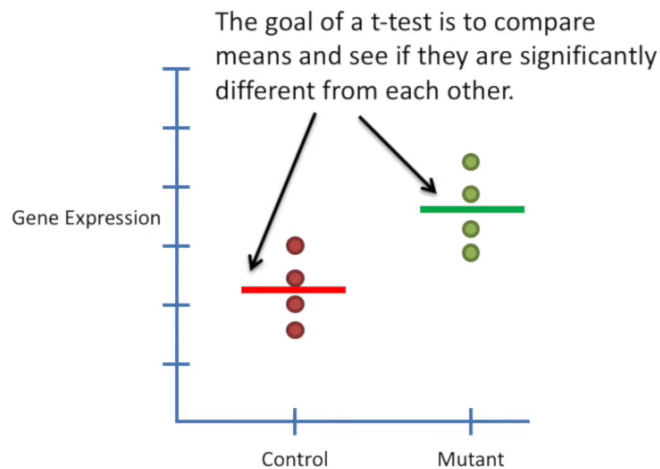
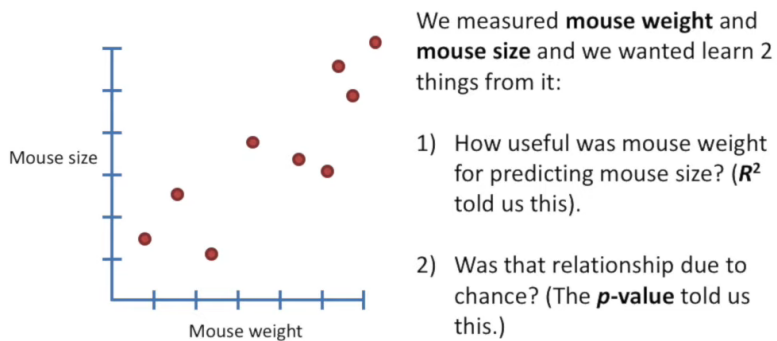


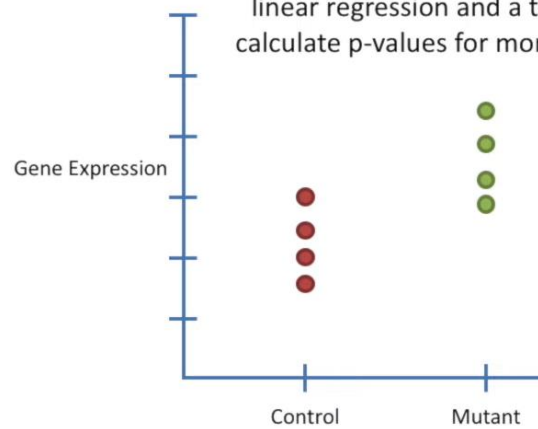
# StatQuest: General Linear Models

## Part 2: t-tests and ANOVA

### Implement the two-sample t-test and ANOVA with linear model method

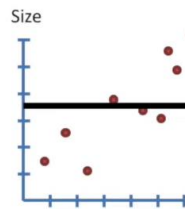


If the same method can calculate  $p$ -values for a linear regression and a  $t$ -test, then we can easily calculate  $p$ -values for more complicated situations.



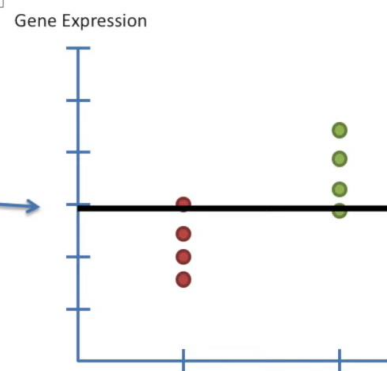
Step 1: Ignore the x-axis and find the overall mean.

Linear Regression



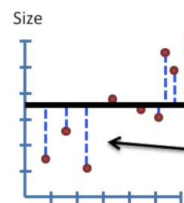
Overall means

t-test



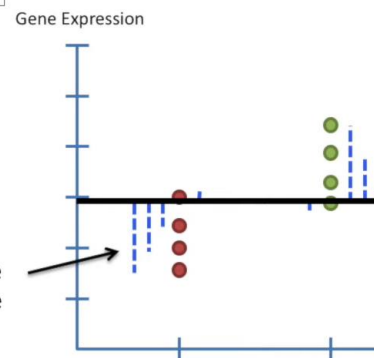
Step 2: Calculate  $SS(\text{mean})$ , the sum of squared residuals around the mean.

