

Stat...

**QUEST!!!**

Stat...

**QUEST!!!**

Stat...

Stat...

**QUEST!!!**

Stat...

**QUEST!!!**

Stat...

**QUEST!!!**

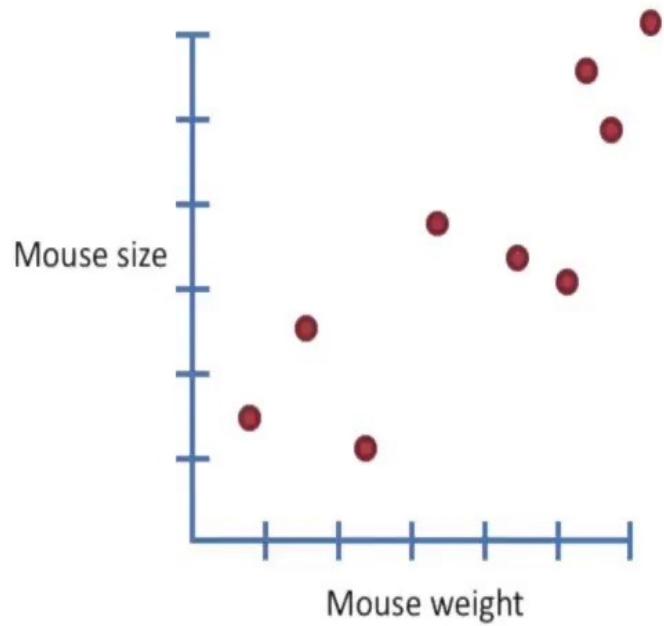
**Yeah!!!**

# StatQuest: General Linear Models

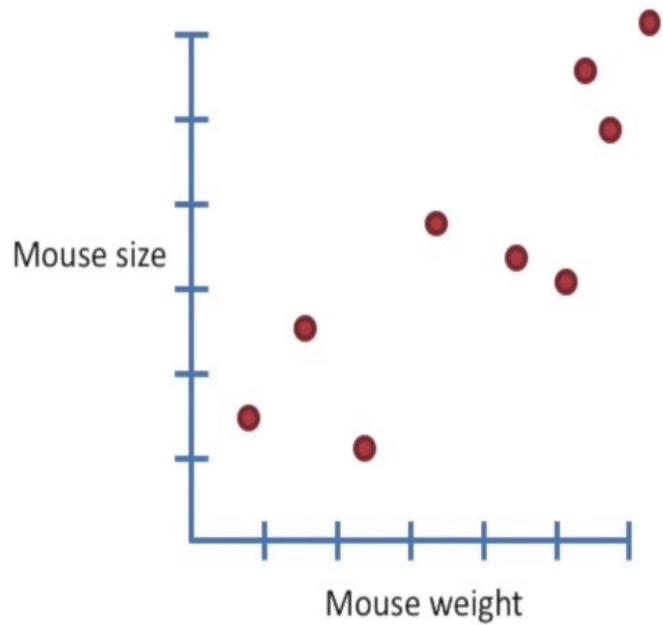
## Part 2:

### t-tests and ANOVA

# Quick Review of Linear Regression...

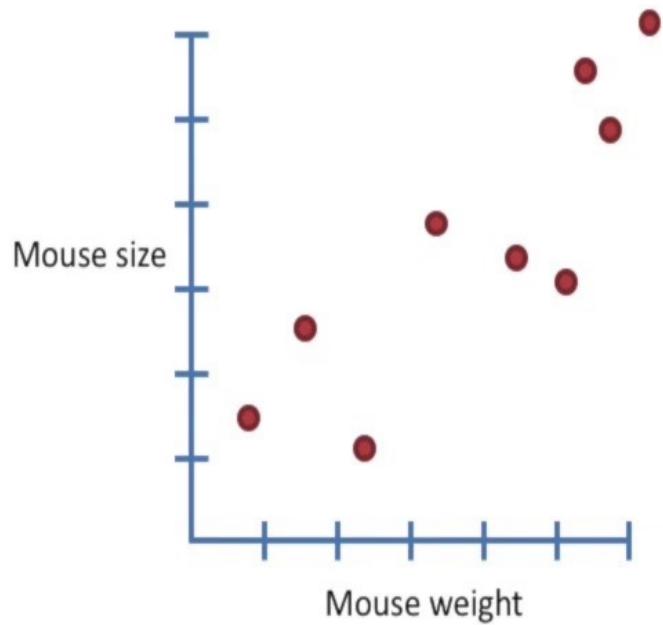


We measured **mouse weight** and **mouse size** and we wanted learn 2 things from it:



We measured **mouse weight** and **mouse size** and we wanted learn 2 things from it:

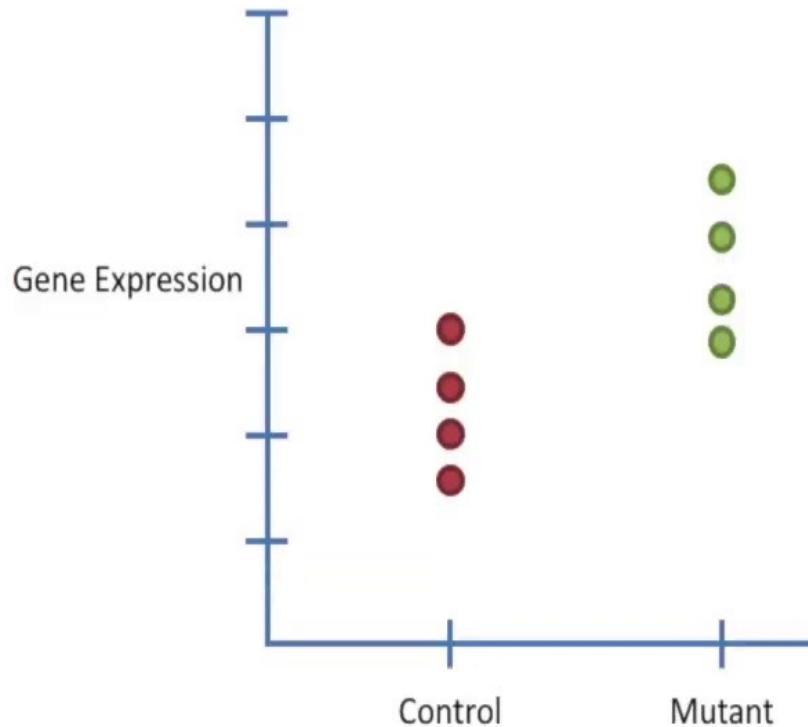
- 1) How useful was mouse weight for predicting mouse size? ( $R^2$  told us this).



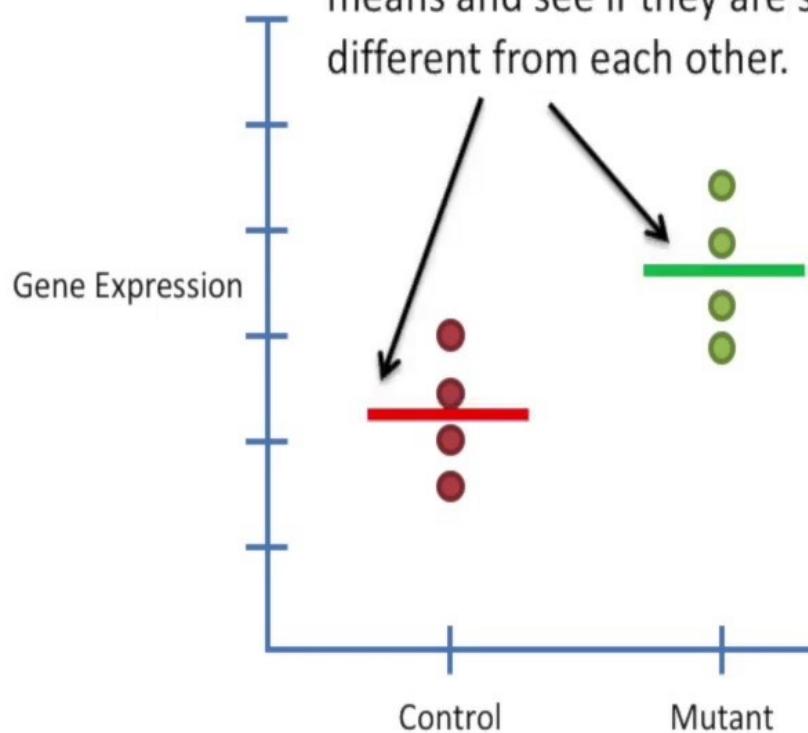
We measured **mouse weight** and **mouse size** and we wanted learn 2 things from it:

- 1) How useful was mouse weight for predicting mouse size? ( $R^2$  told us this).
- 2) Was that relationship due to chance? (The **p-value** told us this.)

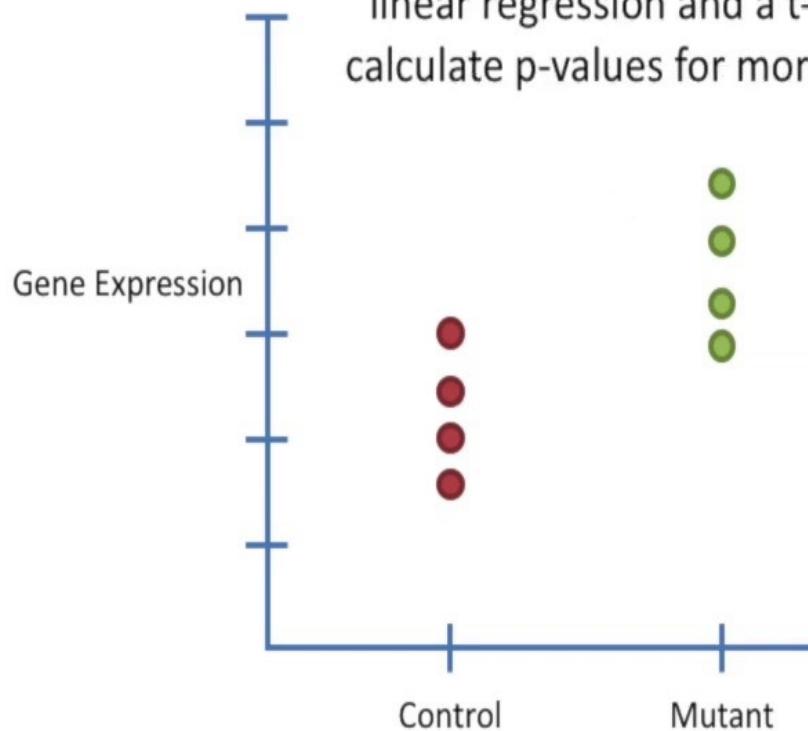
Now let's see if we can apply those concepts to a t-test.



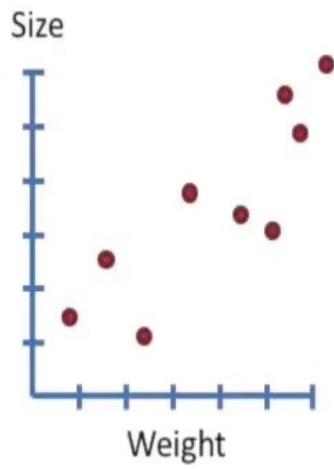
The goal of a t-test is to compare means and see if they are significantly different from each other.



If the same method can calculate  $p$ -values for a linear regression and a t-test, then we can easily calculate  $p$ -values for more complicated situations.

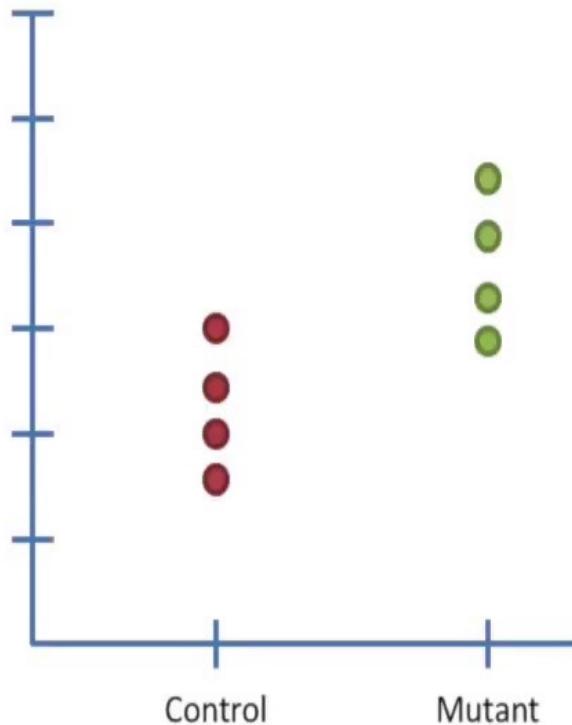


## Linear Regression



t-test

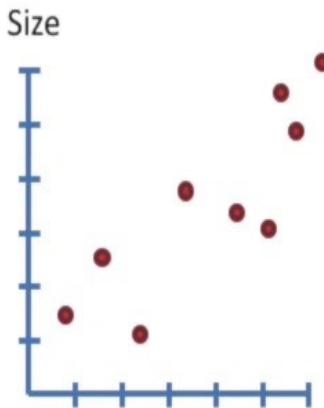
## Gene Expression



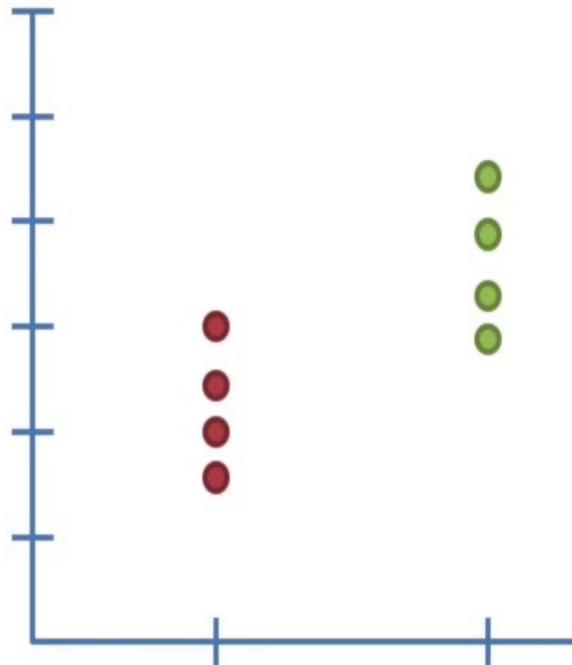
Step 1: Ignore the x-axis and find the overall mean.

t-test

Linear Regression



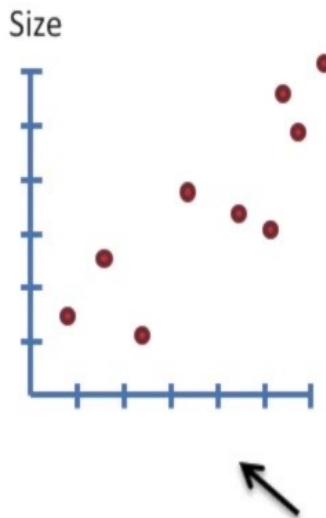
Gene Expression



Step 1: Ignore the x-axis and find the overall mean.

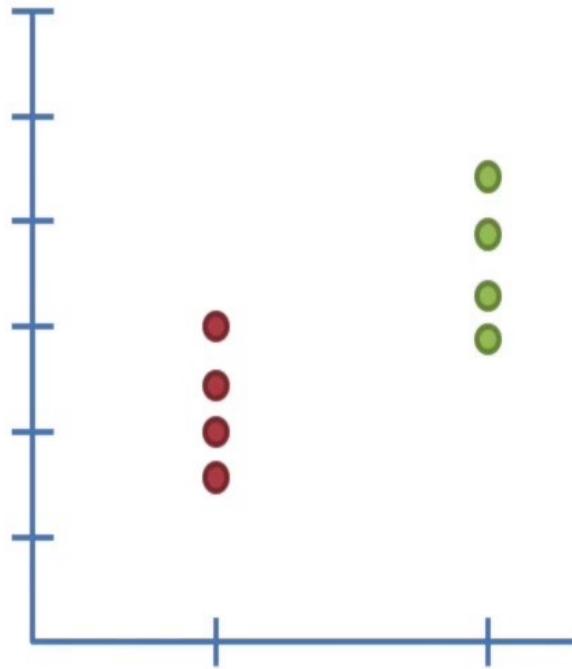
t-test

Linear Regression



To emphasize that we want to focus on the y-axis, I've removed the labels on the x-axis.

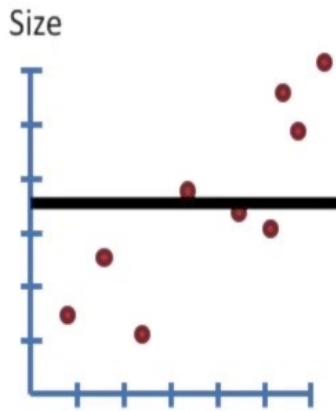
Gene Expression



Step 1: Ignore the x-axis and find the overall mean.

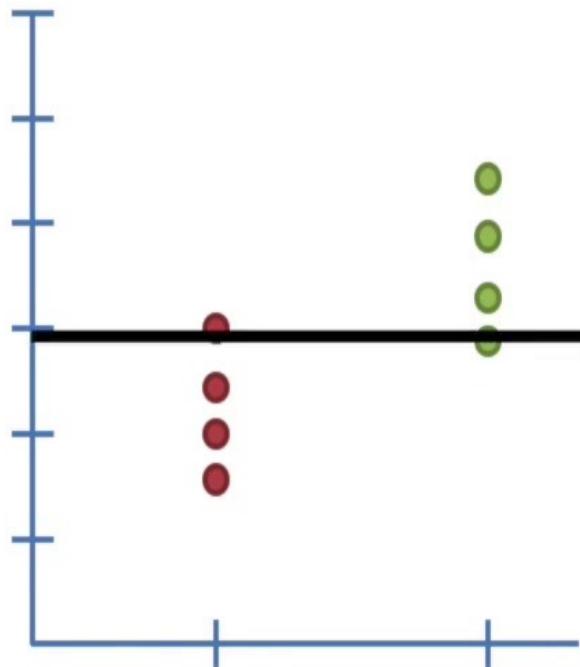
t-test

Linear Regression



Overall means

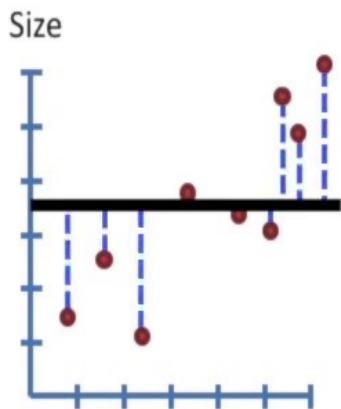
Gene Expression



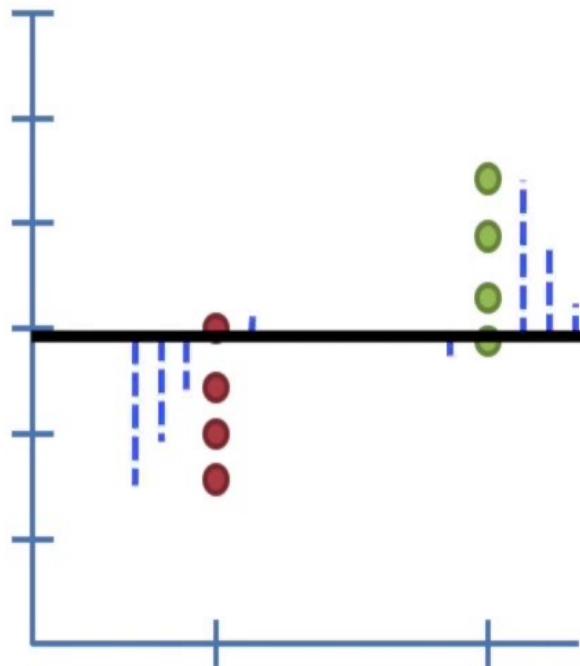
Step 2: Calculate  $SS(\text{mean})$ , the sum of squared residuals around the mean.

t-test

Linear Regression



Gene Expression

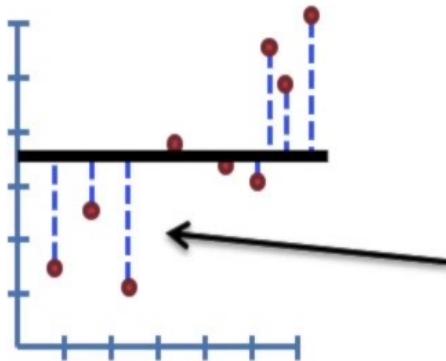


Step 2: Calculate  $SS(\text{mean})$ , the sum of squared residuals around the mean.

t-test

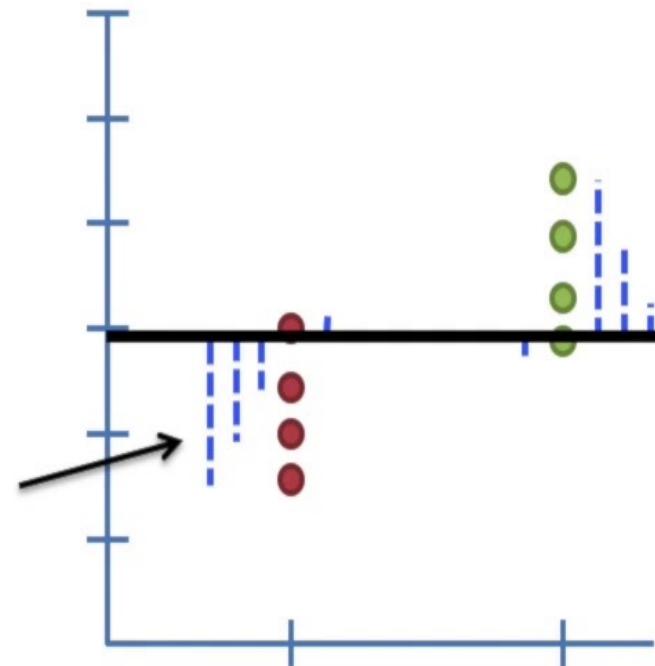
### Linear Regression

Size



These are the residuals, the distance from the data points to the lines.

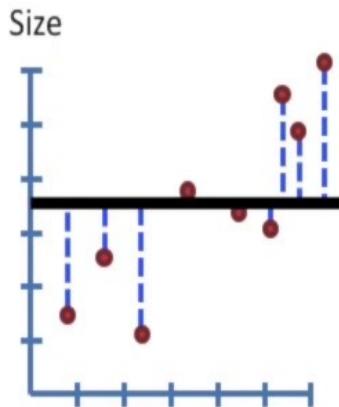
Gene Expression



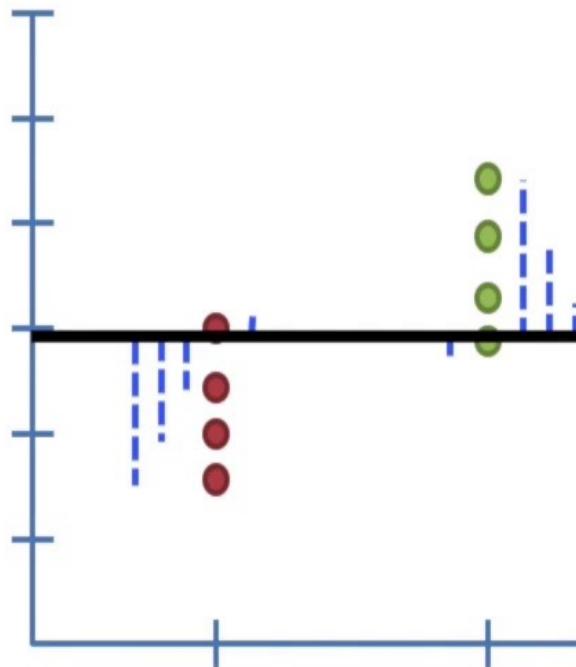
Step 2: Calculate  $SS(\text{mean})$ , the sum of squared residuals around the mean.

t-test

Linear Regression



Gene Expression

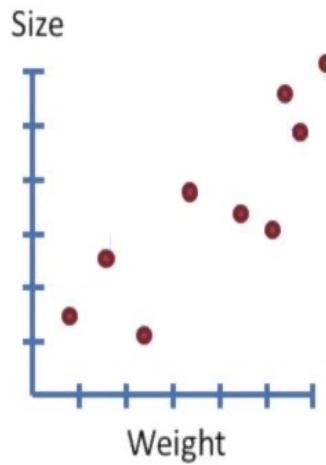


BAM! Calculating  $SS(\text{mean})$  was easy!

