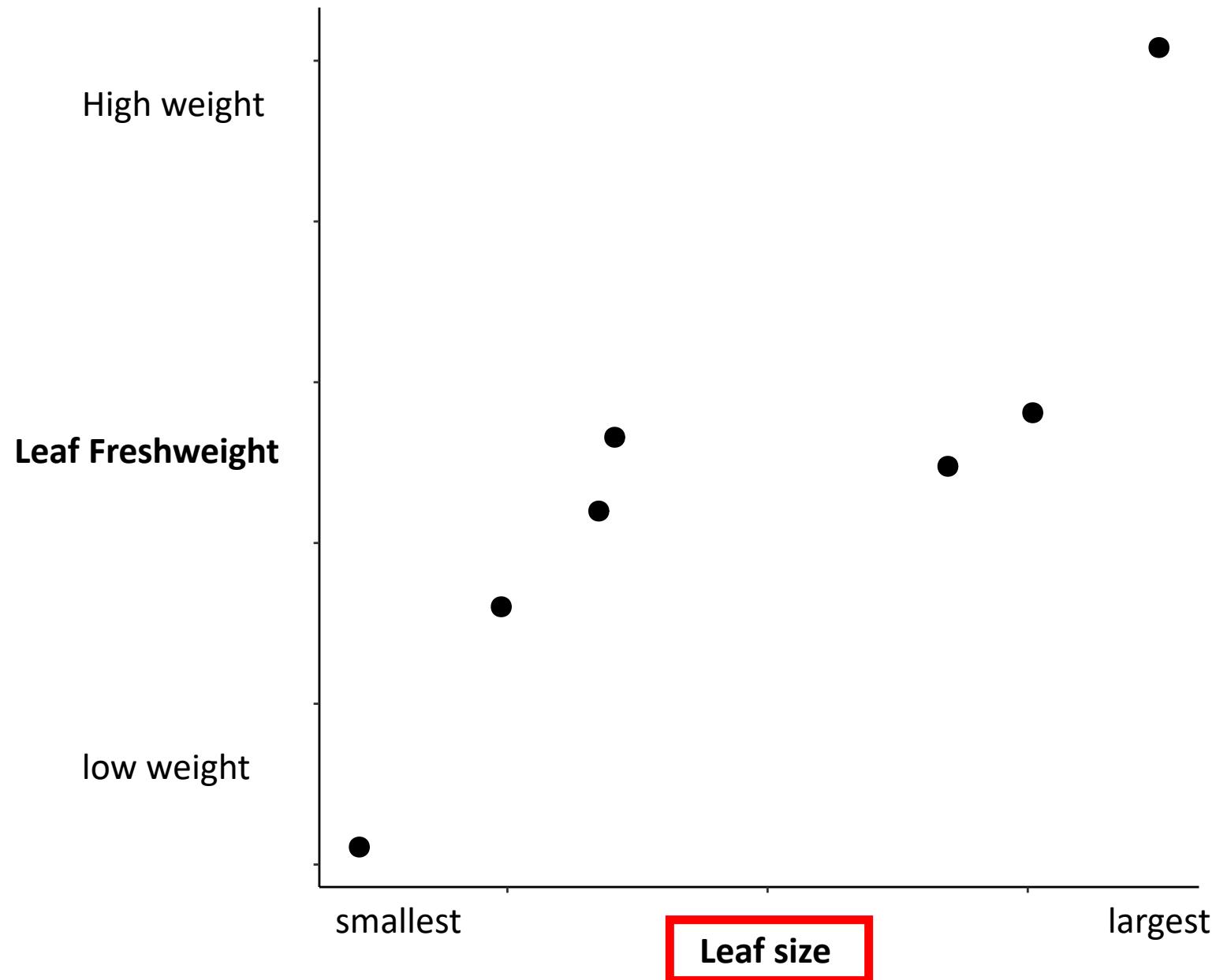
Now, what if, instead of ordering our leafs by their ID#, we ordered them by their size?



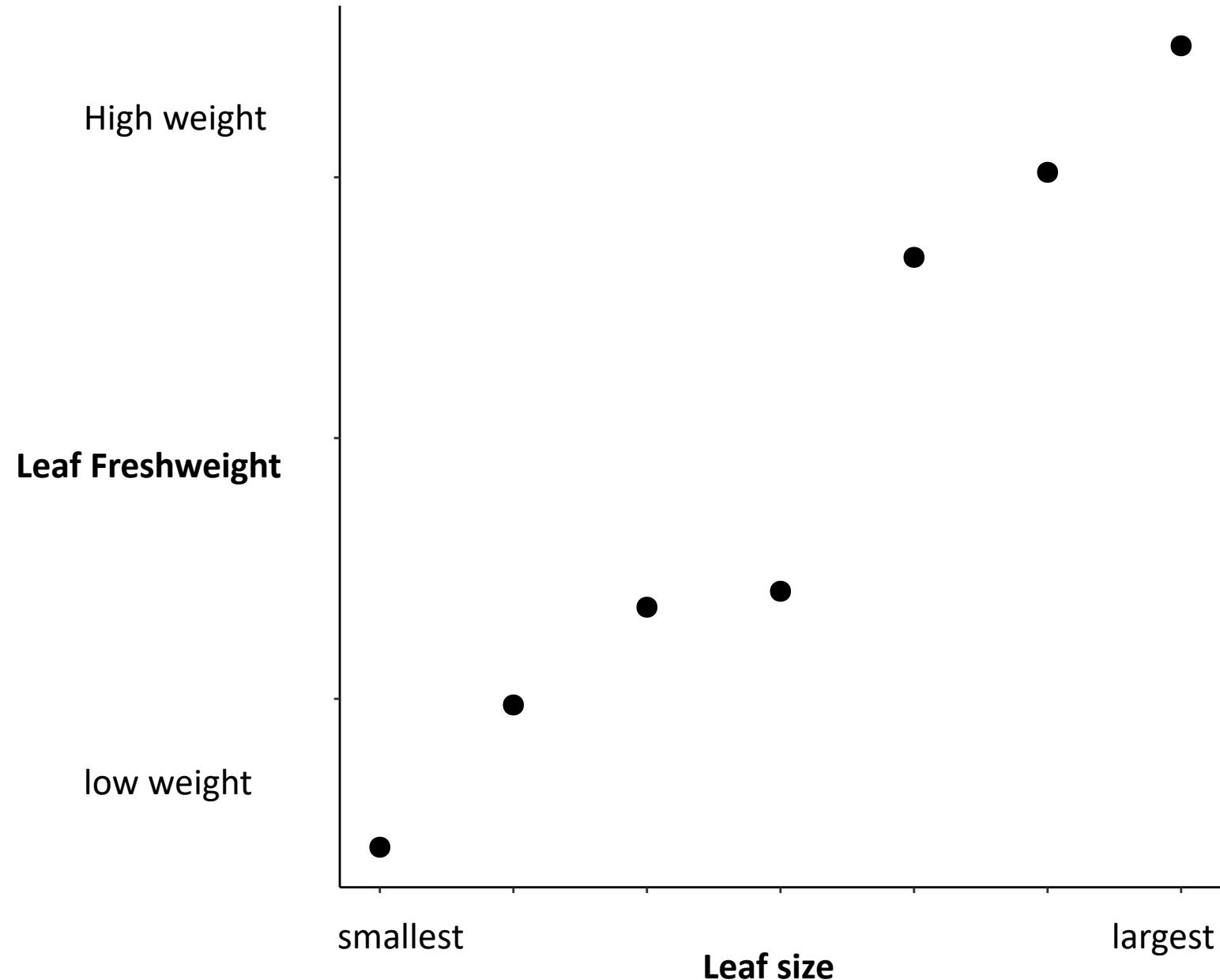
Now, what if, instead of ordering our leafs by their ID#, we ordered them by their size?



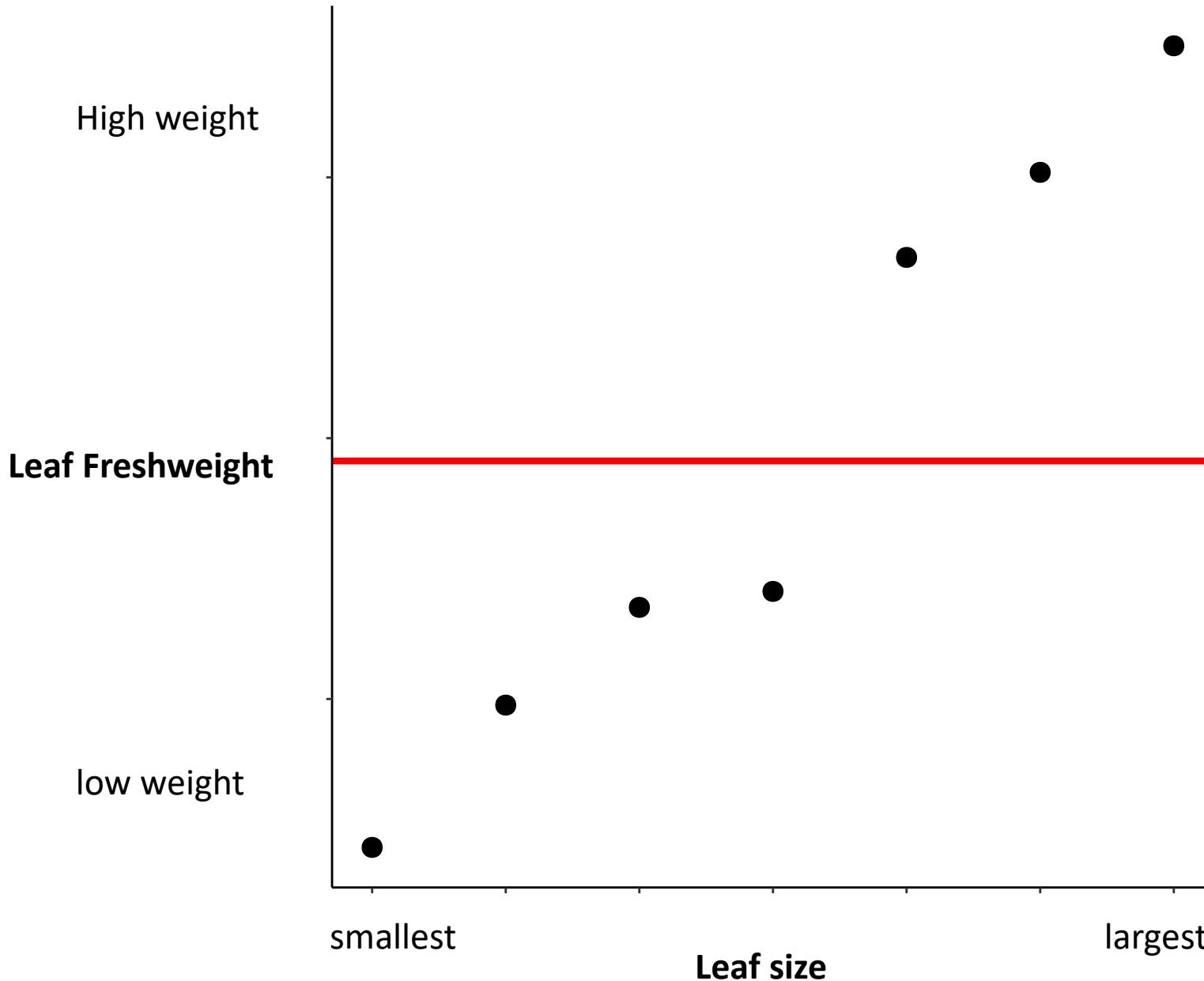
Now, what if, instead of ordering our leafs by their ID#, we ordered them by their size?



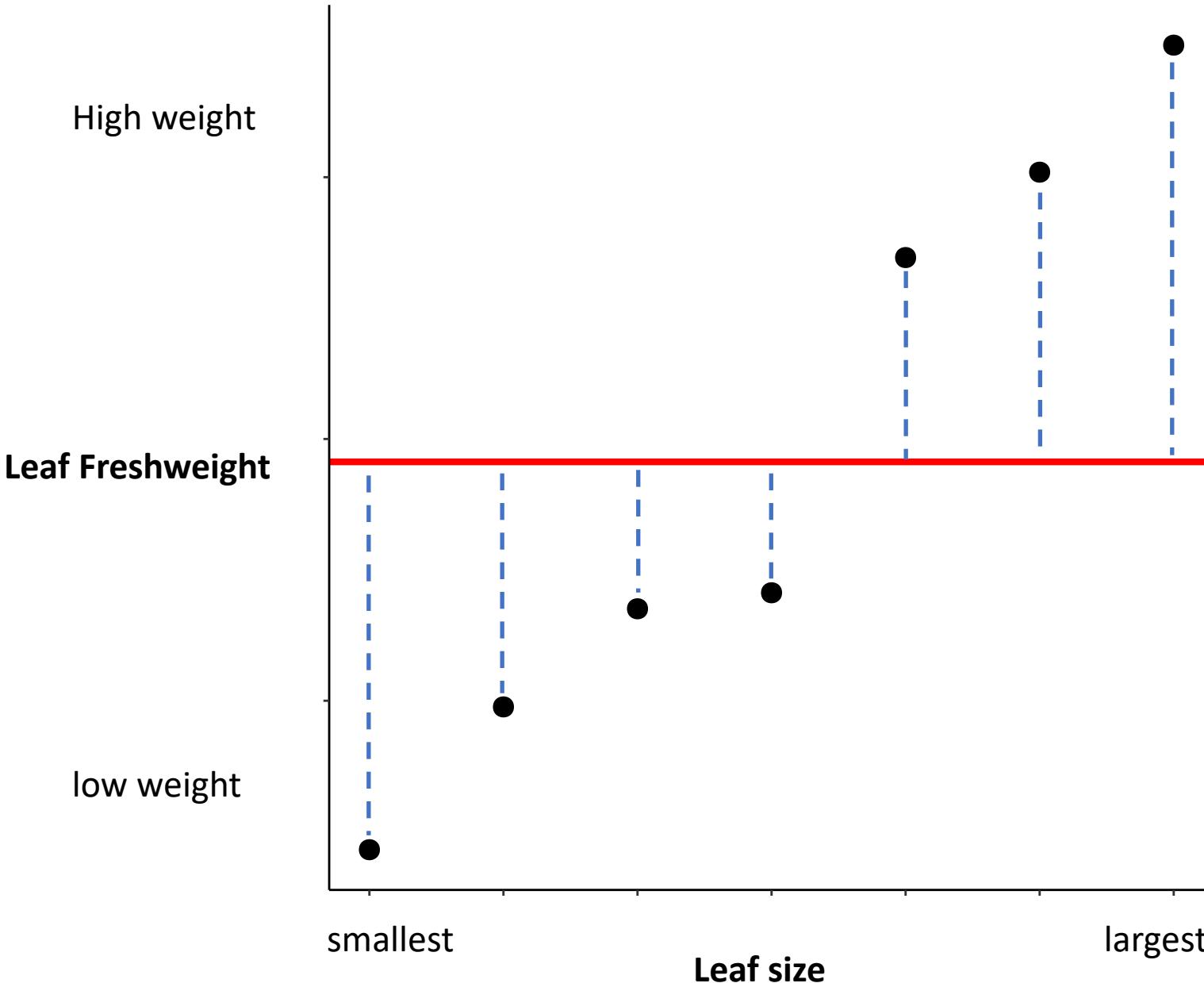
Now, what if, instead of ordering our leafs by their ID#, we ordered them by their size?



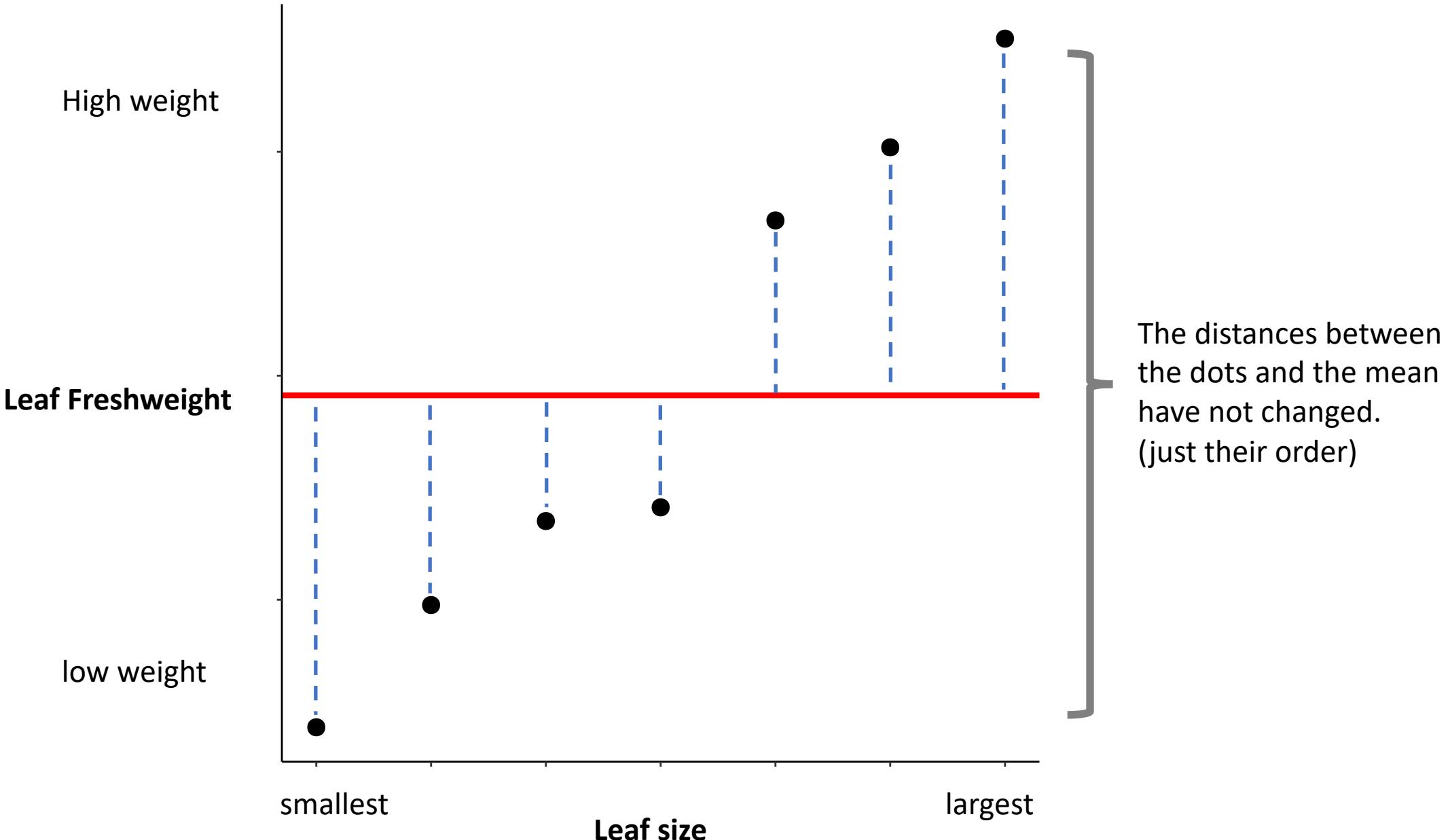
All we have done is reorder the data on the X-axis.  
The mean and variation are the exact same as before.



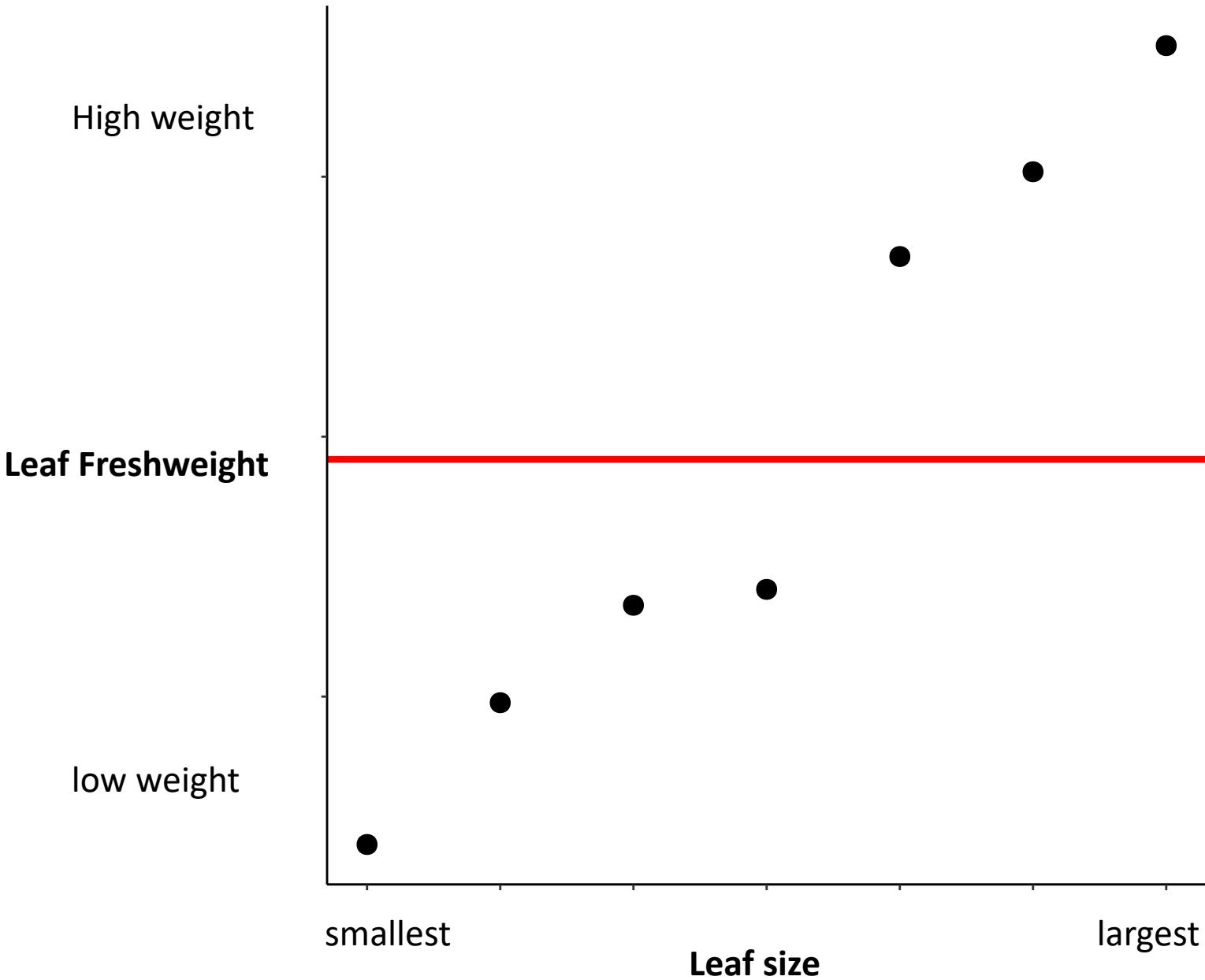
All we have done is reorder the data on the X-axis.  
The mean and variation are the exact same as before.



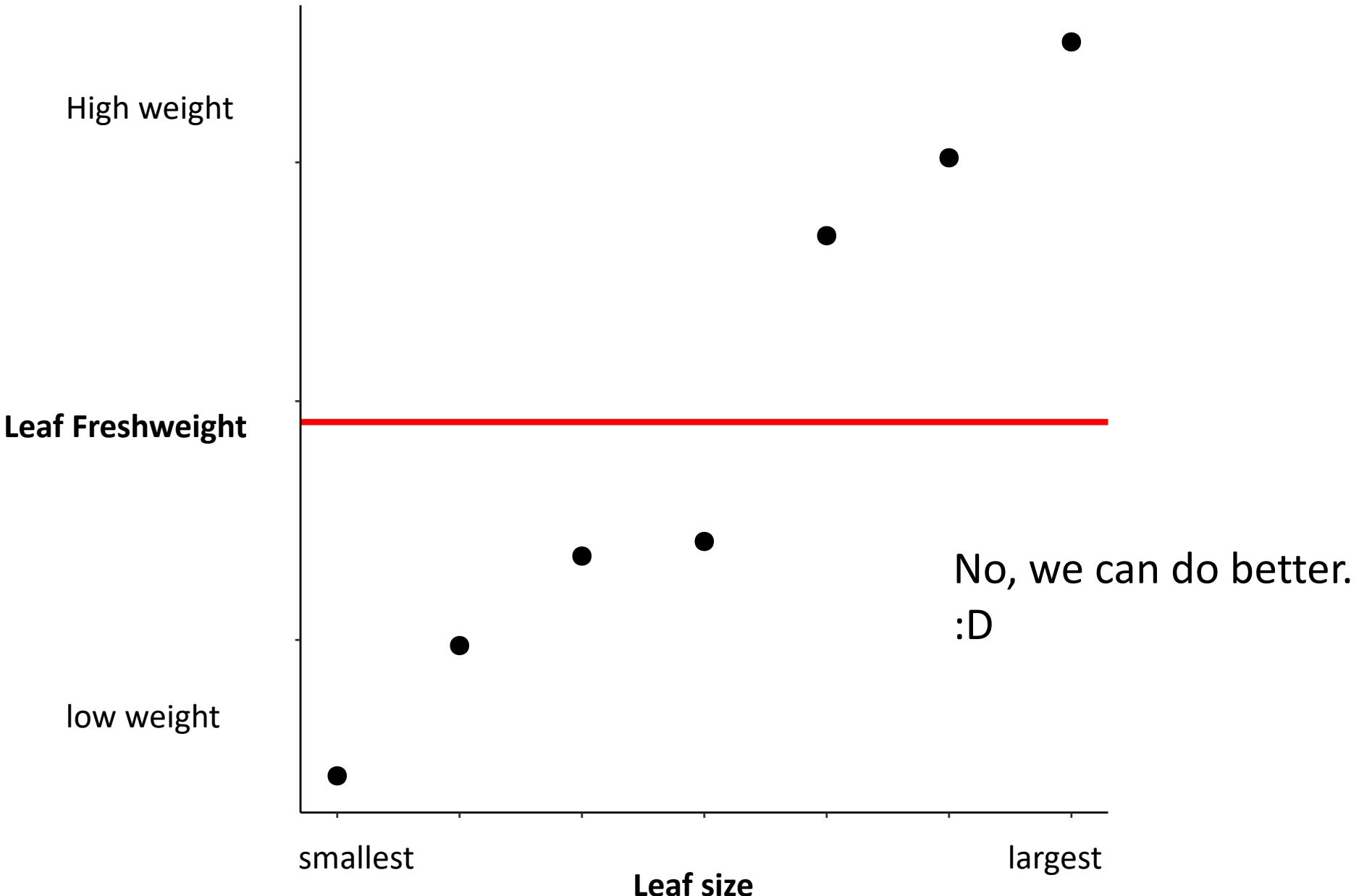
All we have done is reorder the data on the X-axis.  
The mean and variation are the exact same as before.



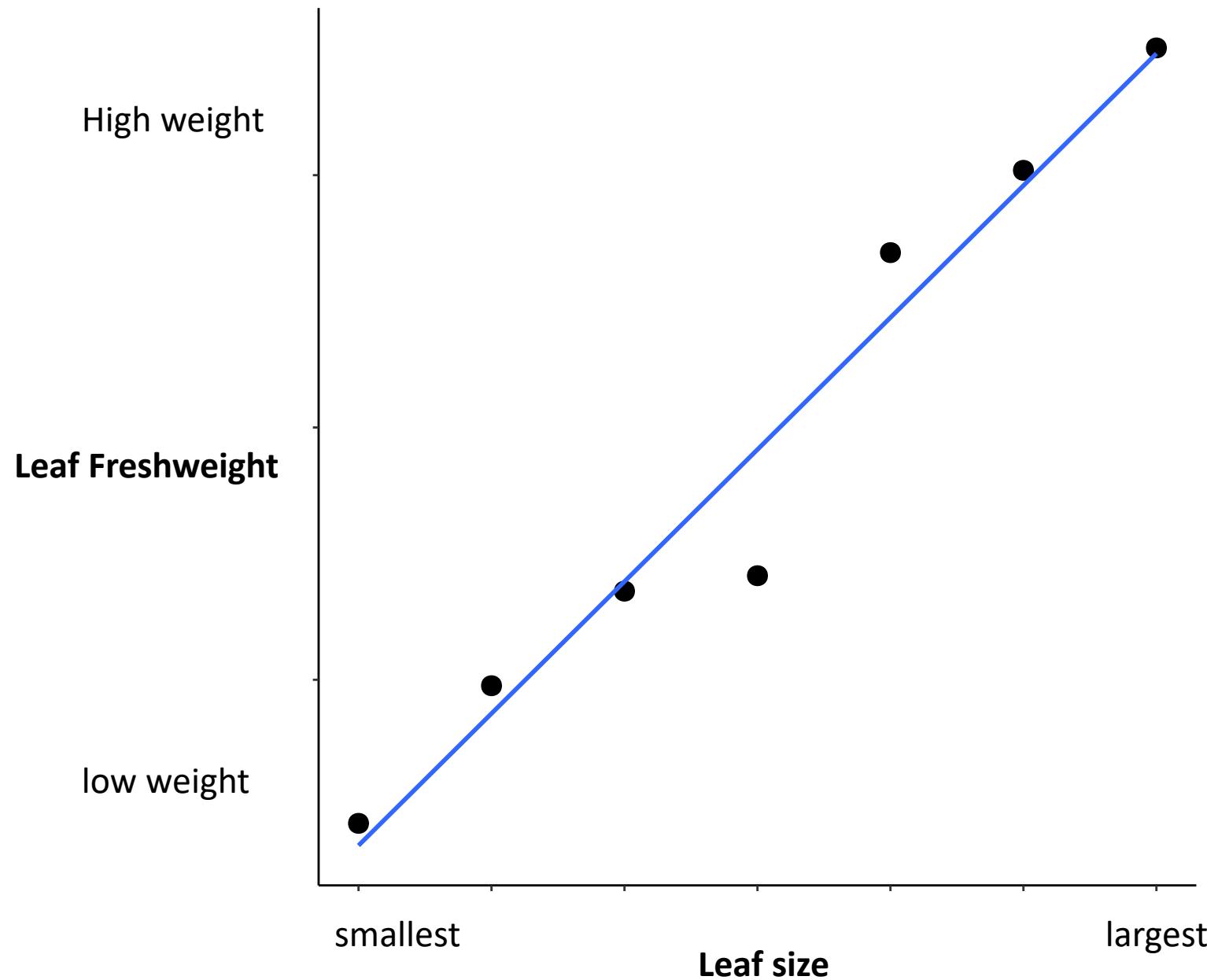
Question: Given that we know the leaf sizes,  
Is the mean weight the best way to predict leaf weight?