Linear Regression



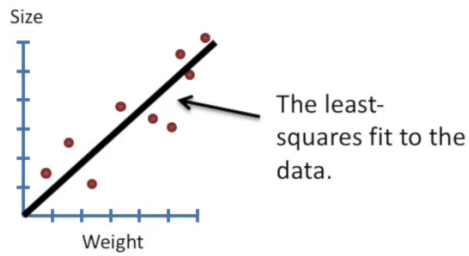
These are the residuals, the distance from the data points to the lines.

t-test

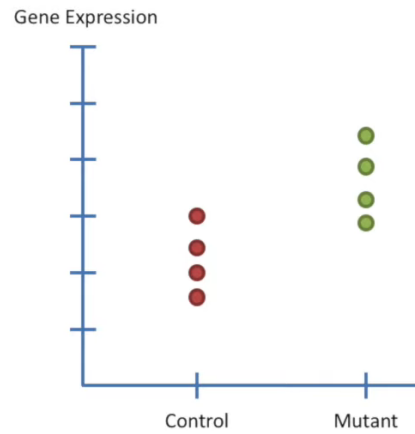


Step 3: Fit a line to the data.

### Linear Regression

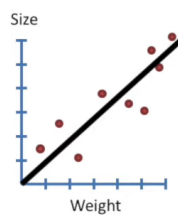


### t-test

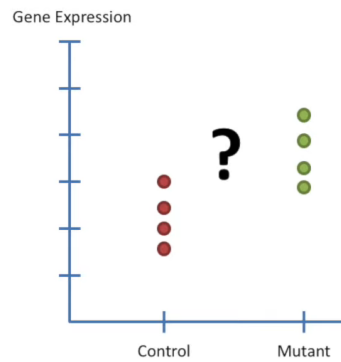


Step 3: Fit a line to the data.

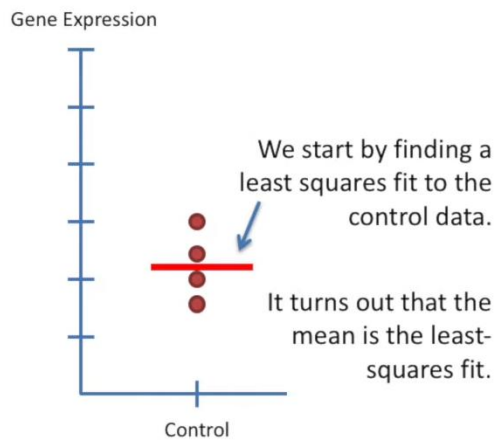
### Linear Regression



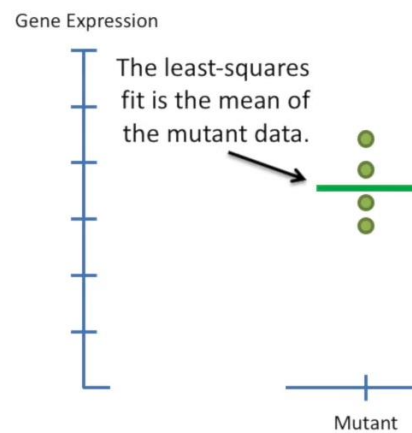
### t-test



### t-test

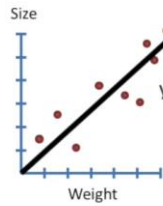


### t-test



Step 3: Fit a line to the data.

Linear Regression



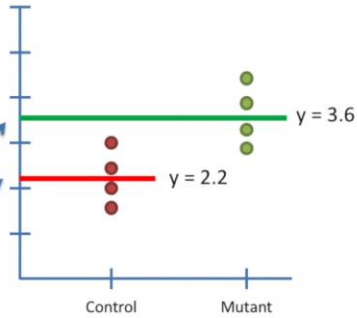
Originally we fit a single line to the data.

$$y = 0.1 + 0.6 \times \text{weight}$$

We have fit two lines to the data.