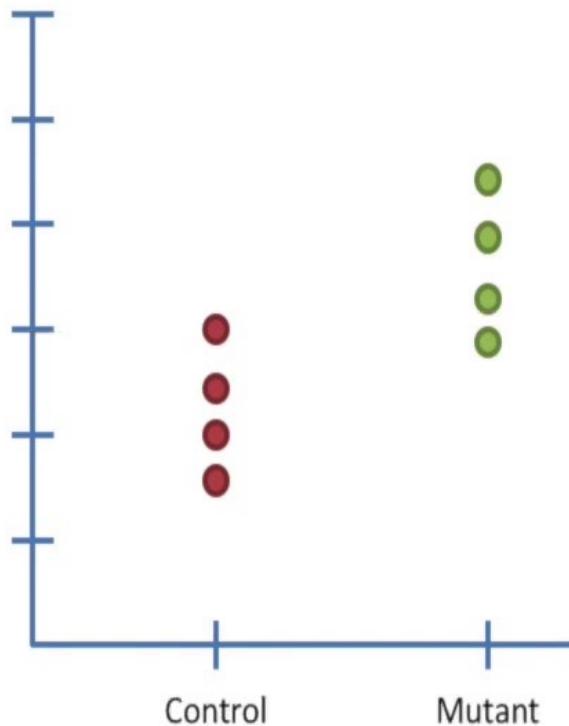
Step 3: Fit a line to the data.

t-test

Linear Regression



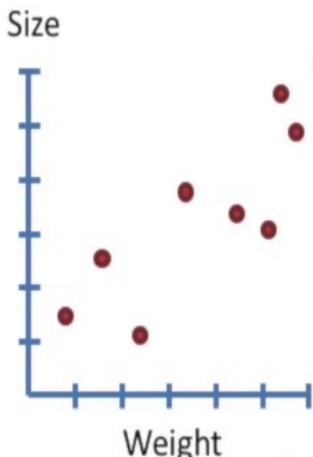
Gene Expression



Step 3: Fit a line to the data.

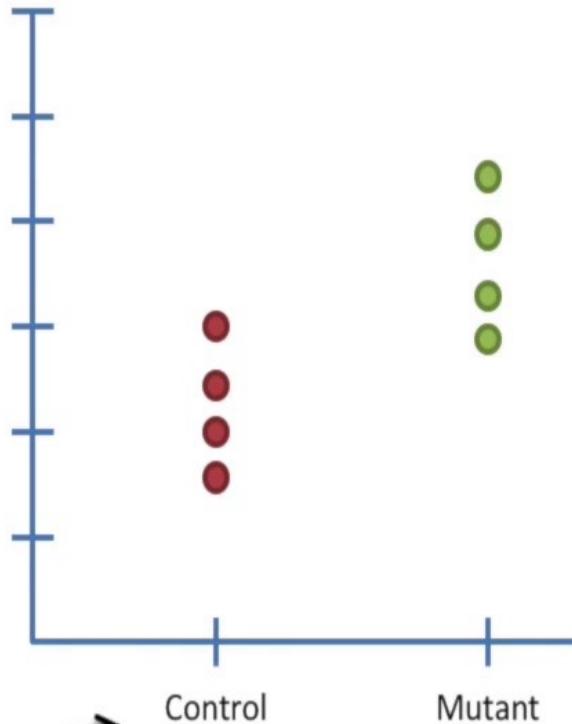
t-test

### Linear Regression



NOTE: This is when we start caring  
about the x-axis again...

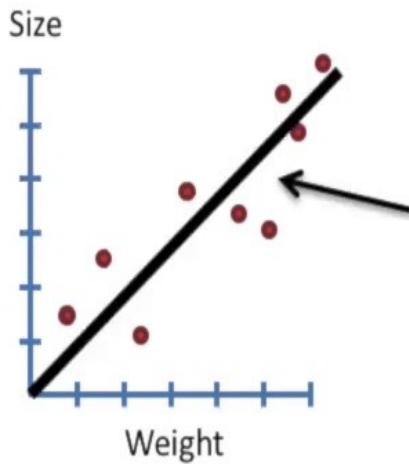
### Gene Expression



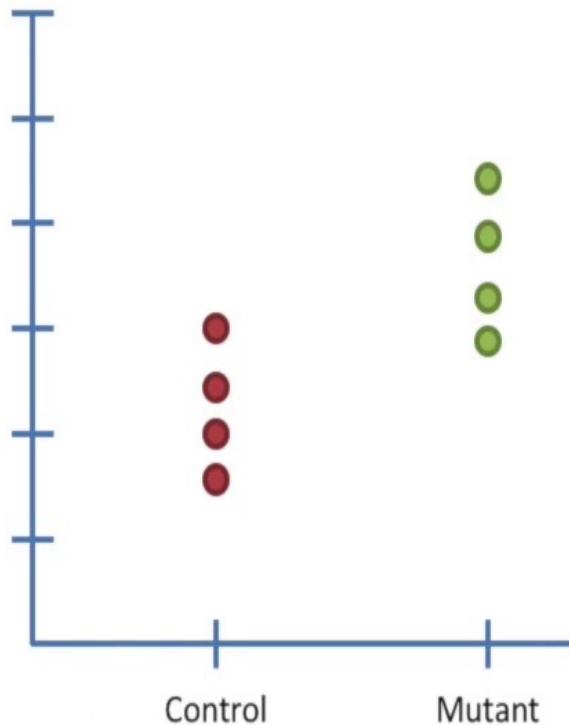
Step 3: Fit a line to the data.

t-test

### Linear Regression



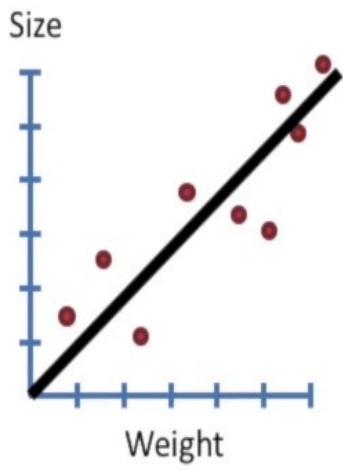
### Gene Expression



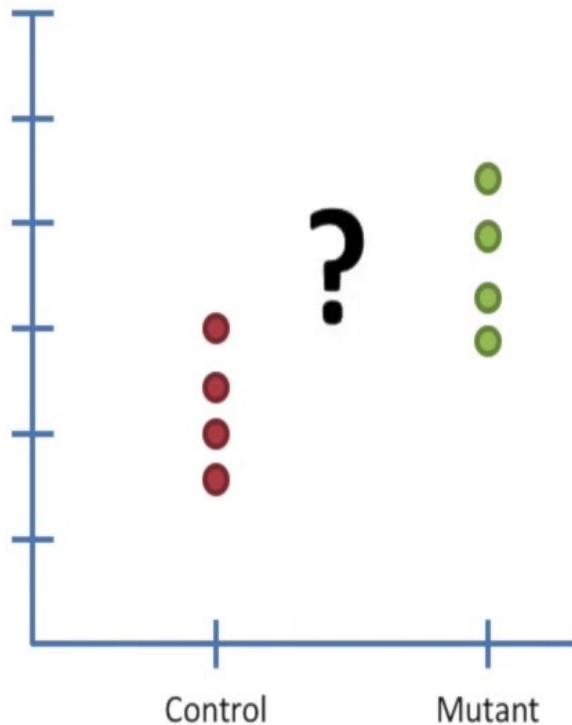
Step 3: Fit a line to the data.

t-test

Linear Regression



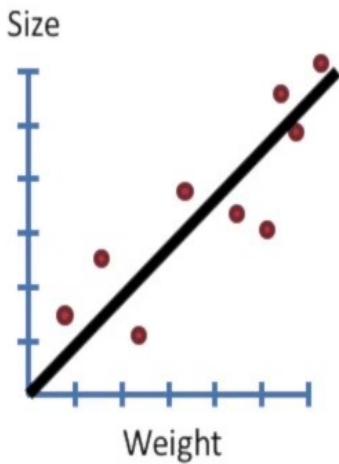
Gene Expression



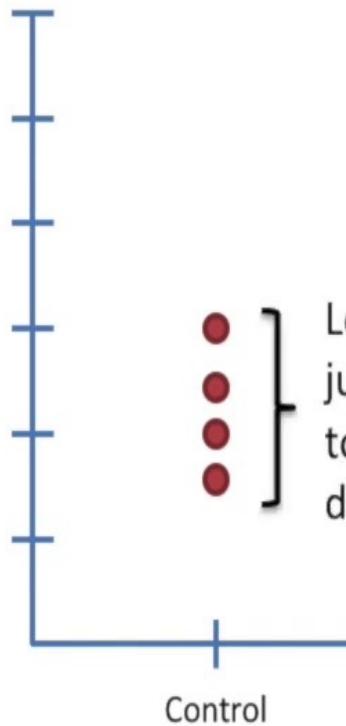
Step 3: Fit a line to the data.

t-test

Linear Regression



Gene Expression

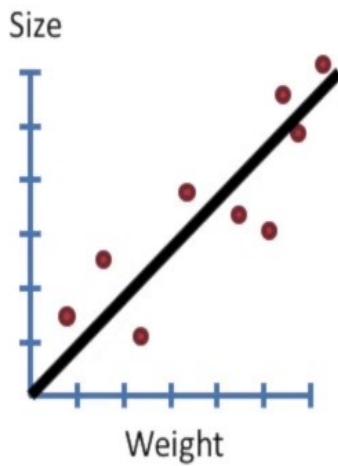


Let's start by  
just fitting a line  
to the control  
data.

Step 3: Fit a line to the data.

t-test

### Linear Regression



### Gene Expression

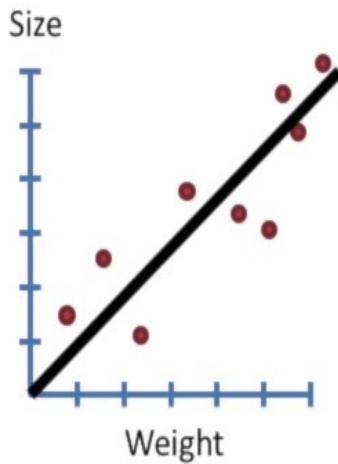


We start by finding a least squares fit to the control data.

Step 3: Fit a line to the data.

t-test

### Linear Regression



### Gene Expression



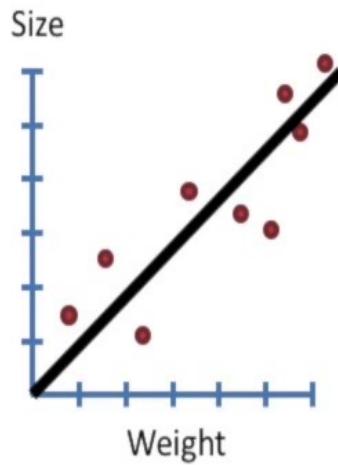
We start by finding a least squares fit to the control data.

It turns out that the mean is the least-squares fit.

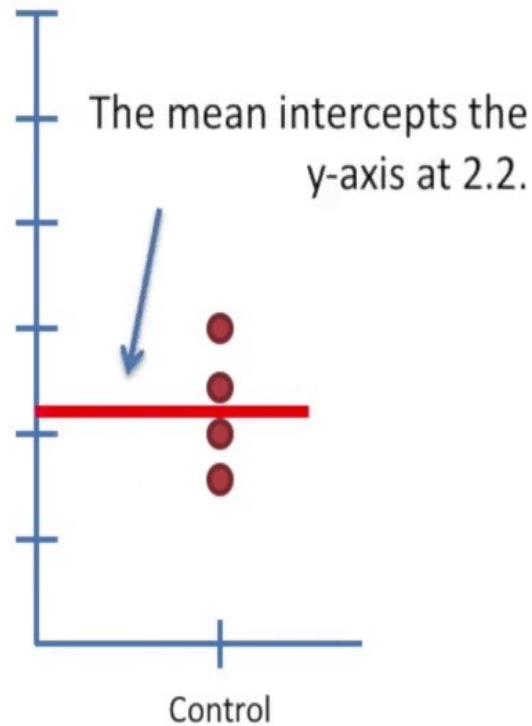
Step 3: Fit a line to the data.

t-test

Linear Regression



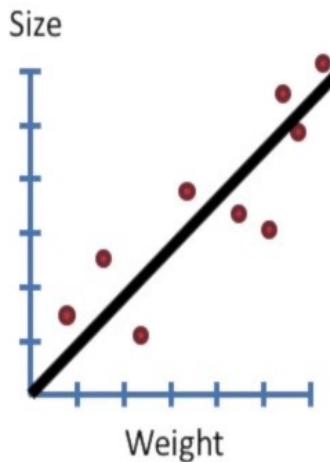
Gene Expression



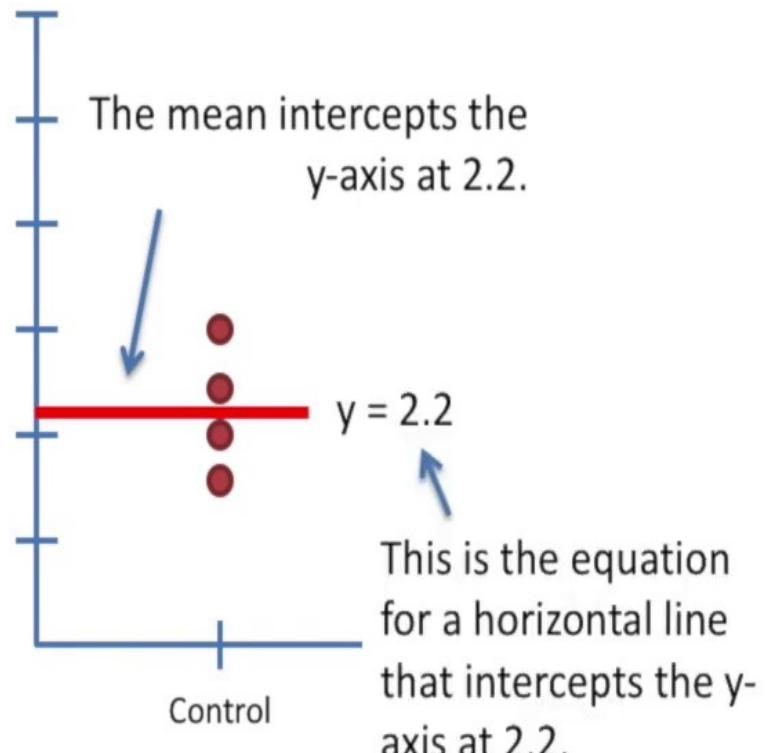
Step 3: Fit a line to the data.

t-test

Linear Regression



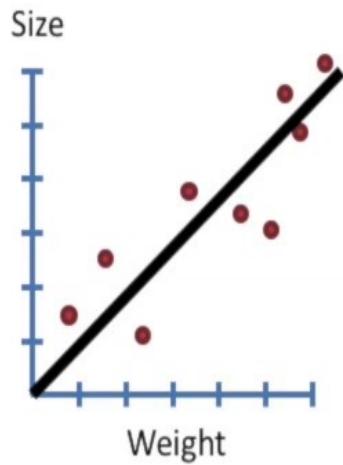
Gene Expression



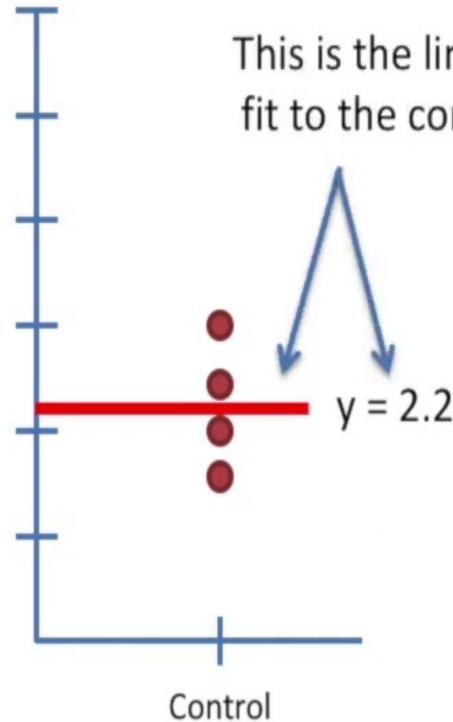
Step 3: Fit a line to the data.

t-test

Linear Regression



Gene Expression

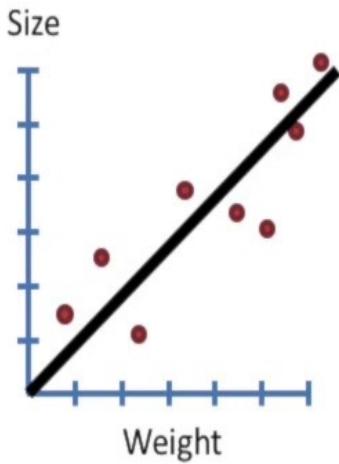


This is the line that we fit to the control data.

Step 3: Fit a line to the data.

t-test

Linear Regression



Gene Expression



Now lets fit a  
line to the  
mutant data.

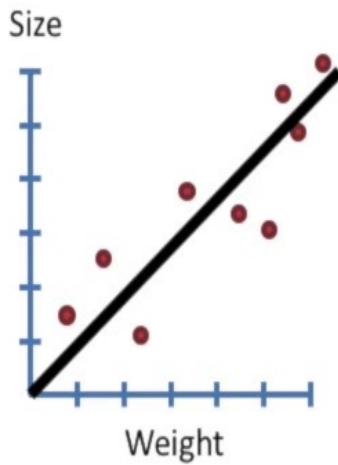


Mutant

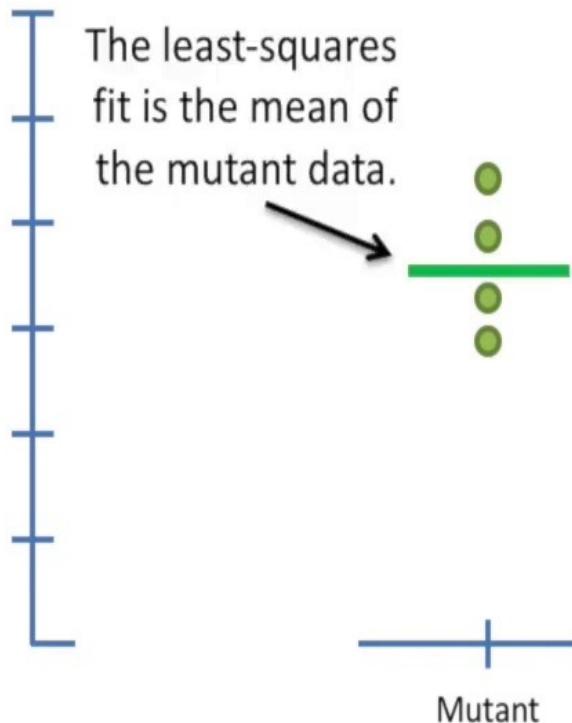
Step 3: Fit a line to the data.

t-test

### Linear Regression



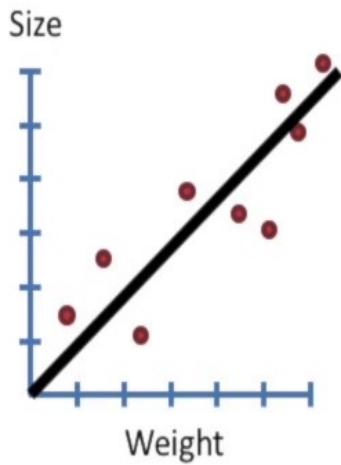
### Gene Expression



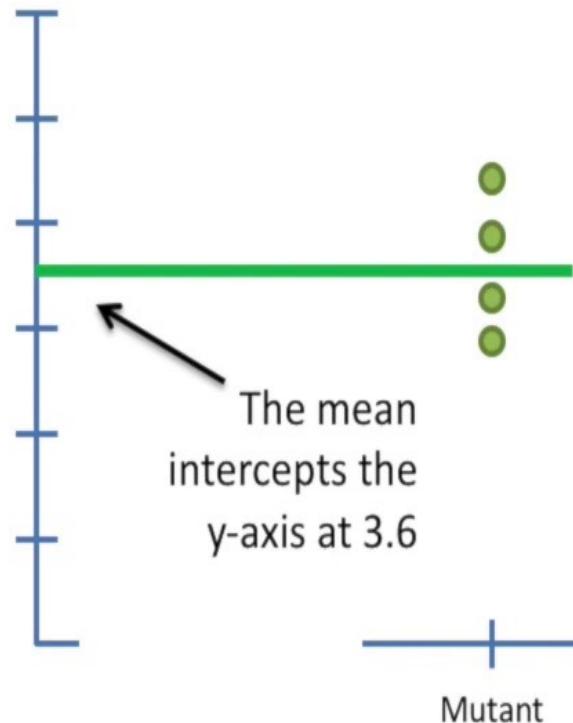
Step 3: Fit a line to the data.

t-test

Linear Regression



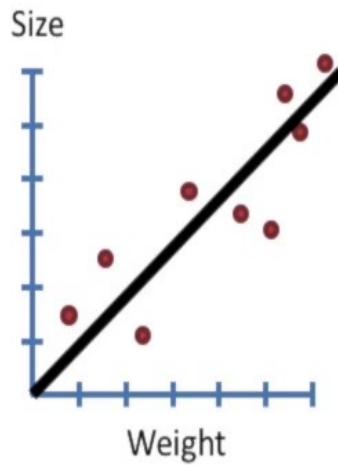
Gene Expression



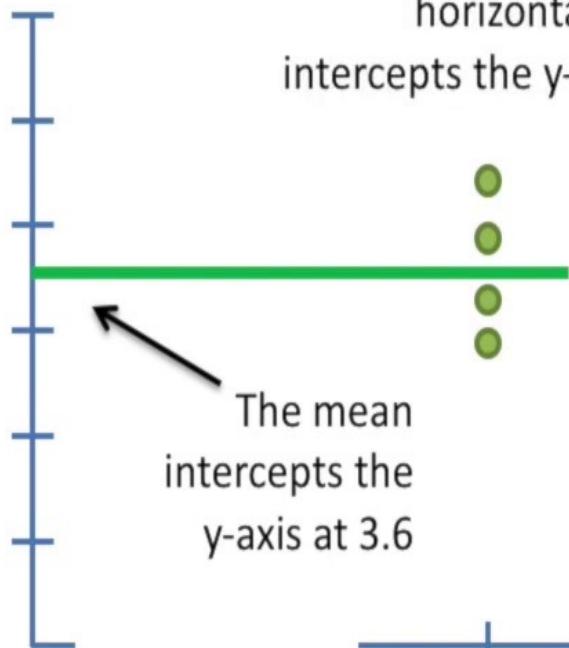
Step 3: Fit a line to the data.

t-test

### Linear Regression



Gene Expression



This is the equation for a horizontal line that intercepts the y-axis at 3.6

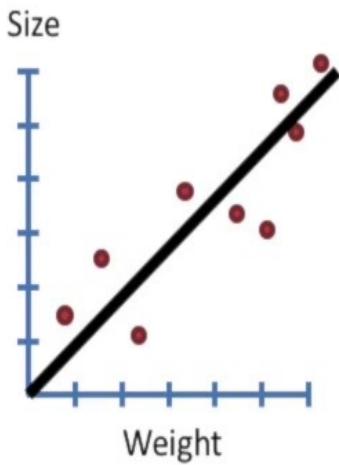
The mean  
intercepts the  
y-axis at 3.6

Mutant

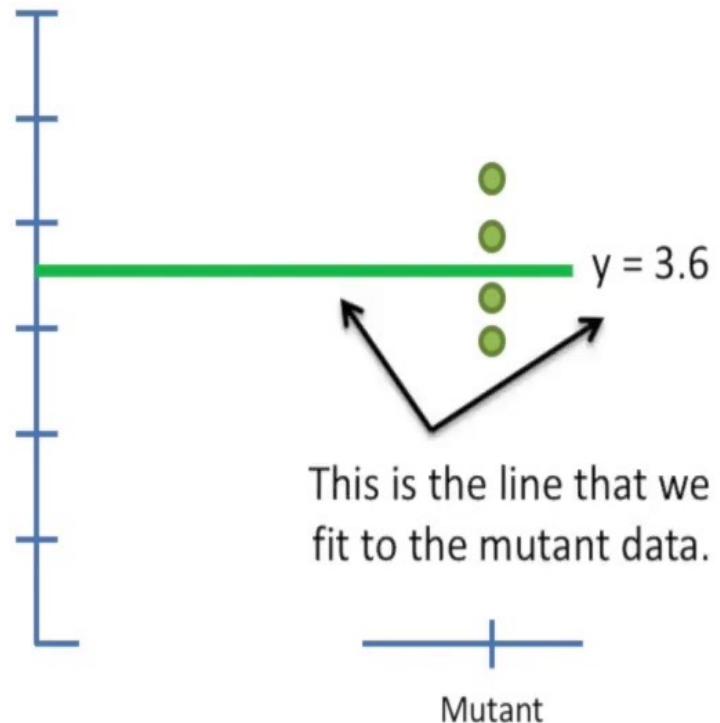
Step 3: Fit a line to the data.

t-test

Linear Regression



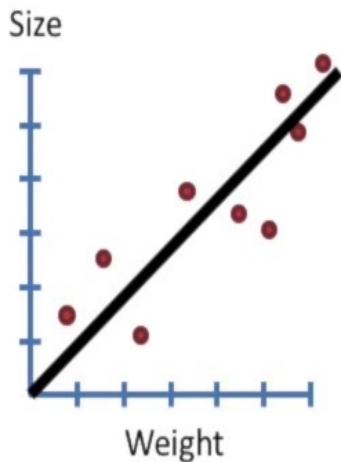
Gene Expression



Step 3: Fit a line to the data.