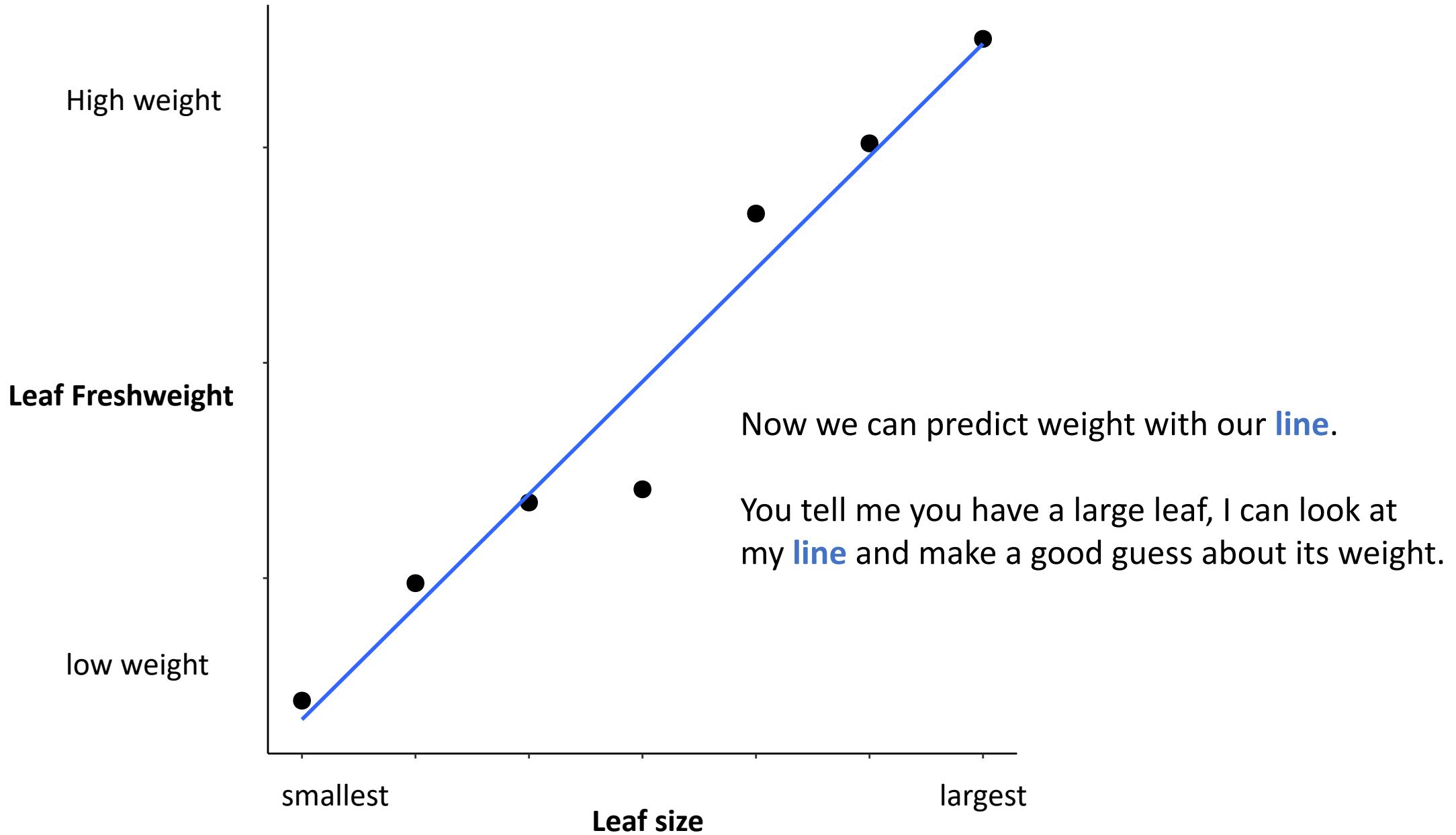


Question: Given that we know the leaf sizes,  
Is the mean weight the best way to predict leaf weight?

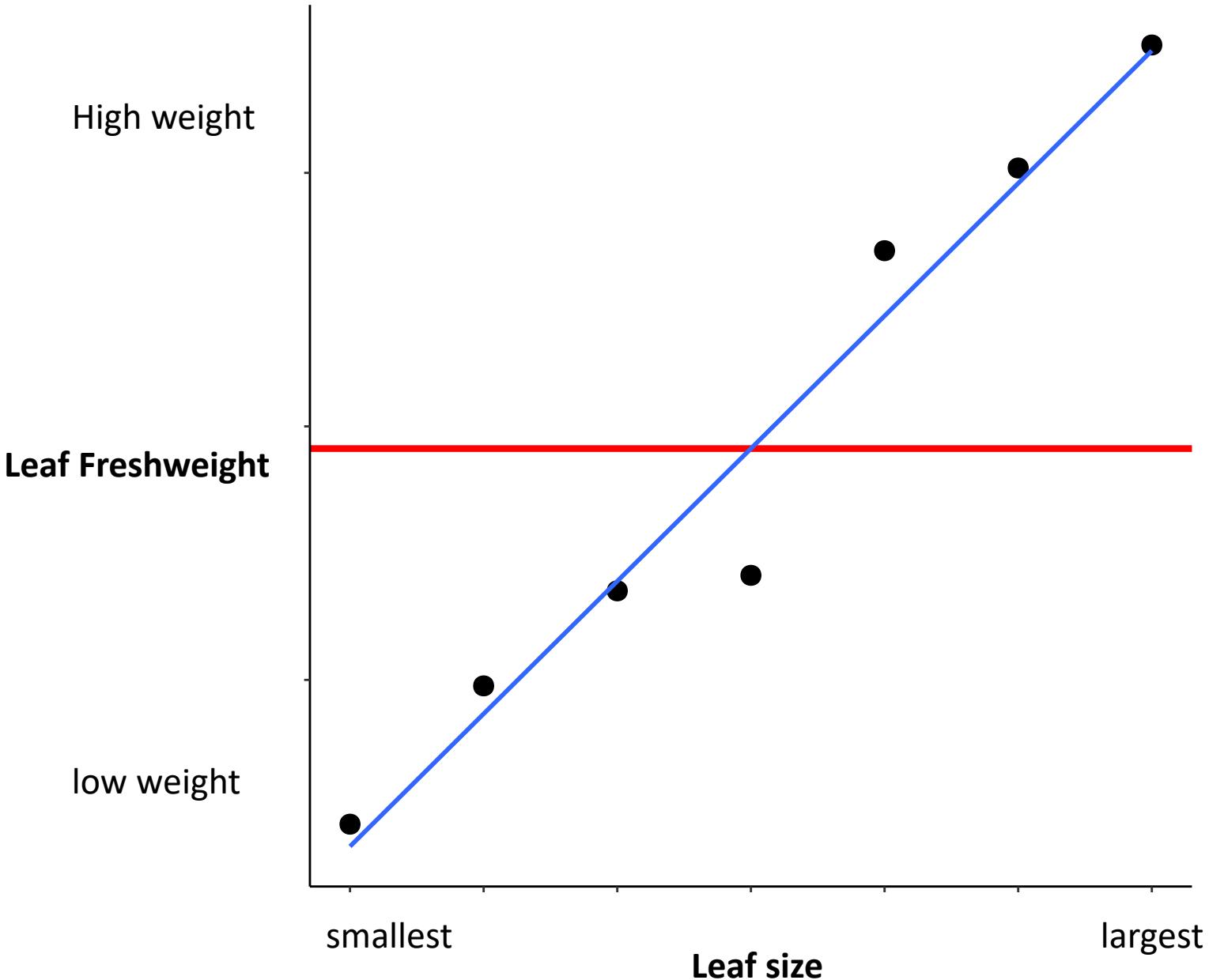


Fit a line to the data...

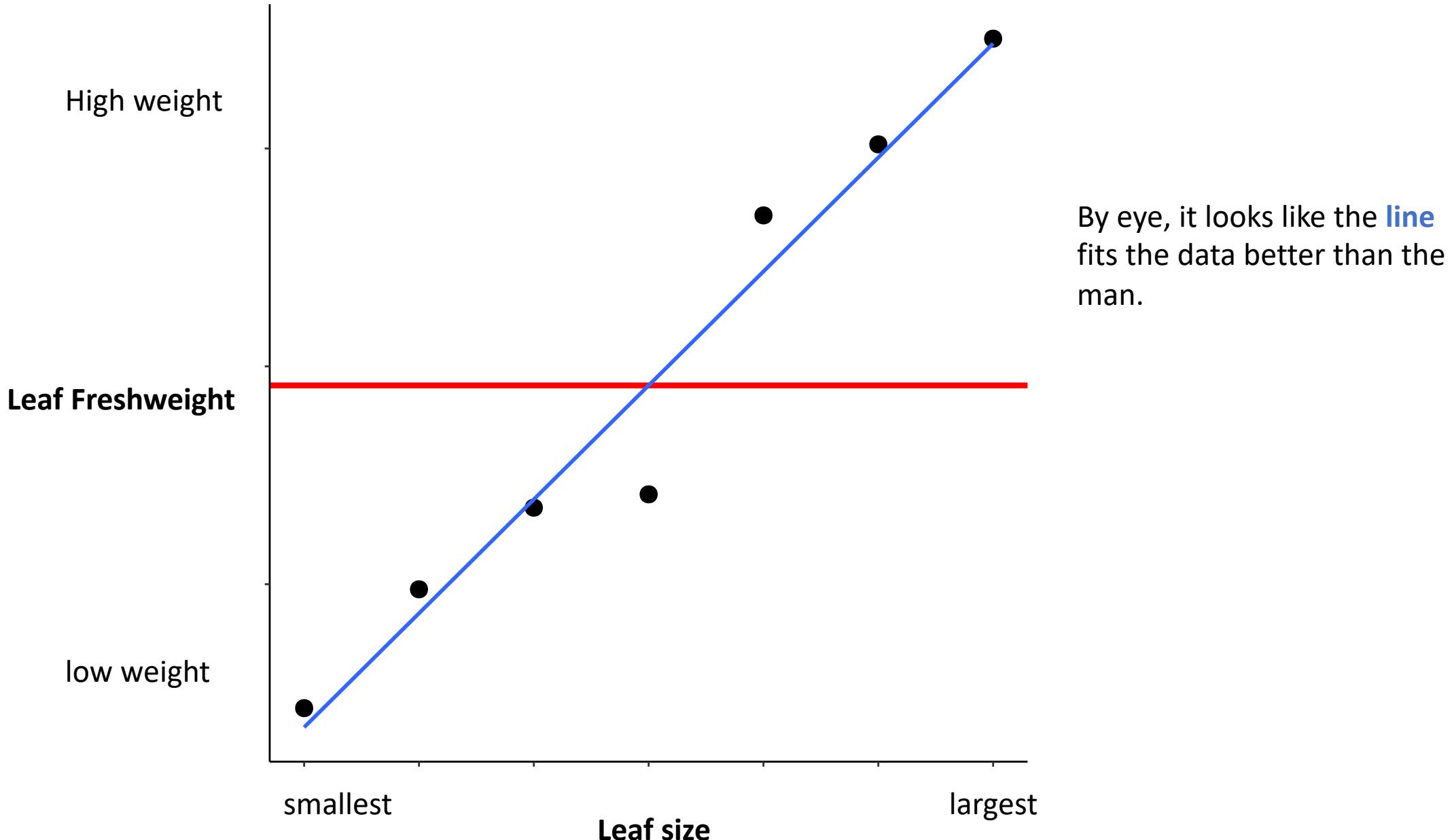




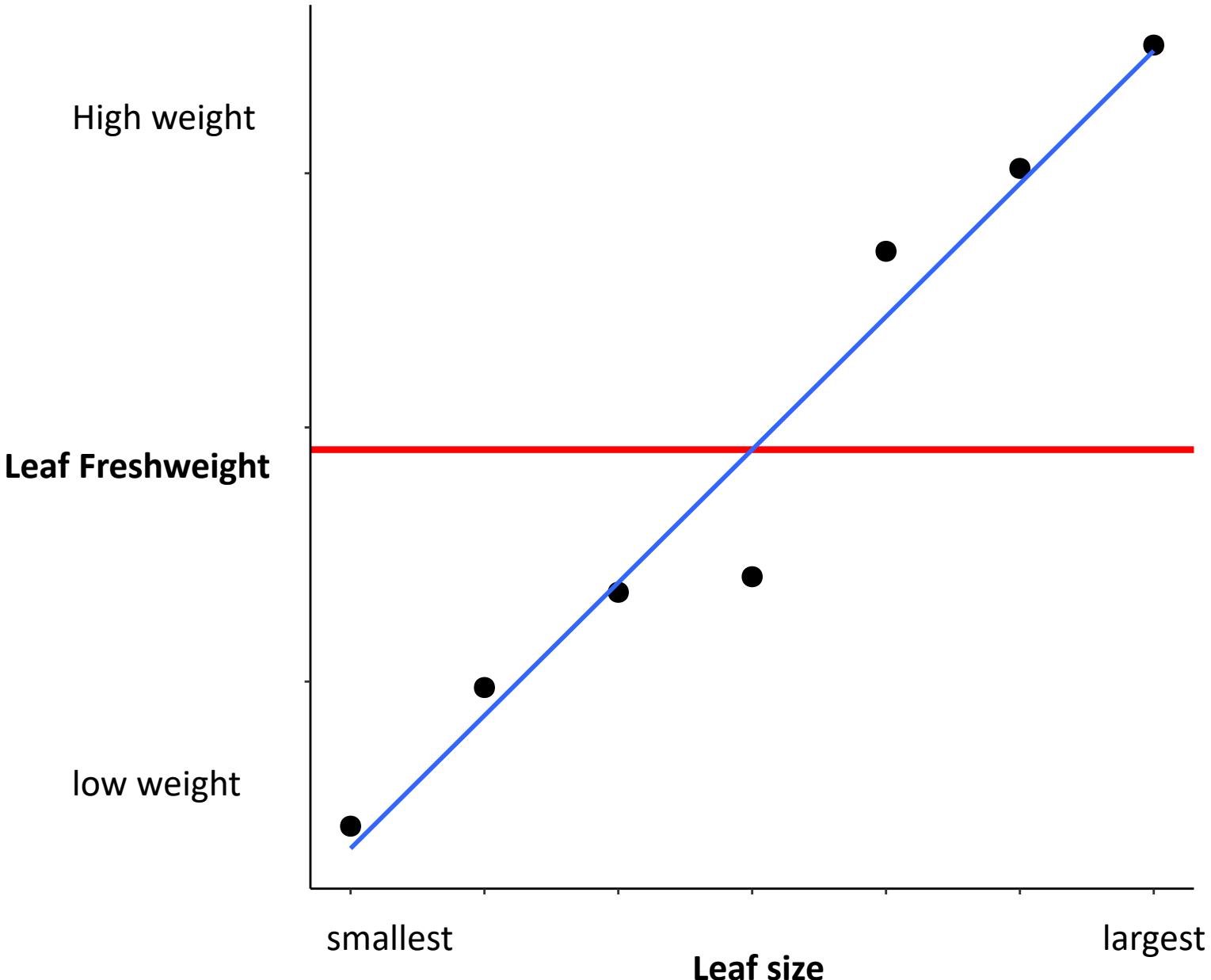
Question: Does the **line** fit the data better than the **mean**?  
If so, how much better?



Question: Does the **line** fit the data better than the **mean**?  
If so, how much better?



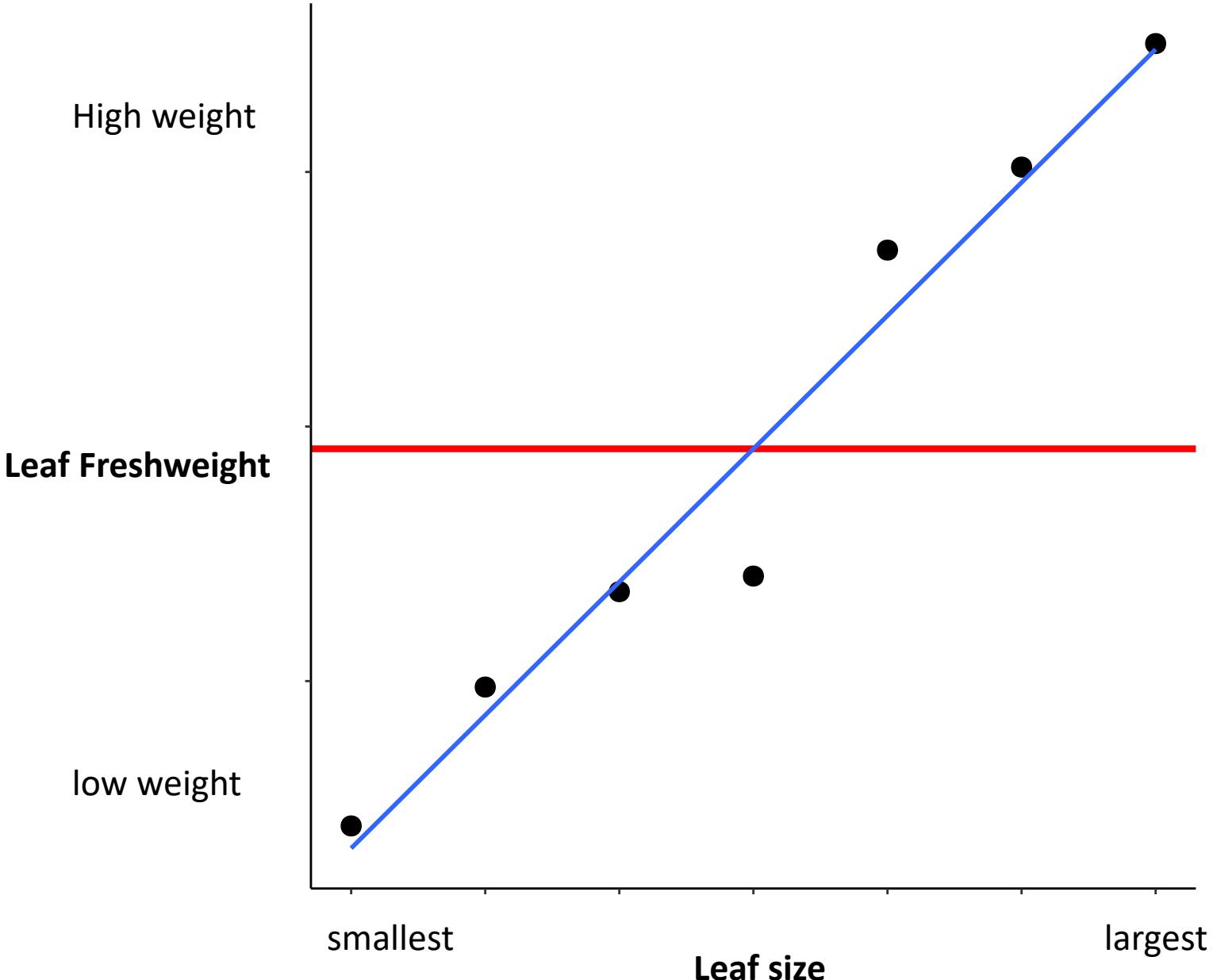
Question: Does the **line** fit the data better than the **mean**?  
If so, how much better?



By eye, it looks like the **line** fits the data better than the **mean**.

How do we quantify that difference?

Question: Does the **line** fit the data better than the **mean**?  
If so, how much better?

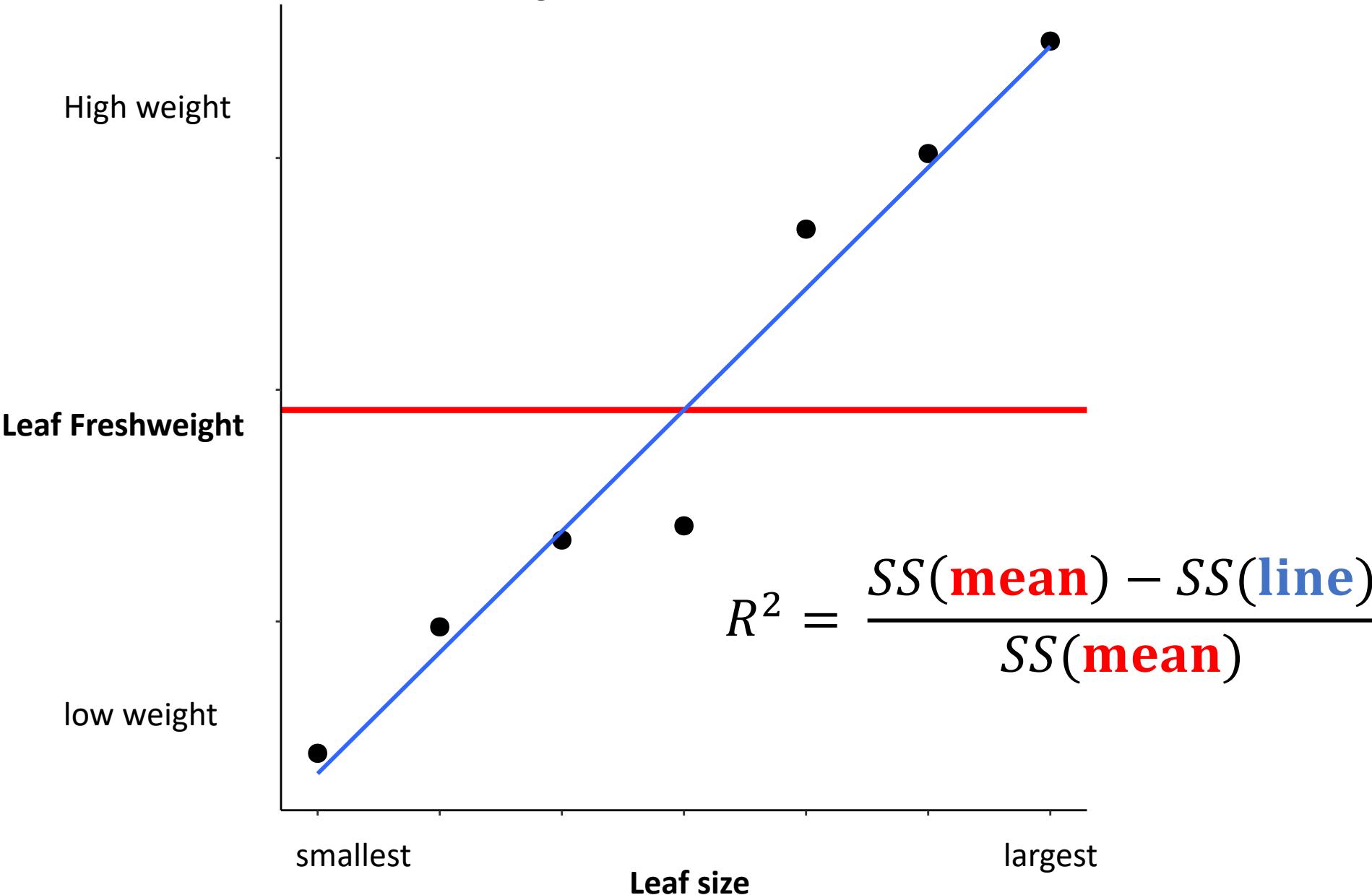


By eye, it looks like the **line** fits the data better than the **mean**.

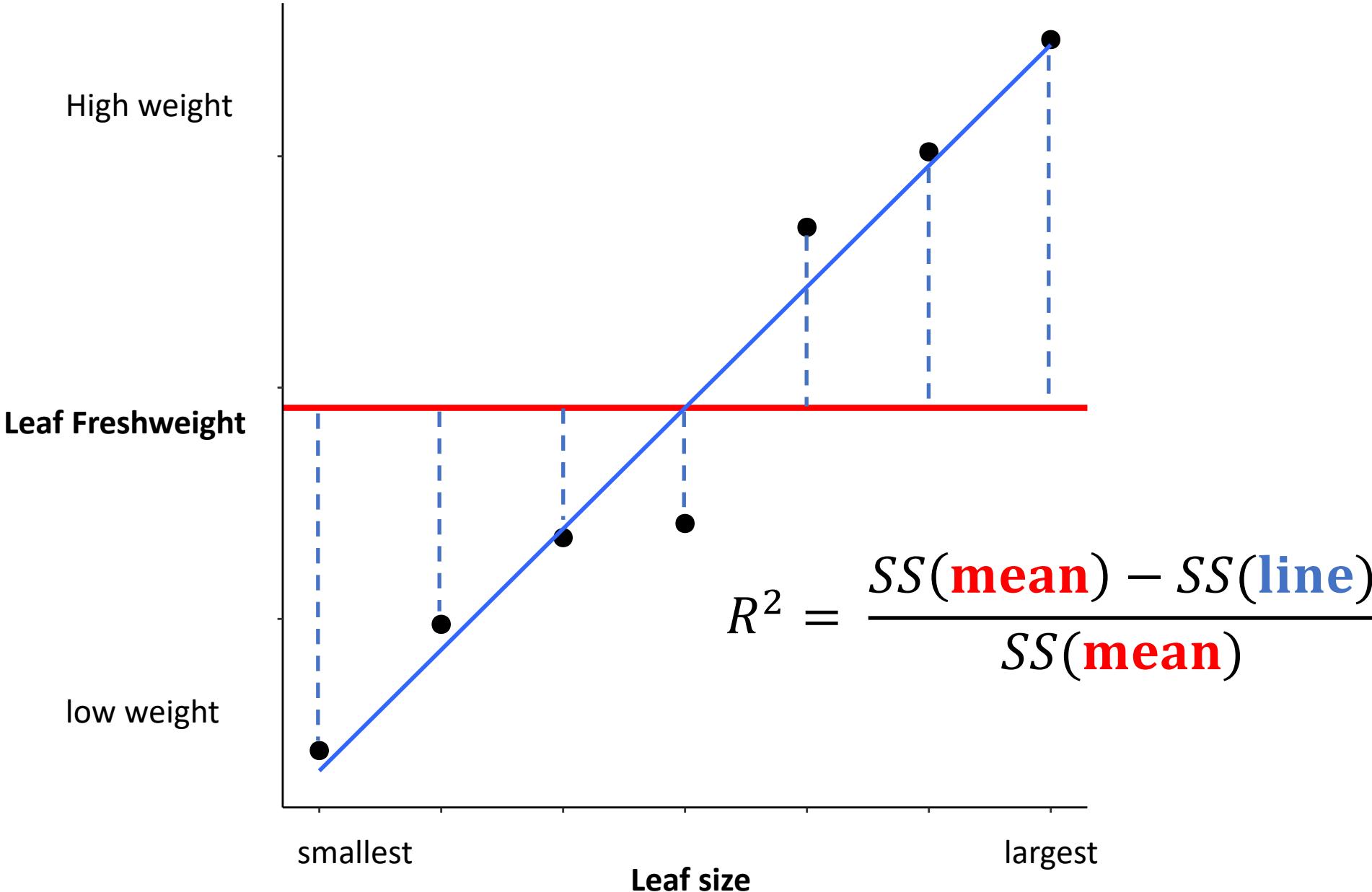
How do we quantify that difference?

$R^2$

Quantifying the difference between the **line** and the mean.  
i.e. Calculating  $R^2$



Quantifying the difference between the **line** and the mean.  
i.e. Calculating  $R^2$



## Quantifying the difference between the **line** and the **mean**. i.e. Calculating $R^2$

Note: the Sum of Squares get quite large quite fast, so for simplicity's sake, we are going to use the Variance instead.

The Variance is just the Mean of the Sum of Squares.  
i.e.:  $SS/n$

Since this is just a scaling factor, the result will be identical.