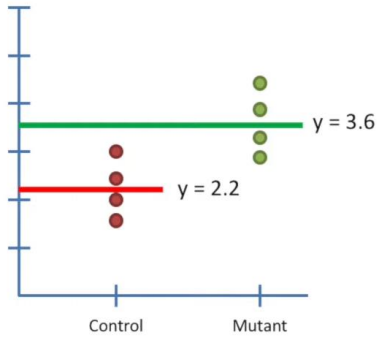
t-test

Gene Expression



t-test

Gene Expression



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

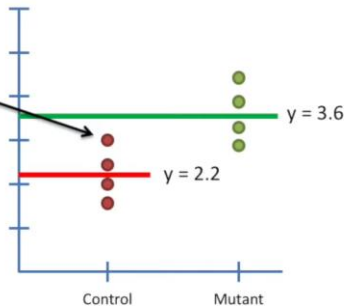
This will make the steps for computing " $F$ " the exact same for the regression and the t-test, which, in turn, means a computer can do it automatically.

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

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

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

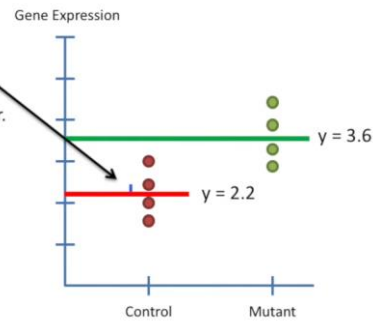
Gene Expression



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

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

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

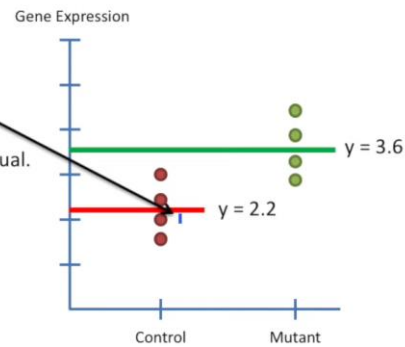


$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

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

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

Again, the only difference is the residual.



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

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

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

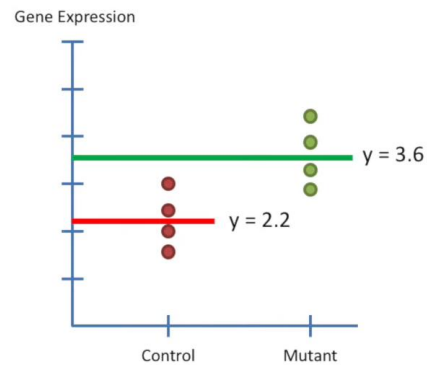
$$y = 1 \times 2.2 + 0 \times 3.6 + \text{the residual}$$

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

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

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

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



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

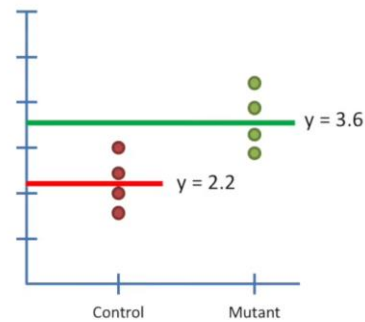
1	0
1	0
1	0
1	0
0	1
0	1
0	1
0	1

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

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.

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



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

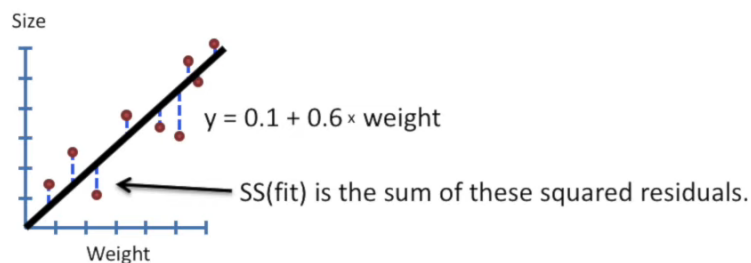
In practice, the role of each column is assumed, and the equation is written out like this:

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

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.