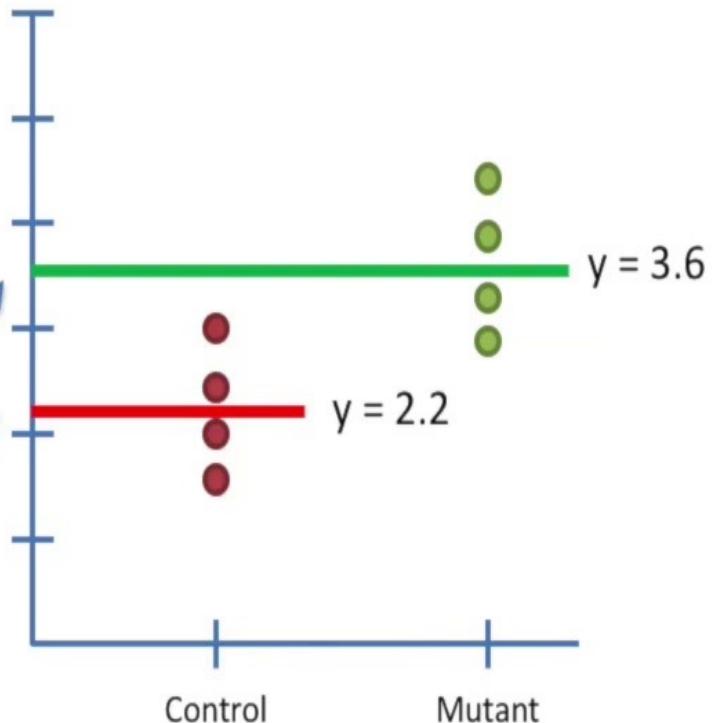
t-test

Linear Regression



We have fit two lines  
to the data.

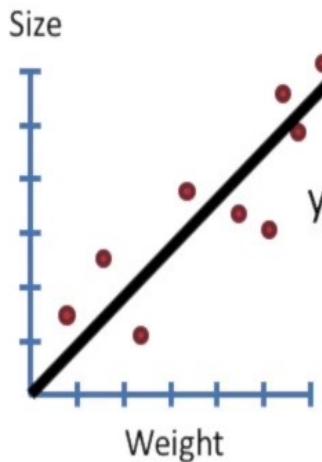
Gene Expression



Step 3: Fit a line to the data.

t-test

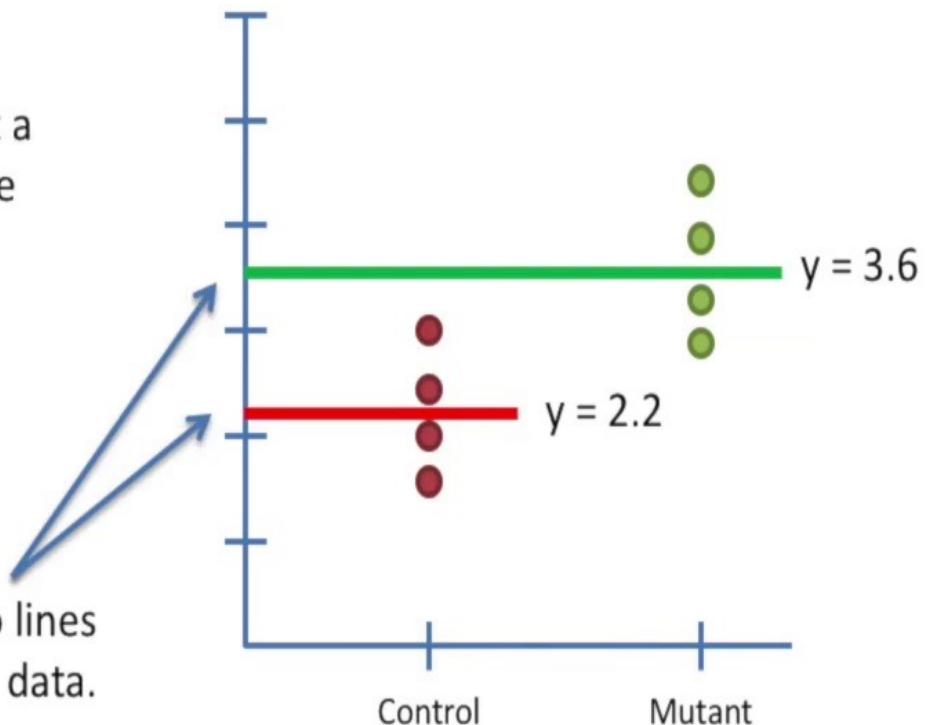
### Linear Regression



Originally we fit a single line to the data.

We have fit two lines to the data.

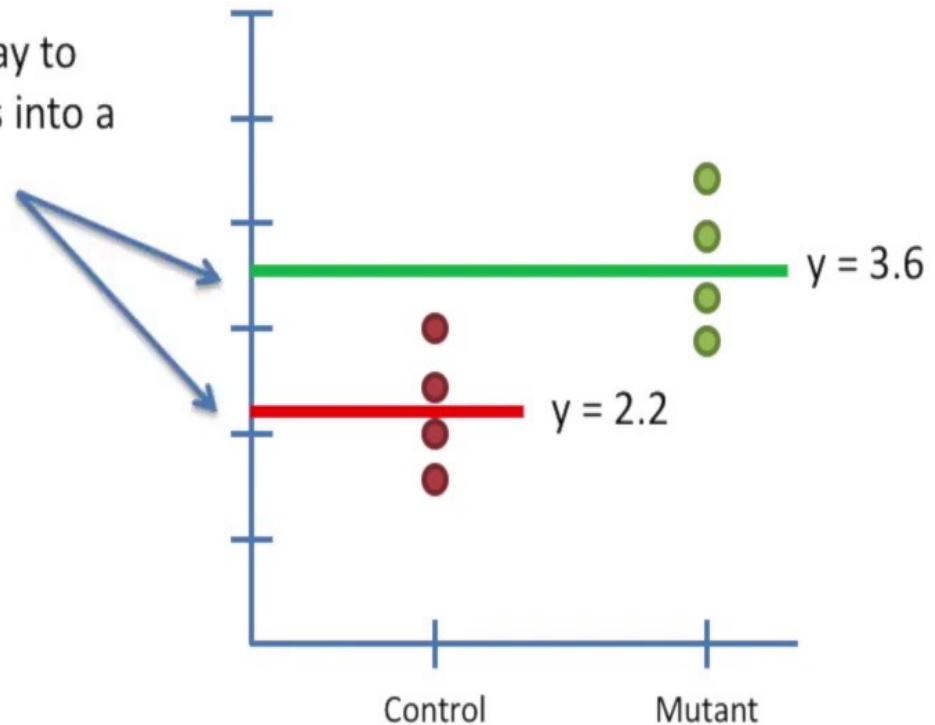
### Gene Expression



t-test

Gene Expression

However, there is a way to combine these two lines into a single equation.

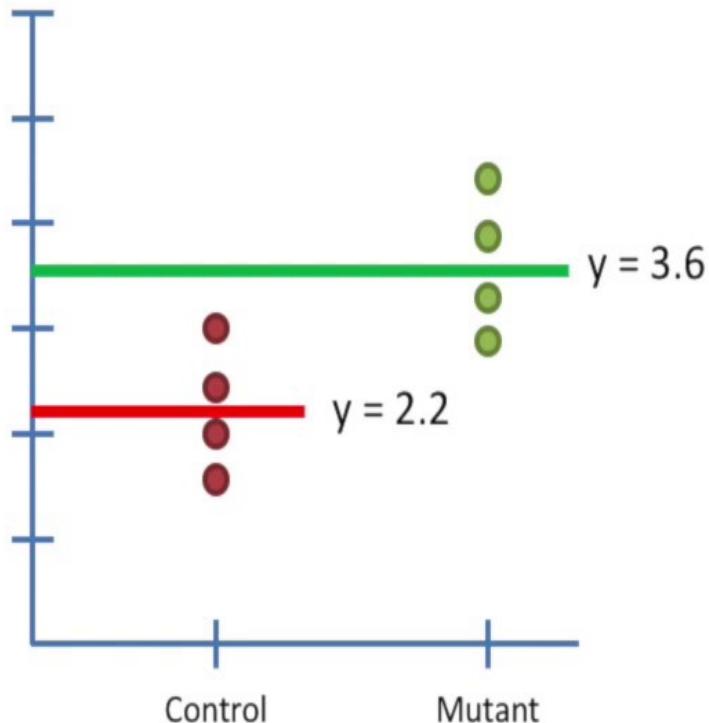


## t-test

Gene Expression

However, there is a way to combine these two lines into a single equation.

This will make the steps for computing “ $F$ ” the exact same for the regression and the t-test, which, in turn, means a computer can do it automatically.



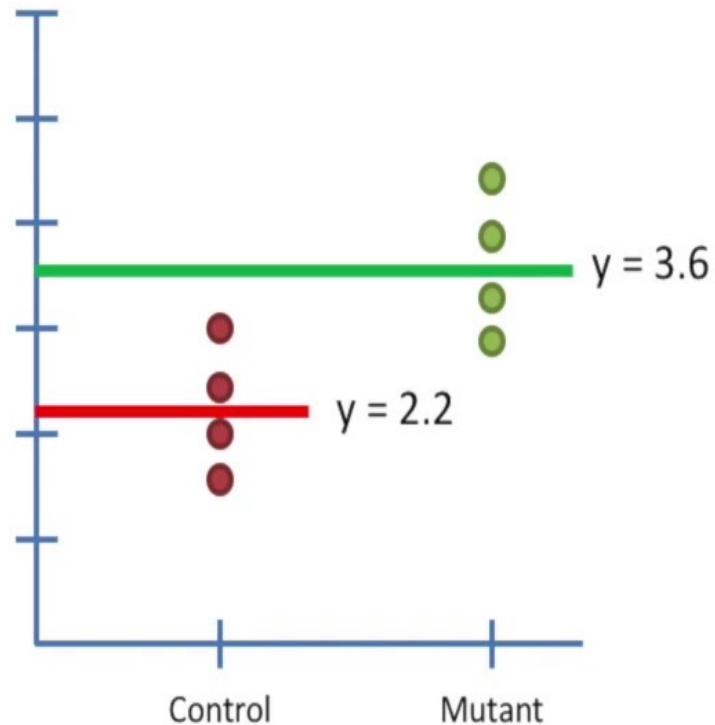
## t-test

### Gene Expression

However, there is a way to combine these two lines into a single equation.

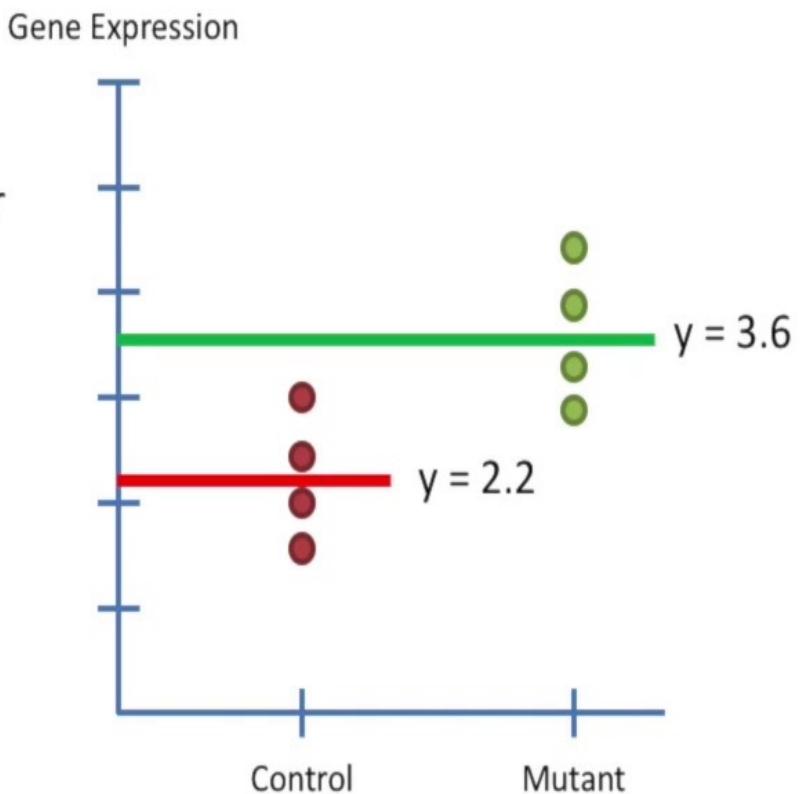
This will make the steps for computing “ $F$ ” the exact same for the regression and the t-test, which, in turn, means a computer can do it automatically.

This is key, because we don’t want to do this by hand, ever....



t-test

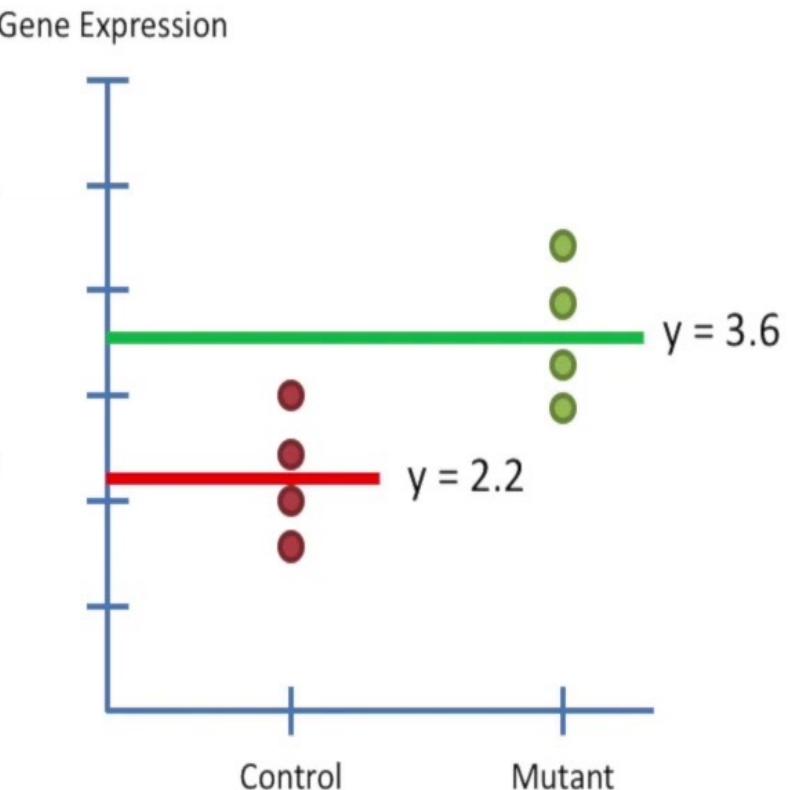
This is going to look weird, but just bear  
with me.



## t-test

This is going to look weird, but just bear with me.

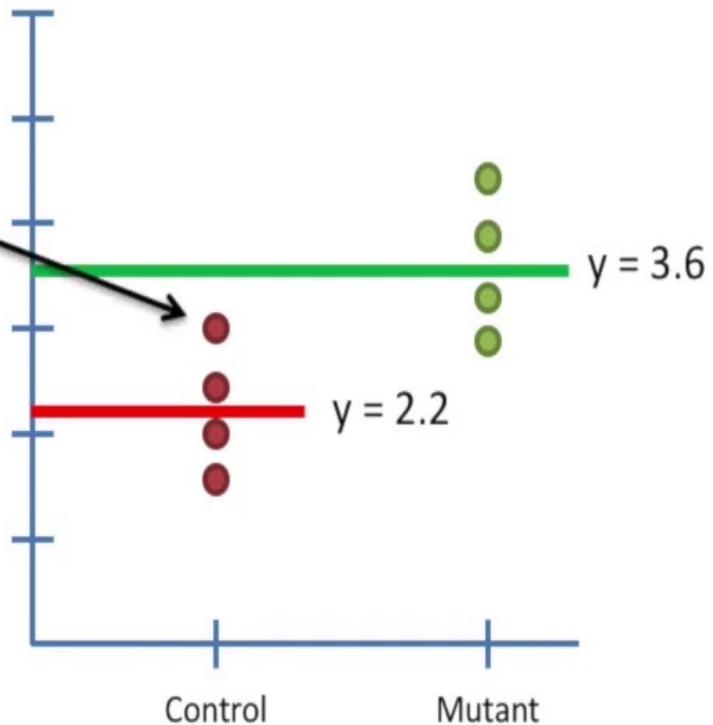
Keep in mind that the goal is to have a flexible way for a computer to solve this, and every other “least-squares” based problem, without having to create a whole new method each time.



$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

This is the equation (which combines both lines) for this point

Gene Expression

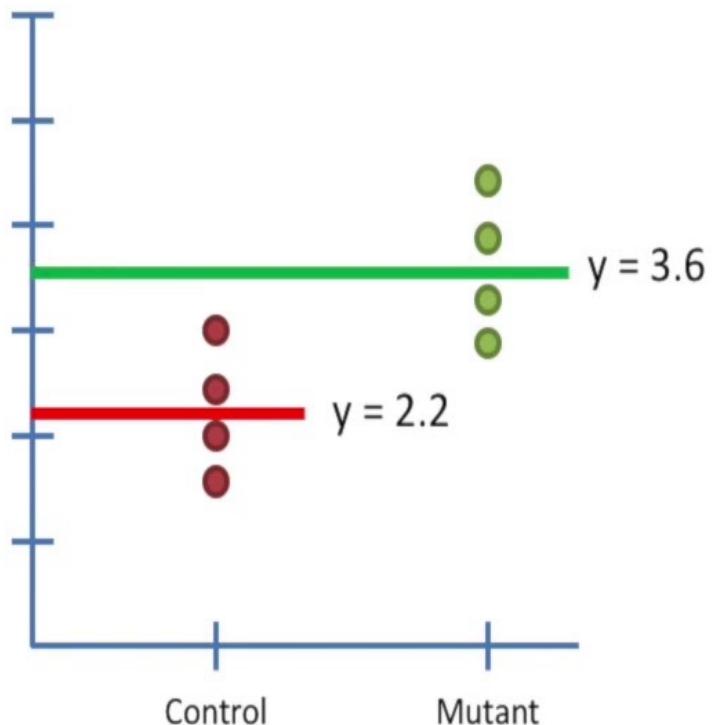


$$y = 1 \times \boxed{2.2} - 0 \times 3.6 + \text{the residual}$$

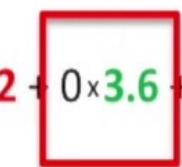


1 times the mean of the **control data**.

Gene Expression

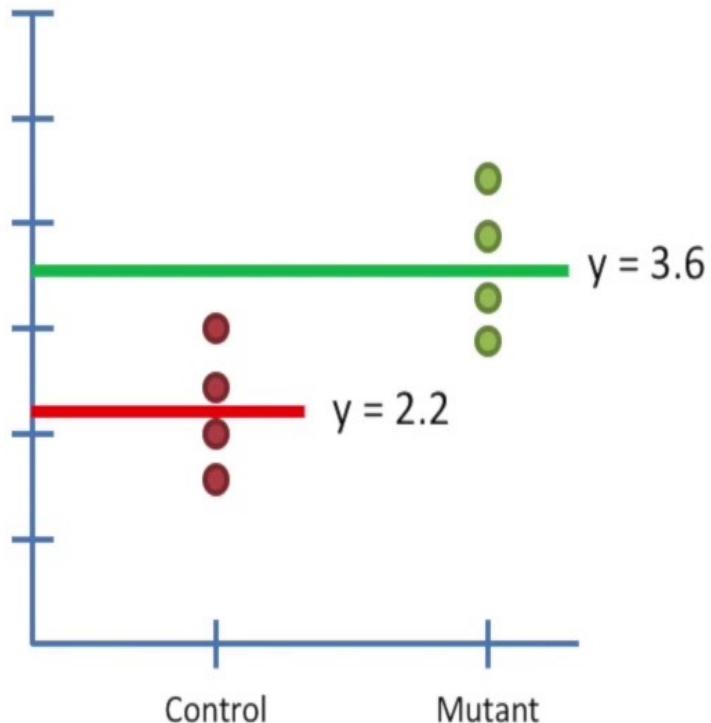


$y = 1 \times 2.2 + 0 \times 3.6$  + the residual



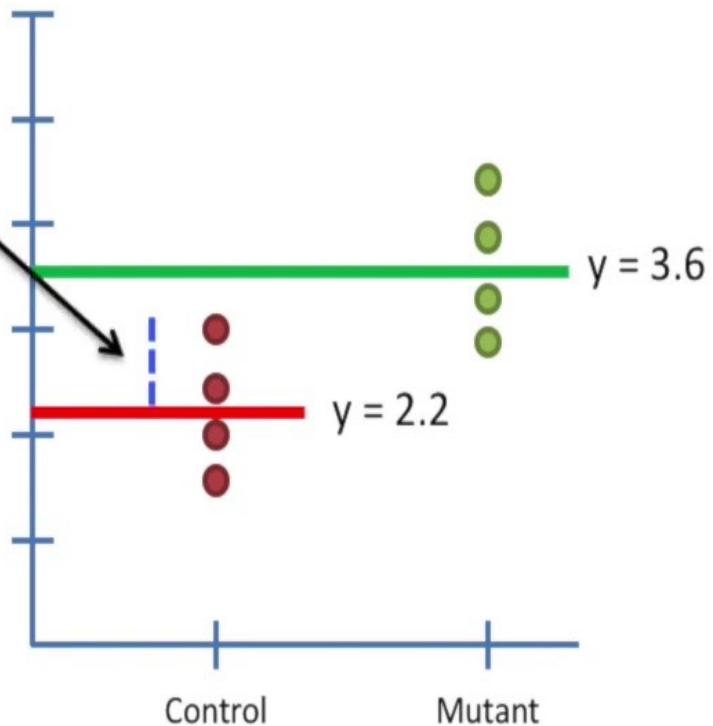
↑  
0 times the mean of the **mutant data**.