Leaf Freshweight

low weight

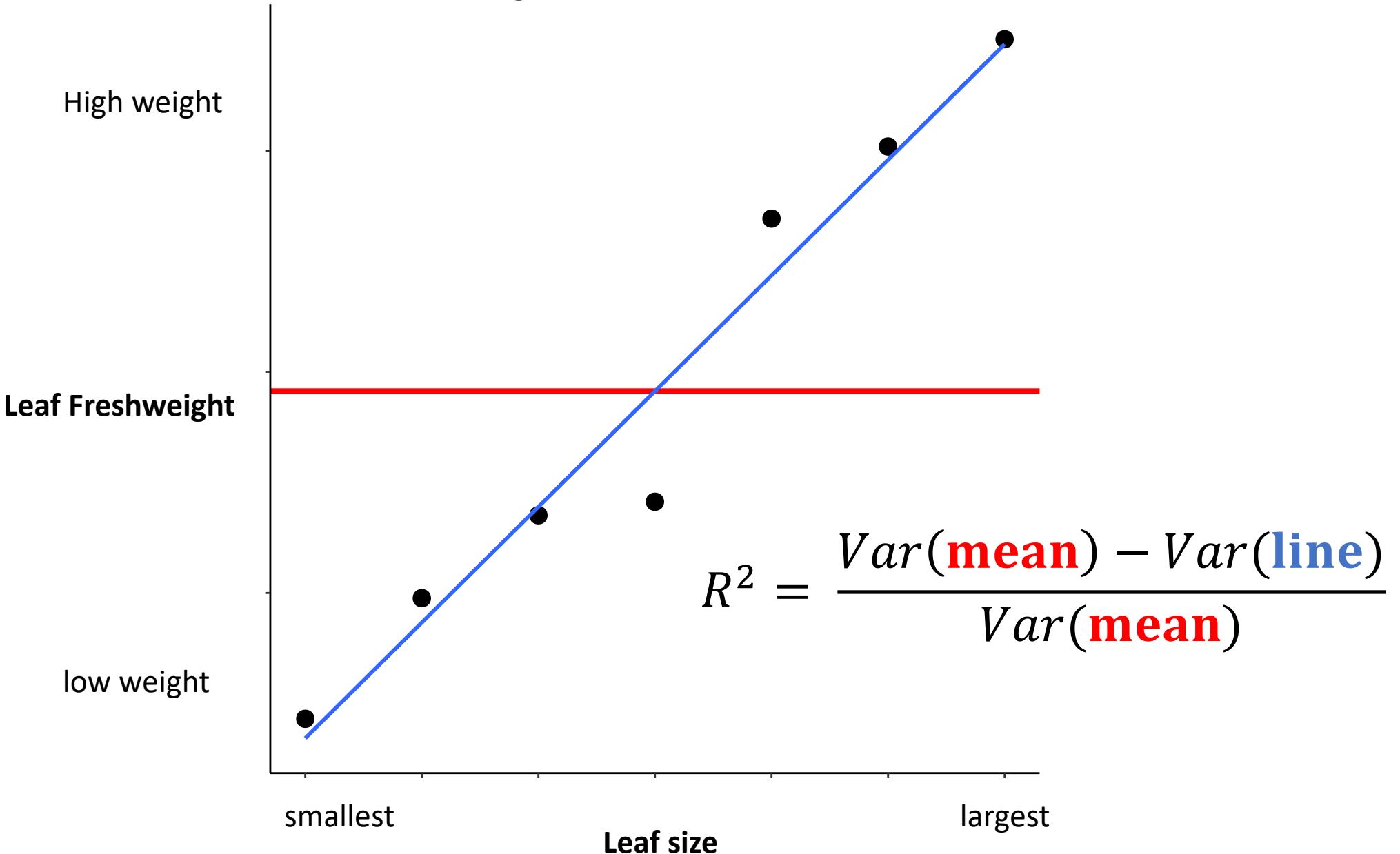
smallest

Leaf size

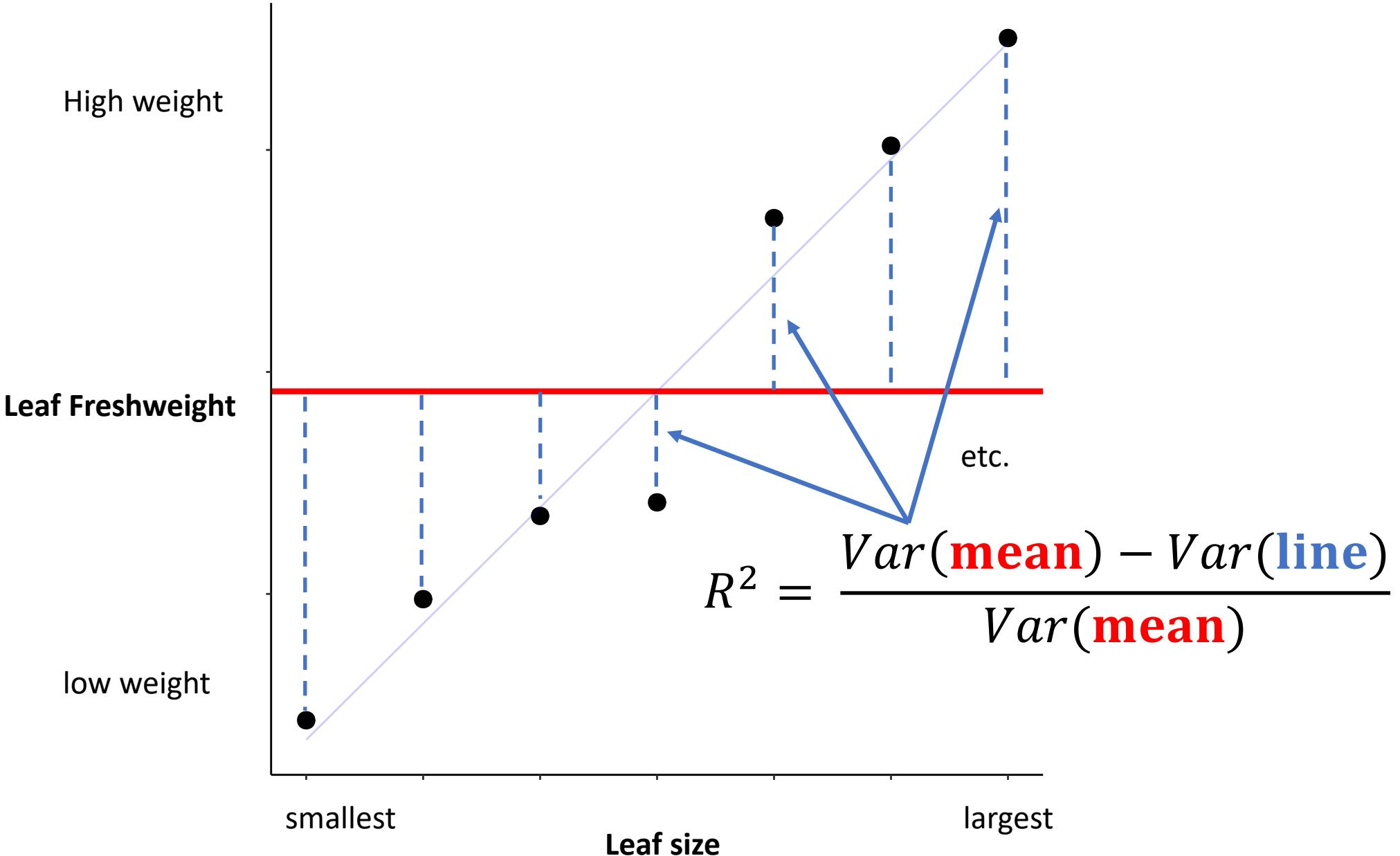
largest

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

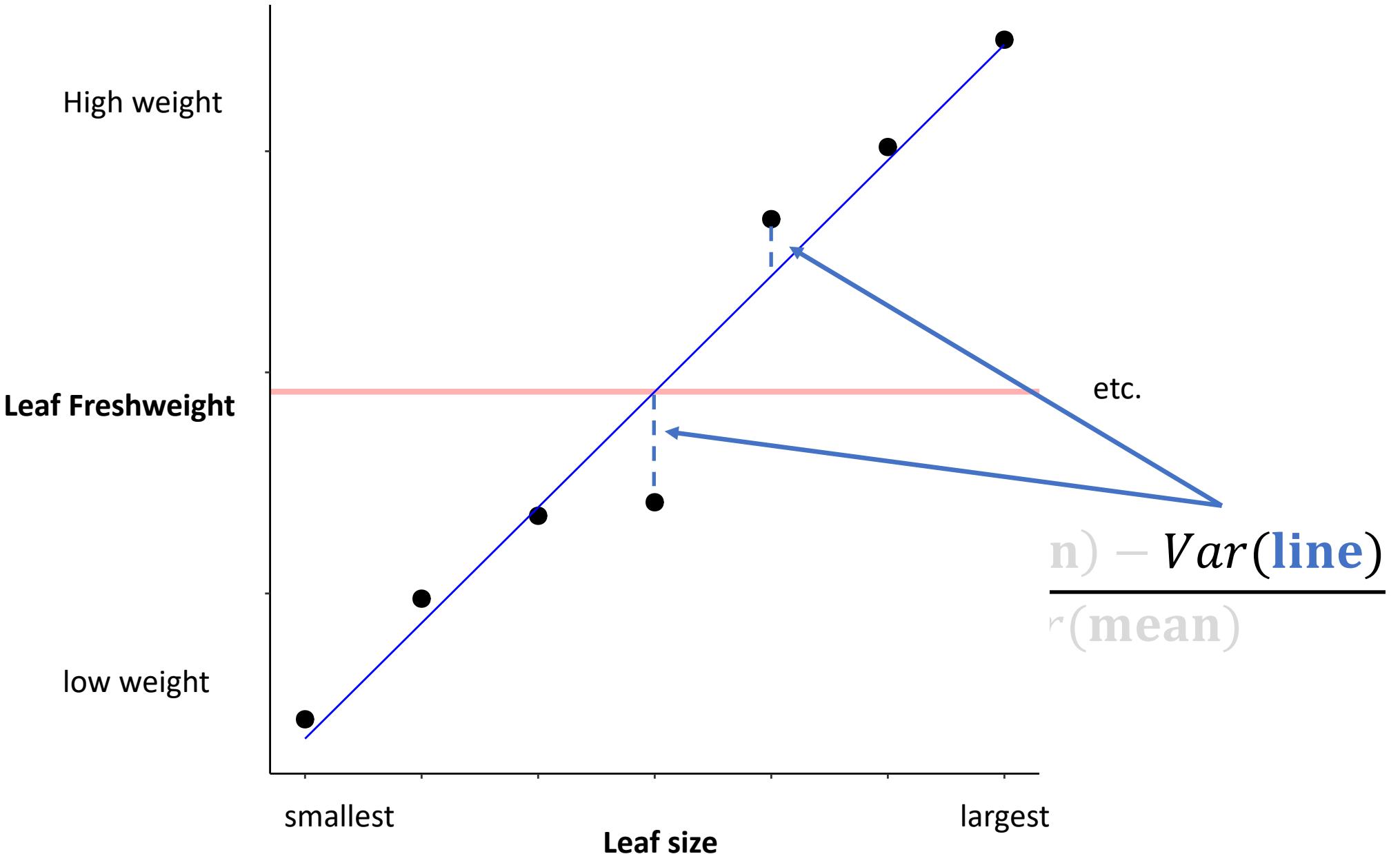
Quantifying the difference between the **line** and the mean.  
i.e. Calculating  $R^2$



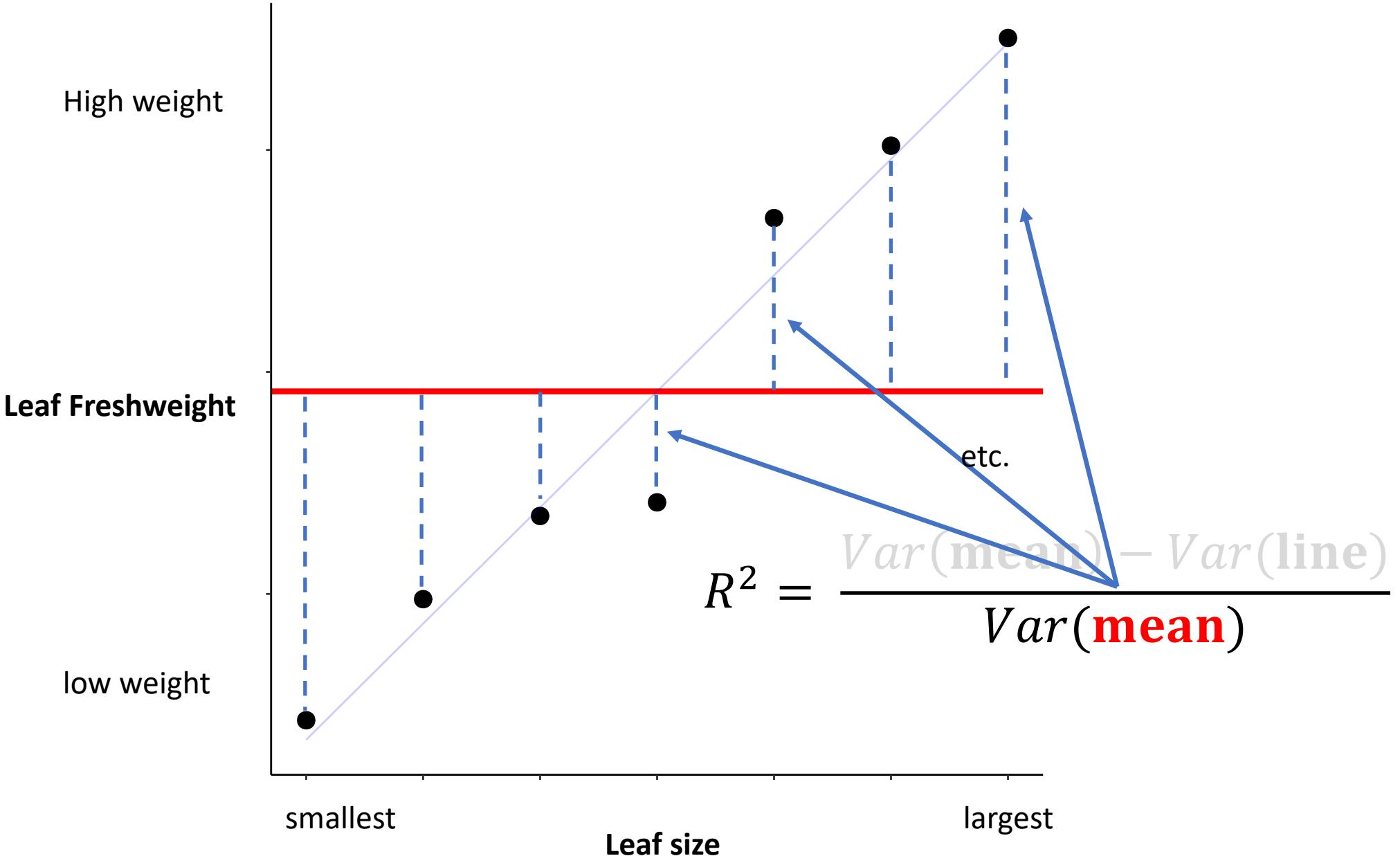
Quantifying the difference between the **line** and the mean.  
i.e. Calculating  $R^2$

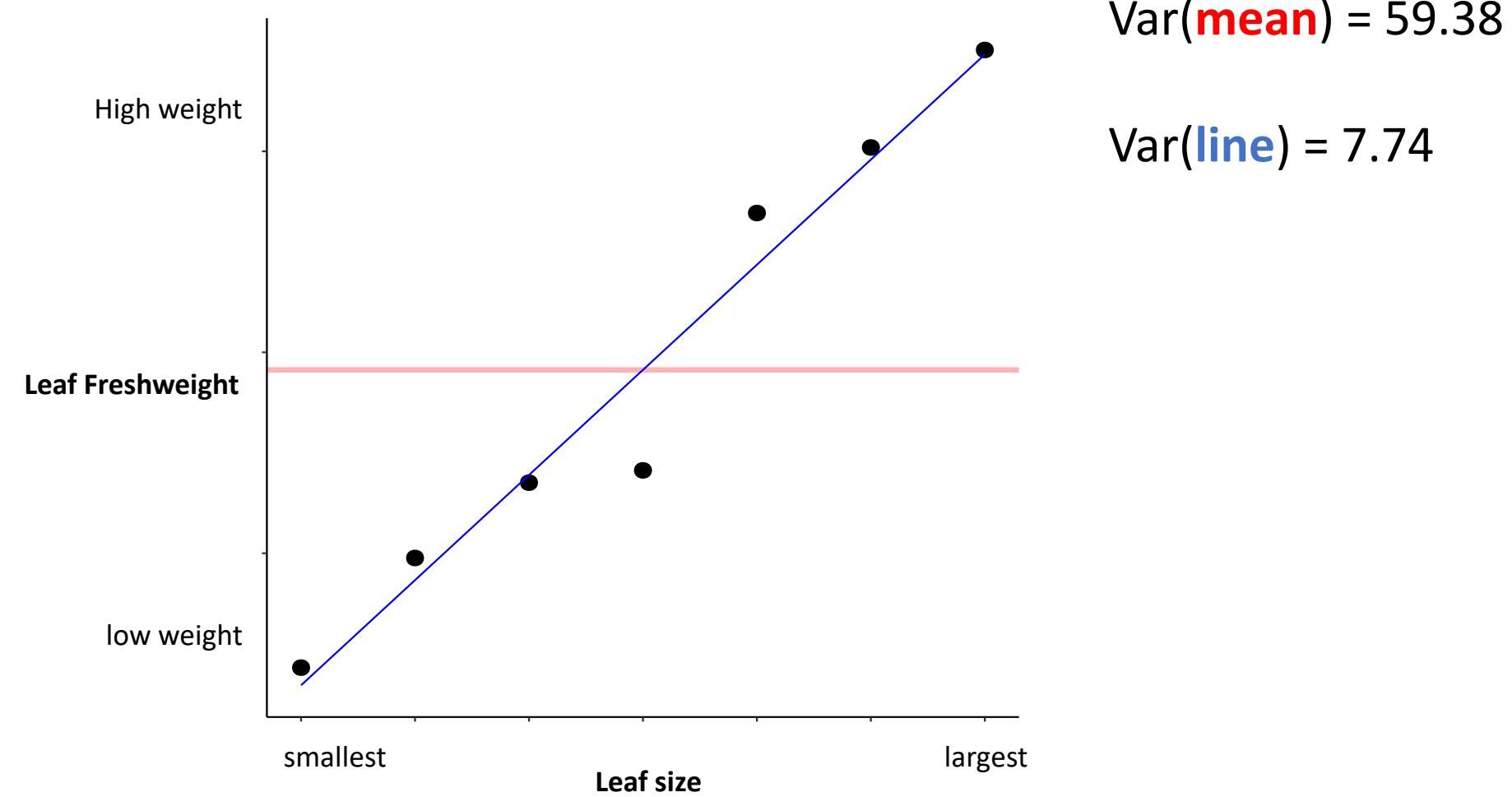


Quantifying the difference between the **line** and the mean.  
i.e. Calculating  $R^2$



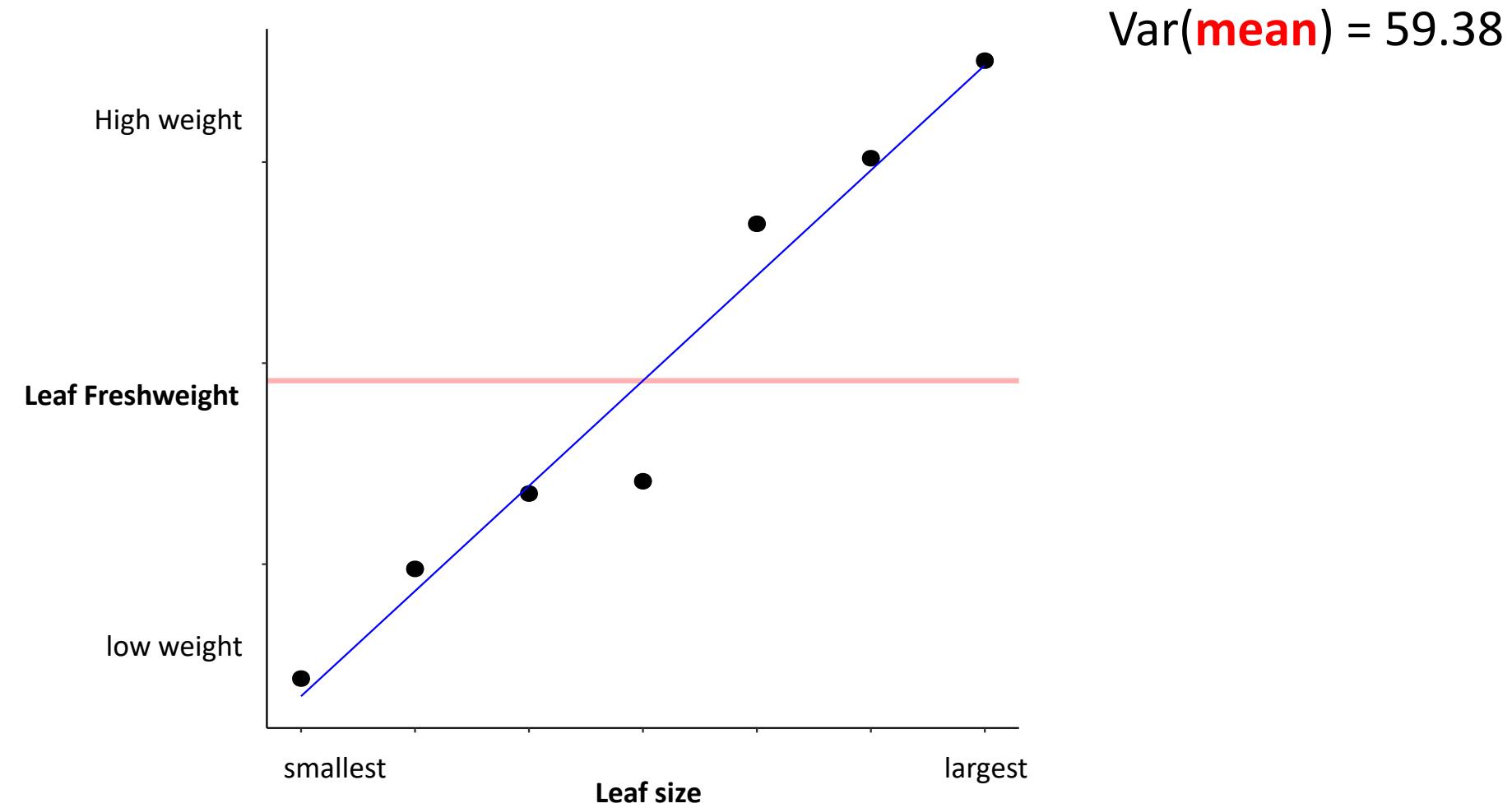
Quantifying the difference between the **line** and the mean.  
i.e. Calculating  $R^2$