Gene Expression



$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$

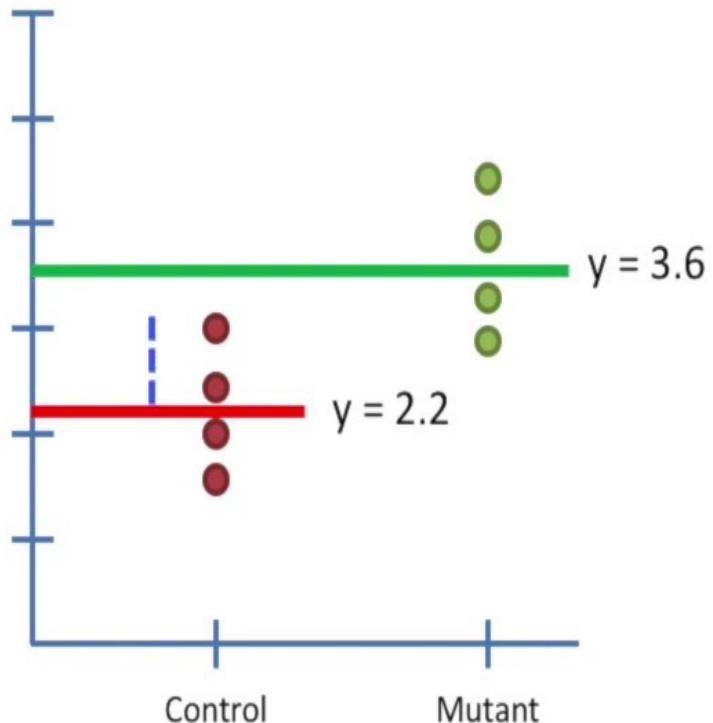
Gene Expression



$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

Yes, this is strange, especially multiplying the mutant mean by 0, but bear with me.

Gene Expression



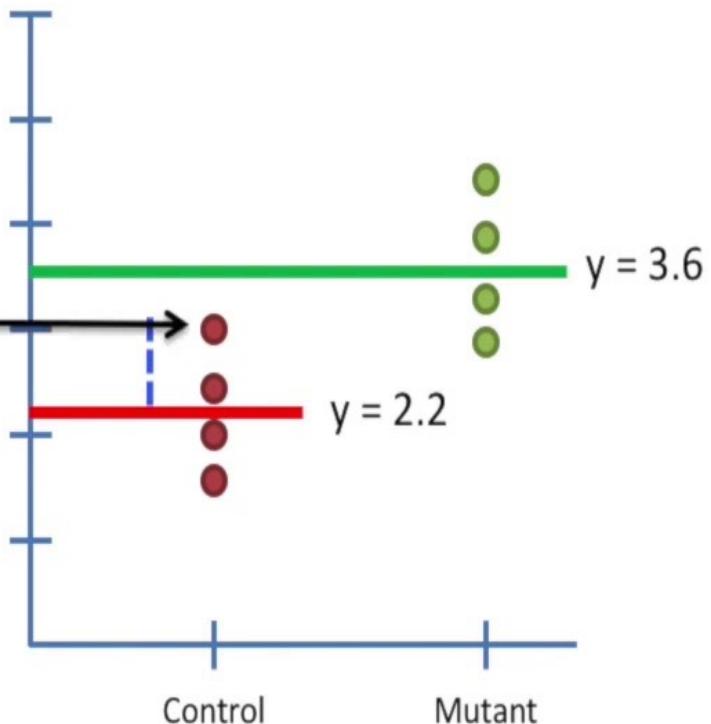
$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

Yes, this is strange, especially multiplying the mutant mean by 0, but bear with me.

If we multiplied things out, the equation for this point would be...

$$y = 2.2 + \text{the residual}$$

Gene Expression



$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

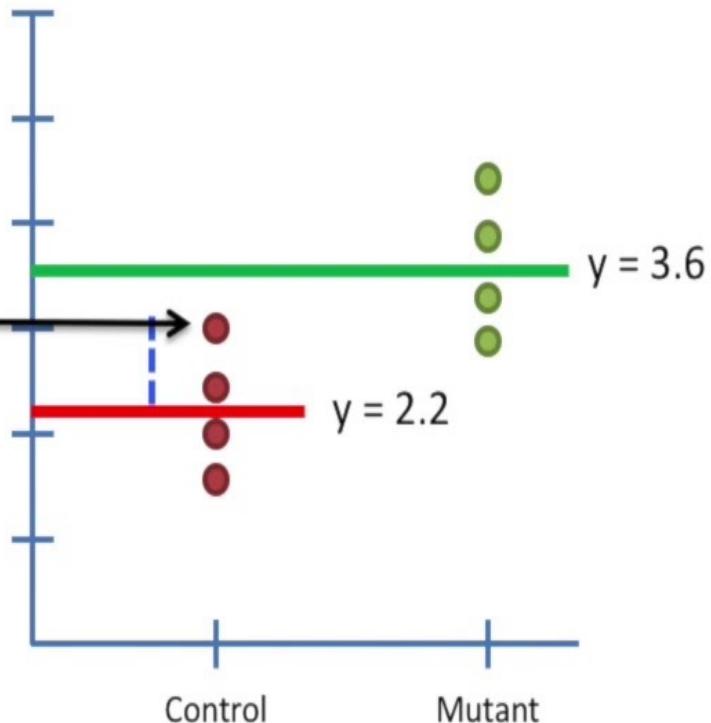
Yes, this is strange, especially multiplying the mutant mean by 0, but bear with me.

If we multiplied things out, the equation for this point would be...

$$y = 2.2 + \text{the residual}$$

... and that sort of makes sense. But just bear with me.

Gene Expression

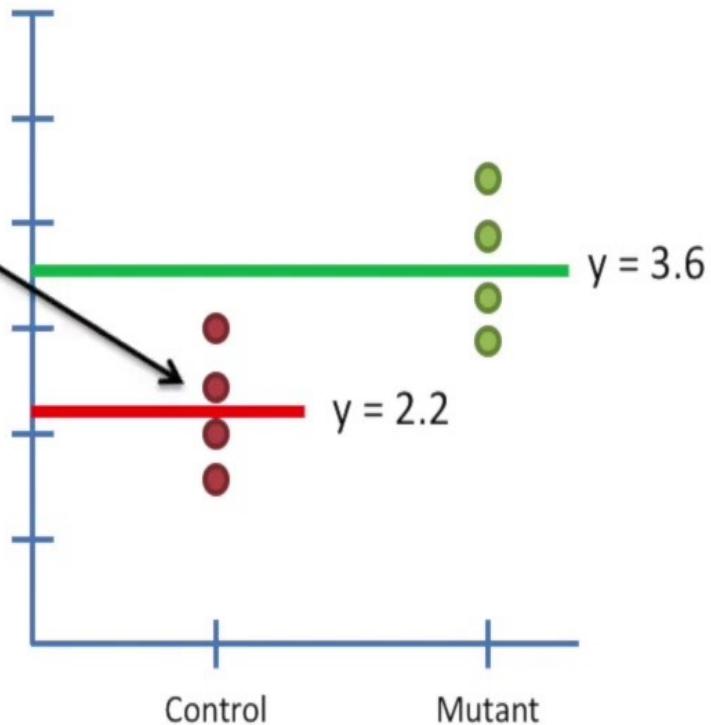


$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

This is the equation for the next point.

Gene Expression

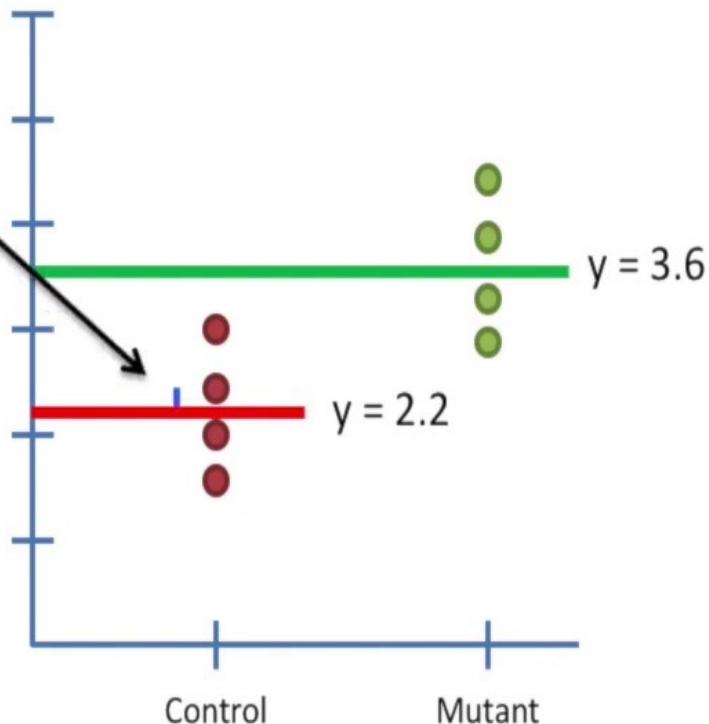


$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

The only difference is the residual. This one is smaller.

Gene Expression



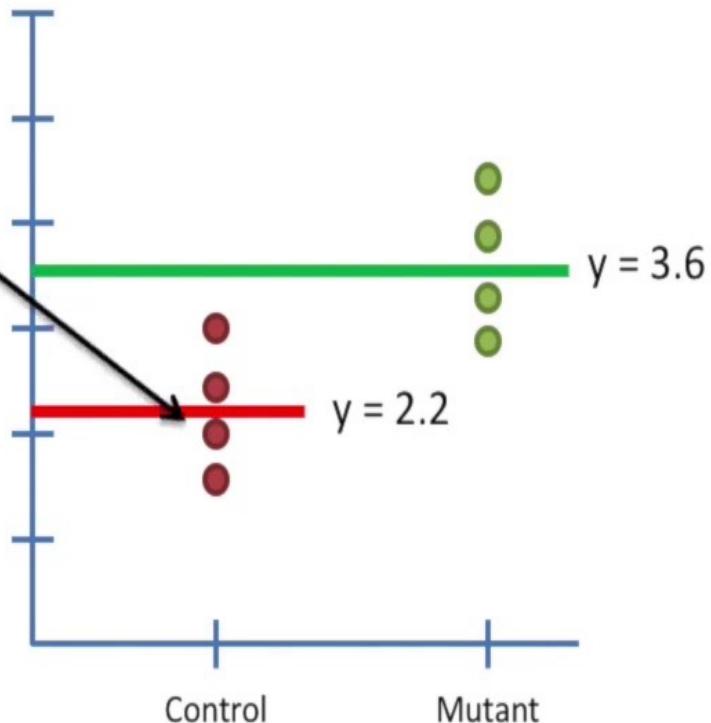
$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$

$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$

$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$

This is the equation for the next point.

Gene Expression



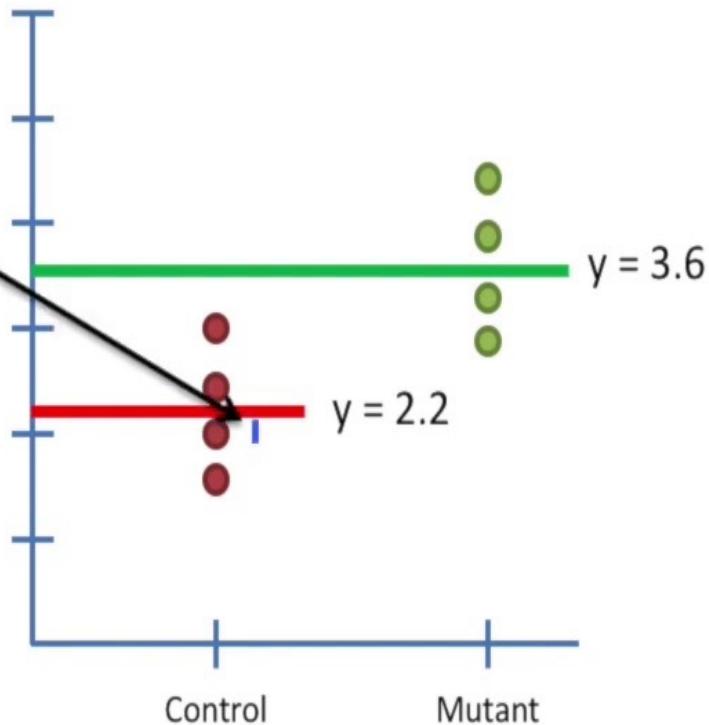
$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

Again, the only difference is the residual.

Gene Expression



$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$

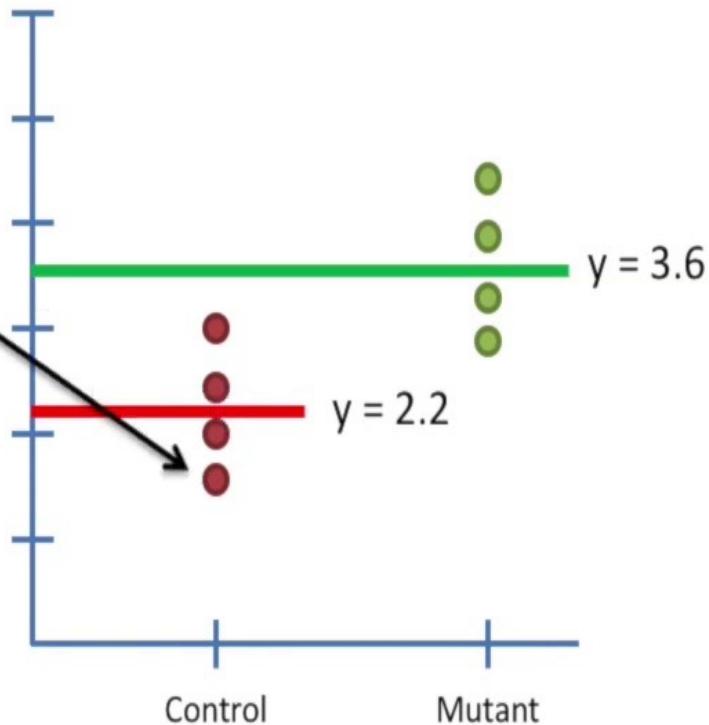
$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$

$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$

$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$

This is the equation for the  
next point.

Gene Expression



$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

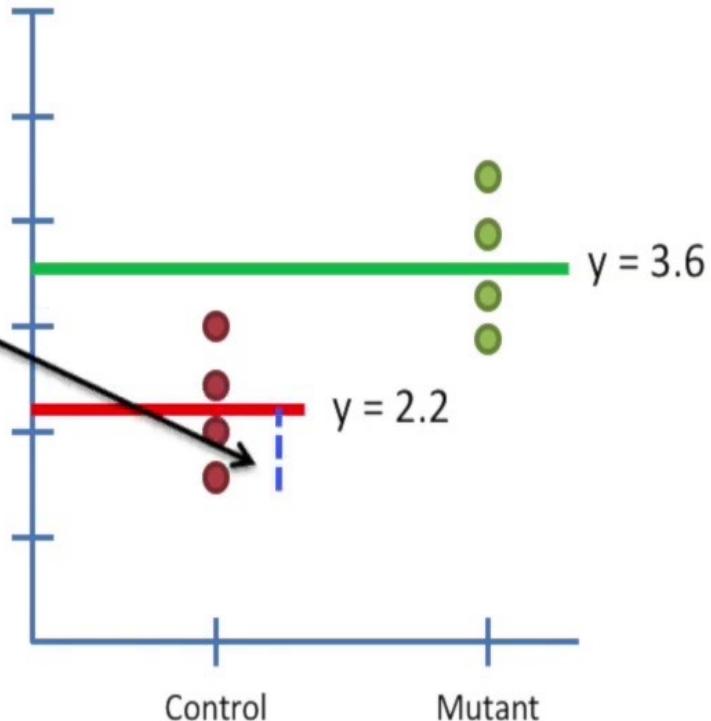
$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

Again, the only difference is the residual.

Gene Expression



$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

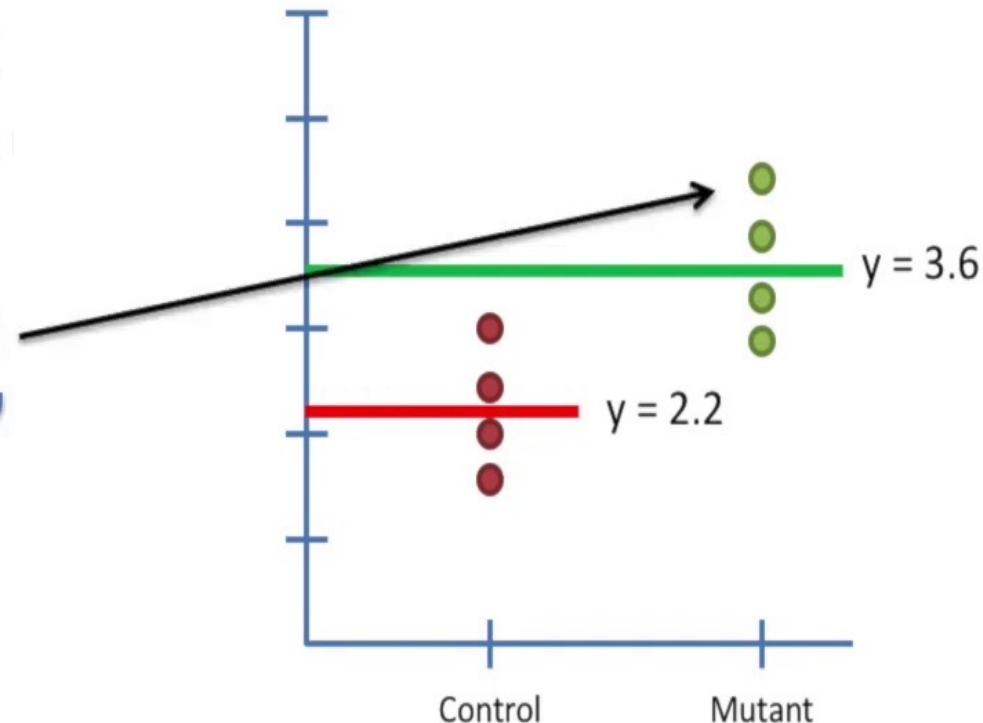
$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 0 \times 2.2 + 1 \times 3.6 + \text{the residual}$$

This is the equation for the first point in the mutant dataset.

Gene Expression



$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

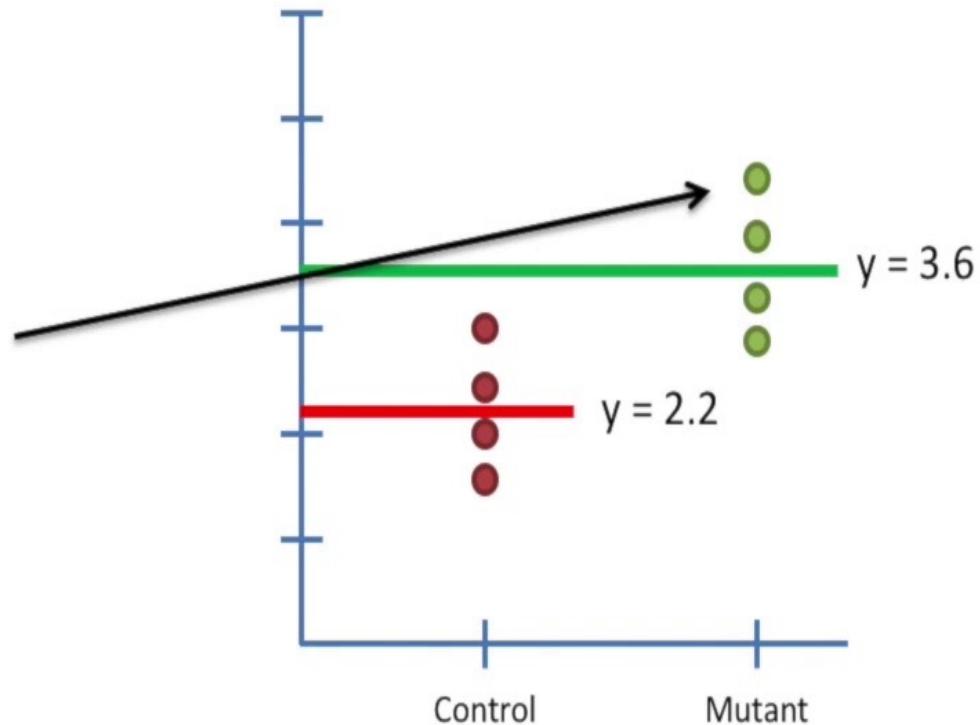
$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 0 \times 2.2 - 1 \times 3.6 + \text{the residual}$$

↑  
Now we are  
multiplying the **control**  
mean by 0...

Gene Expression



$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

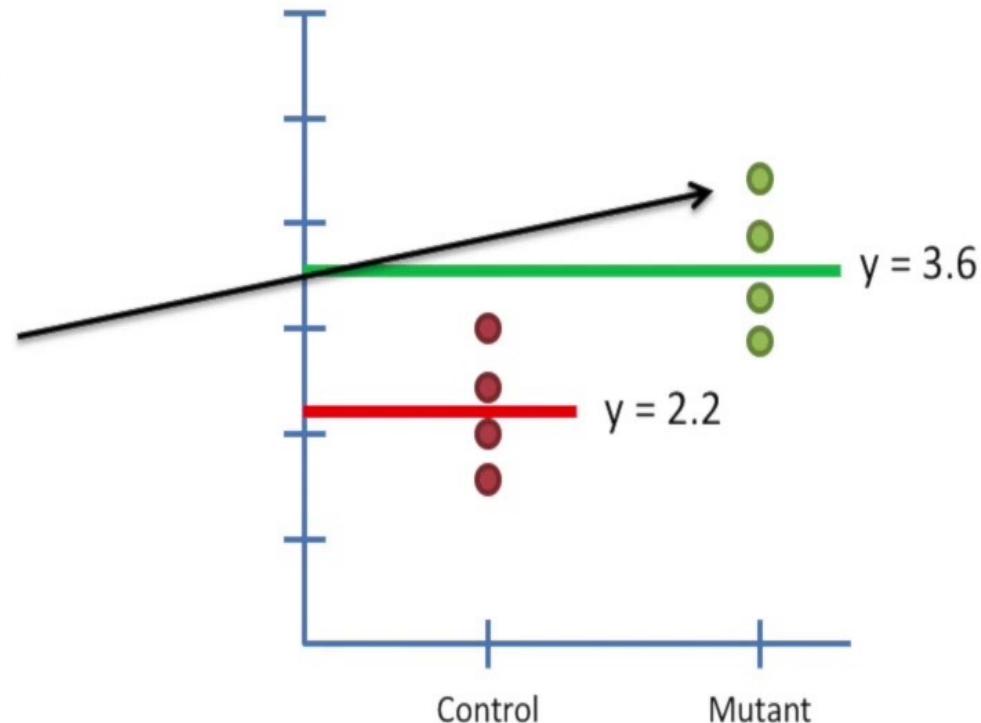
$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 0 \times 2.2 + 1 \times 3.6 - \text{the residual}$$

...and multiplying the  
**mutant** mean by 1.