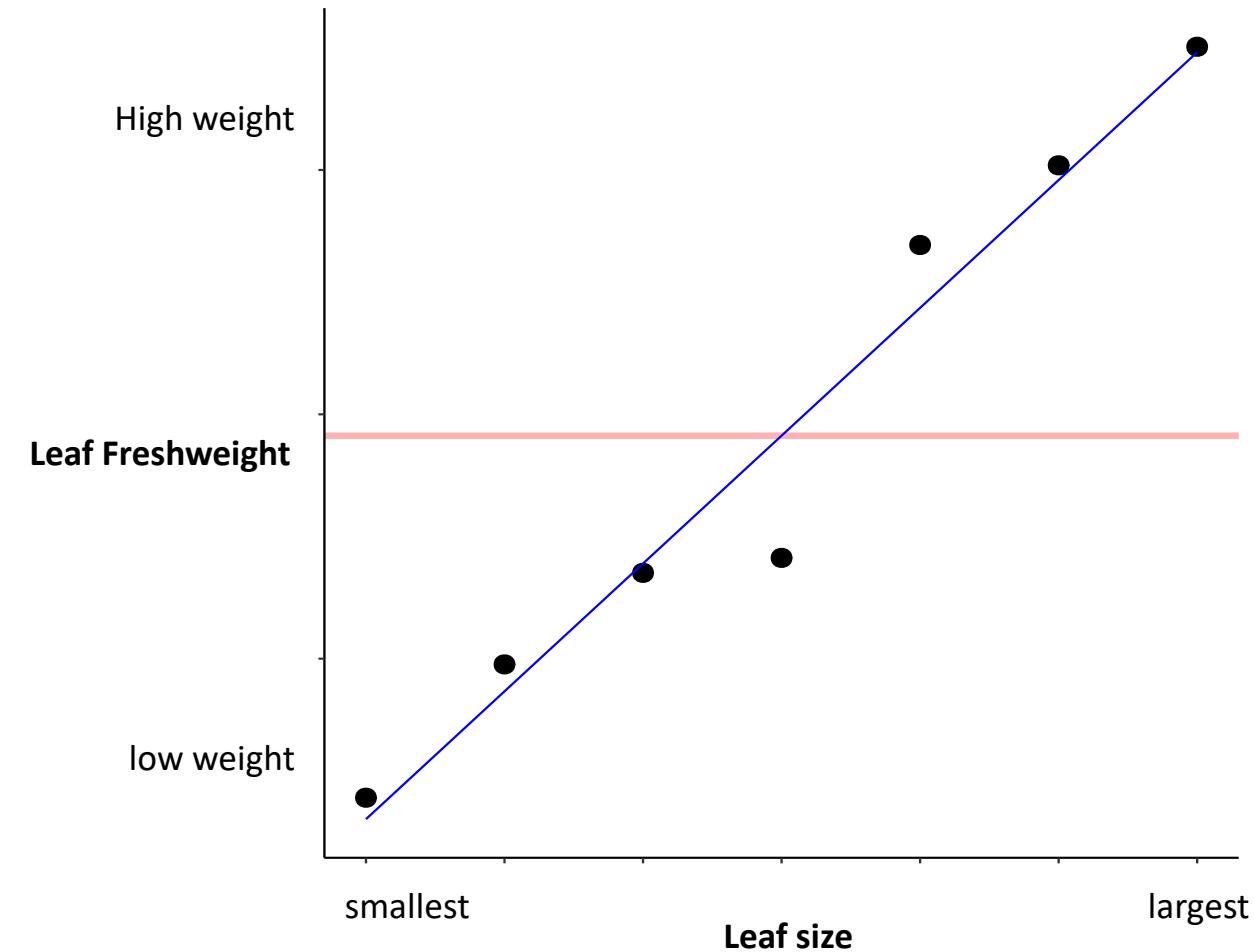


$\text{Var}(\text{mean}) = 59.38$

$\text{Var}(\text{line}) = 7.74$



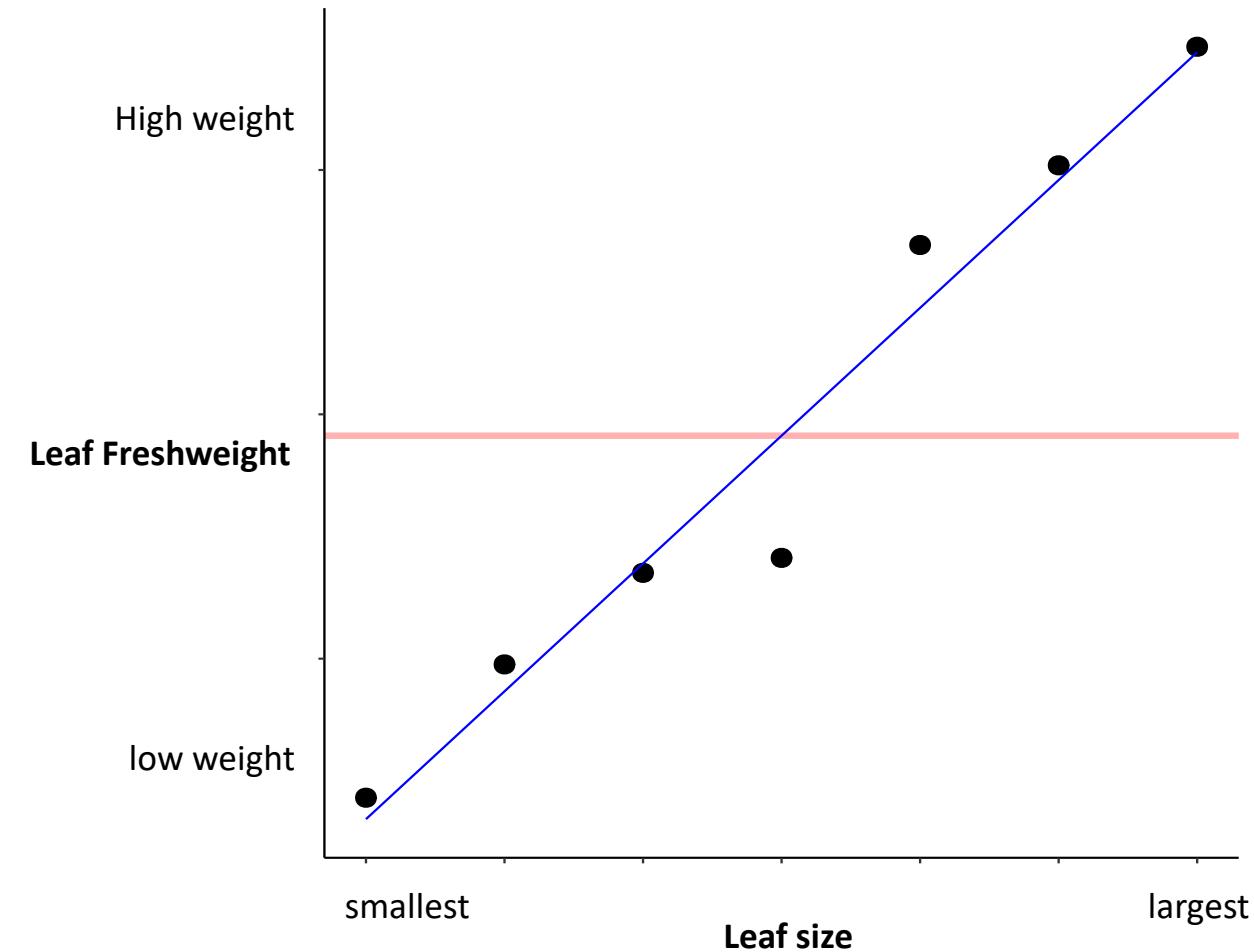
$$\text{Var}(\text{mean}) = 59.38$$



$$\text{Var}(\text{mean}) = 59$$

$$\text{Var}(\text{line}) = 8$$

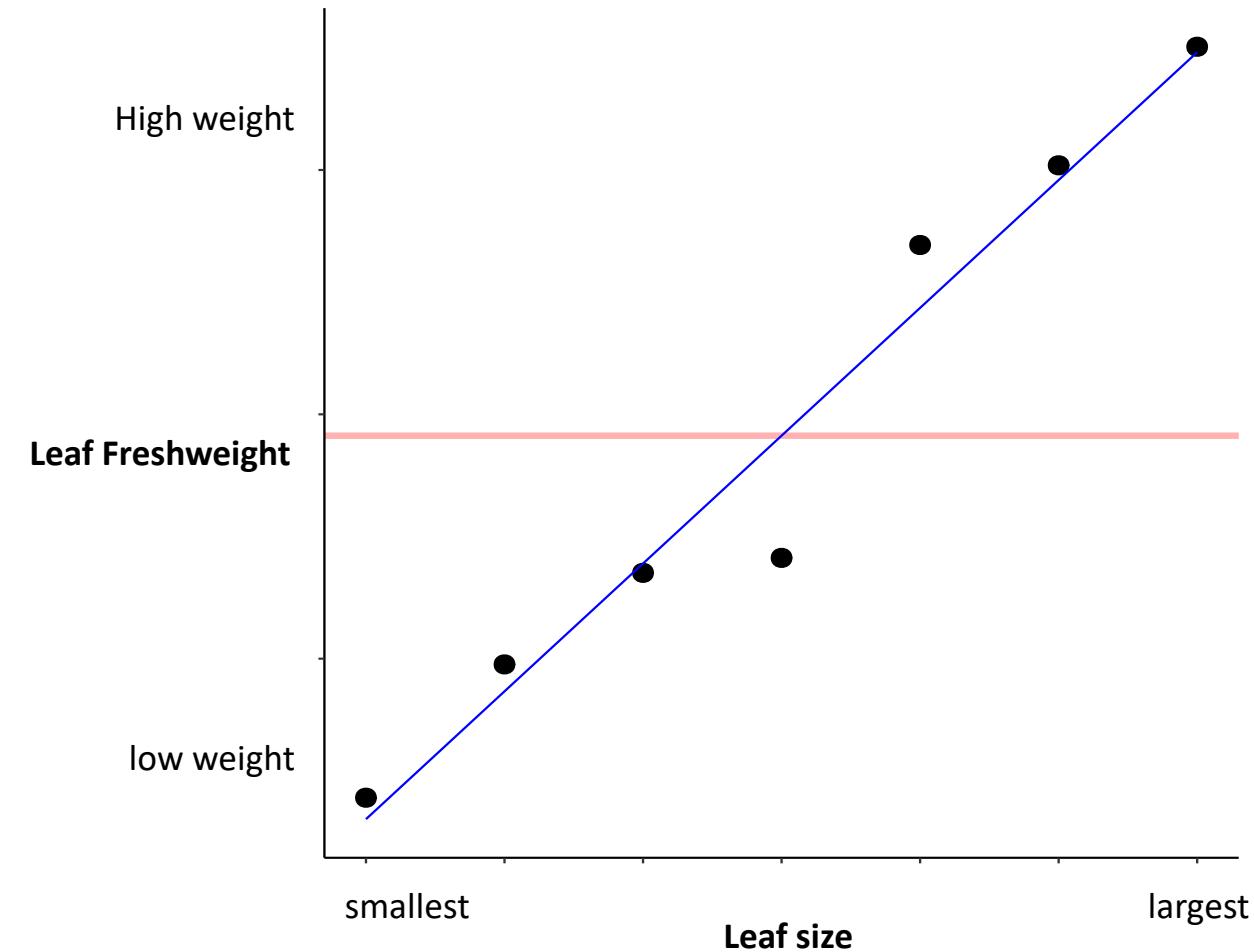
$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{line})}{\text{Var}(\text{mean})}$$



$$\text{Var}(\text{mean}) = 59$$

$$\text{Var}(\text{line}) = 8$$

$$R^2 = \frac{59 - 8}{59}$$

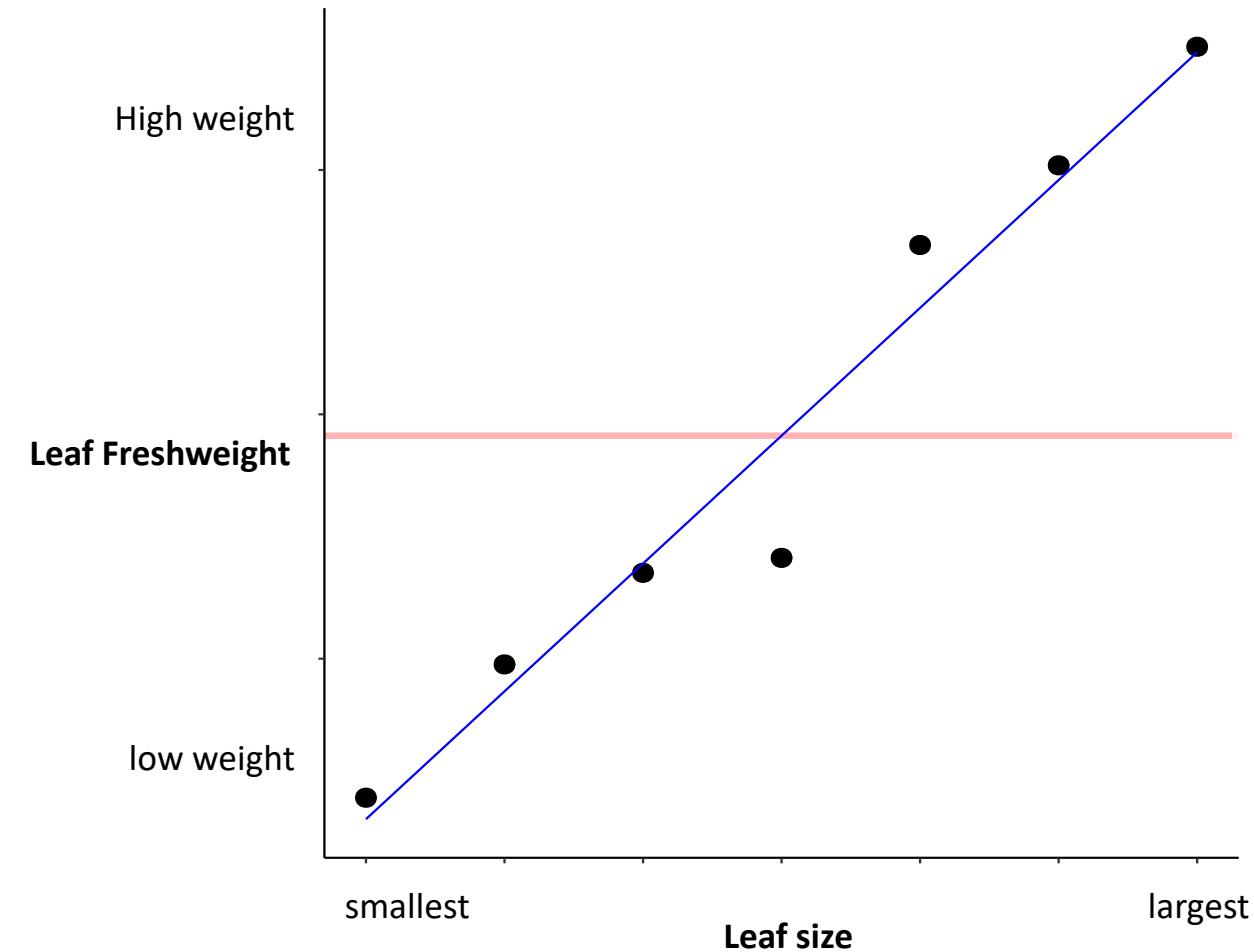


$$\text{Var}(\text{mean}) = 59$$

$$\text{Var}(\text{line}) = 8$$

$$R^2 = \frac{59 - 8}{59}$$

$$R^2 = \frac{51}{59} = 0.86 = 86\%$$



Var(**mean**) = 59

Var(**line**) = 8

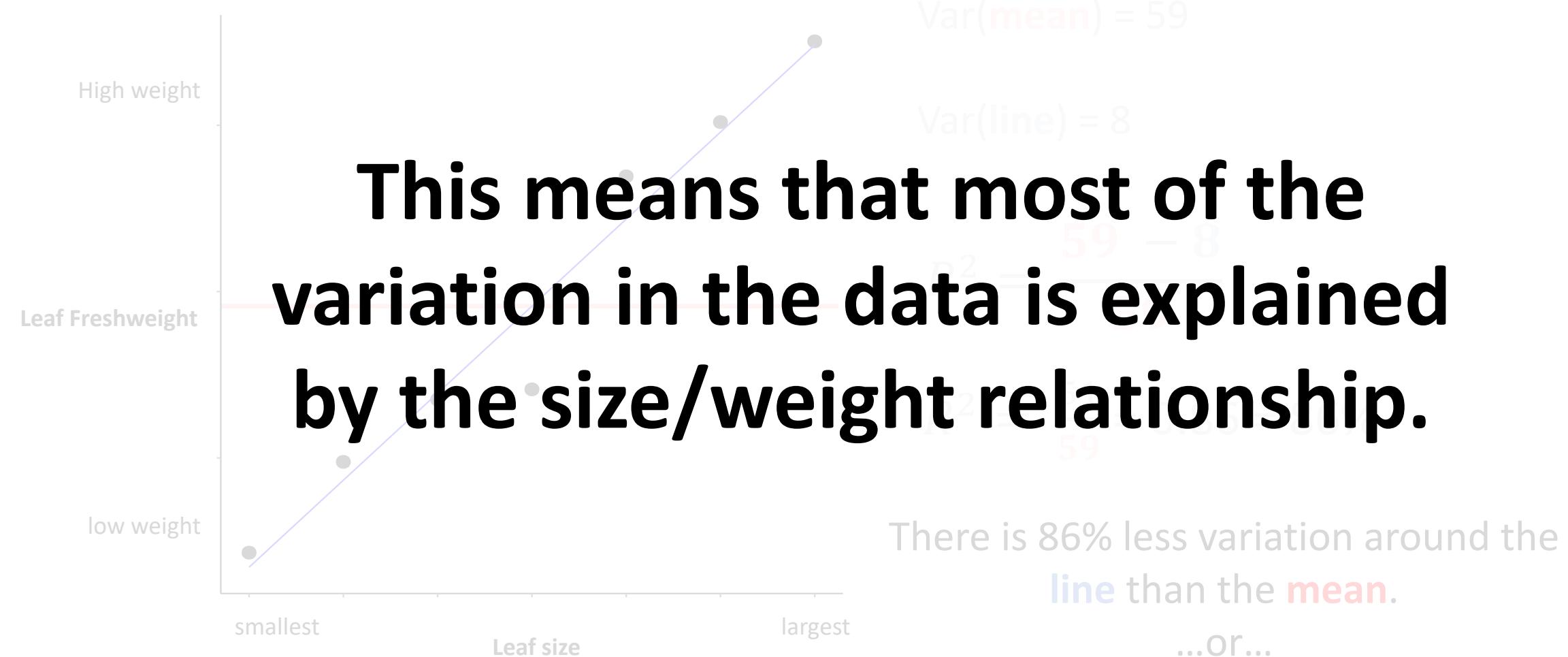
$$R^2 = \frac{59 - 8}{59}$$

$$R^2 = \frac{51}{59} = 0.86 = 86\%$$

There is 86% less variation around the **line** than the **mean**.

...or...

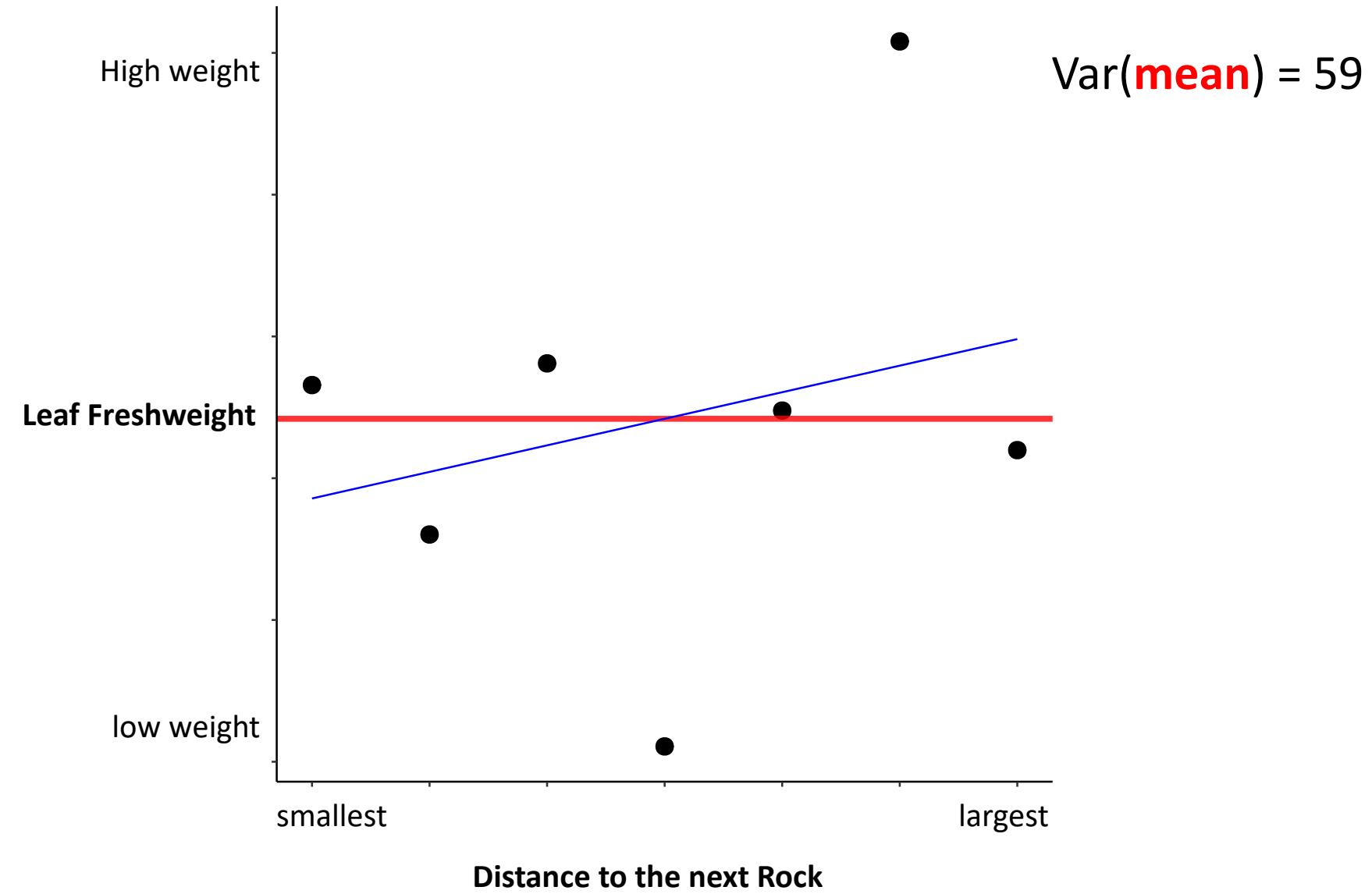
The size/weight relationship accounts for 81% of the variation

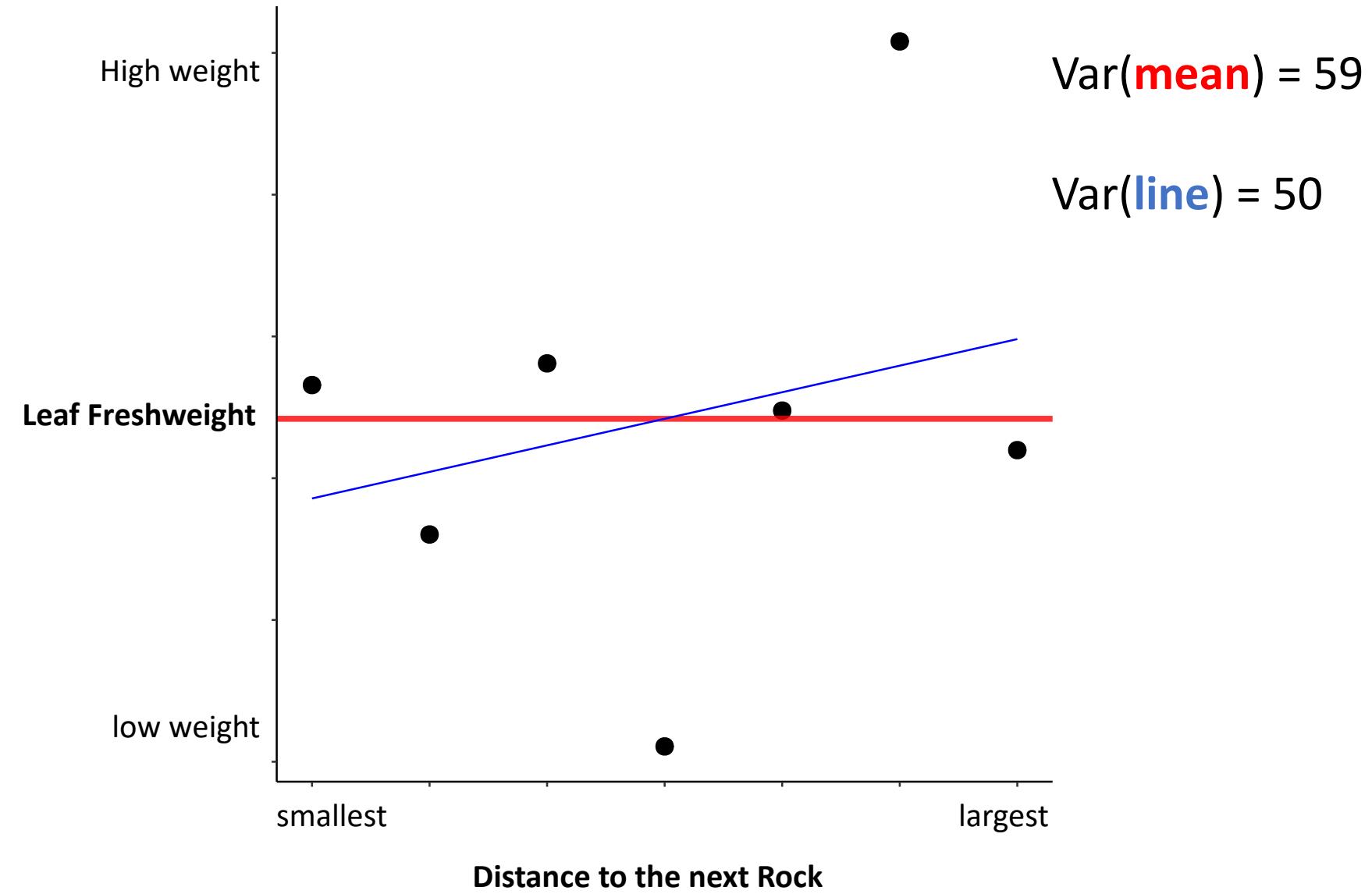


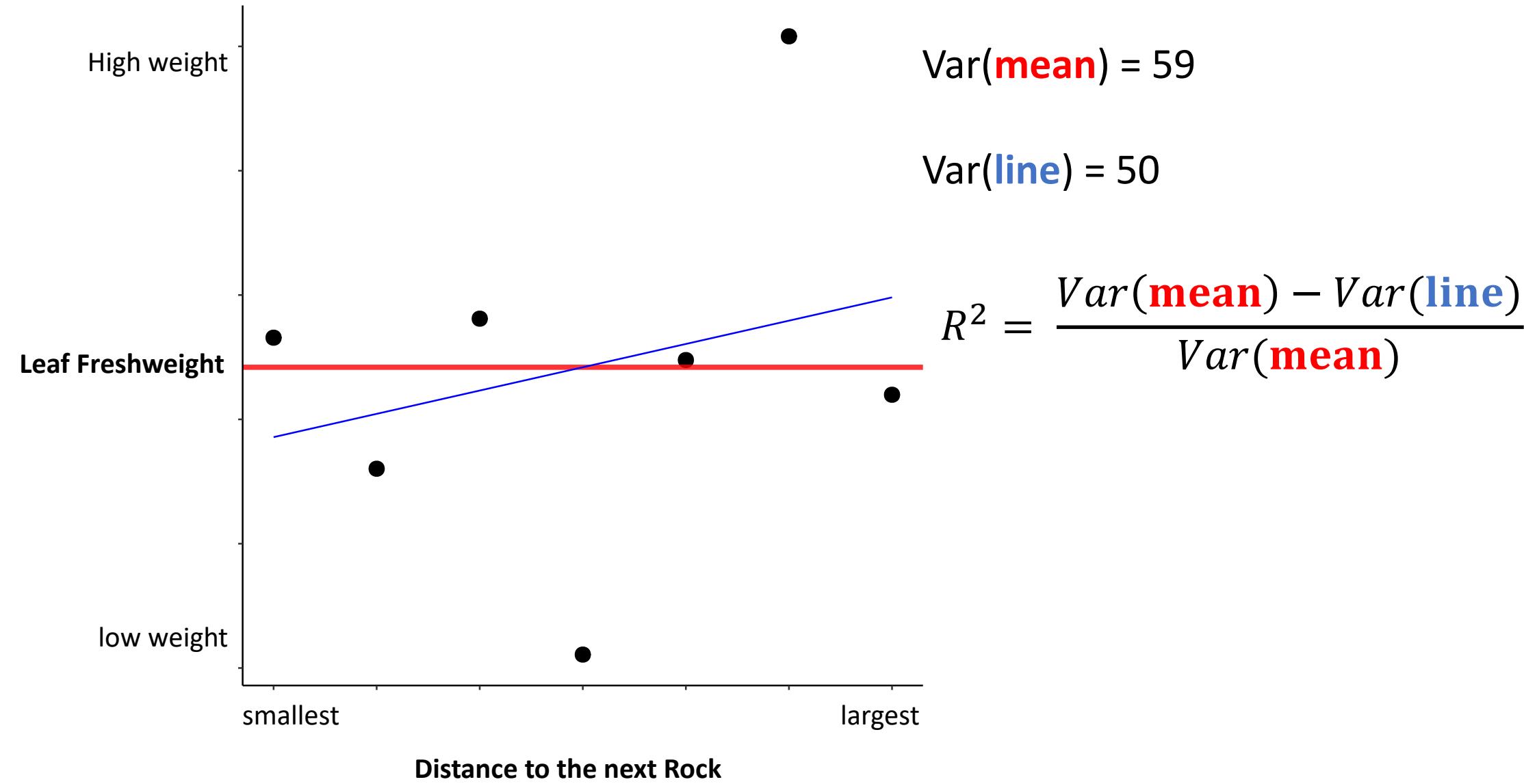
There is 86% less variation around the line than the mean.

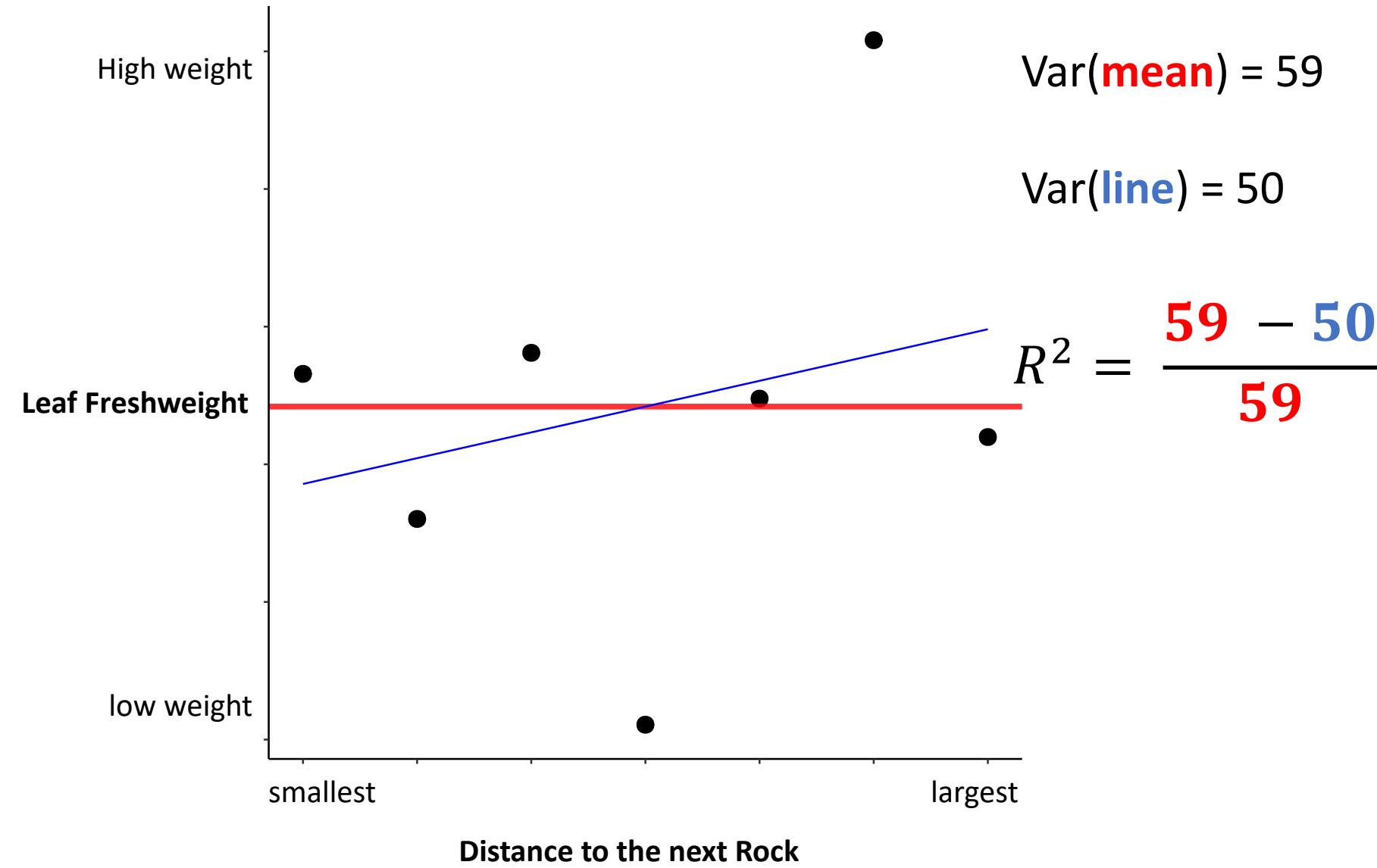
...or...

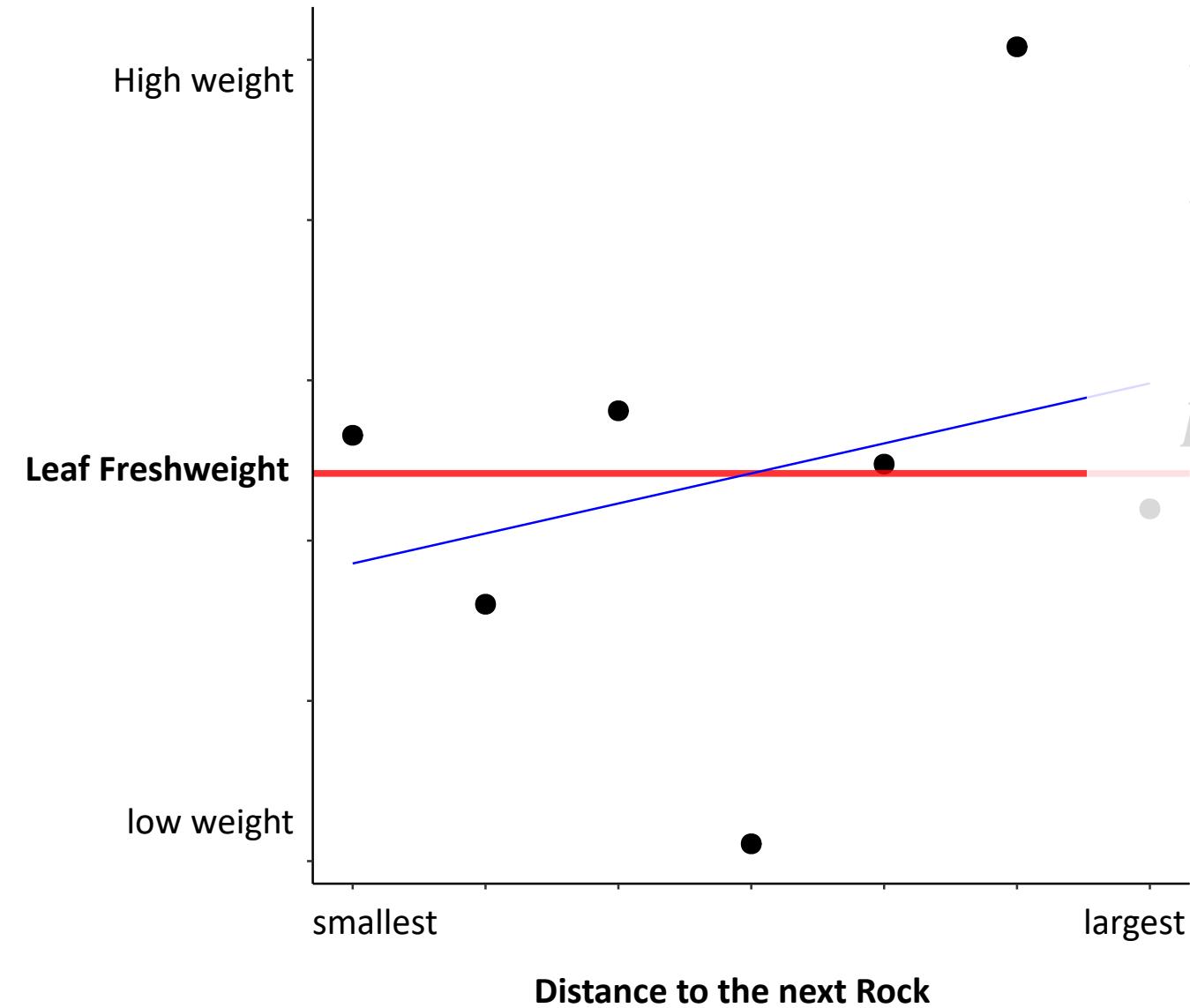
The size/weight relationship accounts for 81% of the variation











Var(**mean**) = 59

Var(**line**) = 50

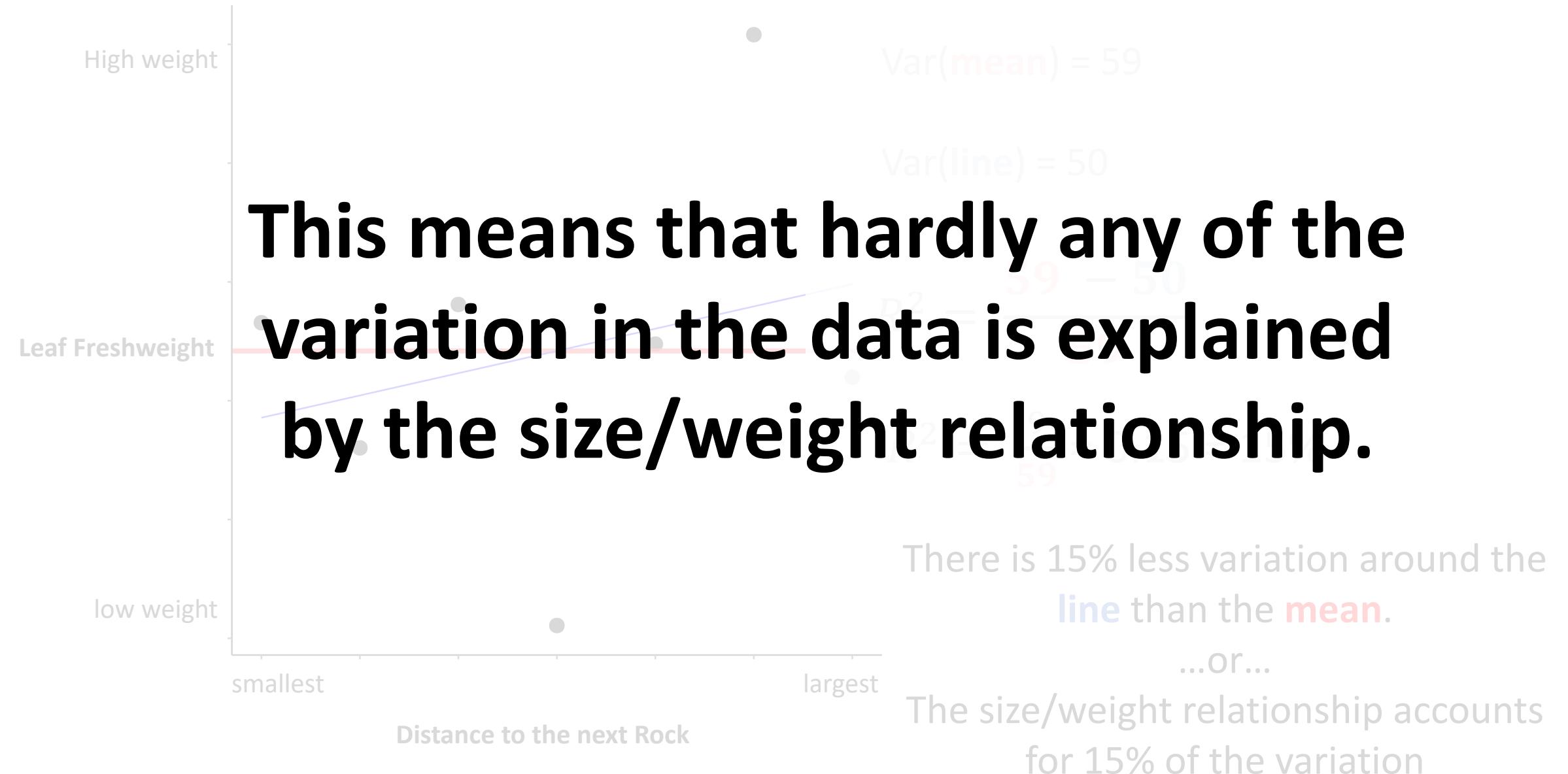
$$R^2 = \frac{59 - 50}{59}$$

$$R^2 = \frac{9}{59} = 0.15 = 15\%$$

There is 15% less variation around the **line** than the **mean**.

...or...

The size/weight relationship accounts for 15% of the variation



- Now, when someone says..
    - “the statistically significant  $R^2$  was 0.9”
  - You can think to yourself..
    - “Very good! The relationship between the two variables explains 90% of the variation in the data!”
- 
- And when someone else says...
    - “The statistically significant  $R^2$  was 0.01”
  - You can think to yourself..
    - “Who cares if that relationship is significant, it only accounts for 1% of the variation in the data.”
    - “Something else must explain the remaining 99%”

- Now, when someone says..
  - “the statistically significant  $R^2$  was 0.9”
- You can think to yourself..
  - “Very good! The relationship between the two variables explains 90% of the variation in the data!”
- And when someone else says...
  - “The statistically significant  $R^2$  was 0.01”
- You can think to yourself..
  - “Who cares if that relationship is significant, it only accounts for 1% of the variation in the data.”
  - “Something else must explain the remaining 99%”

# So, what about R? How is it related to $R^2$ ?

- $R^2$  is just the square of R.
- Now, when someone says...
  - “The statistically significant R was 0.9”
- You can think to yourself..
  - “0.9 times 0.9 = 0.81. Very good! The relationship between the two variables explains 81% of the variation in the data!”

I prefer  $R^2$  over plain R, because it is easier to interpret:

How much better is  $R=0.7$  than  $R=0.5$ ?

I prefer  $R^2$  over plain R, because it is easier to interpret:

How much better is  $R=0.7$  than  $R=0.5$ ?

If we convert R to  $R^2$ , we see that:

$R^2=0.7^2=0.5$       50% of the original variation is explained

$R^2=0.5^2=0.25$       25% of the original variation is explained

That said,  $R^2$  does not indicate the direction of the correlation, because it can never be negative.

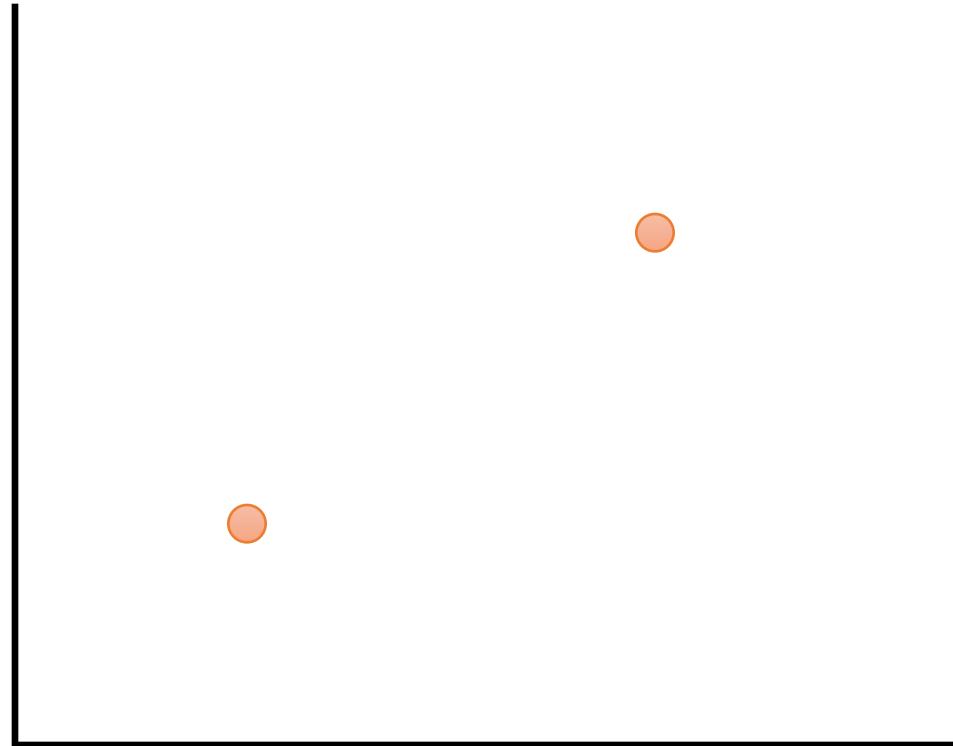
If the direction of the correlation isn't obvious, you can say:

“the two variables were positively (or negatively) correlated with  $R^2=..$ ”

# Main Ideas for $R^2$

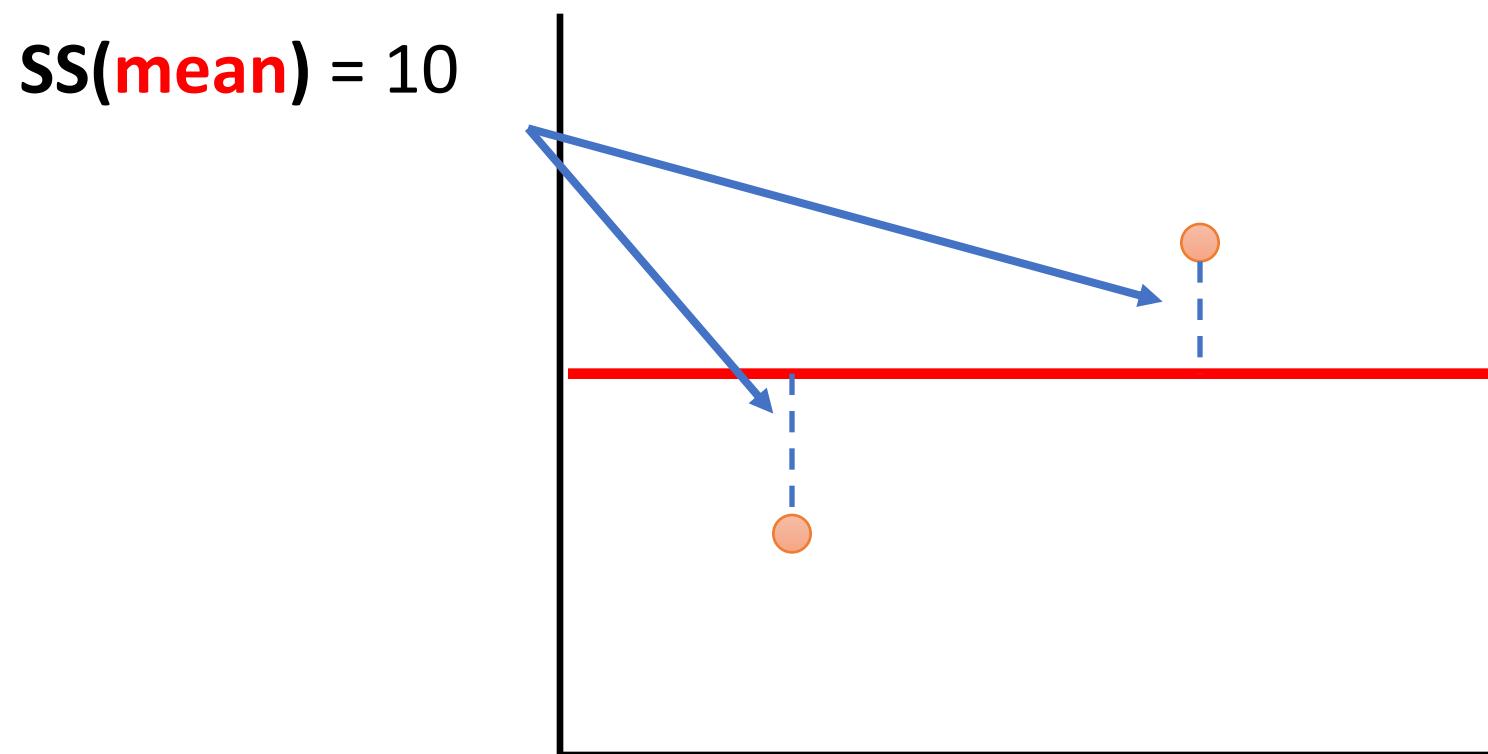
- $R^2$  is the percentage of variation explained by the relationship between two variables.
- If someone gives you a value for plain old R, square it!

So, this is nice and all, but:



What if we had just 2 measurements?

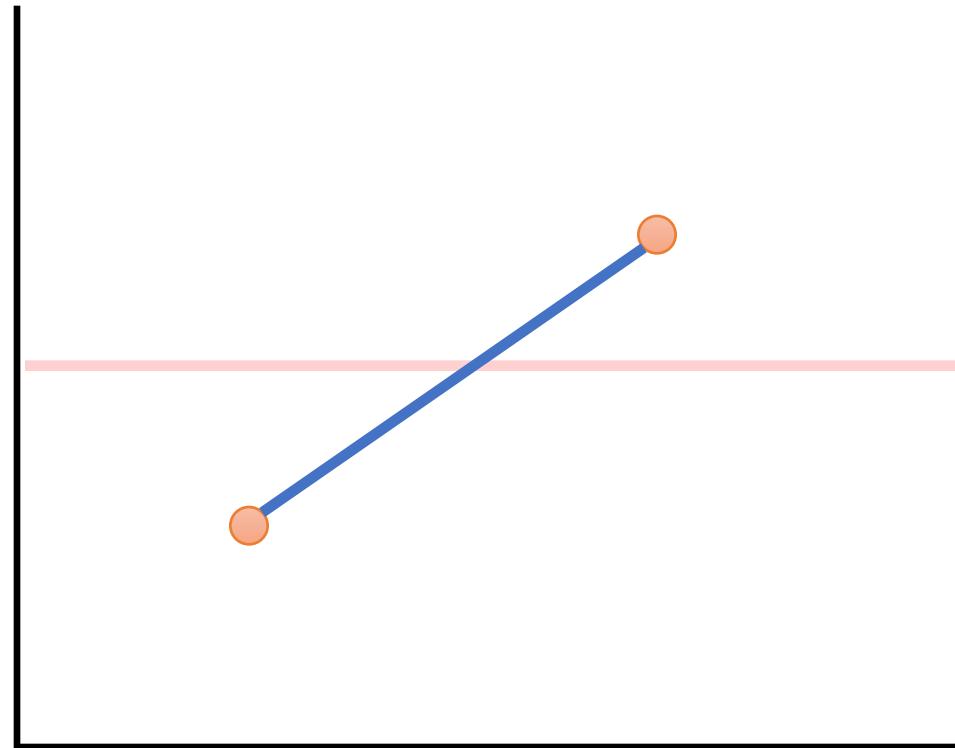
So, this is nice and all, but:



So, this is nice and all, but:

$\text{SS}(\text{mean}) = 10$

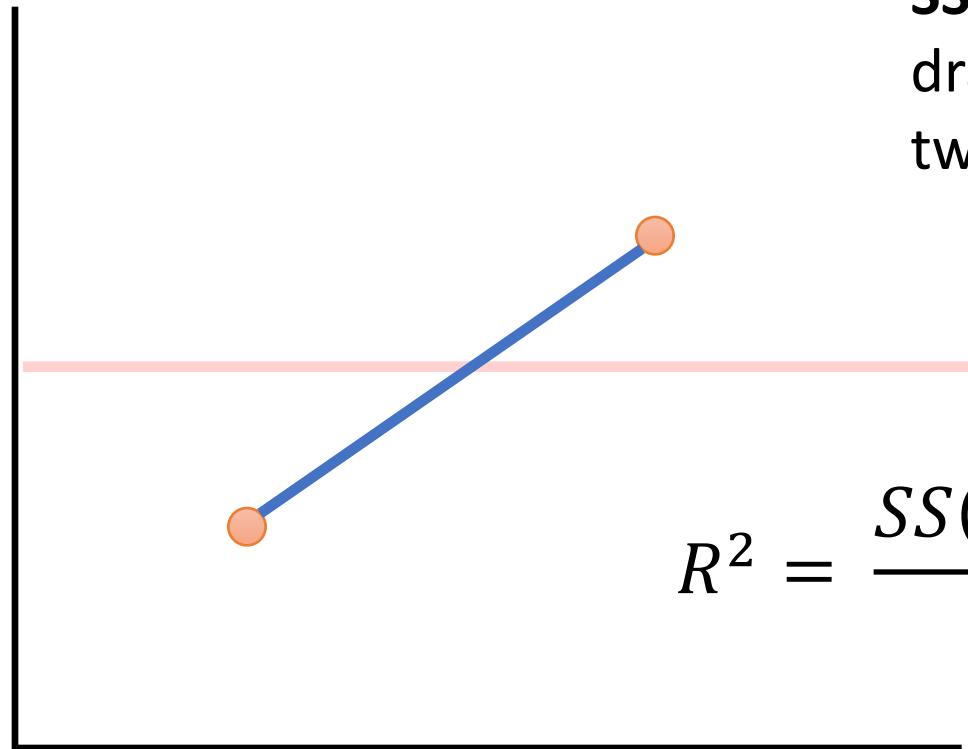
$\text{SS}(\text{fit}) = 0$



So, this is nice and all, but:

**SS(mean) = 10**

**SS(fit) = 0**



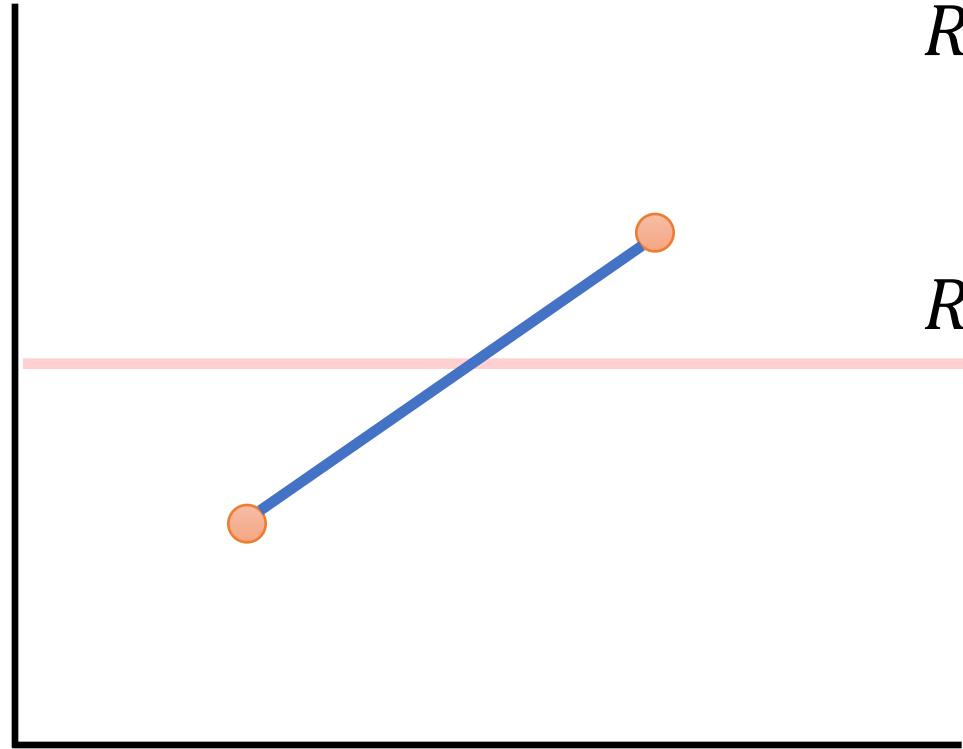
**SS(fit) = 0**, because you can always draw a straight line to connect **any** two dots

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

So, this is nice and all, but:

$$SS(\text{mean}) = 10$$

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



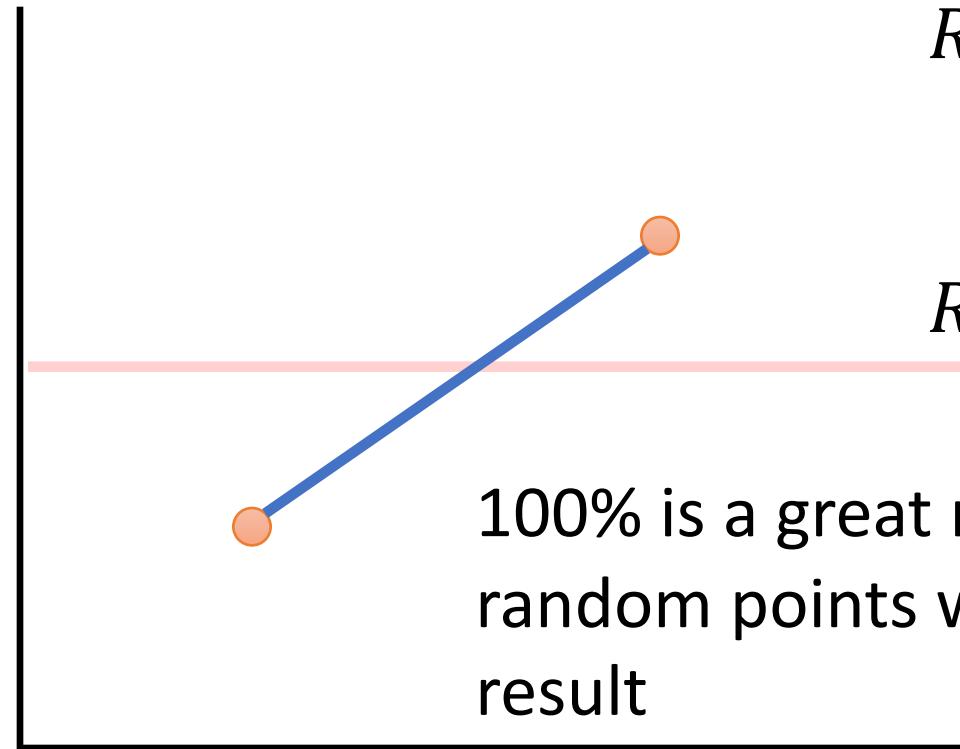
$$R^2 = \frac{SS(\text{mean}) - SS(\text{fit})}{SS(\text{mean})}$$

$$R^2 = \frac{100 - 0}{100} = 1 = 100\%$$

So, this is nice and all, but:

$$SS(\text{mean}) = 10$$

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

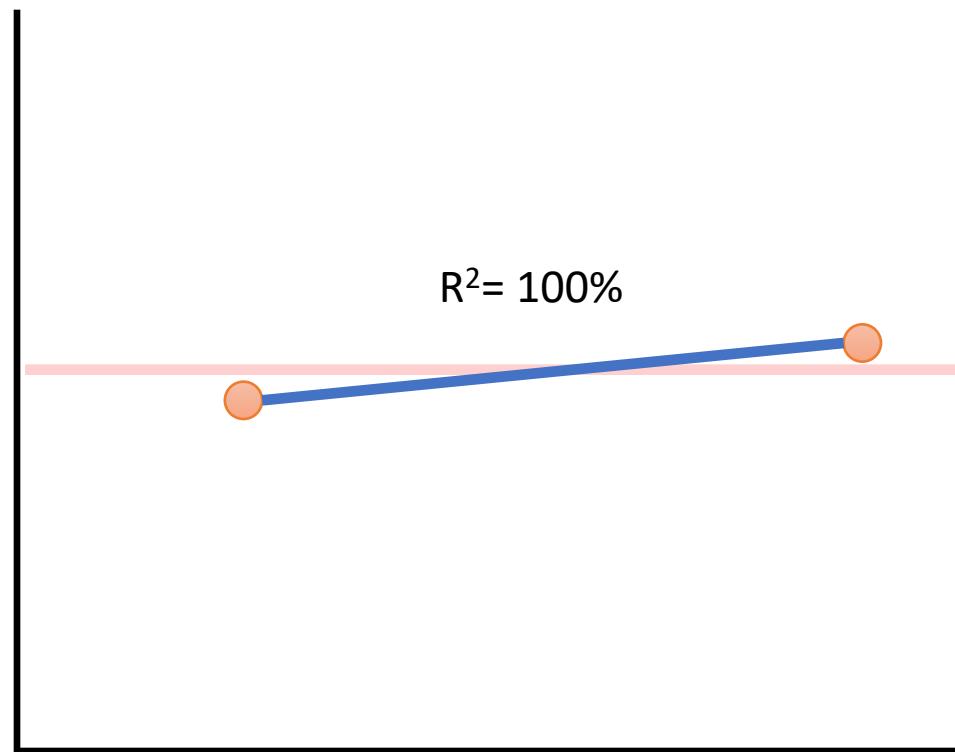


$$R^2 = \frac{SS(\text{mean}) - SS(\text{line})}{SS(\text{mean})}$$

$$R^2 = \frac{100 - 0}{100} = 1 = 100\%$$

100% is a great number, but any two random points will yield the same result

We need a way to determine if the  $R^2$  is statistically significant.



We need a way to determine if the  $R^2$  is statistically significant.

We need a p-value.