Gene Expression



$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

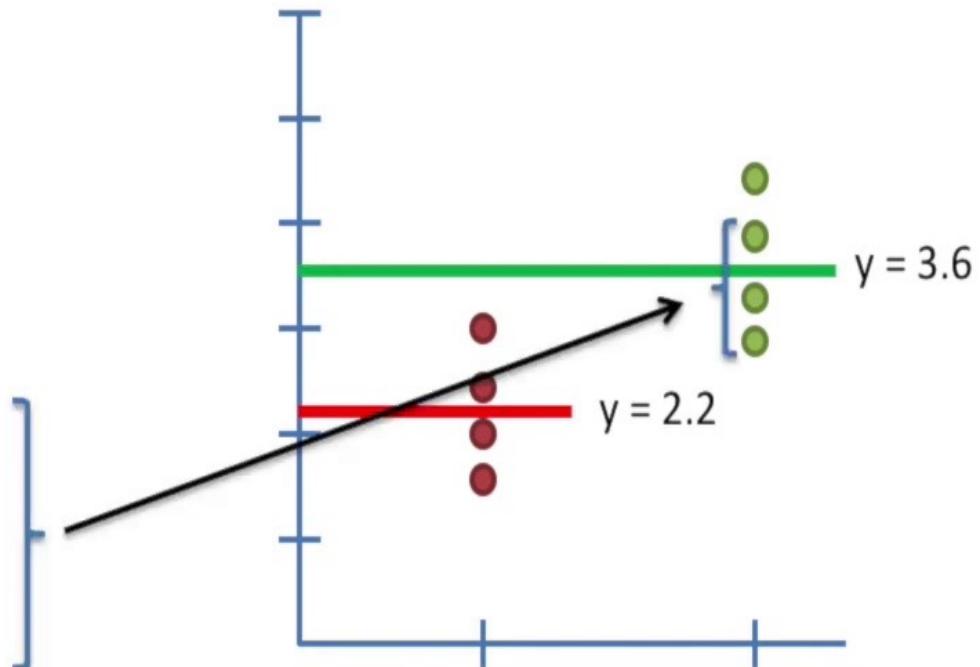
$$y = 0 \times 2.2 + 1 \times 3.6 + \text{the residual}$$

$$y = 0 \times 2.2 + 1 \times 3.6 + \text{the residual}$$

$$y = 0 \times 2.2 + 1 \times 3.6 + \text{the residual}$$

$$y = 0 \times 2.2 + 1 \times 3.6 + \text{the residual}$$

Gene Expression



These are the  
equations for the  
remaining points.

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

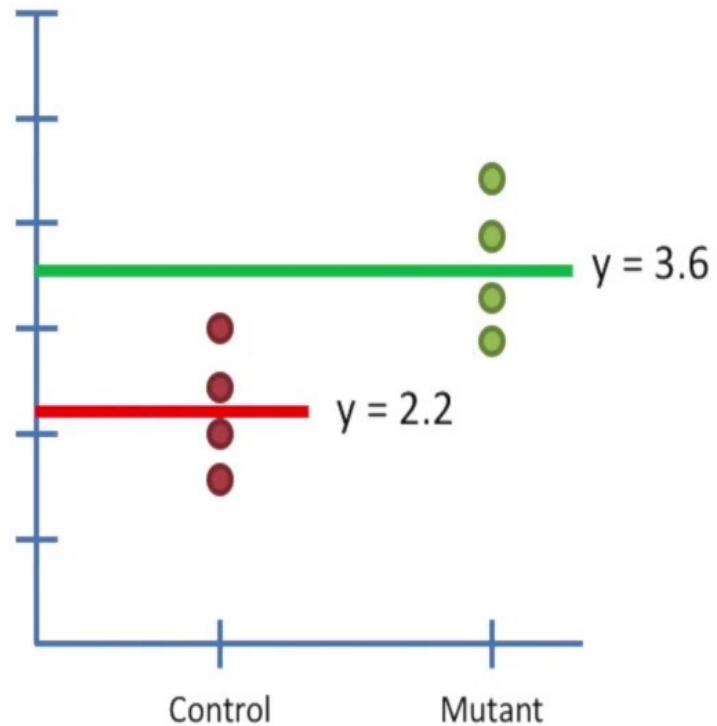
$$y = 0 \times 2.2 + 1 \times 3.6 + \text{the residual}$$

$$y = 0 \times 2.2 + 1 \times 3.6 + \text{the residual}$$

$$y = 0 \times 2.2 + 1 \times 3.6 + \text{the residual}$$

$$y = 0 \times 2.2 + 1 \times 3.6 + \text{the residual}$$

Gene Expression



Now let's focus on the 1's and 0's.

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

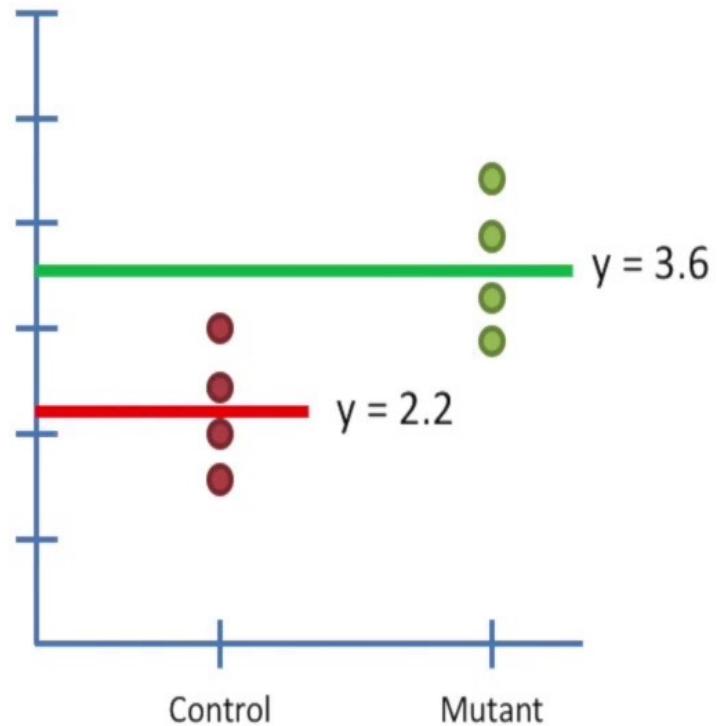
$$y = 0 \times 2.2 + 1 \times 3.6 + \text{the residual}$$

$$y = 0 \times 2.2 + 1 \times 3.6 + \text{the residual}$$

$$y = 0 \times 2.2 + 1 \times 3.6 + \text{the residual}$$

$$y = 0 \times 2.2 + 1 \times 3.6 + \text{the residual}$$

Gene Expression

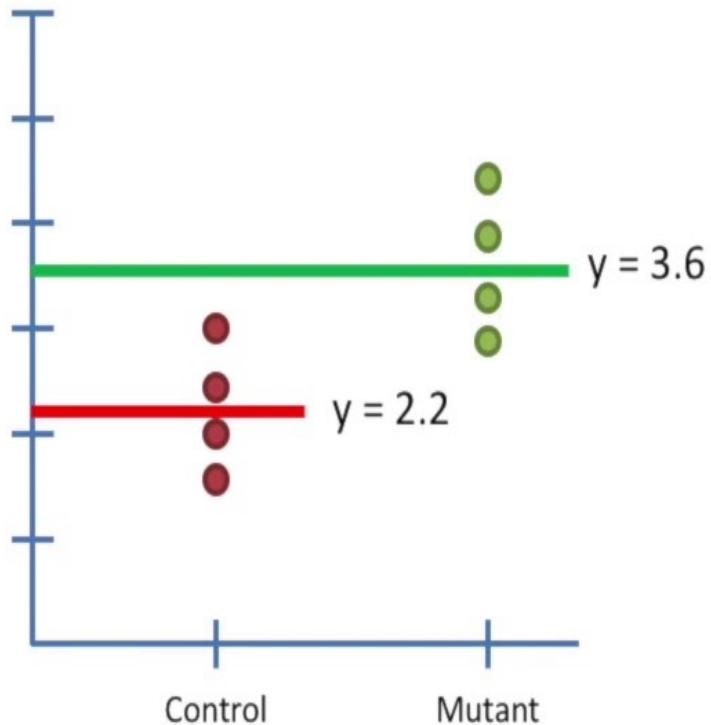


They function like on/off switches for the two means.

1	0
1	0
1	0
1	0
0	1
0	1
0	1
0	1

When we isolate the 1s and 0s, they form a matrix called a “design matrix”.

Gene Expression

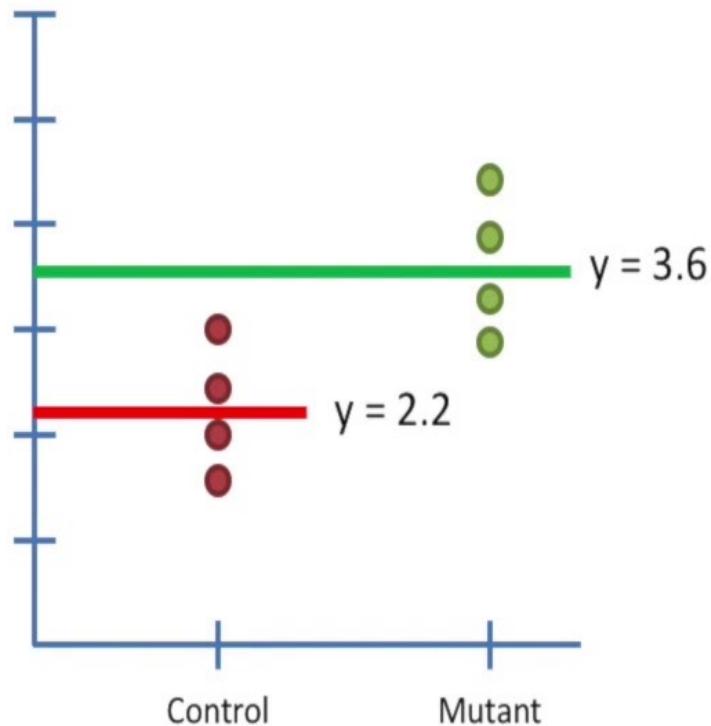


1	0
1	0
1	0
1	0
0	1
0	1
0	1
0	1

The design matrix can be combined with an abstract version of the equation to represent a “fit” to the data.



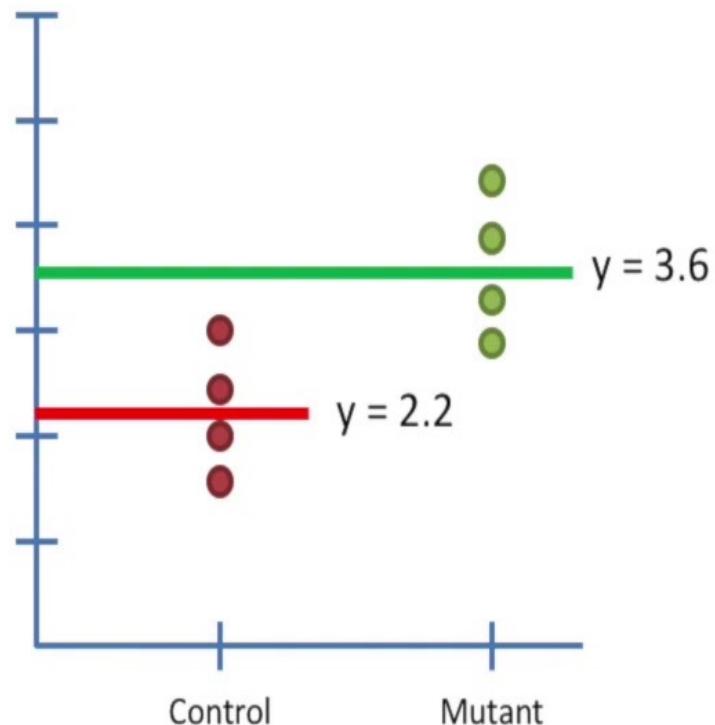
$$y = \text{column1} \times 2.2 + \text{column2} \times 3.6$$



1	0
1	0
1	0
1	0
0	1
0	1
0	1
0	1

$$y = \text{column1} \times 2.2 + \text{column2} \times 3.6$$

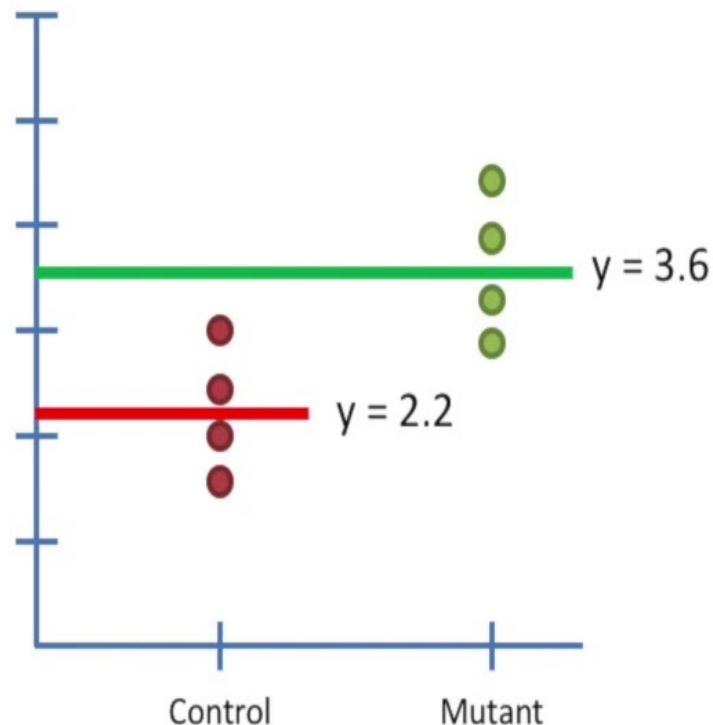
Column1 turns the **control** mean “on” or “off”.



1	0
1	0
1	0
1	0
0	1
0	1
0	1
0	1

$$y = \text{column1} \times 2.2 + \text{column2} \times 3.6$$

Column2 turns the **mutant** mean “on” or “off”.

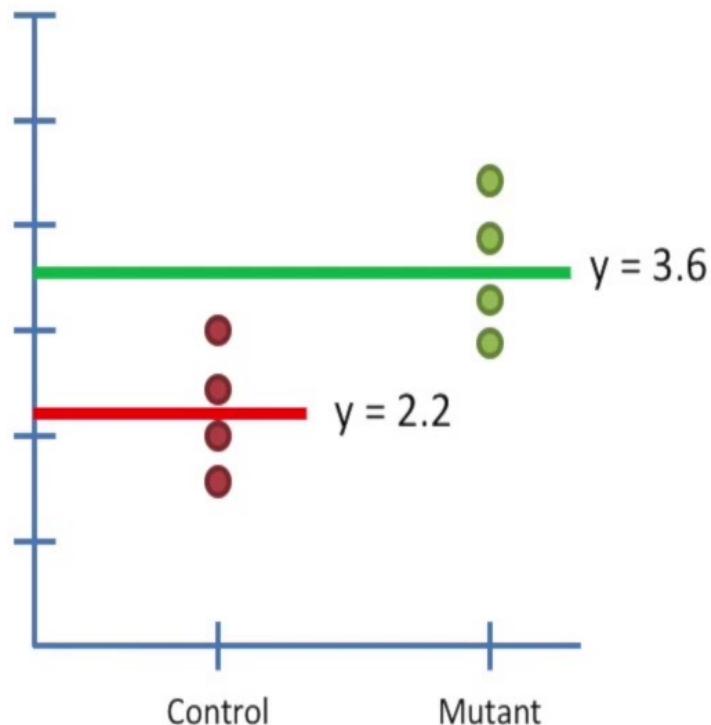


1	0
1	0
1	0
1	0
0	1
0	1
0	1
0	1

$$y = \text{column1} \times 2.2 + \text{column2} \times 3.6$$

In practice, the role of each column is assumed, and the equation is written out like this:

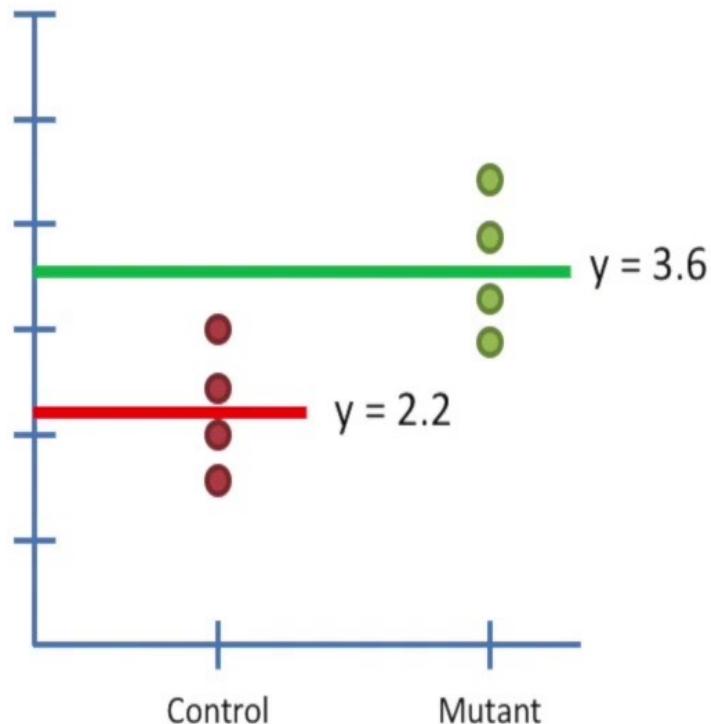
$$y = \text{mean}_{\text{control}} + \text{mean}_{\text{mutant}}$$



1	0
1	0
1	0
1	0
0	1
0	1
0	1
0	1

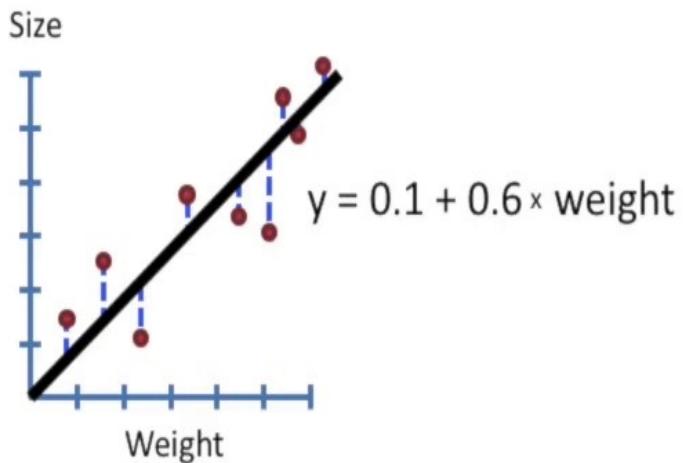
Now that we have the “fit” for the control and mutant data down to a single equation (plus design matrix). We can move on to calculating  $F$  and the p-value.

$$y = \text{mean}_{\text{control}} + \text{mean}_{\text{mutant}}$$



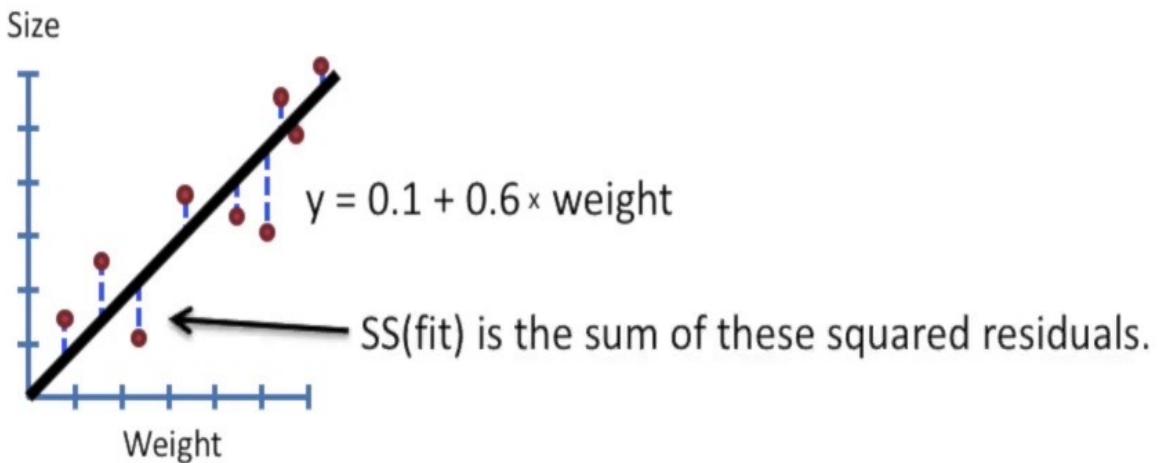
Step 4: Calculate  $SS(\text{fit})$ , the sum of squares of the residuals around the fitted line(s)

## Linear Regression



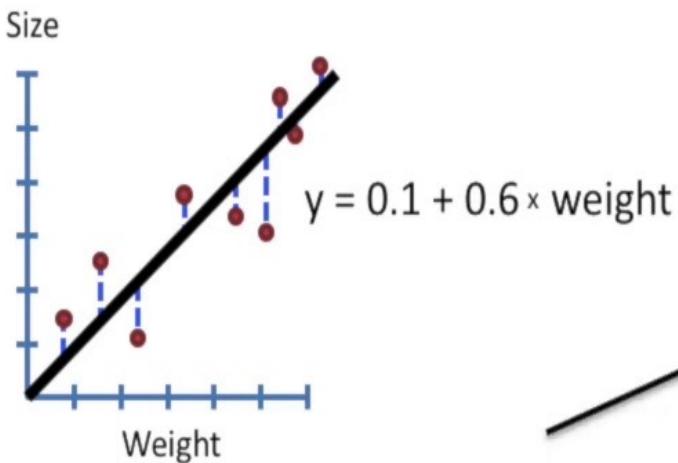
Step 4: Calculate  $SS(\text{fit})$ , the sum of squares of the residuals around the fitted line(s)

## Linear Regression



Step 4: Calculate  $SS(\text{fit})$ , the sum of squares of the residuals around the fitted line(s)

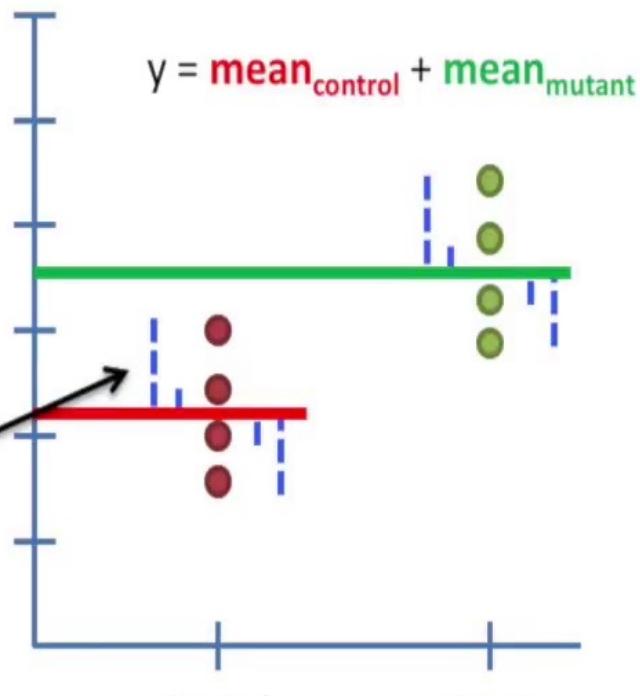
### Linear Regression



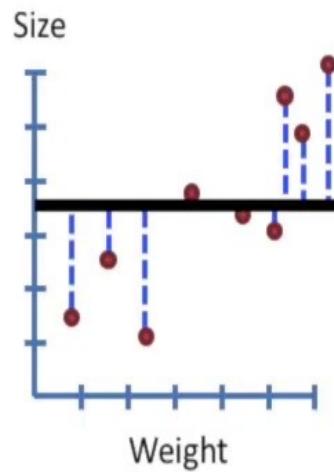
$SS(\text{fit})$  for the t-test is the sum of these squared residuals.