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

$$\frac{SS(\text{leaf weight}) - SS(\text{after taking size into account })}{SS(\text{leaf weight})}$$

$$\frac{\text{Variation in leaf weight explained by size}}{\text{Variation in leaf weight without taking size into account}}$$

$$R^2 = \frac{\text{Variation in leaf weight explained by size}}{\text{Variation in leaf weight without taking size into account}}$$

The p-value for  $R^2$  comes from something called “F”

$$R^2 = \frac{\text{Variation in leaf weight explained by size}}{\text{Variation in leaf weight without taking size into account}}$$

$$F = \frac{\text{Variation in leaf weight explained by size}}{\text{Variation in leaf weight not explained by weight}}$$

The p-value for  $R^2$  comes from something called “F”

$$R^2 = \frac{\text{Variation in leaf weight explained by size}}{\text{Variation in leaf weight without taking size into account}}$$

The numerators are  
the same

$$F = \frac{\text{Variation in leaf weight explained by size}}{\text{Variation in leaf weight not explained by weight}}$$

The p-value for  $R^2$  comes from something called “F”

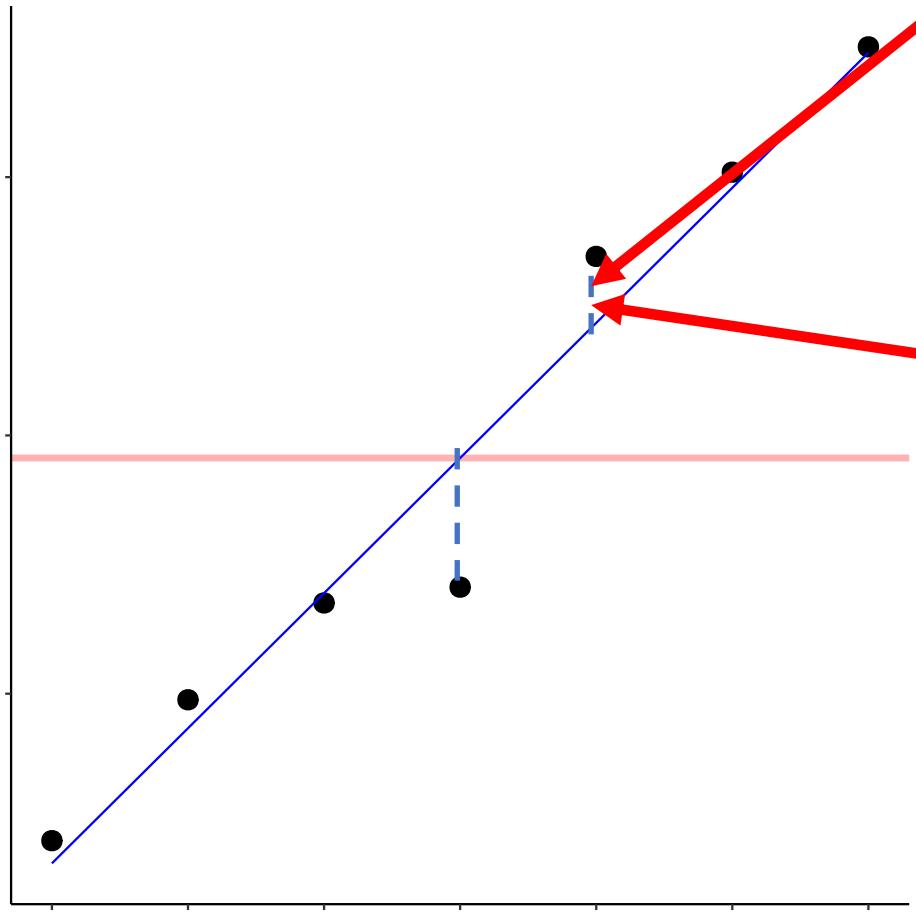
$$R^2 = \frac{\text{Variation in leaf weight explained by size}}{\text{Variation in leaf weight without taking size into account}}$$

The denominator is slightly different

$$F = \frac{\text{Variation in leaf weight explained by size}}{\text{Variation in leaf weight not explained by size}}$$

The p-value for  $R^2$  comes from something called “F”

$$F = \frac{\text{Variation in leaf weight explained by size}}{\text{Variation in leaf weight not explained by size}}$$



These dotted lines (residuals) represent the variation that remains after fitting the line.  
**This is the variation that is not explained by size.**

Now let's look at the underlying mathematics

Now let's look at the underlying mathematics

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

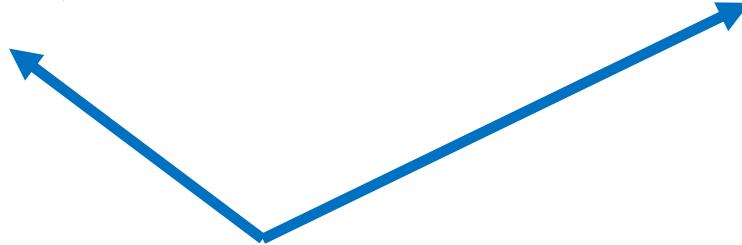
Now let's look at the underlying mathematics

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

$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$

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

$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$



The “meat” of these equations are very similar and rely on the same “sum of squares”

$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$



These numbers over here are the “degrees of freedom”

$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$



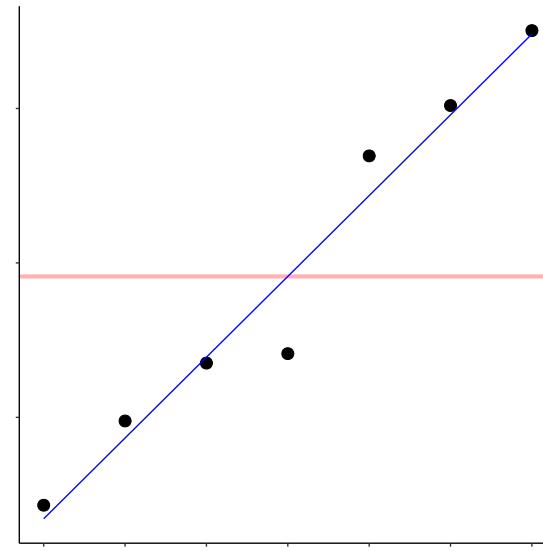
These numbers over here are the “degrees of freedom”

They turn the sums of squares into variances

$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$



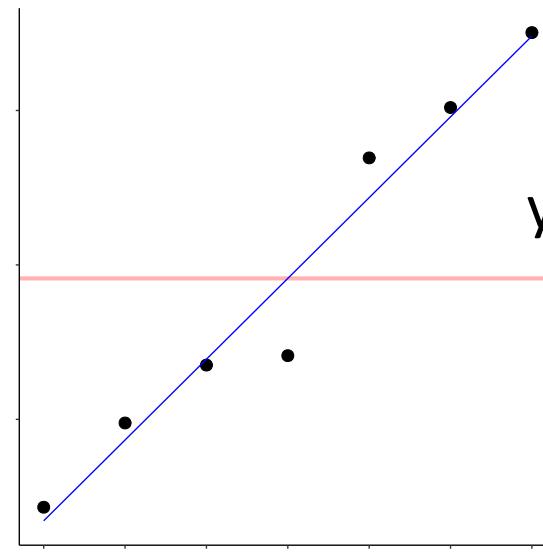
Maybe we can dedicate a whole Seminar to degrees of freedom. For now, let's see if we can get an intuitive feel for what they are doing here.



$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$



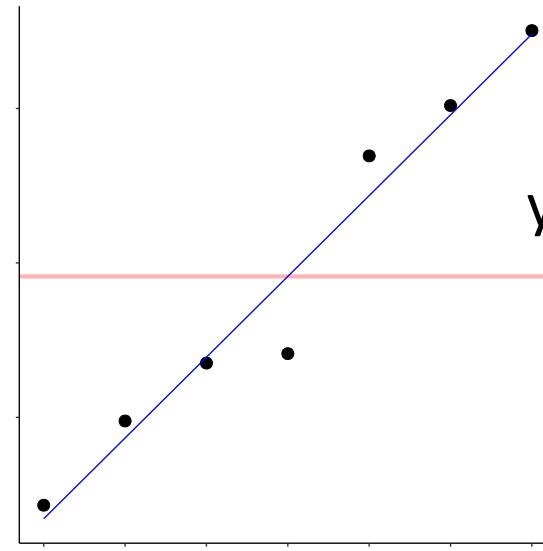
$P_{fit}$  is the number of parameters in the fit line



$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$



$p_{fit}$  is the number of parameters in the fit line



$$y = y\text{-intercept} + \text{slope } x$$

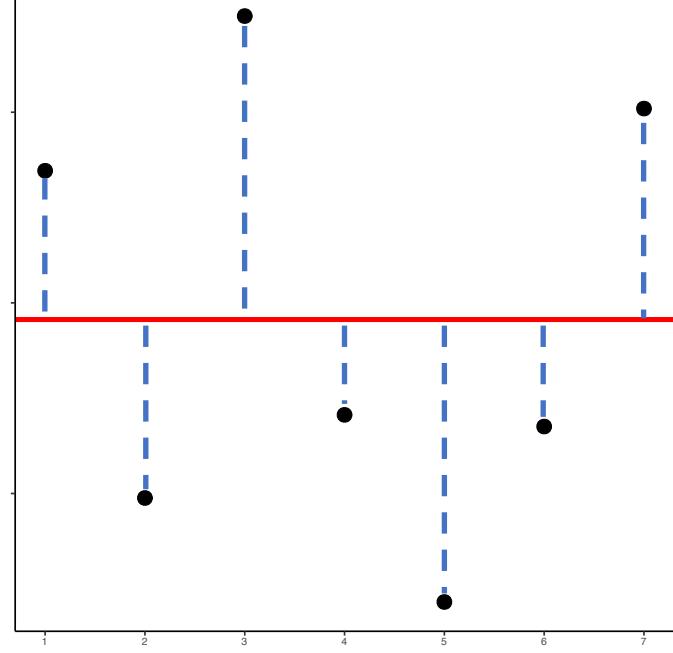
2 parameters

$$p_{fit} = 2$$

$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$

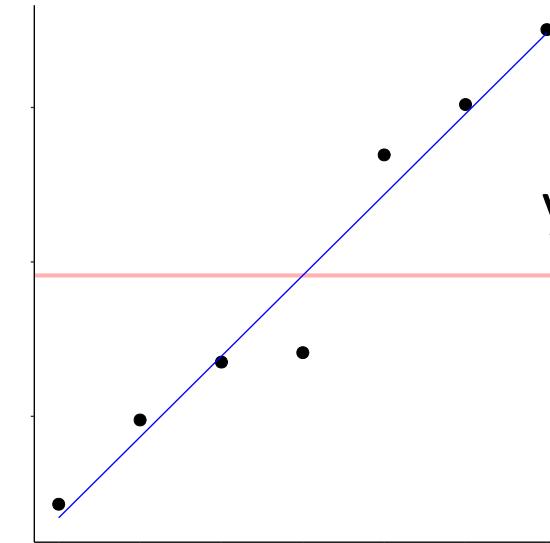


$p_{fit}$  is the number of parameters in the fit line



$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$

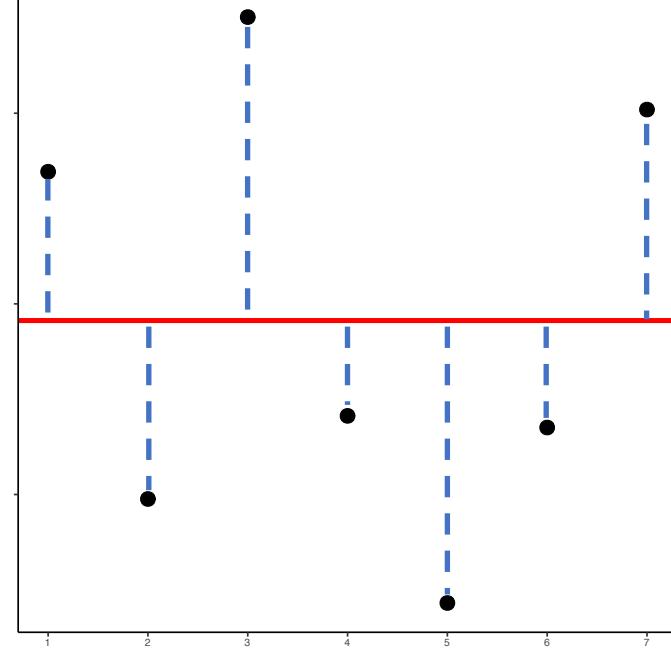
$p_{fit}$  is the number of parameters in the **fit** line  
 $p_{mean}$  is the number of parameters in the **mean** line



$y = y\text{-intercept} + \text{slope } x$

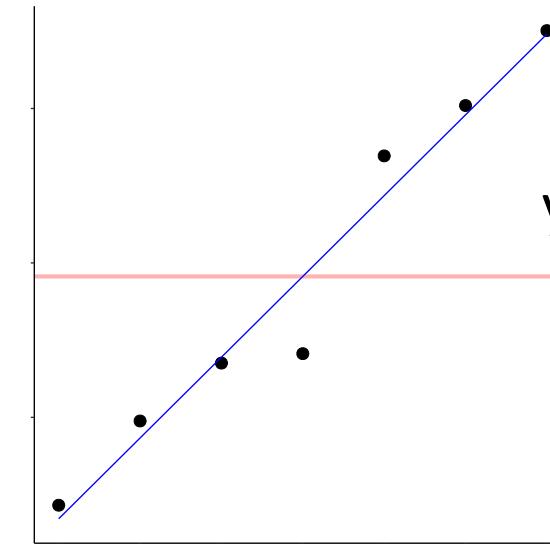
2 parameters

$p_{fit} = 2$



$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$

$y = y\text{-intercept}$



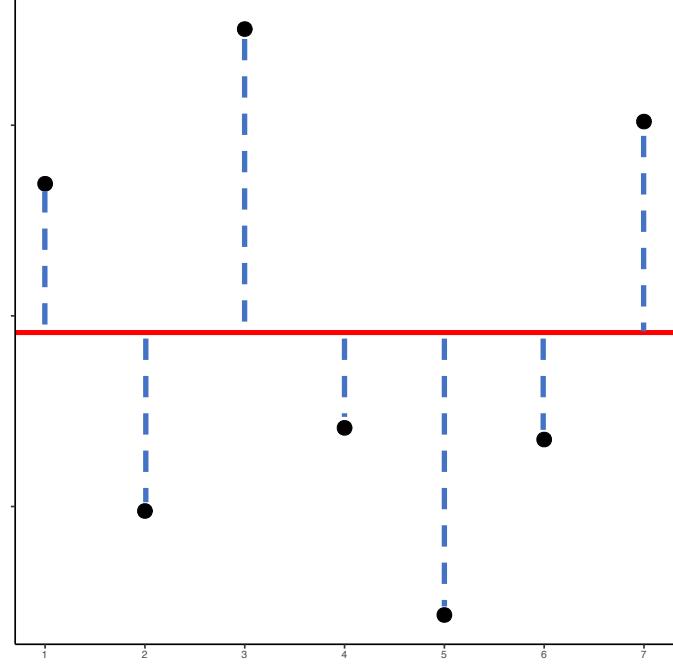
$y = y\text{-intercept} + \text{slope } x$

2 parameters

$p_{fit} = 2$

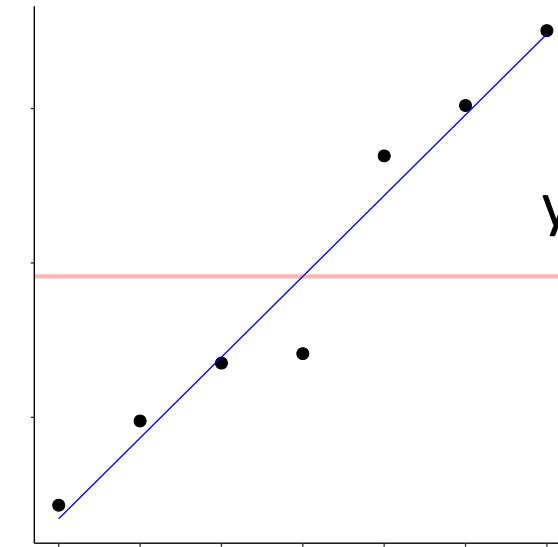
$p_{fit}$  is the number of parameters in the **fit** line

$p_{mean}$  is the number of parameters in the **mean** line



$y = y\text{-intercept}$   
1 parameters  
 $p_{mean} = 1$

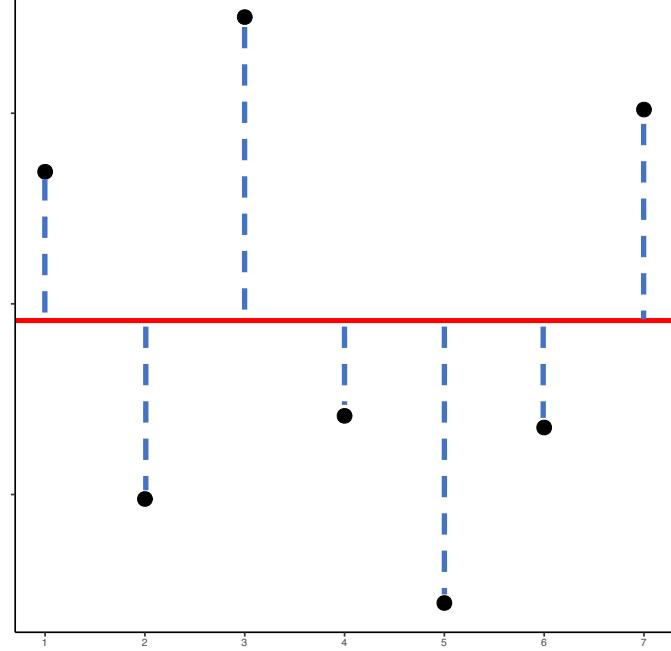
$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$



$y = y\text{-intercept} + \text{slope } x$   
2 parameters  
 $p_{fit} = 2$

Both equations have a parameter for the y-intercept.

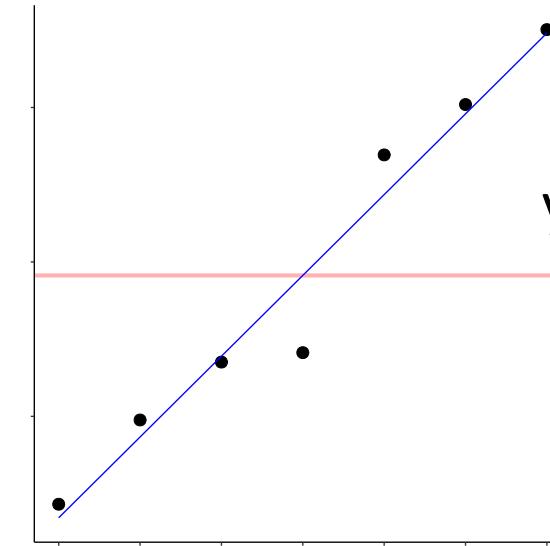
However, the **fit** line has one extra parameter, the slope.  
In our example, this slope is the relationship between size and weight.



$y = y\text{-intercept}$   
1 parameters  
 $p_{mean} = 1$

$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$

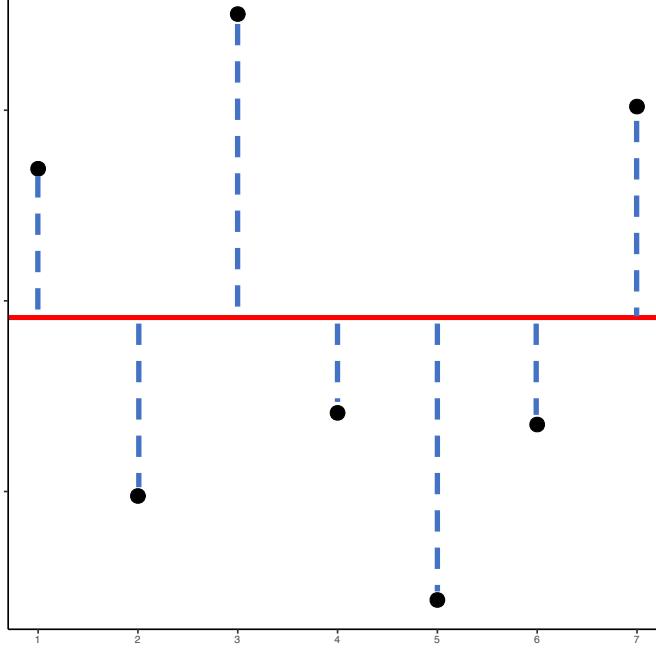
$$(p_{fit} - p_{mean}) = (2-1) = 1$$



$y = y\text{-intercept} + \text{slope } x$

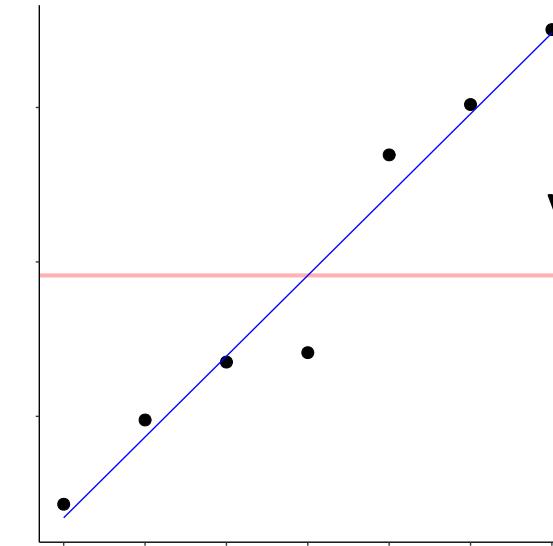
2 parameters

$$p_{fit} = 2$$



$y = y\text{-intercept}$   
1 parameters  
 $p_{mean} = 1$

$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$

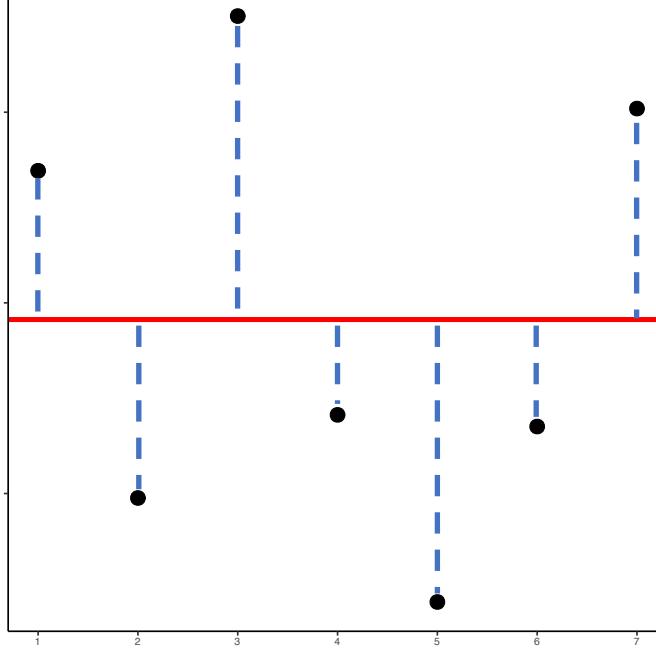


$y = y\text{-intercept} + \text{slope } x$

2 parameters

$p_{fit} = 2$

Thus, the numerator is the variance explained by the extra parameter. In our example, that's the variance in leaf weight explained by leaf size.



$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$

$y = y\text{-intercept}$

1 parameters

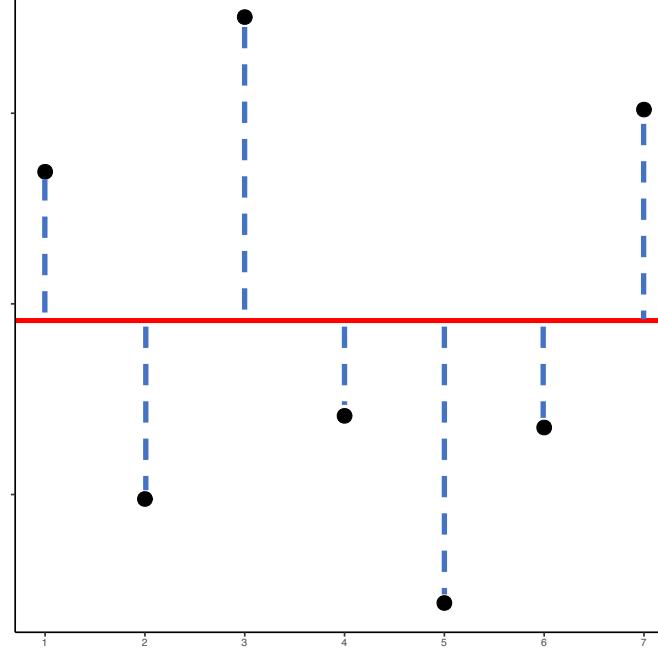
$p_{mean} = 1$

$y = y\text{-intercept} + \text{slope } x + \text{slope } z$

3 parameters

$p_{fit} = 3$

If we had a third variable, the amount of water in the leaf, we'd end up with this equation.



$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$

$y = y\text{-intercept}$

1 parameters

$p_{mean} = 1$

$y = y\text{-intercept} + \text{slope } x + \text{slope } z$

3 parameters

$p_{fit} = 3$

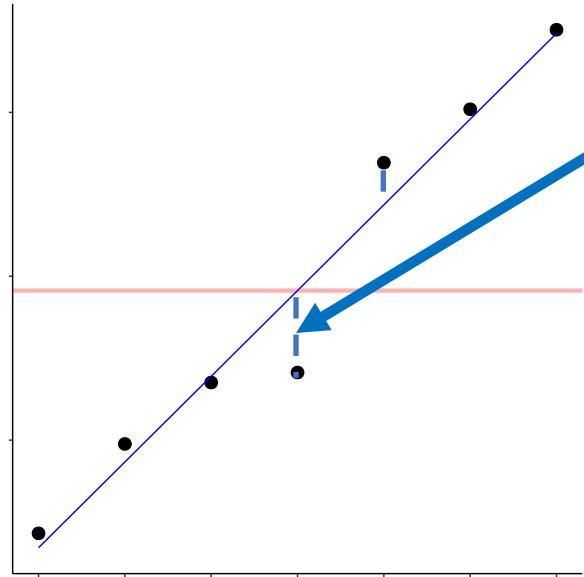
With the fancier **fit**, the numerator is the variance in leaf weight explained by leaf size and water content.

Let's talk about the denominator

$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$

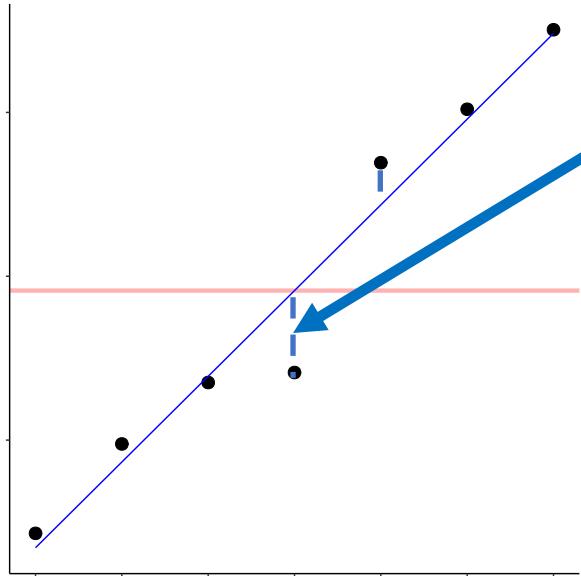
The variation in leaf weight not explained by the fit.

$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$



The variation in leaf weight not explained by the fit.

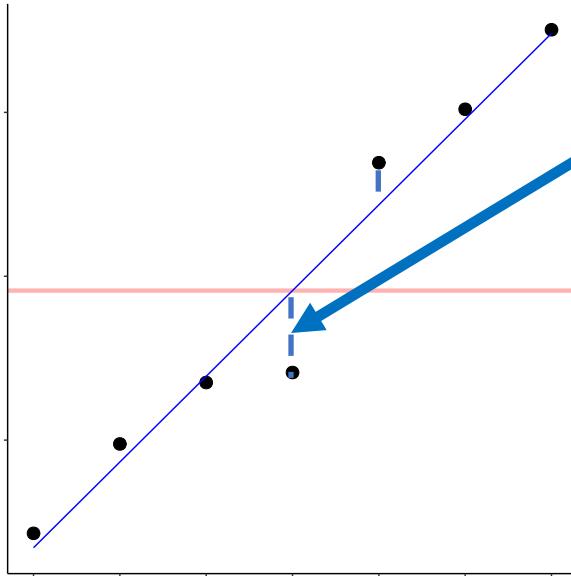
$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$



Why divide  $SS(\text{fit})$  by  $n - p_{fit}$  instead of just  $n$ ?

The variation in leaf weight not explained by the fit.

$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{\text{fit}} - p_{\text{mean}})}{SS(\text{fit})/(n - p_{\text{fit}})}$$



Why divide  $SS(\text{fit})$  by  $n - p_{\text{fit}}$  instead of just  $n$ ?

Intuitively, the more parameters you have in your equation, the more data you need to estimate them. For example, you only need two points to estimate a line, but you need 3 points to estimate a plane.

If the **fit** is good, then ...

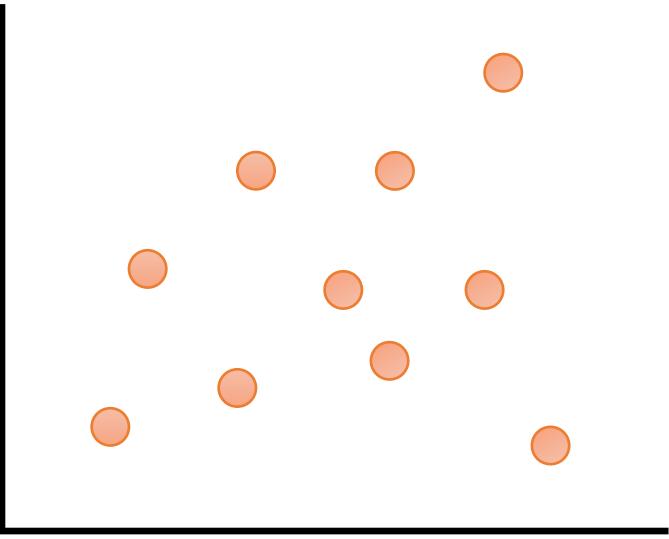
$$F = \frac{\text{Variation in leaf weight explained by size}}{\text{Variation in leaf weight not explained by size}} \rightarrow \frac{\text{large number}}{\text{small number}}$$

If the **fit** is good, then ...

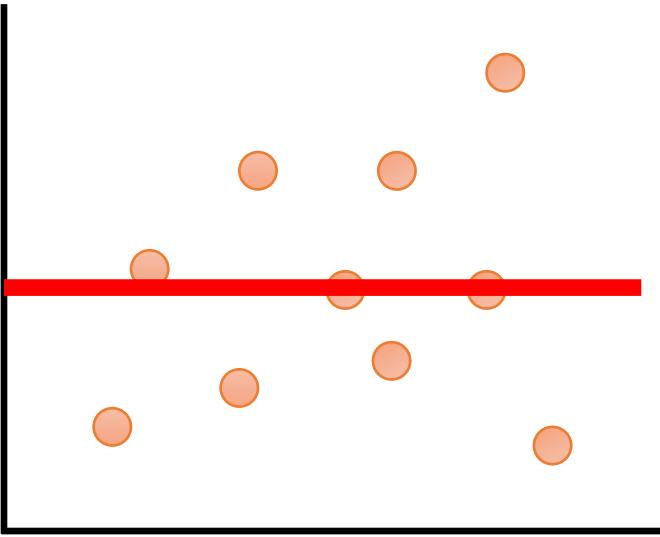
$$F = \frac{\text{Variation in leaf weight explained by size}}{\text{Variation in leaf weight not explained by size}} \rightarrow \frac{\text{large number}}{\text{small number}}$$

$F = \text{really large number}$

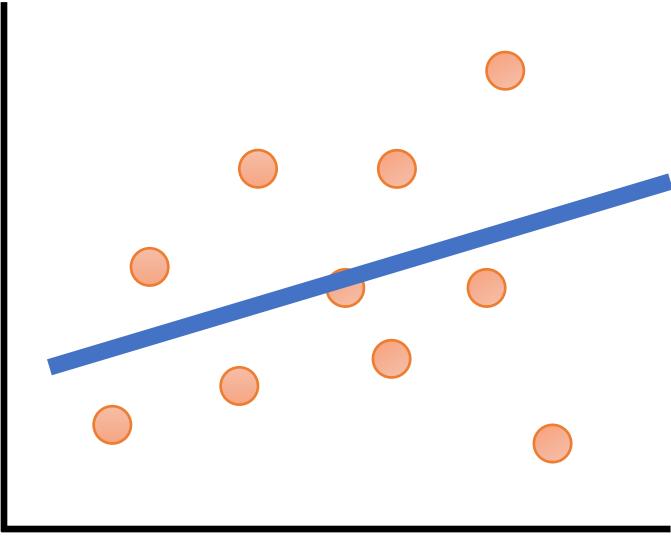
How do we turn this into a p-value?



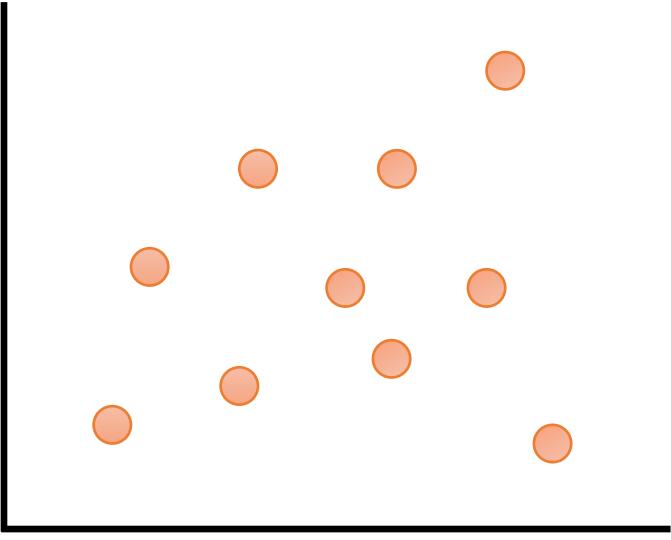
Generate a set of random data...



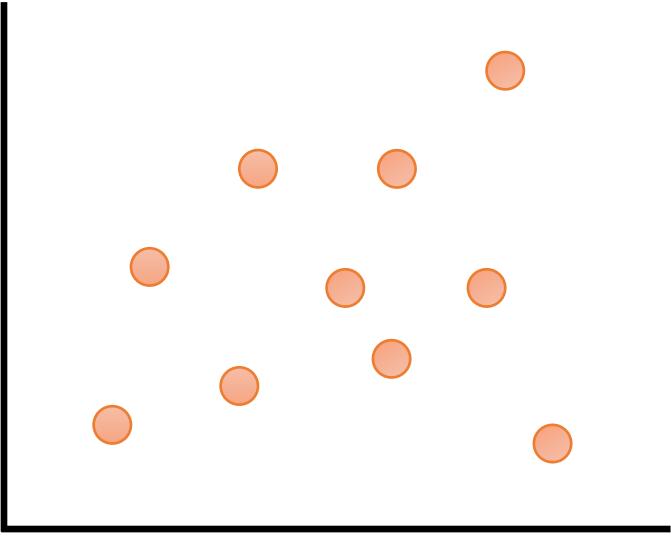
...calculate the **mean** and  $SS(\text{mean})$



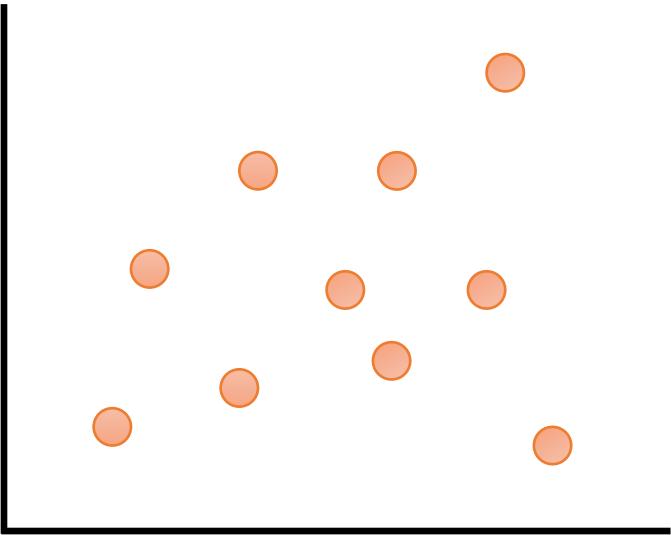
...calculate the **fit** and  $SS(\text{fit})$ ...



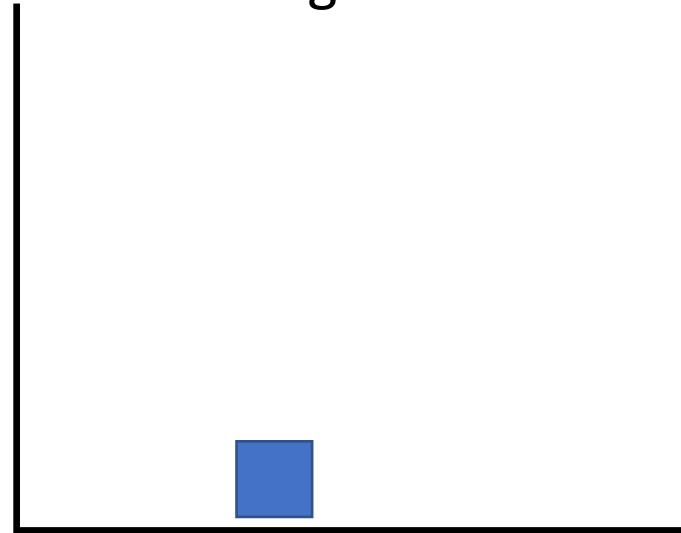
$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$



$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})} \rightarrow F = 2$$



Histogram



$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})} \rightarrow F = 2$$

## Histogram

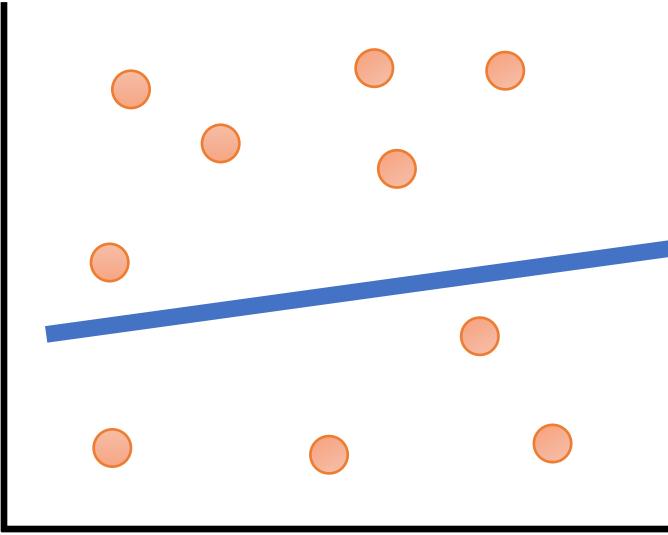
Generate another set  
of random data...



## Histogram

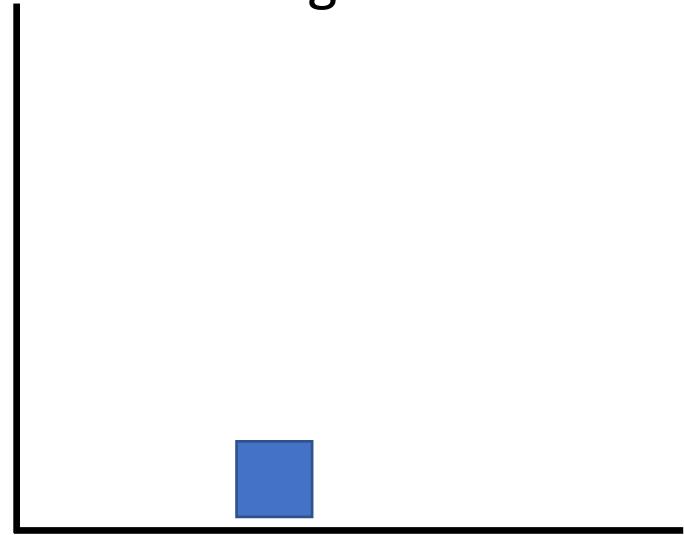
...calculate the **mean**  
and  $SS(\text{mean})$ ...

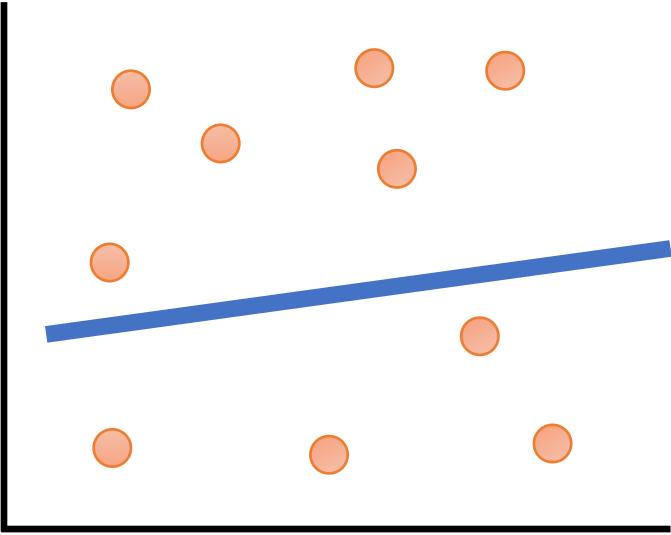




...calculate the **fit** and  
SS(**fit**)...

Histogram

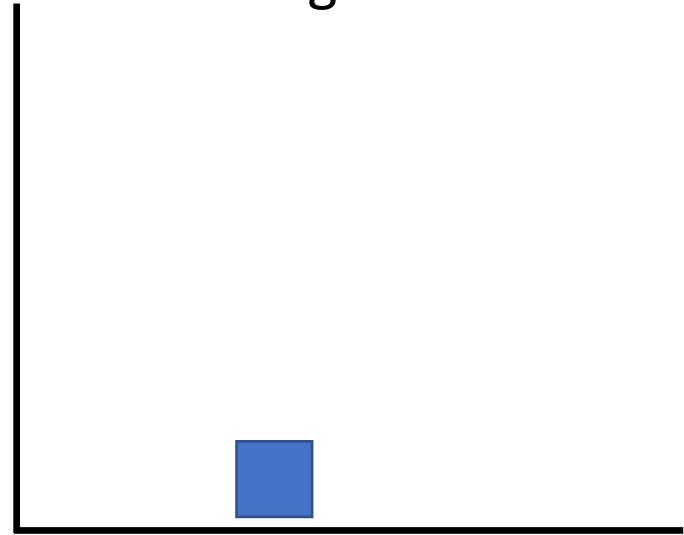


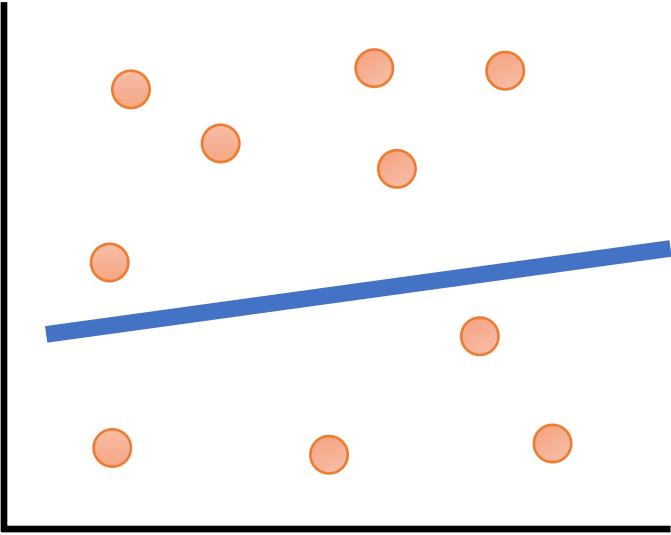


...calculate the **fit** and  
SS(**fit**)...

$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})}$$

Histogram

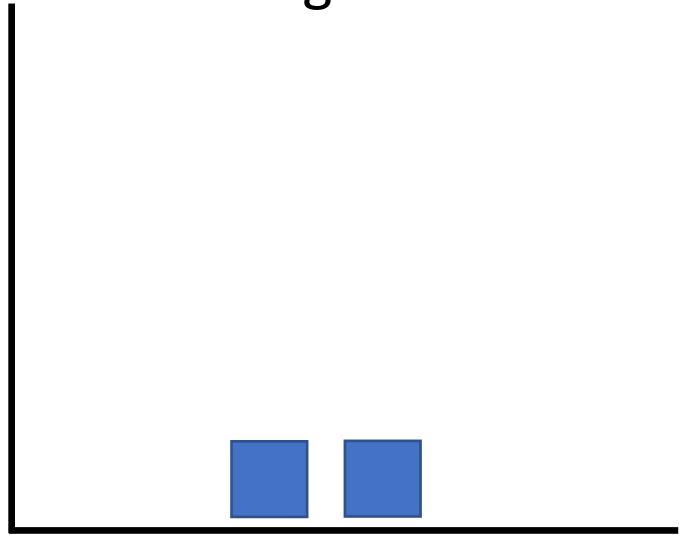


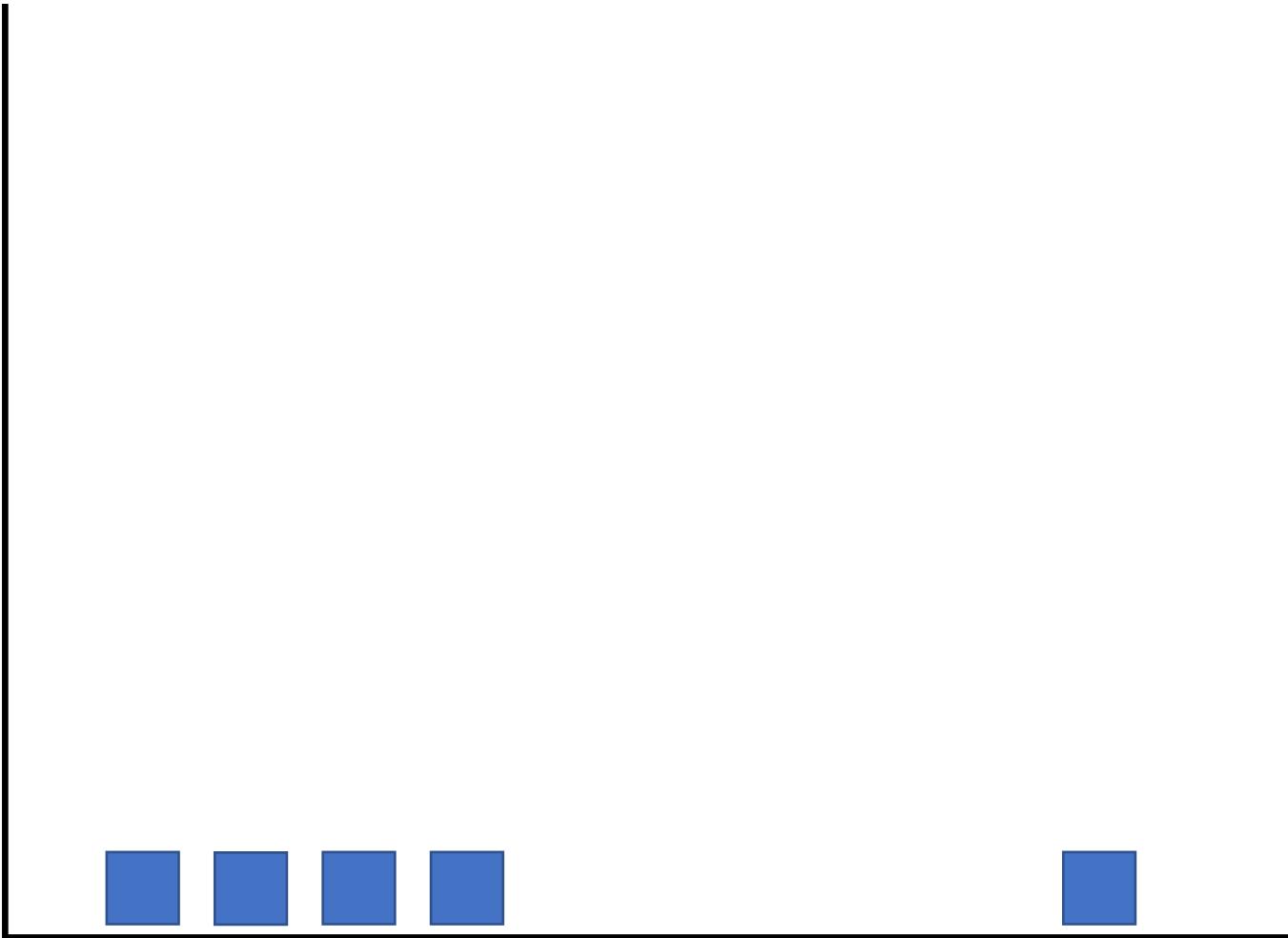


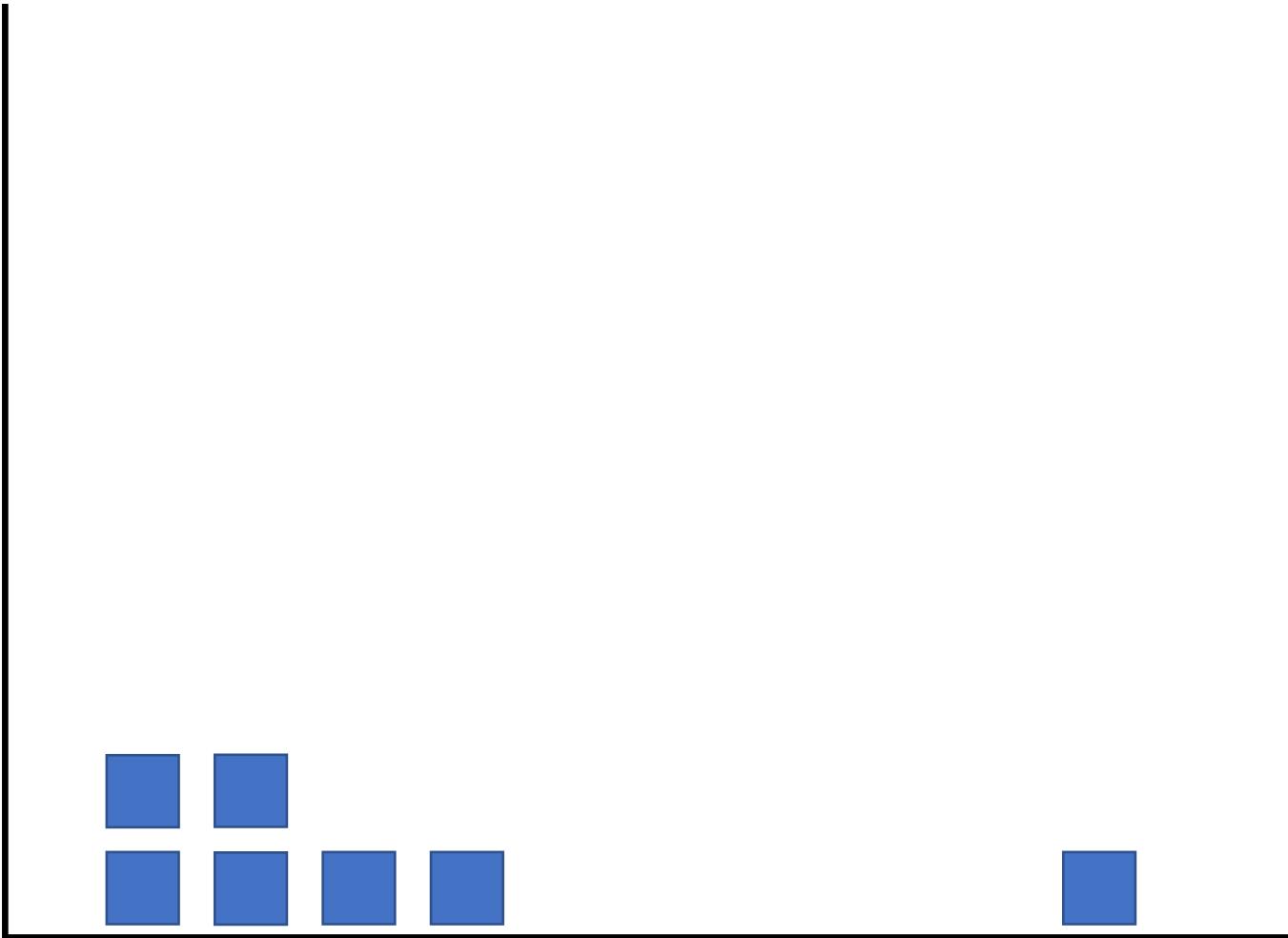
...calculate the **fit** and  
SS(**fit**)...

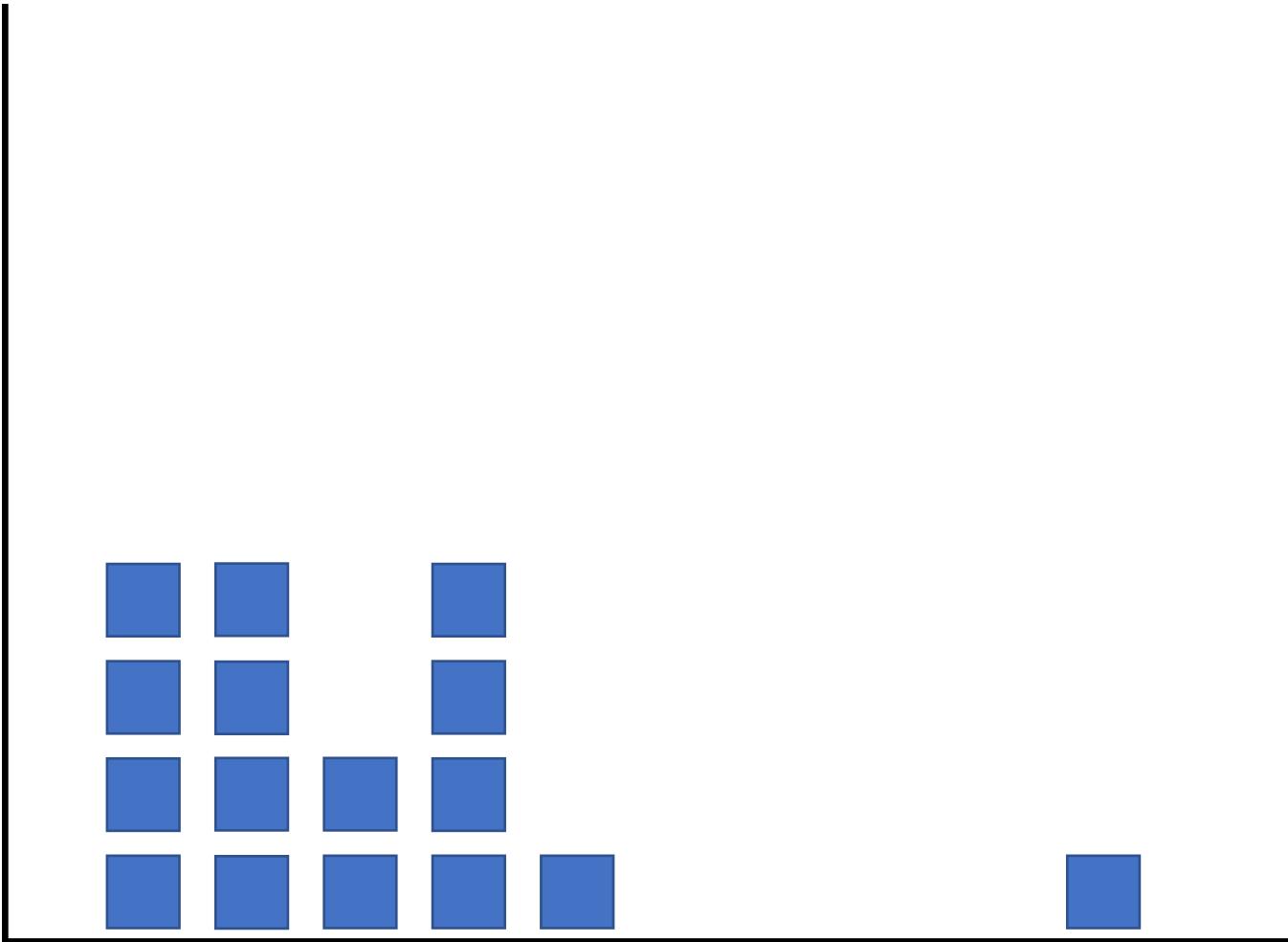
$$F = \frac{SS(\text{mean}) - SS(\text{fit})/(p_{fit} - p_{mean})}{SS(\text{fit})/(n - p_{fit})} \rightarrow F = 3$$