### t-test

#### Gene Expression

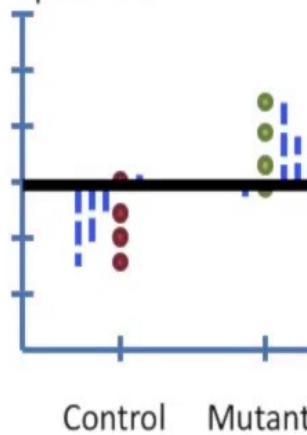


## Linear Regression



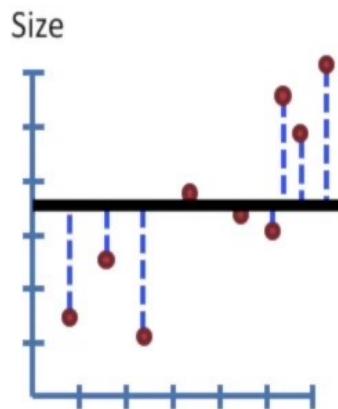
## t-test

### Gene Expression

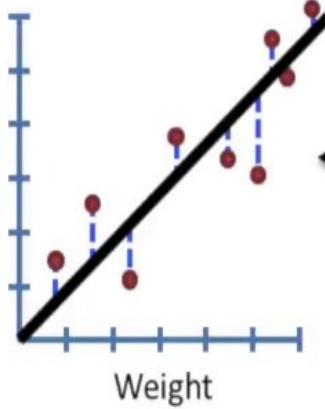


$SS(\text{mean})$

## Linear Regression



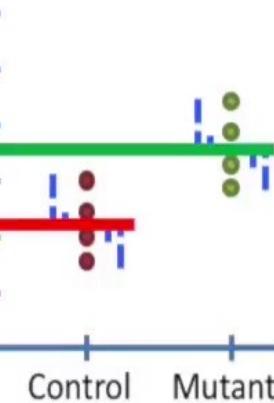
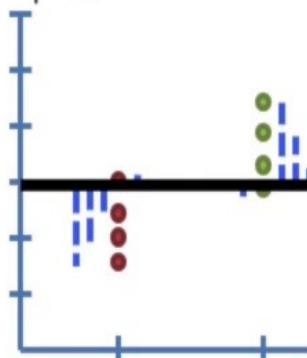
$SS(\text{mean})$



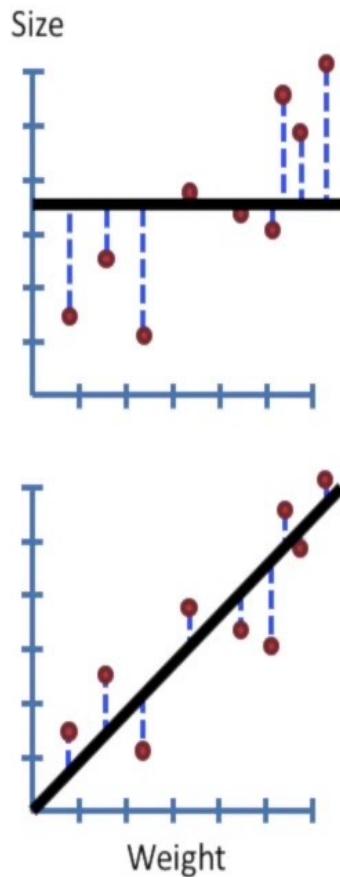
$SS(\text{fit})$

## t-test

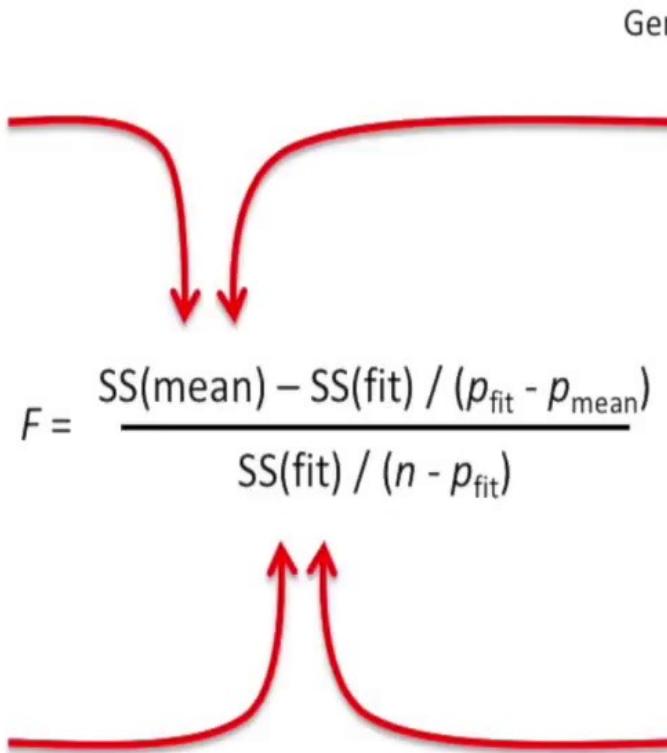
### Gene Expression



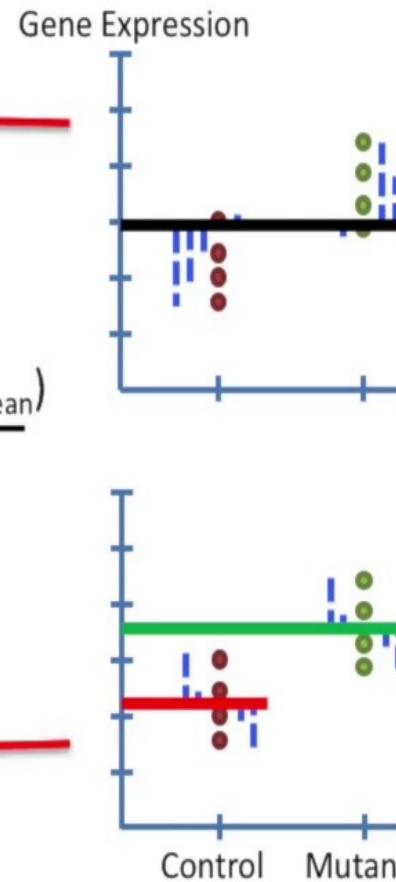
## Linear Regression



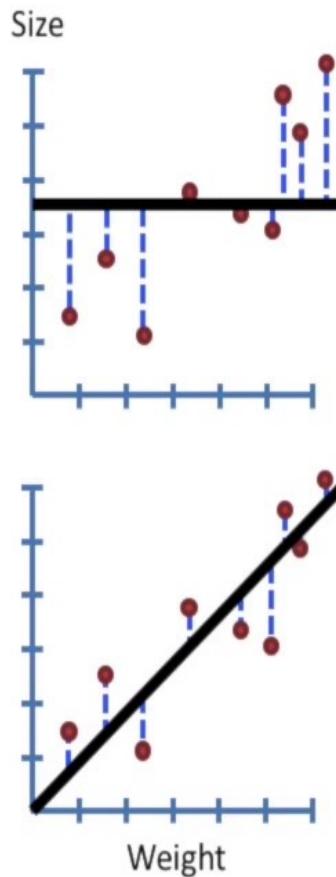
$$F = \frac{SS(\text{mean}) - SS(\text{fit}) / (p_{\text{fit}} - p_{\text{mean}})}{SS(\text{fit}) / (n - p_{\text{fit}})}$$



## t-test



## Linear Regression



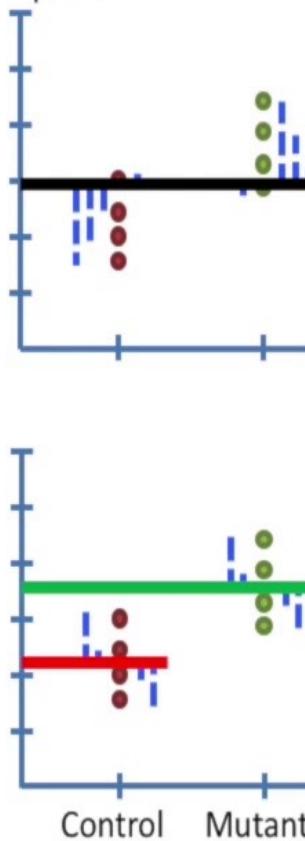
$y = \text{mean mouse size}$

$p_{\text{mean}} = 1$

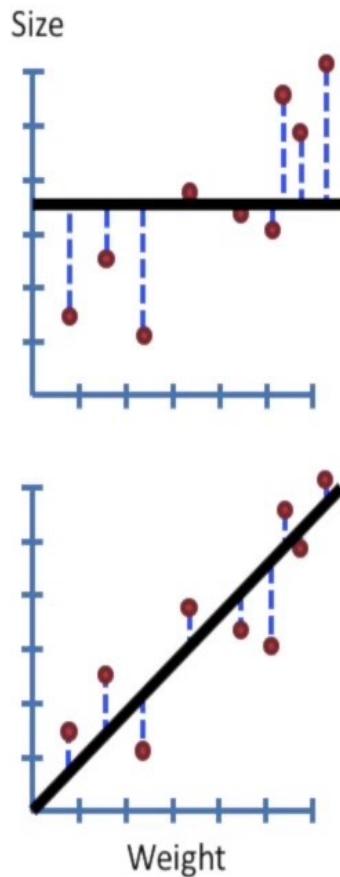
$$F = \frac{\text{SS}(\text{mean}) - \text{SS}(\text{fit}) / (p_{\text{fit}} - p_{\text{mean}})}{\text{SS}(\text{fit}) / (n - p_{\text{fit}})}$$

## t-test

### Gene Expression



## Linear Regression



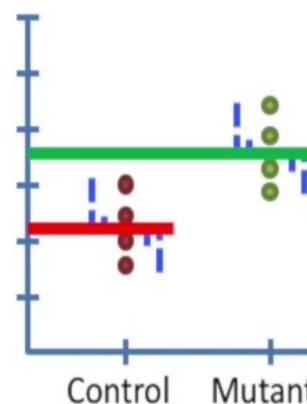
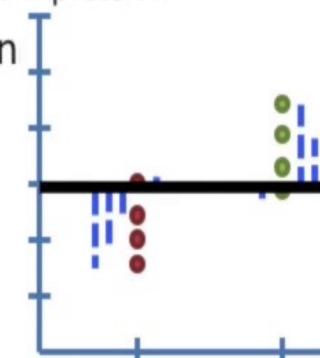
$$F = \frac{SS(\text{mean}) - SS(\text{fit}) / (p_{\text{fit}} - p_{\text{mean}})}{SS(\text{fit}) / (n - p_{\text{fit}})}$$

Gene Expression

$y = \text{mean gene expression}$

$p_{\text{mean}} = 1$

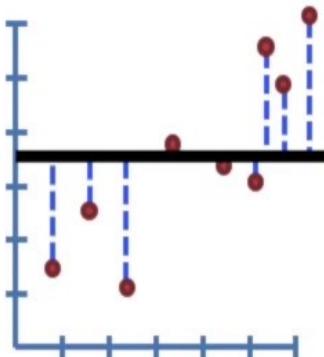
A red arrow starts at the mean value on the y-axis and points towards the fitted regression line, indicating the difference between the mean and the fit.



## t-test

## Linear Regression

Size



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

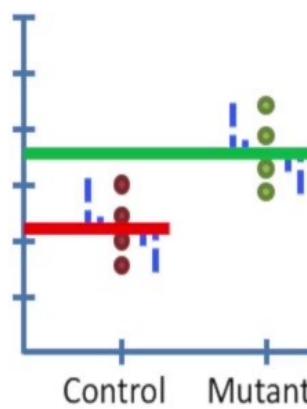
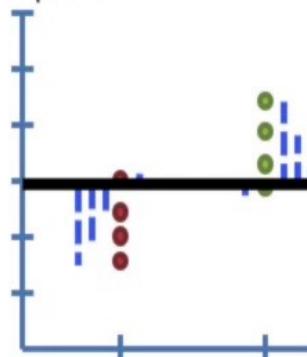
$$y = \text{intercept} + \text{slope}$$

$$p_{\text{fit}} = 2$$

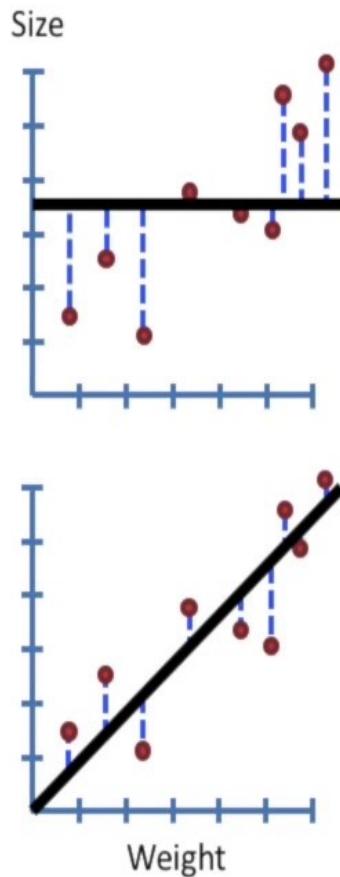
Weight

## t-test

Gene Expression



## Linear Regression



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

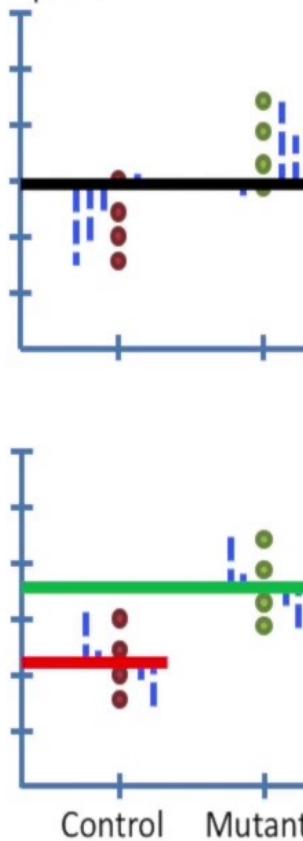
$y = \text{mean}_{\text{control}} + \text{mean}_{\text{mutant}}$

A red wavy line with arrows at both ends, representing a fluctuating signal or noise.

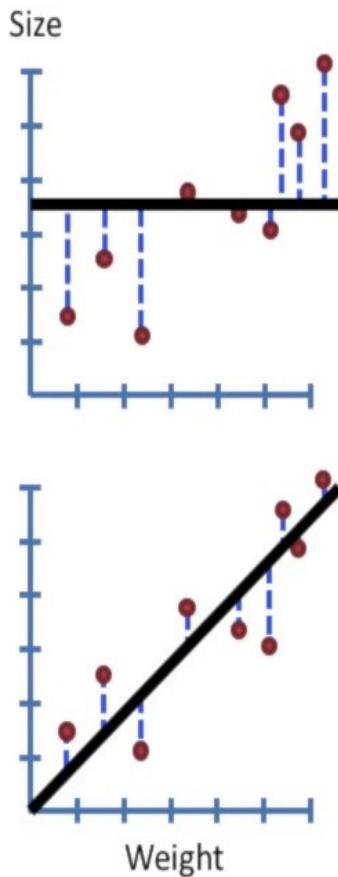
$$p_{\text{fit}} = 2$$

## t-test

### Gene Expression



## Linear Regression

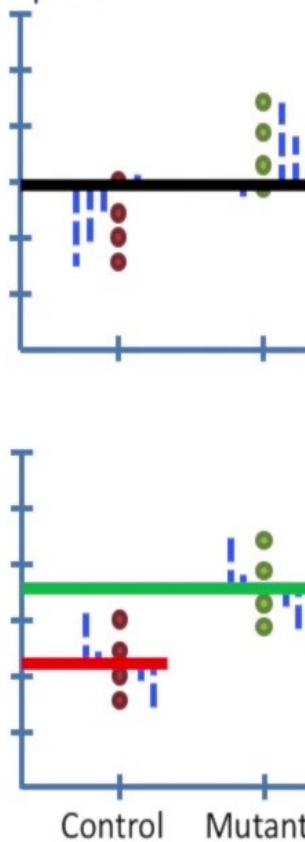


Now we can calculate a p-value  
for the t-test!

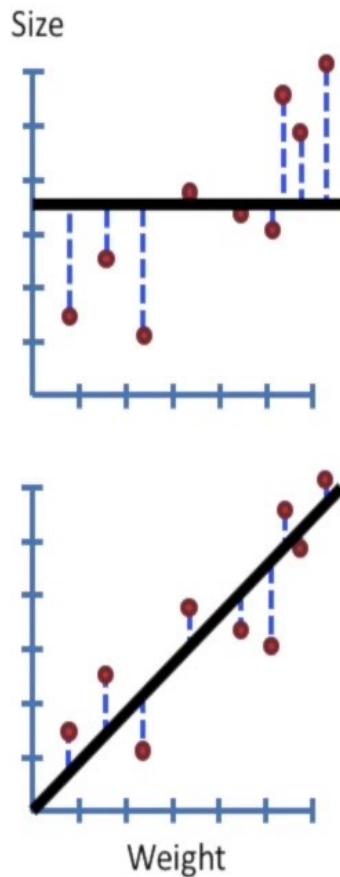
$$F = \frac{SS(\text{mean}) - SS(\text{fit}) / (p_{\text{fit}} - p_{\text{mean}})}{SS(\text{fit}) / (n - p_{\text{fit}})}$$

## t-test

### Gene Expression



## Linear Regression

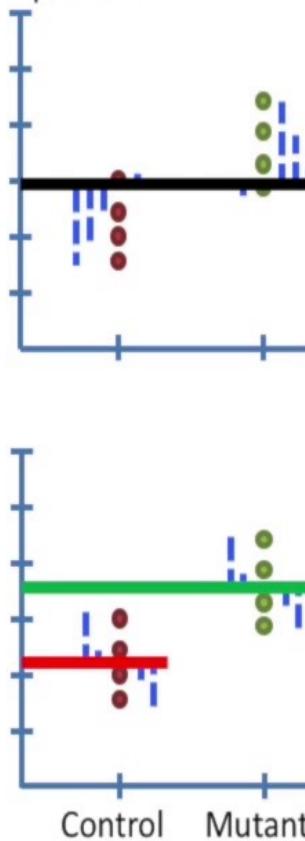


BAM!!!

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

## t-test

### Gene Expression



Let's review what we've done so far..

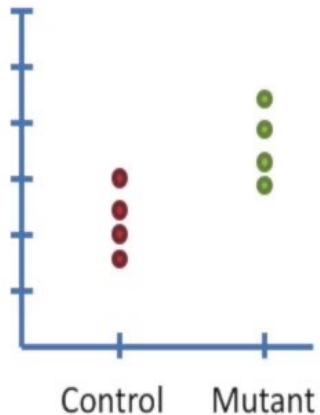
The original data.

Gene expression

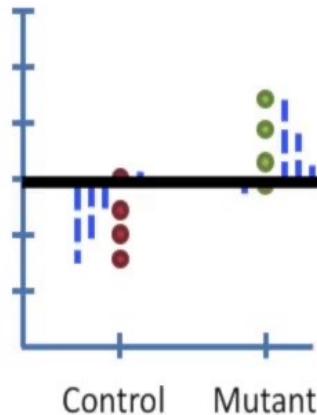


The original data.

Gene expression



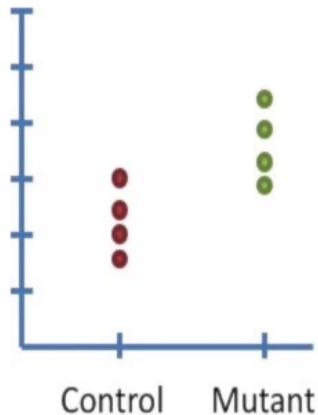
Calculate SS(mean)



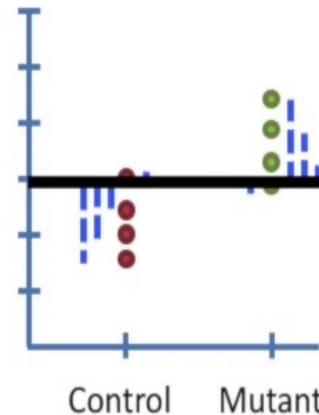
$y$  = overall mean

The original data.

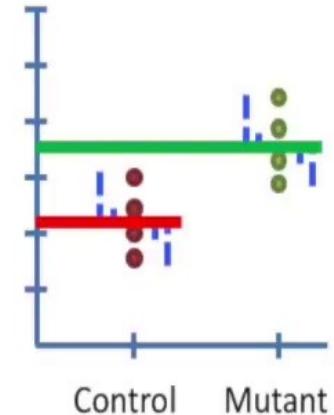
Gene expression



Calculate SS(mean)



Calculate SS(fit)

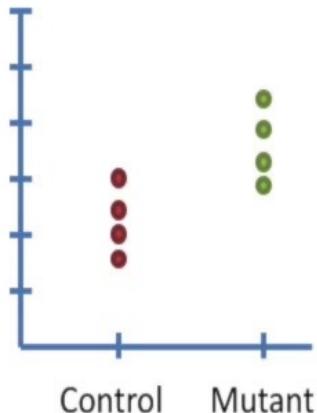


$$y = \text{overall mean}$$

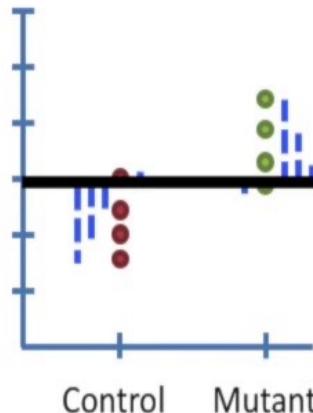
$$y = \text{mean}_{\text{control}} + \text{mean}_{\text{mutant}}$$

The original data.

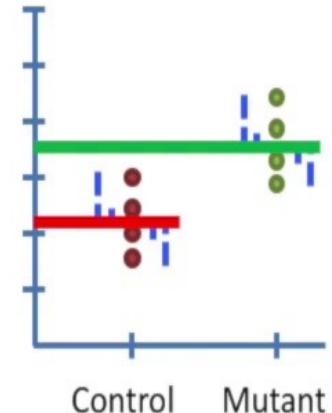
Gene expression



Calculate SS(mean)



Calculate SS(fit)



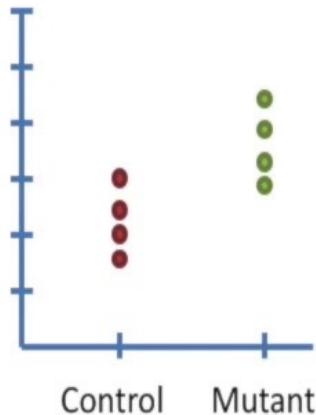
$$y = \text{overall mean}$$

$$y = \text{mean}_{\text{control}} + \text{mean}_{\text{mutant}}$$

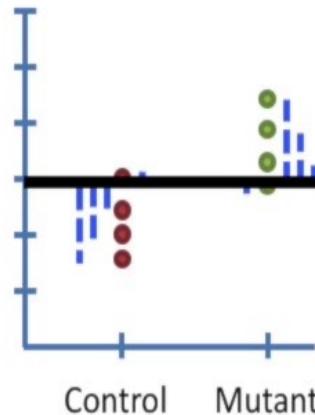
$$F = \frac{\text{SS}(\text{mean}) - \text{SS}(\text{fit}) / (p_{\text{fit}} - p_{\text{mean}})}{\text{SS}(\text{fit}) / (n - p_{\text{fit}})}$$

The original data.

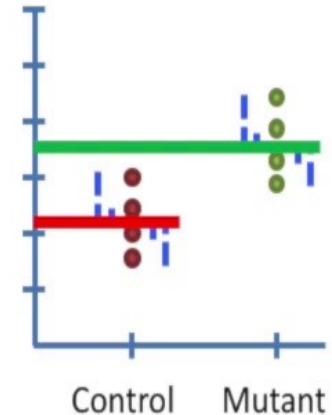
Gene expression



Calculate SS(mean)



Calculate SS(fit)



$y$  = overall mean

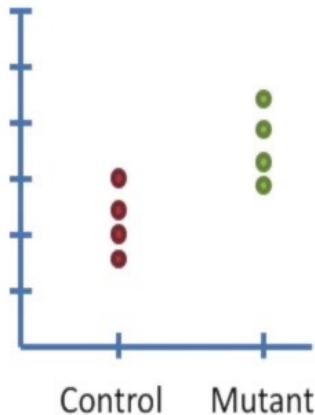
$y$  =  $\text{mean}_{\text{control}}$  +  $\text{mean}_{\text{mutant}}$

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

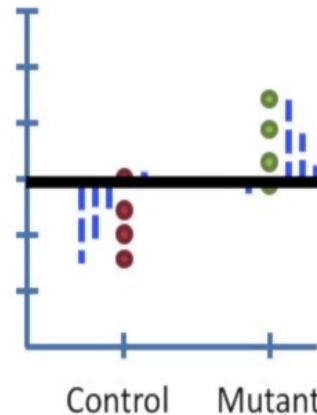
$p_{\text{mean}} = 1$

The original data.