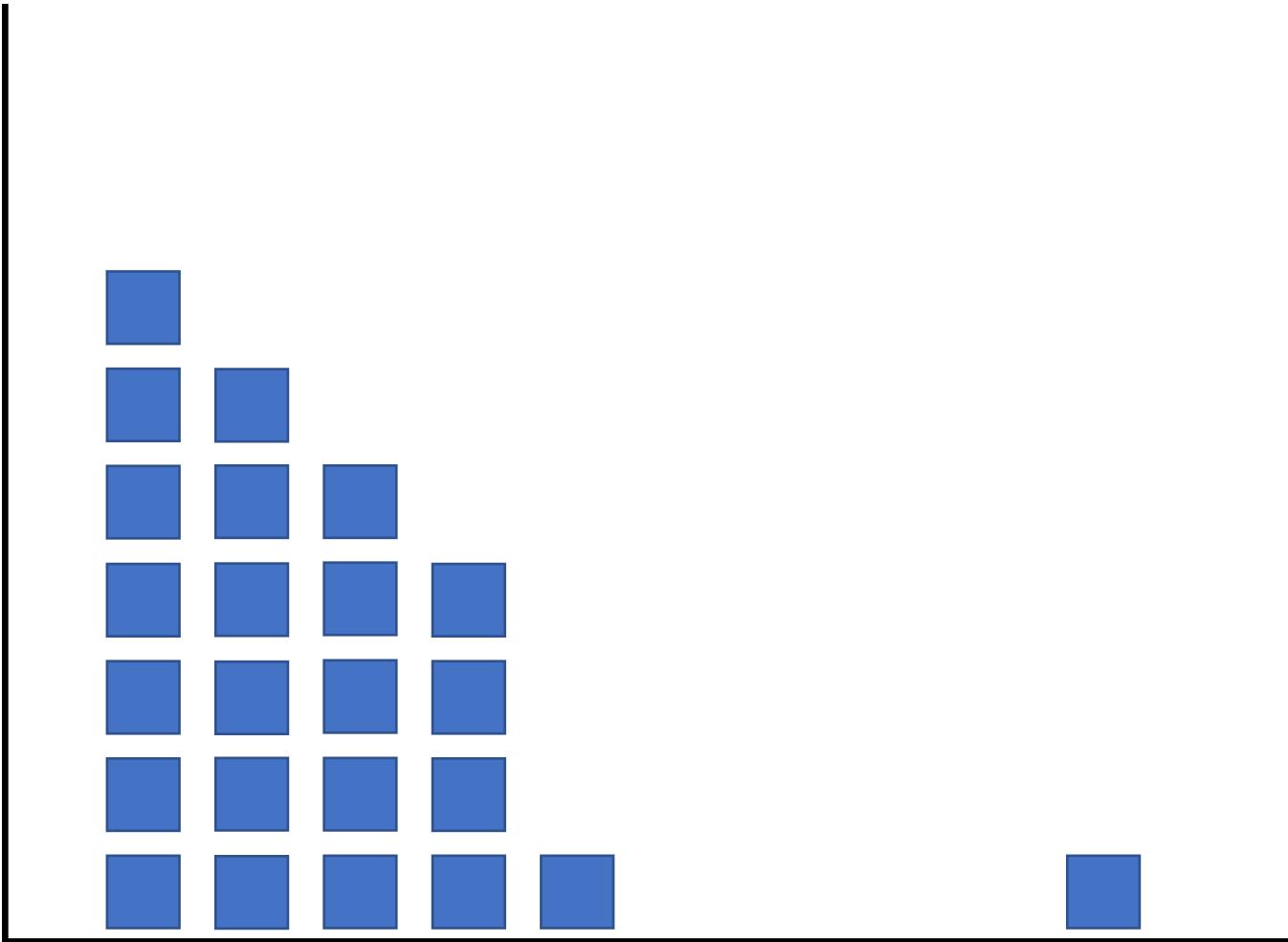
Histogram

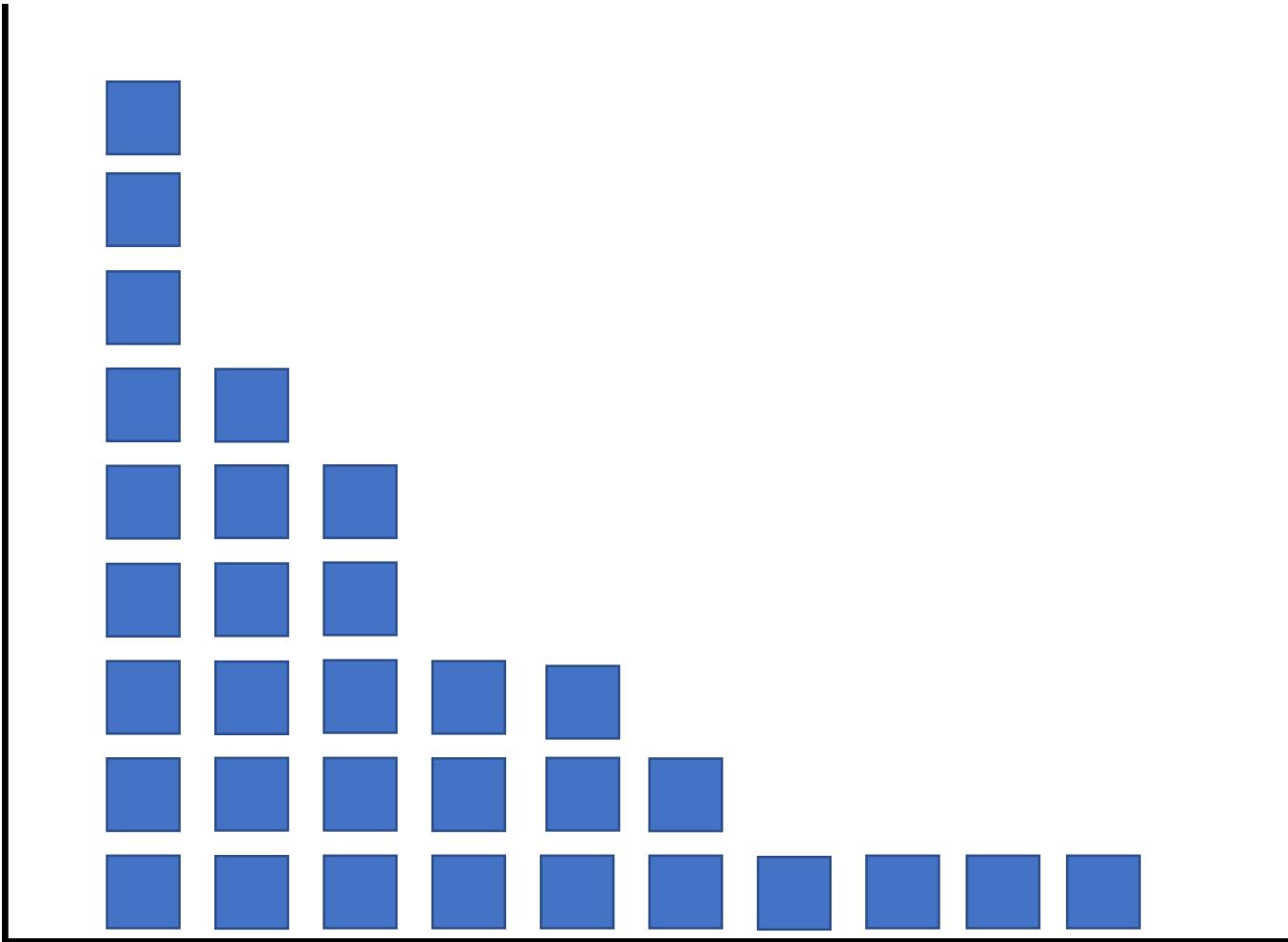




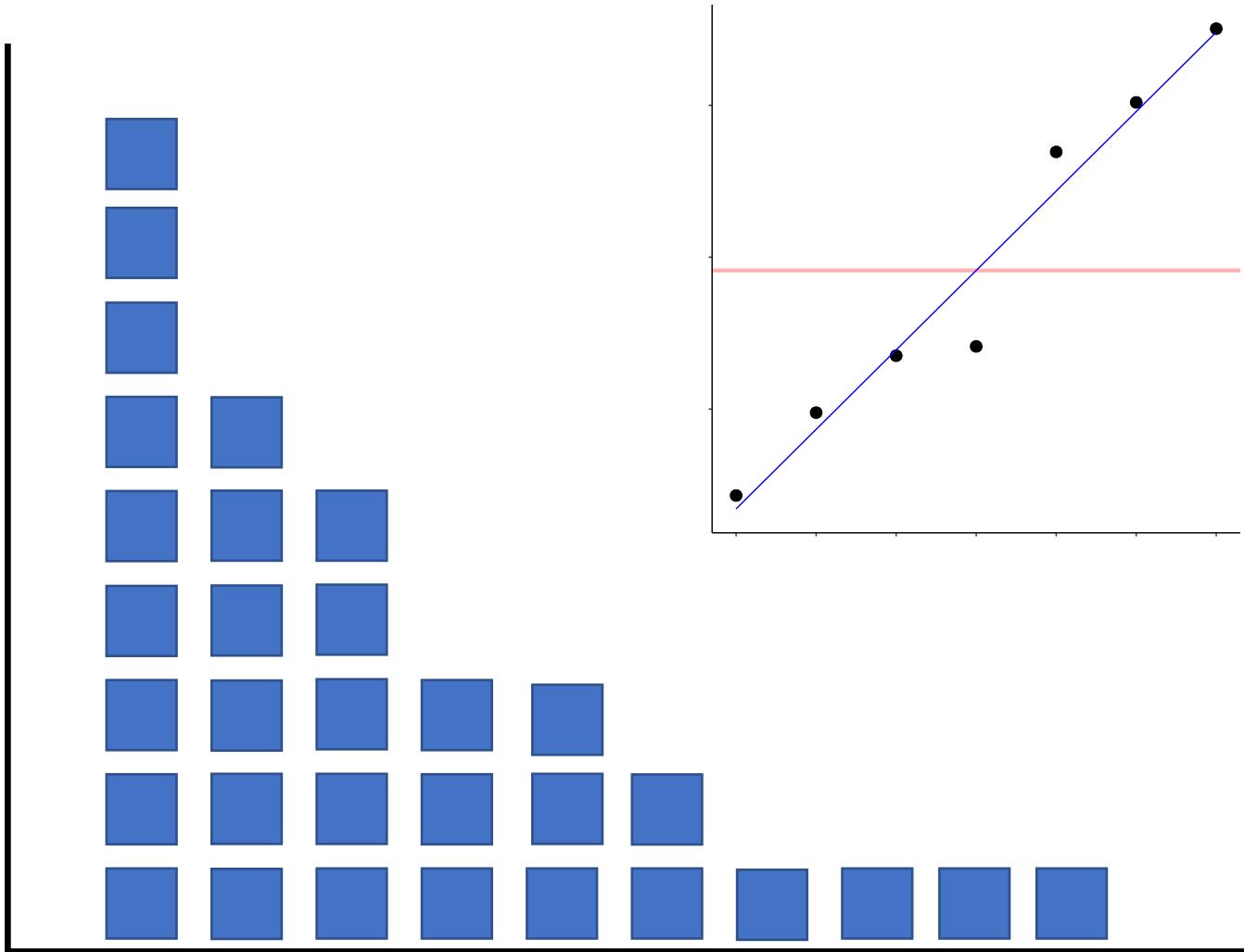




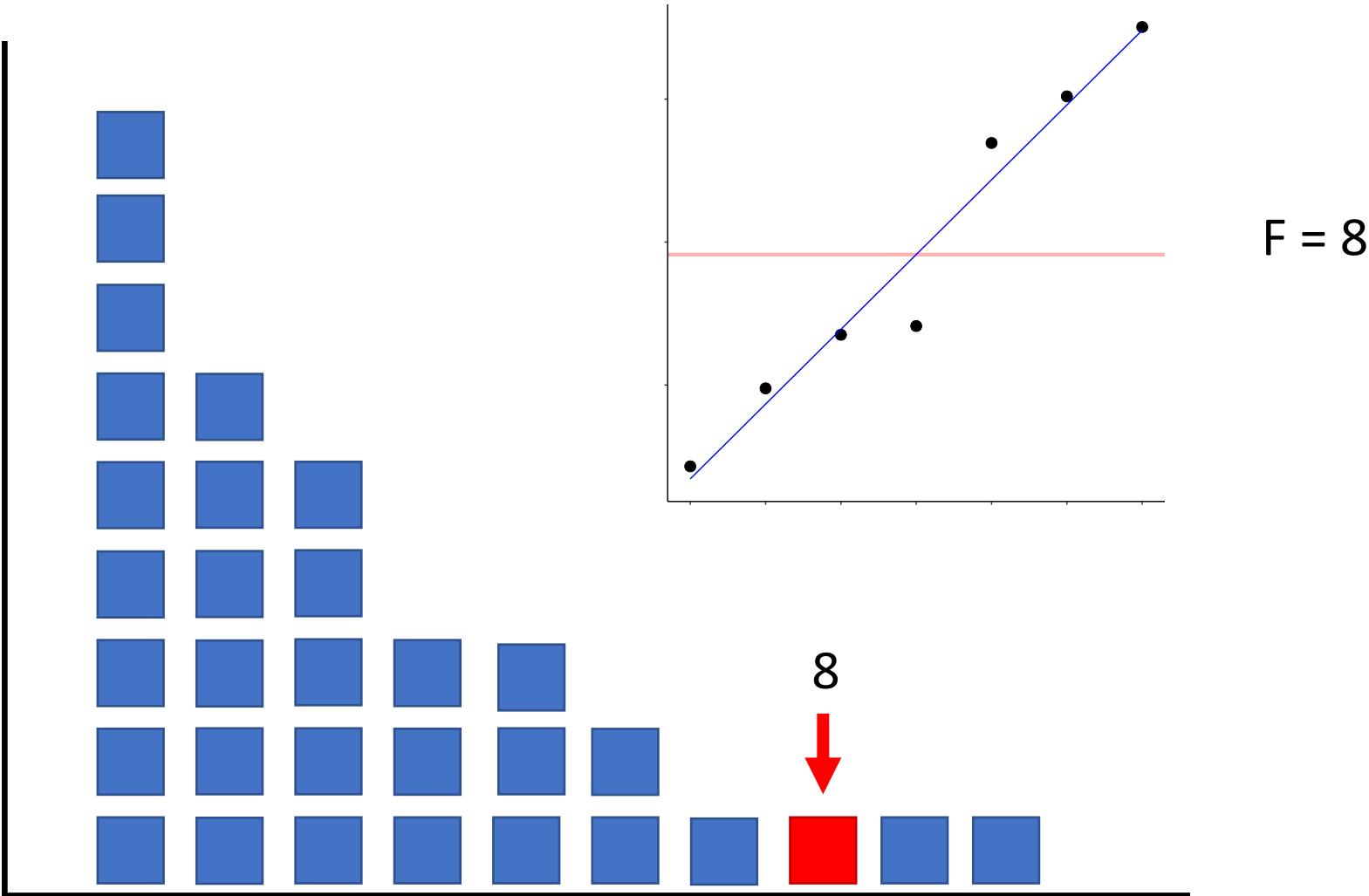


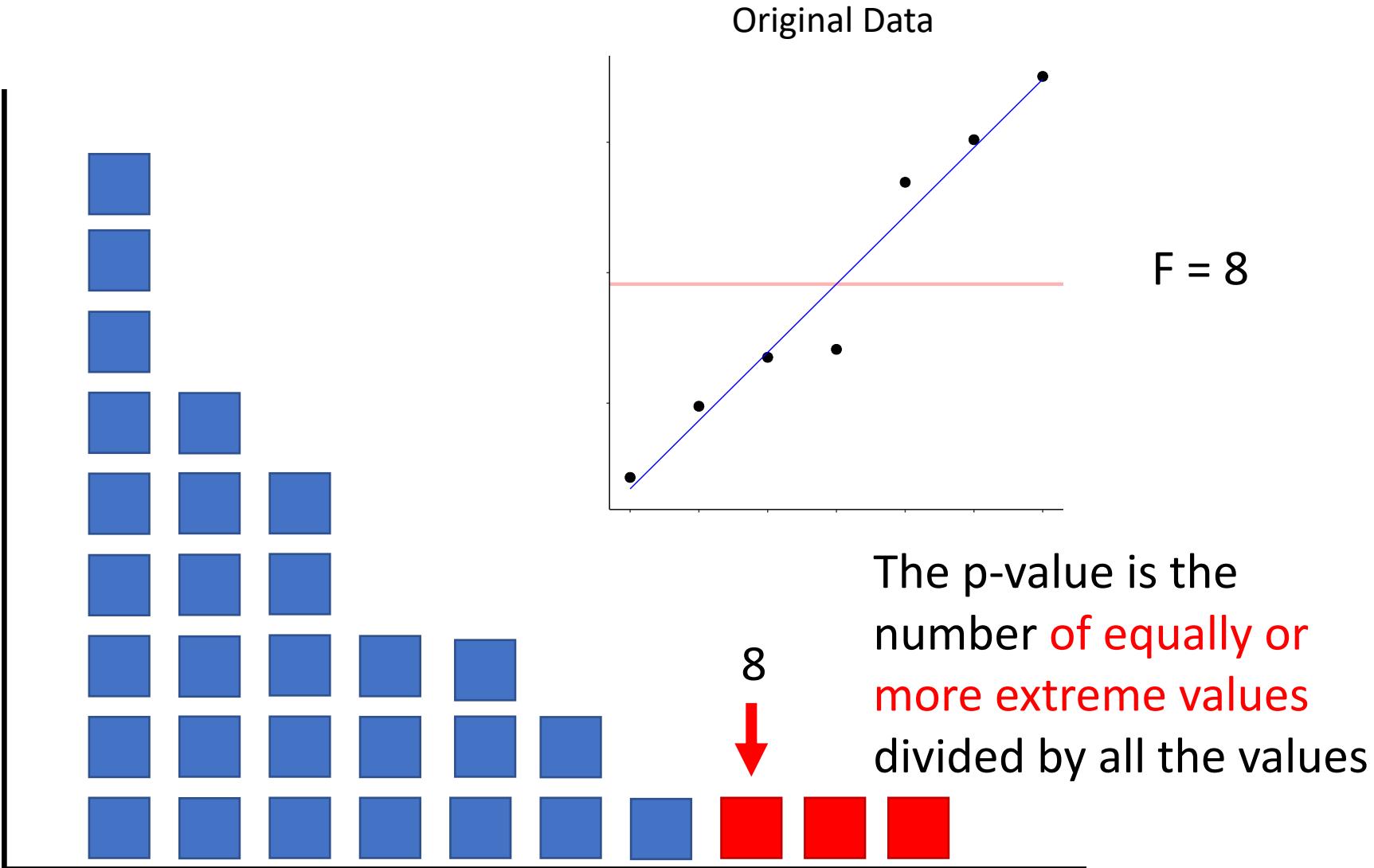


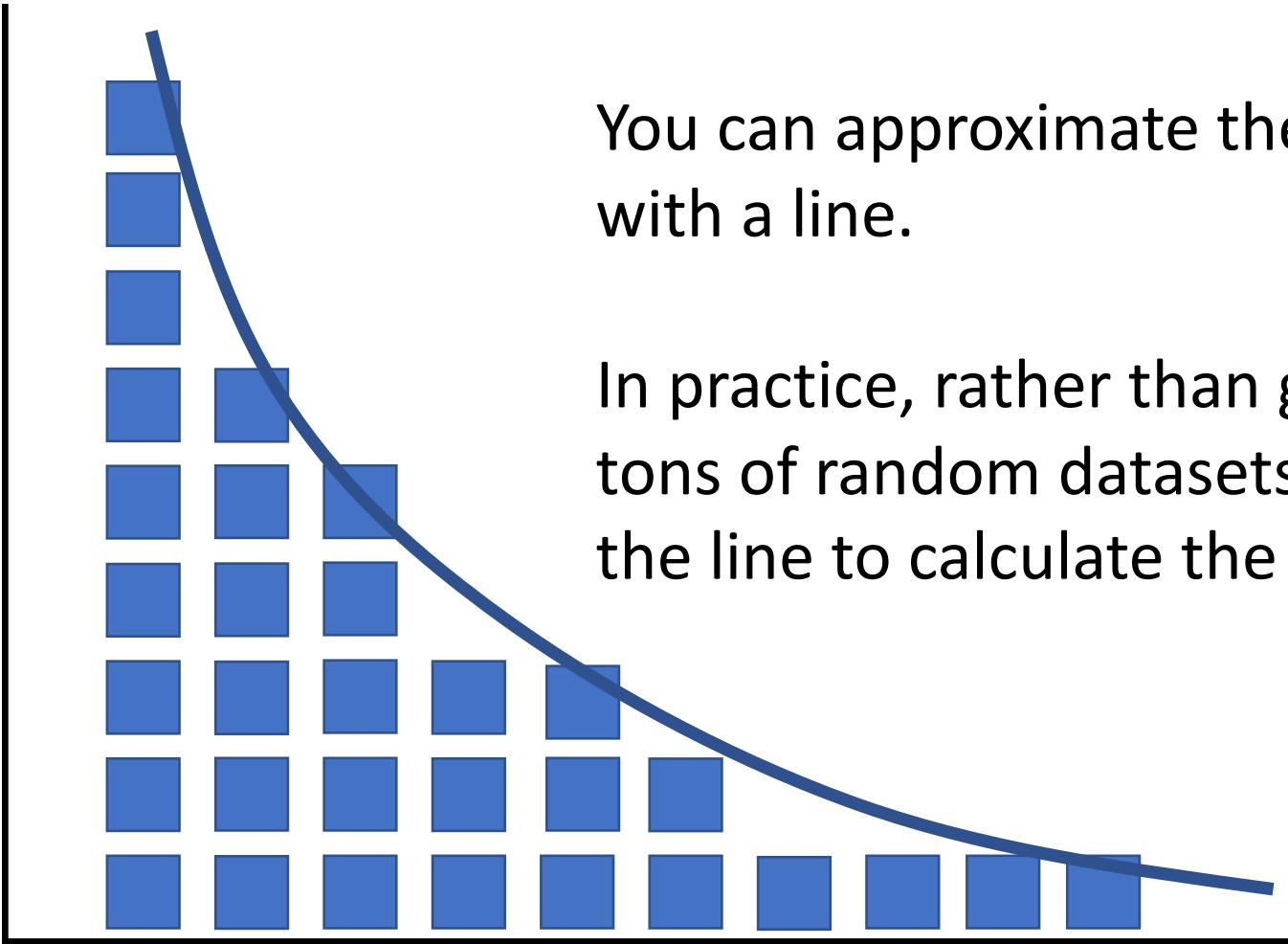
Original Data



Original Data



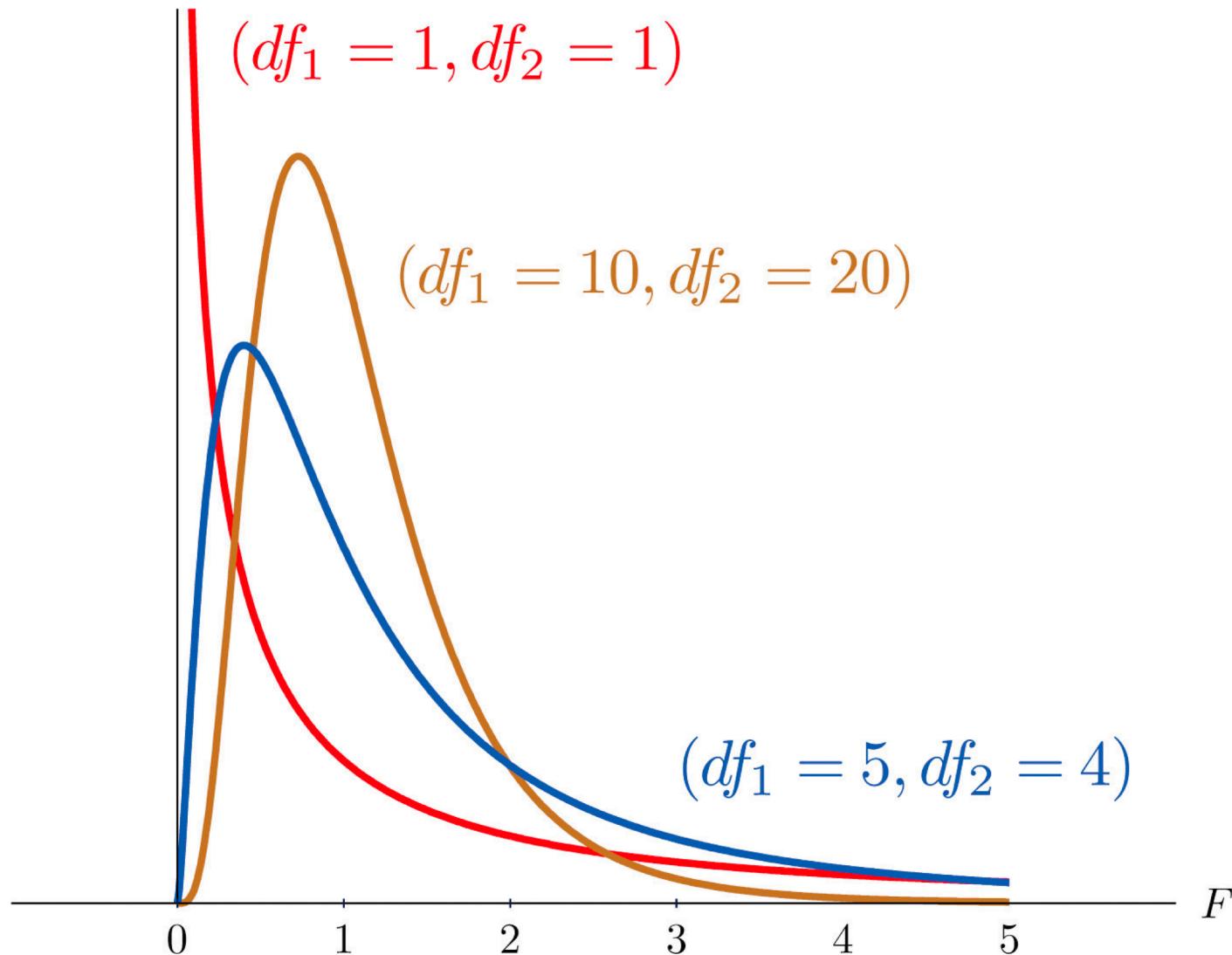


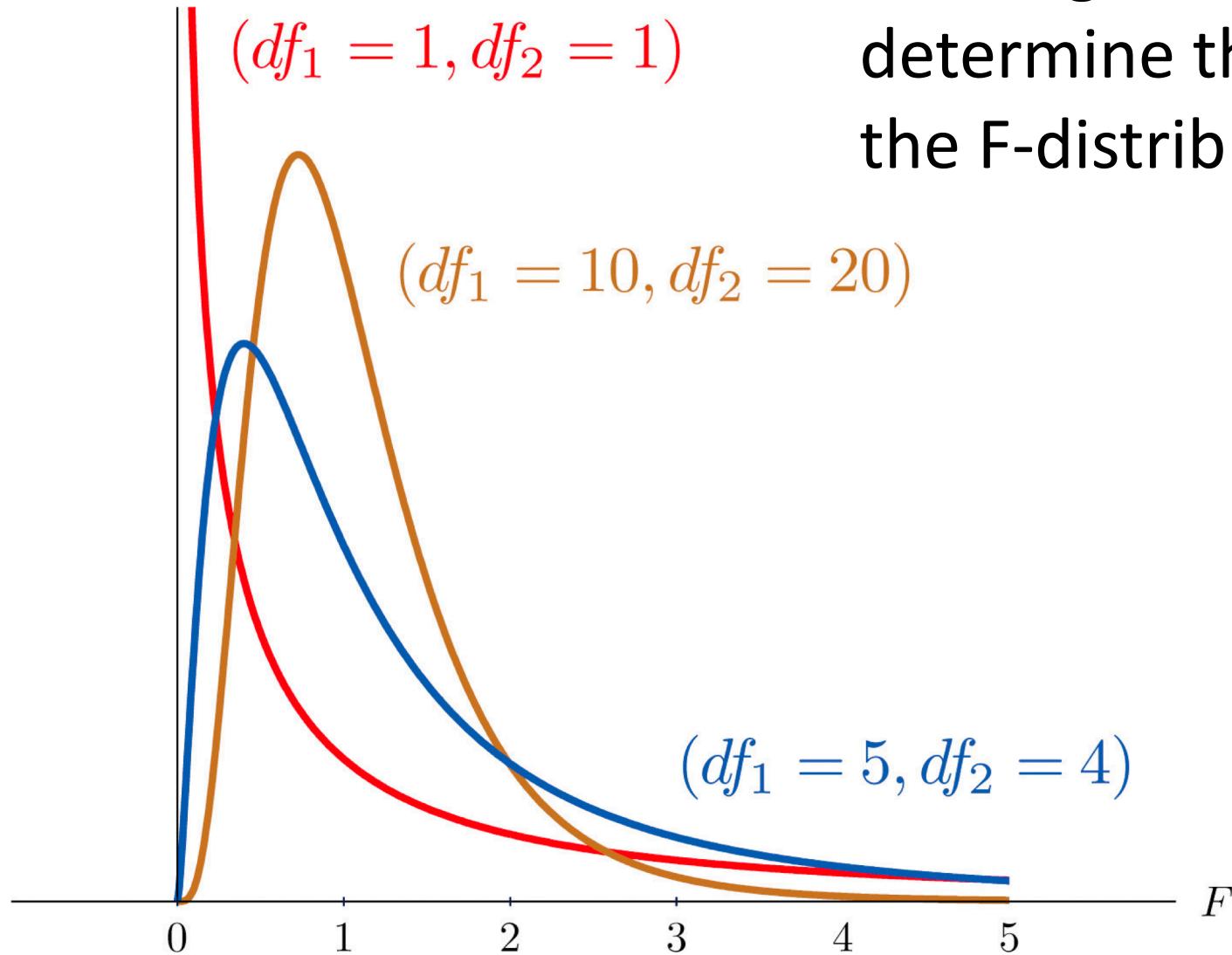


You can approximate the histogram with a line.

In practice, rather than generating tons of random datasets, people use the line to calculate the p-value

## Here are some F-distributions

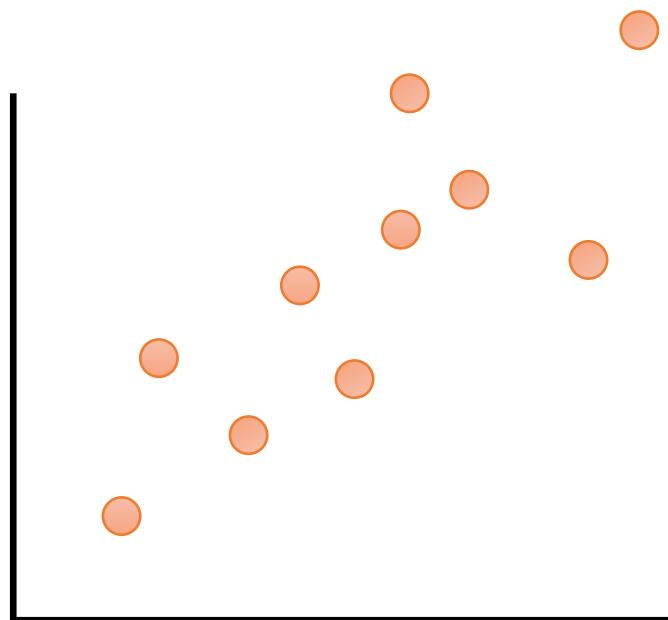




The degrees of freedom  
determine the shape of  
the F-distribution

# Main ideas

# Given some data that you think are related...



Linear regression:

- 1) Quantifies the relationship in the data ( $R^2$ )
- 2) Determines how reliable that relationship is (this is the p-value that we calculate with F)

You need both to have an interesting result!