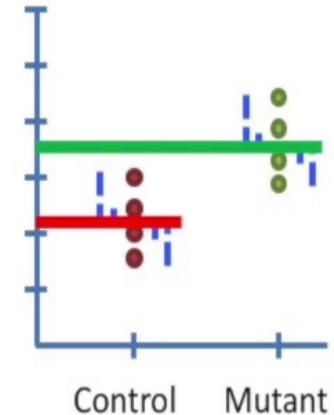
Gene expression



Calculate SS(mean)



Calculate SS(fit)



$$y = \text{overall mean}$$

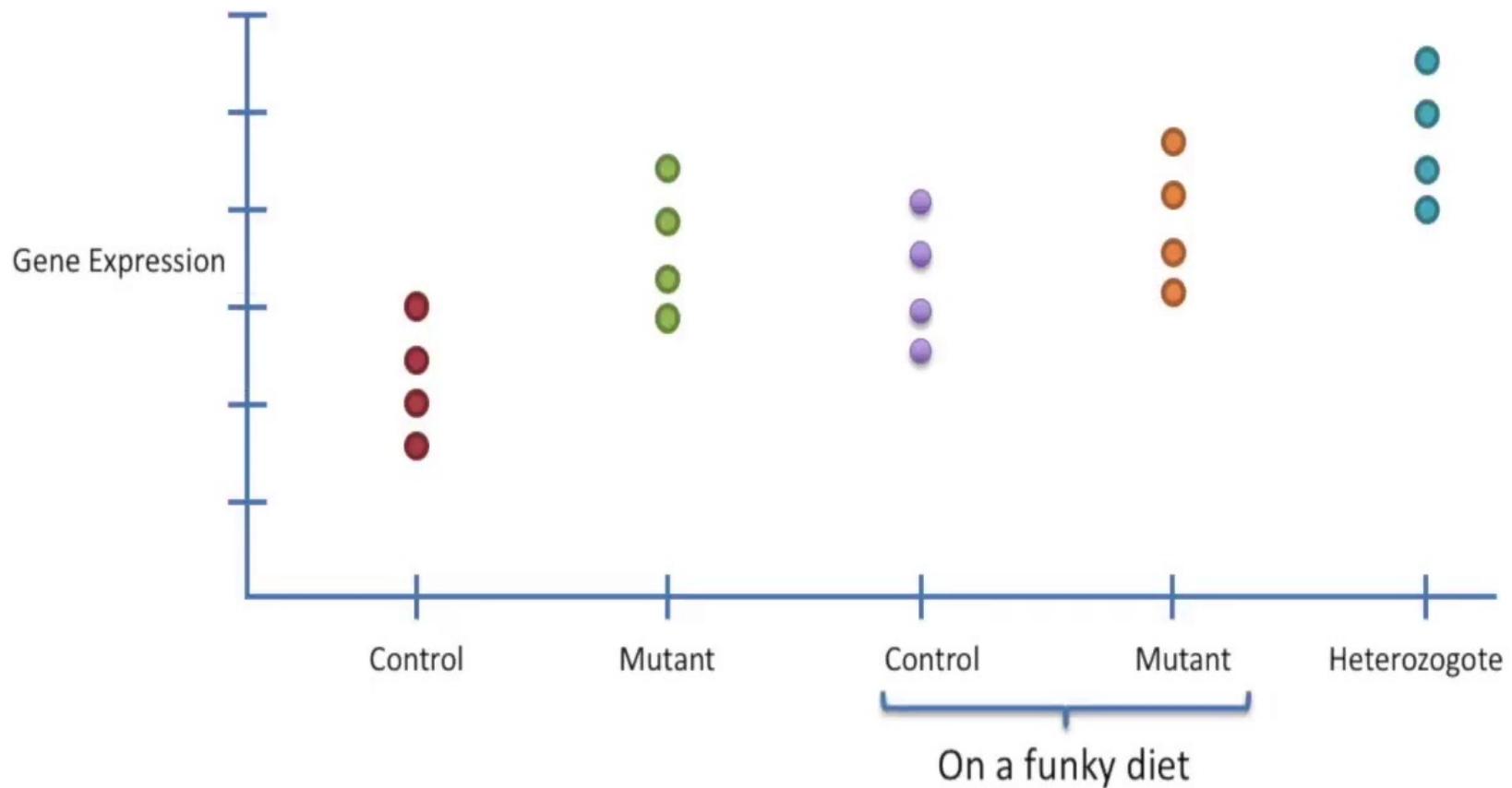
$$F = \frac{\text{SS}(\text{mean}) - \text{SS}(\text{fit}) / (p_{\text{fit}} - p_{\text{mean}})}{\text{SS}(\text{fit}) / (n - p_{\text{fit}})}$$

$$p_{\text{mean}} = 1$$

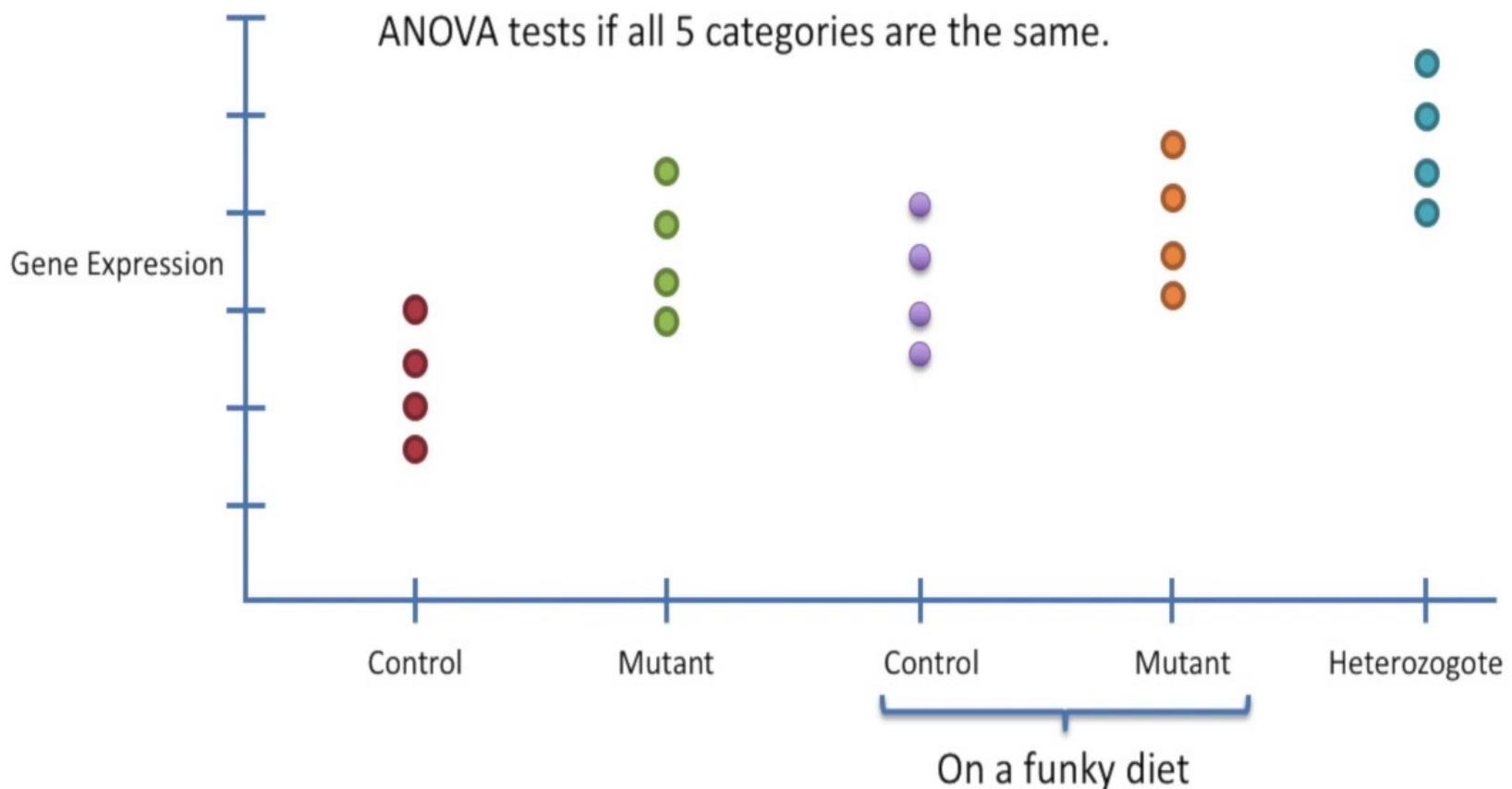
$$y = \text{mean}_{\text{control}} + \text{mean}_{\text{mutant}}$$

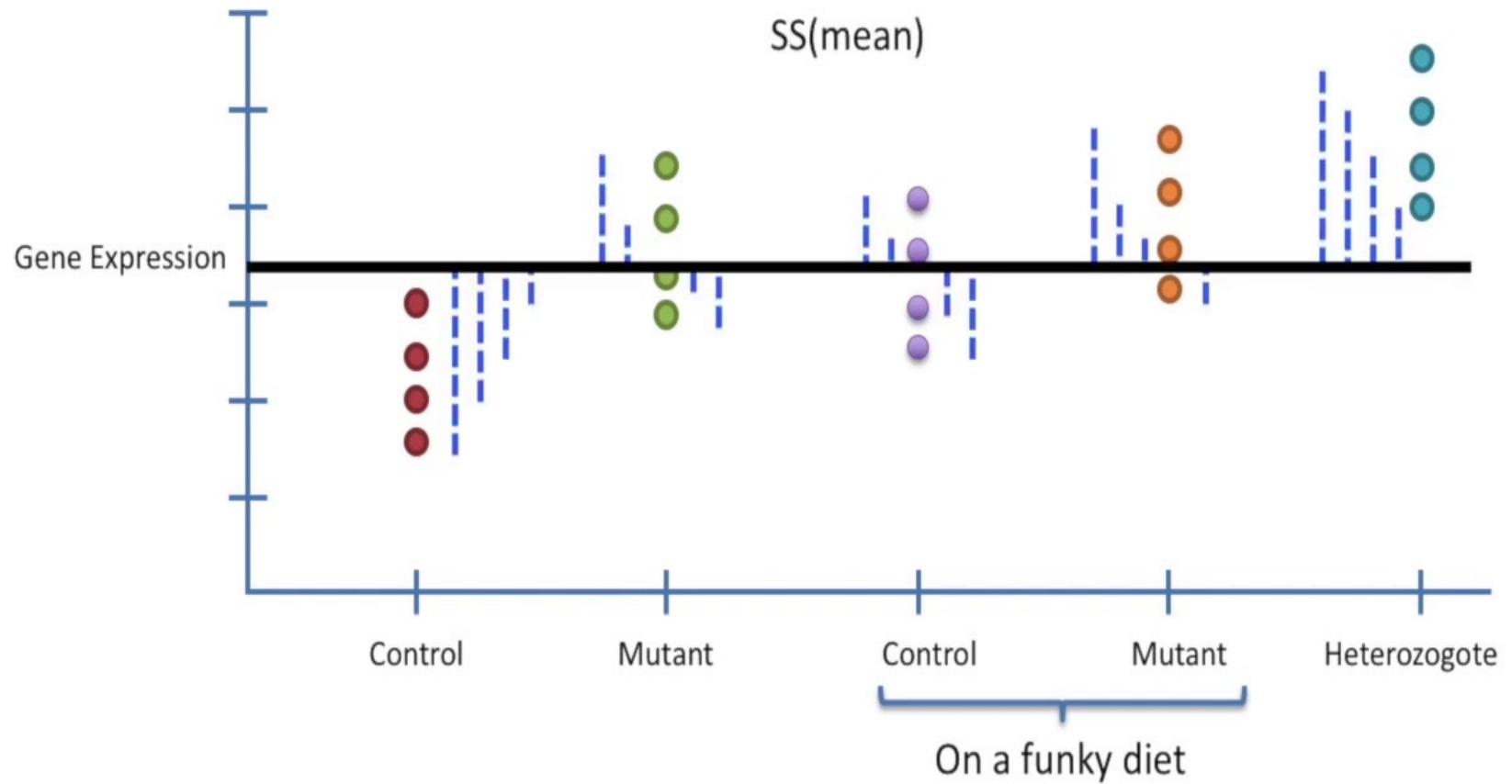
$$p_{\text{fit}} = 2$$

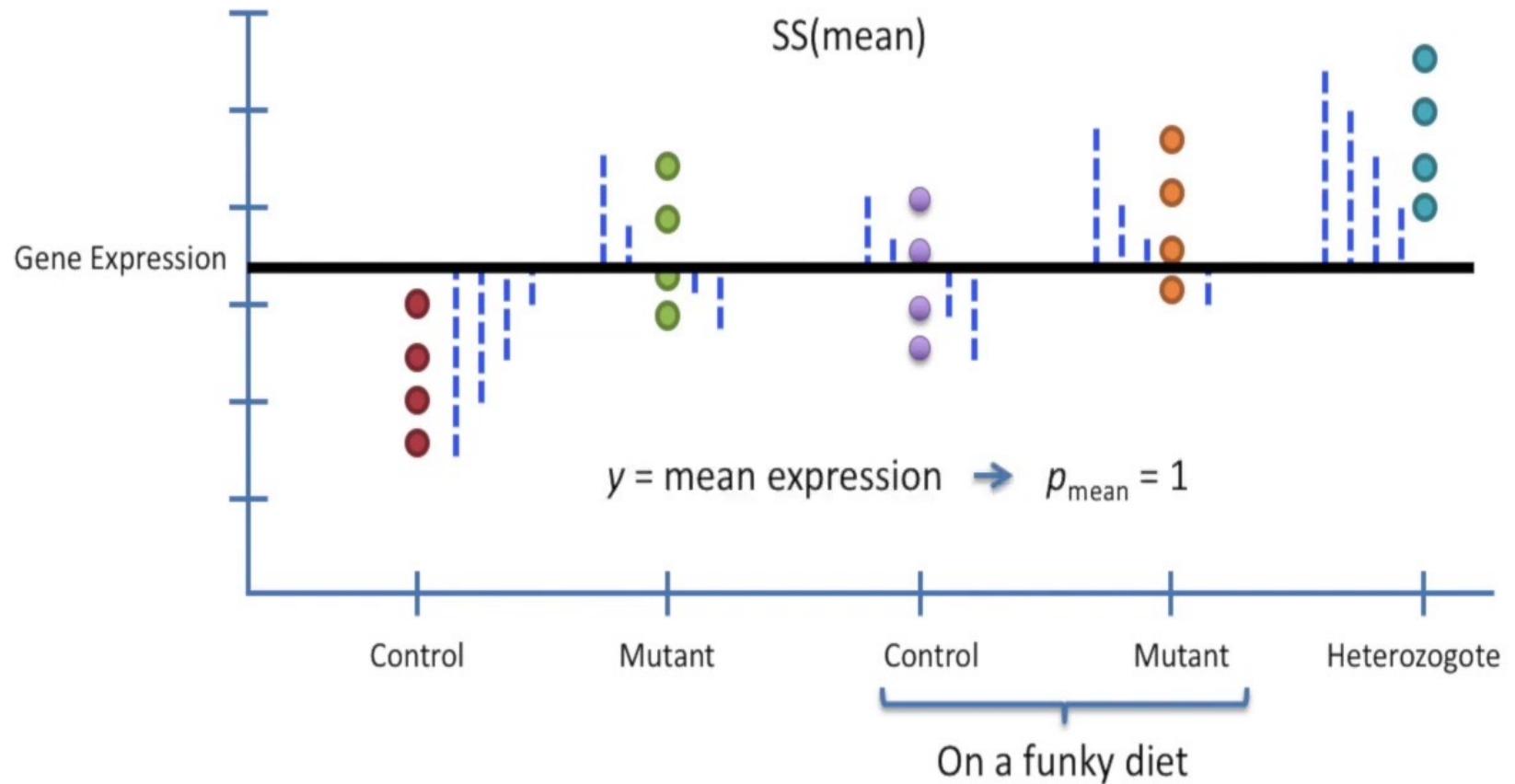
Now let's do an ANOVA!!

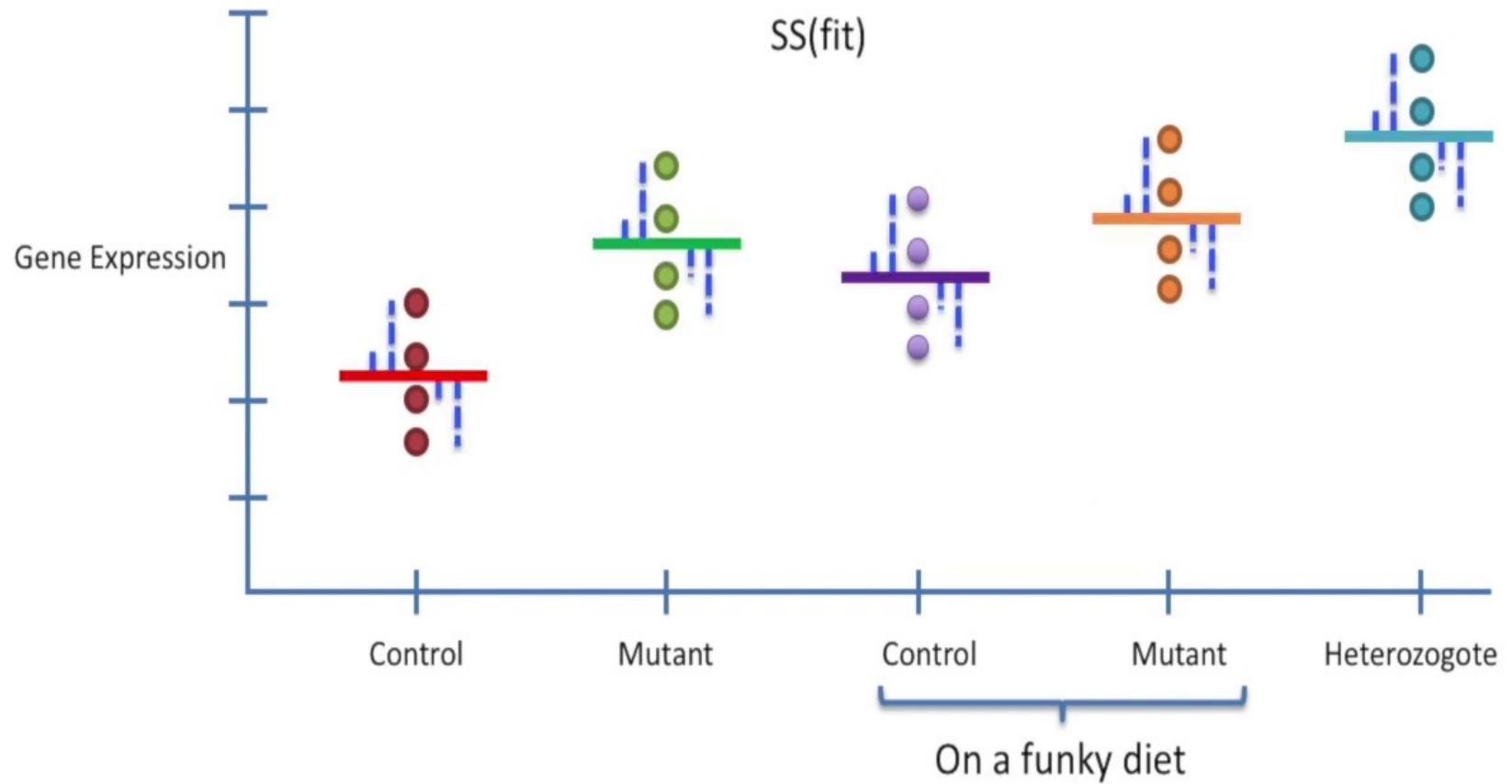


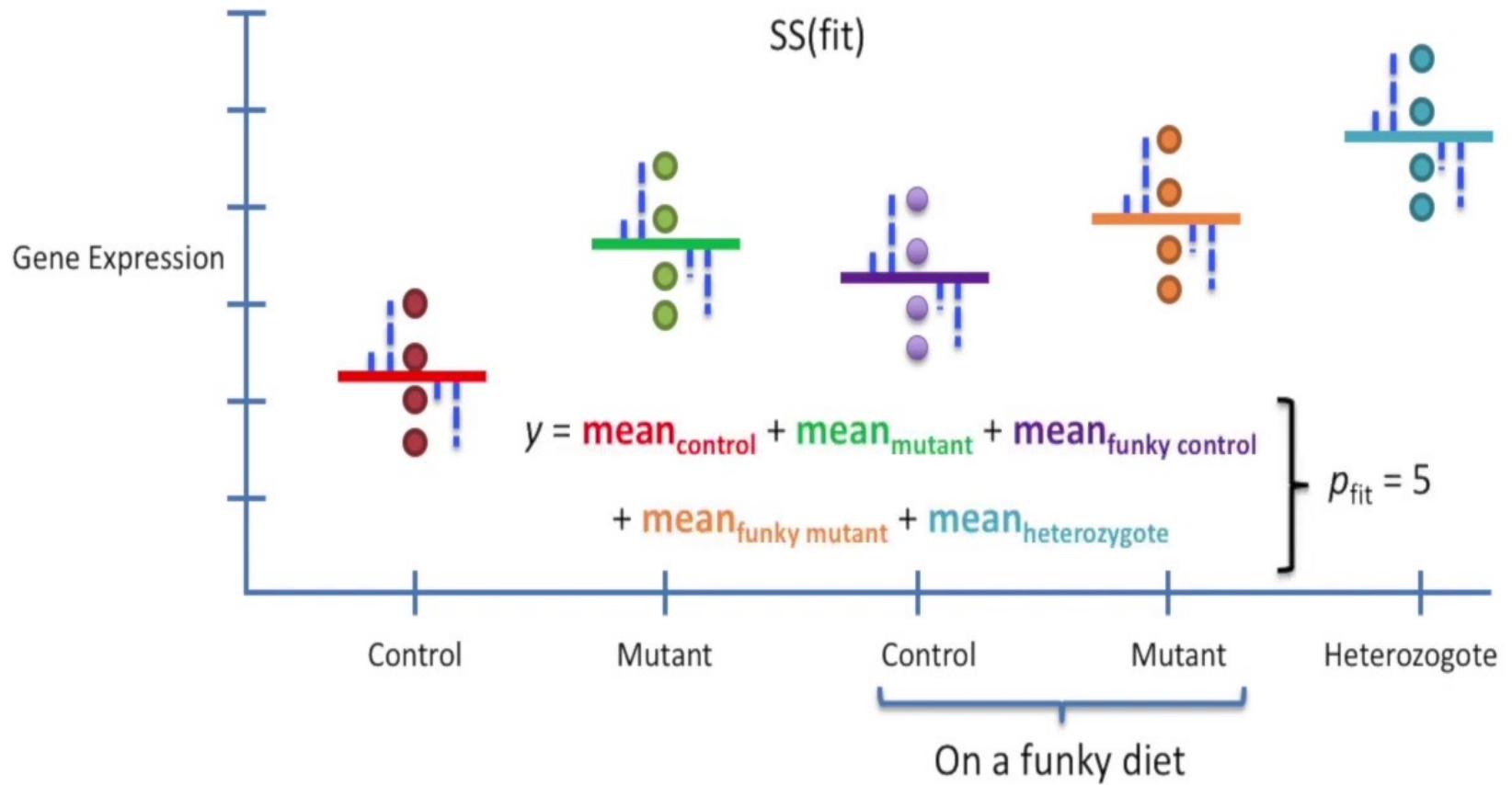
Now let's do an ANOVA!!





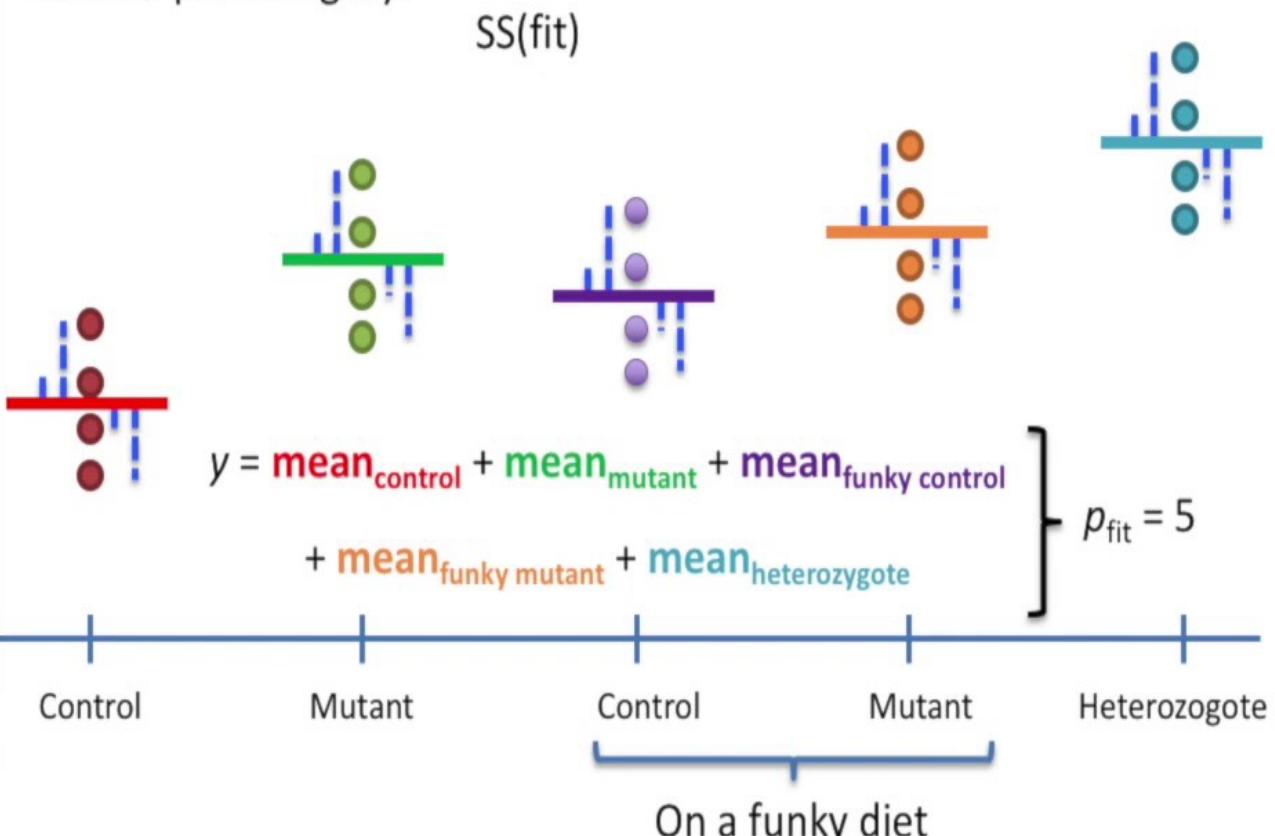


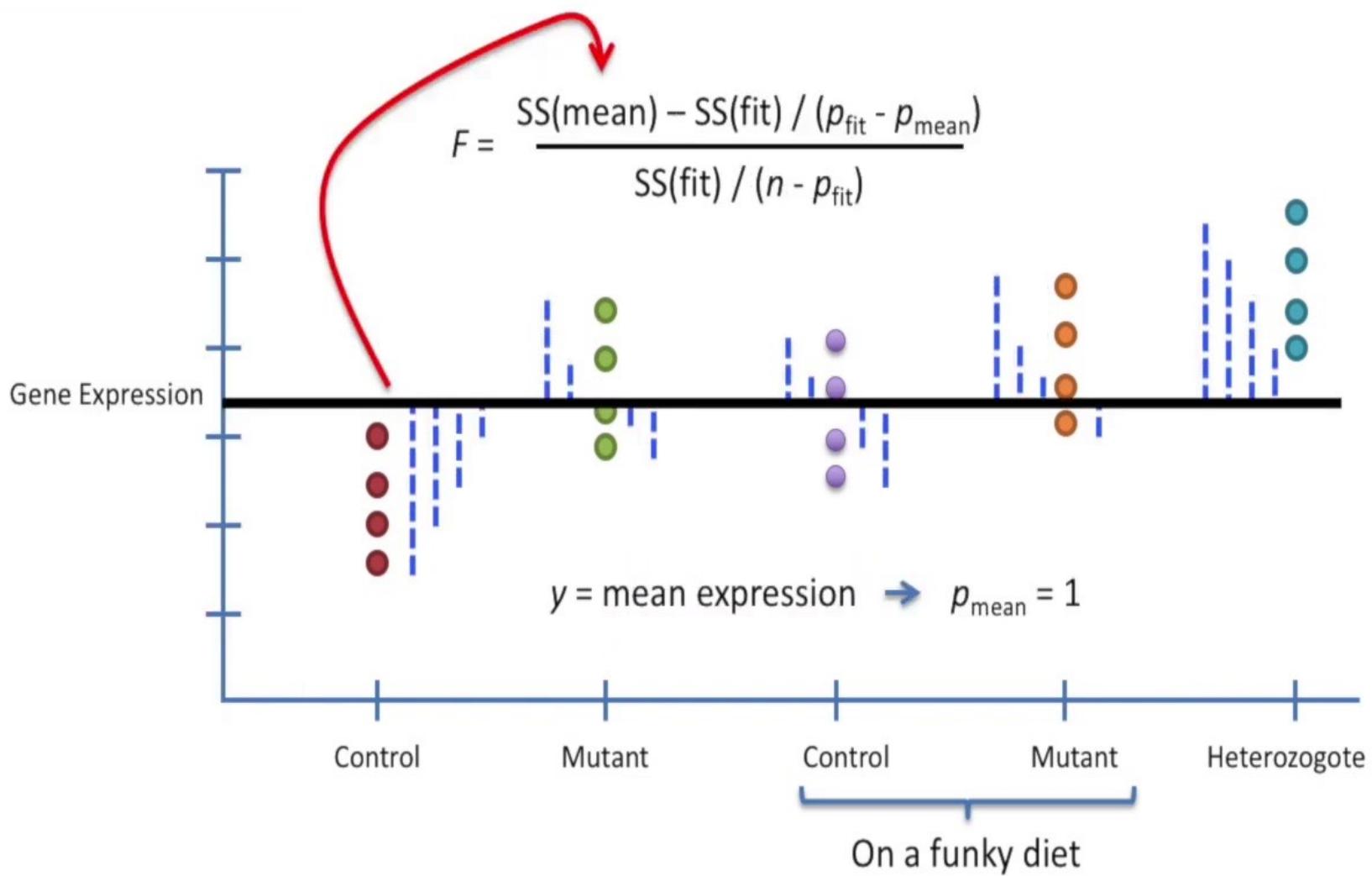


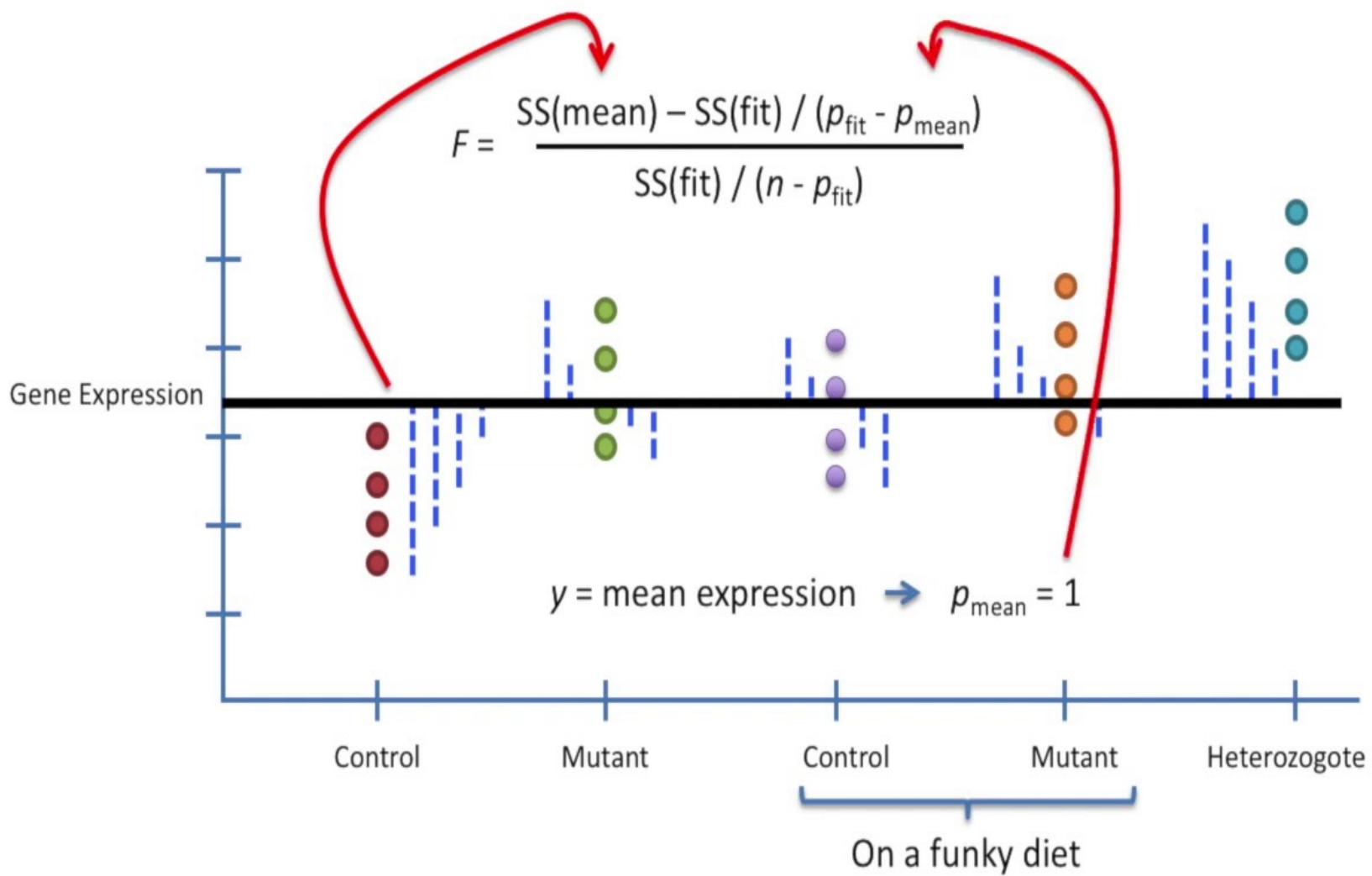


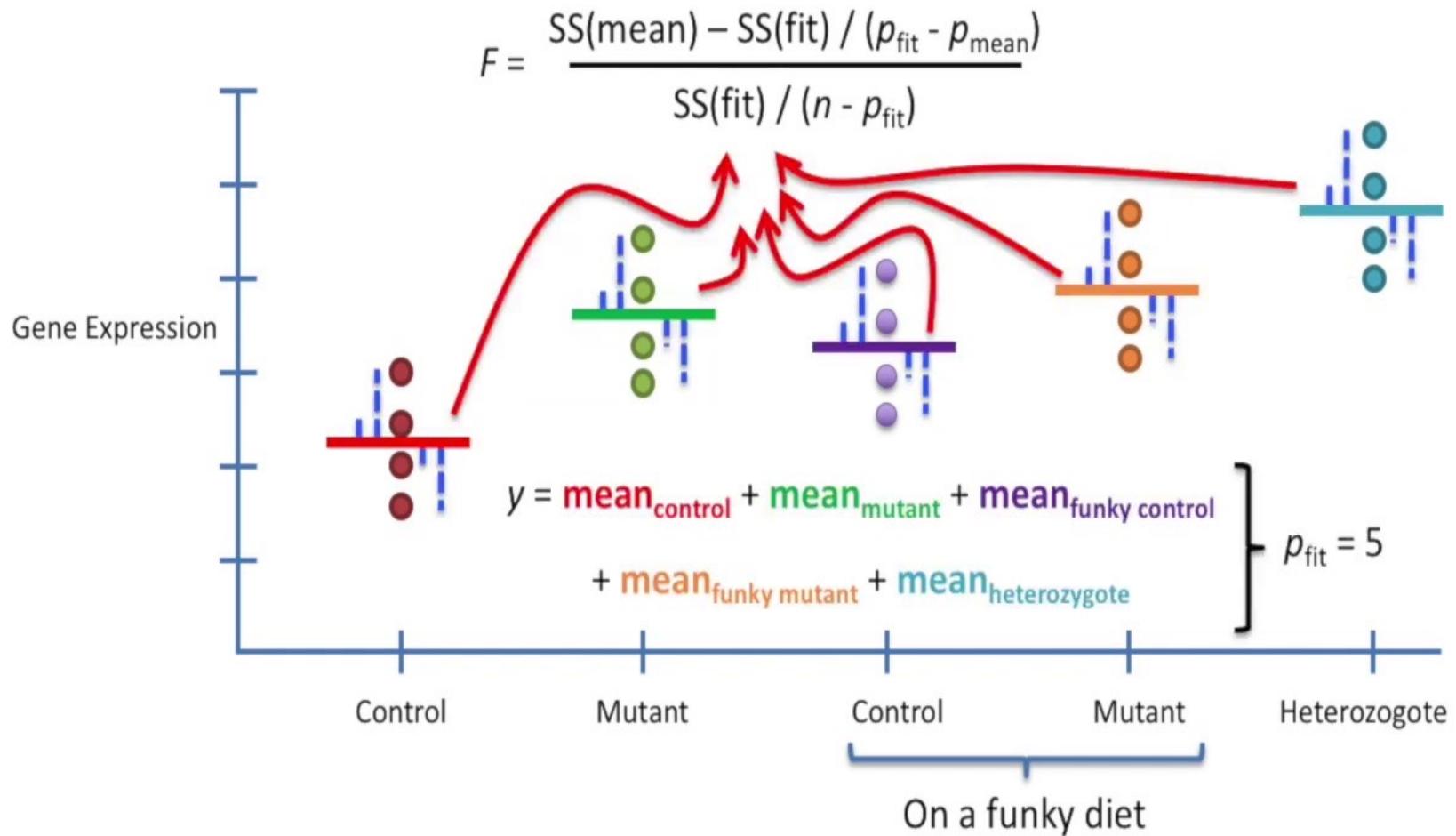
1	0	0	0	0
1	0	0	0	0
1	0	0	0	0
1	0	0	0	0
0	1	0	0	0
0	1	0	0	0
0	1	0	0	0
0	1	0	0	0
0	0	1	0	0
0	0	1	0	0
0	0	1	0	0
0	0	1	0	0
0	0	0	1	0
0	0	0	1	0
0	0	0	1	0
0	0	0	0	1
0	0	0	0	1
0	0	0	0	1

Here's what the design matrix looks like - one column per category.

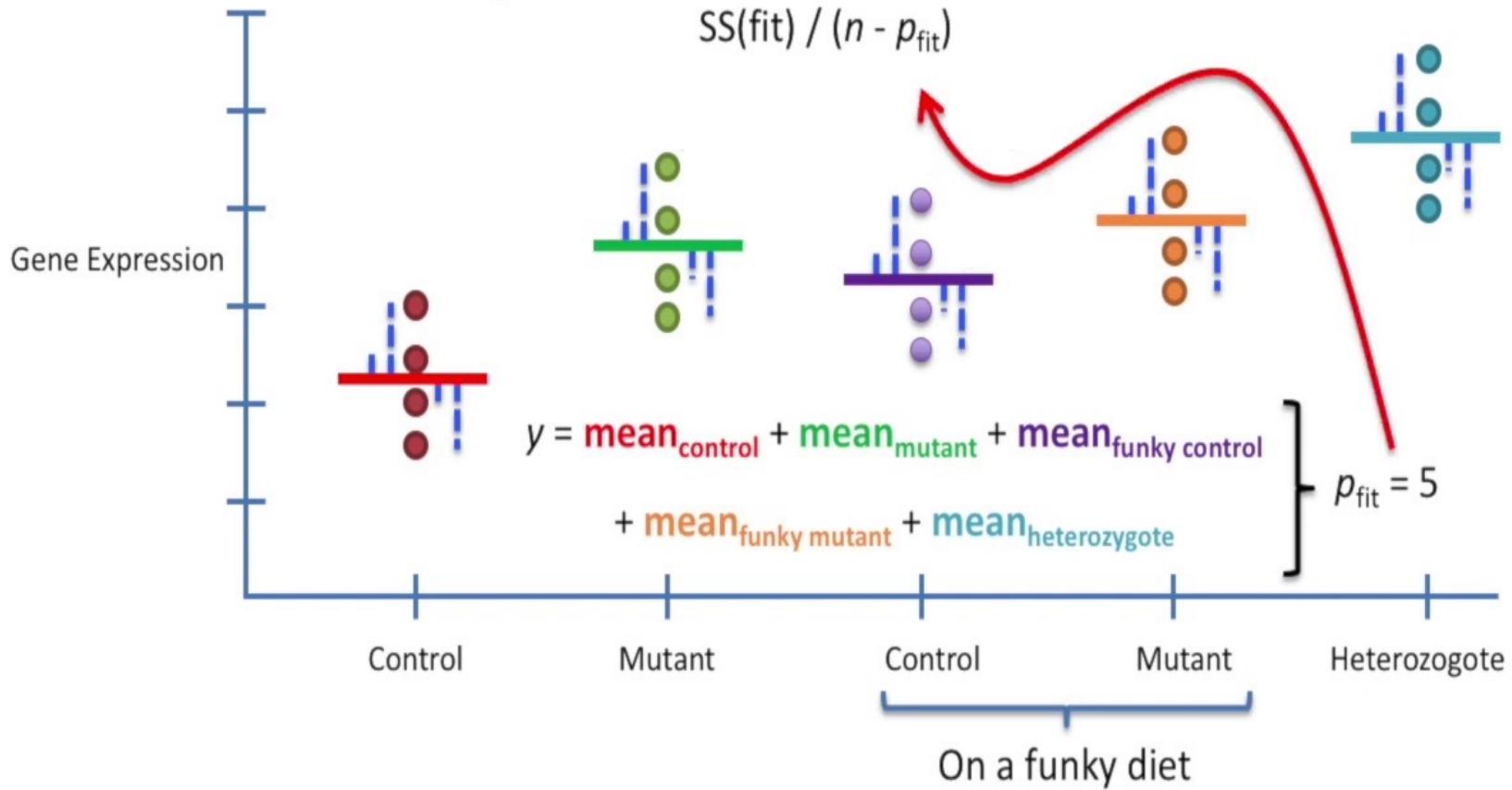


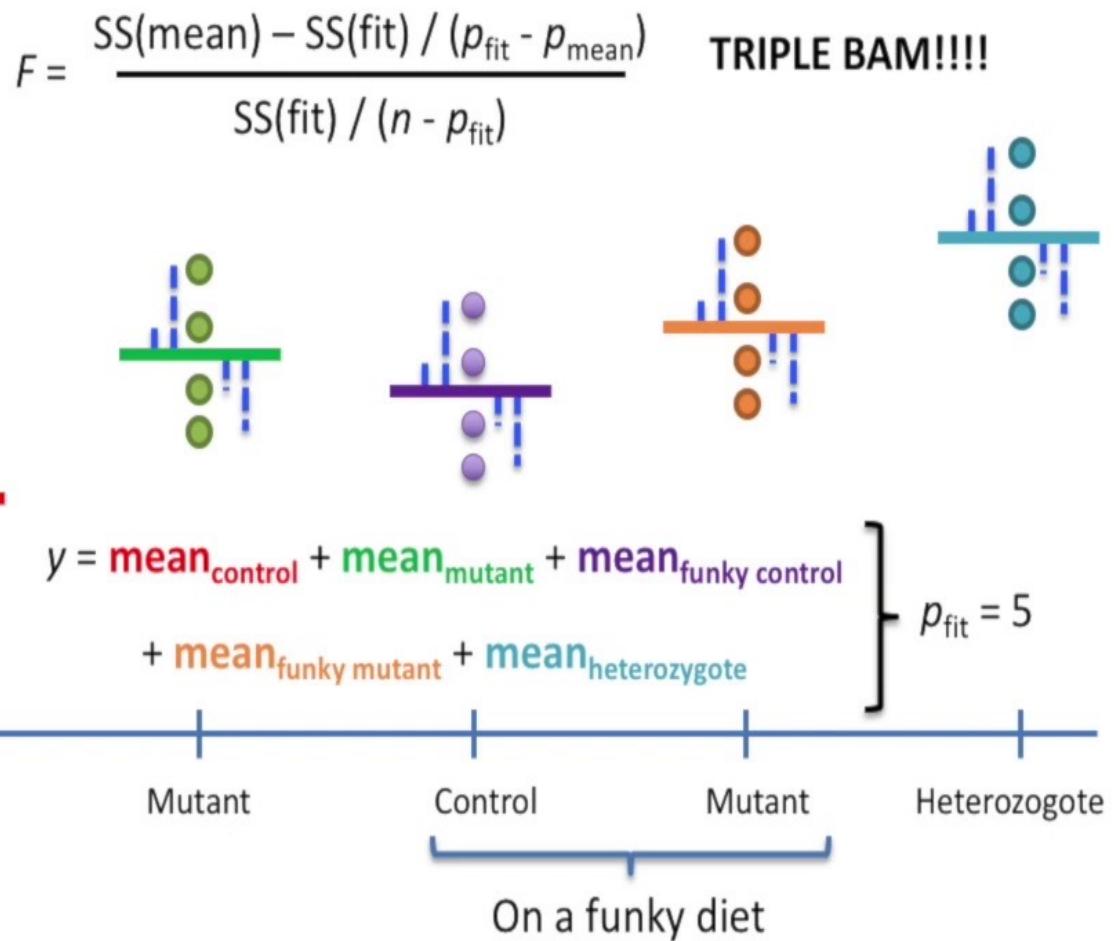






$$F = \frac{\frac{SS(\text{mean}) - SS(\text{fit})}{(p_{\text{fit}} - p_{\text{mean}})}}{\frac{SS(\text{fit})}{(n - p_{\text{fit}})}}$$





## One last important detail before we're done...

The design matrices that I've shown you are not the standard design matrices used for doing t-tests and ANOVA.

# One last important detail before we're done...

The design matrices that I've shown you are not the standard design matrices used for doing t-tests and ANOVA.

1	0
1	0
1	0
1	0
0	1
0	1
0	1
0	1

This is what we used for the t-test in this StatQuest...

$$y = \text{mean}_{\text{control}} + \text{mean}_{\text{mutant}}$$

# One last important detail before we're done...

1	0
1	0
1	0
1	0
0	1
0	1
0	1
0	1

The design matrices that I've shown you are not the standard design matrices used for doing t-tests and ANOVA.

This is what we used for the t-test in this StatQuest...

$$y = \text{mean}_{\text{control}} + \text{mean}_{\text{mutant}}$$

...but this is a more common design matrix. We'll talk about this one and other, more elaborate designs in the next StatQuest.

1	0
1	0
1	0
1	0
1	1
1	1
1	1
1	1

**The End!!!!**