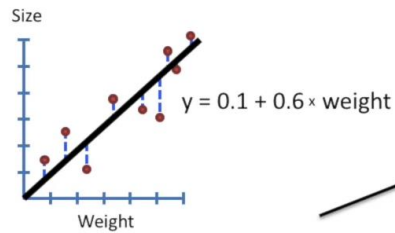
Step 4: Calculate  $SS(\text{fit})$ , the sum of squares of the residuals around the fitted line(s)

## Linear Regression



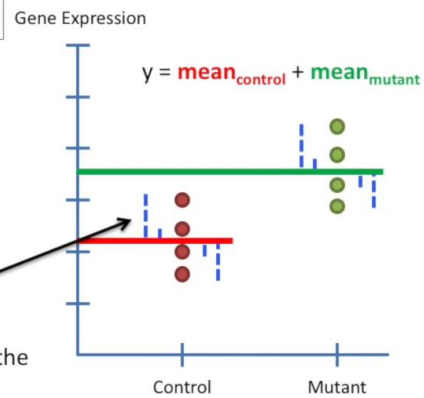
Step 4: Calculate  $SS(\text{fit})$ , the sum of squares of the residuals around the fitted line(s)

Linear Regression

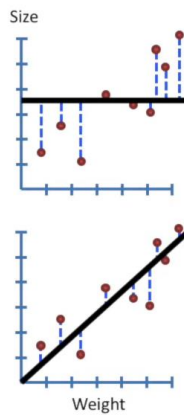


$SS(\text{fit})$  for the t-test is the sum of these squared residuals.

t-test

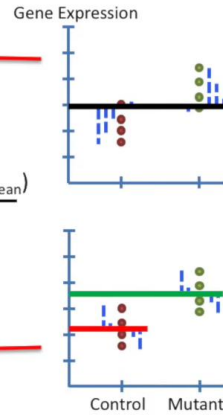


Linear Regression



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

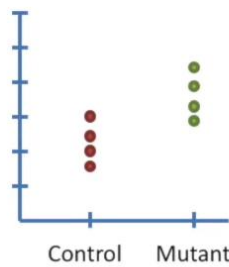
t-test



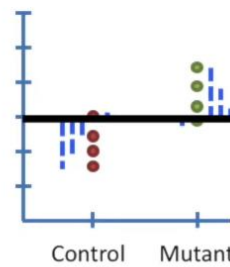
## Summary and ANOVA

The original data.

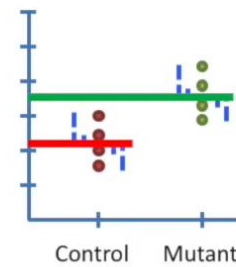
Gene expression



Calculate SS(mean)



Calculate SS(fit)



$y = \text{overall mean}$

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

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

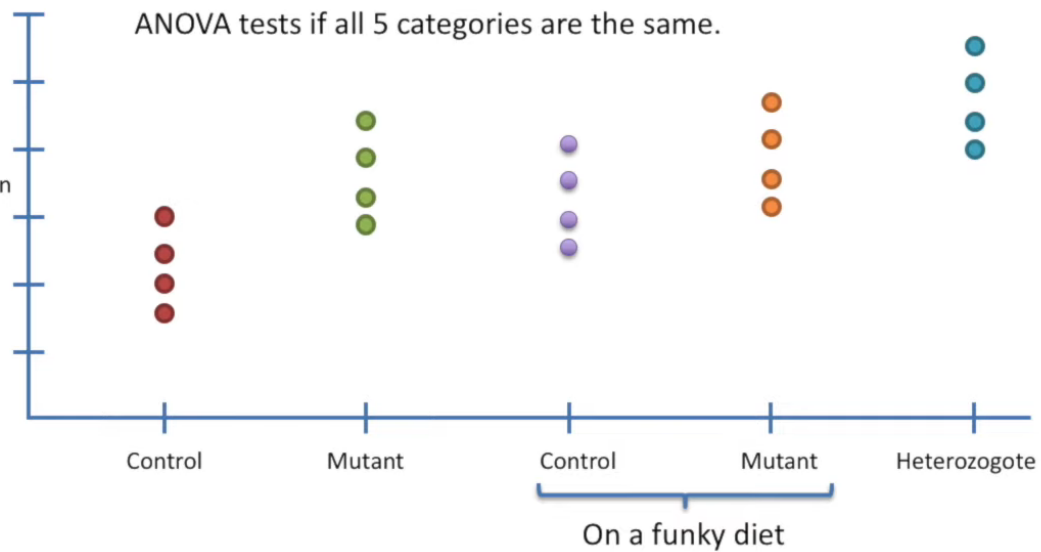
$p_{\text{mean}} = 1$

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

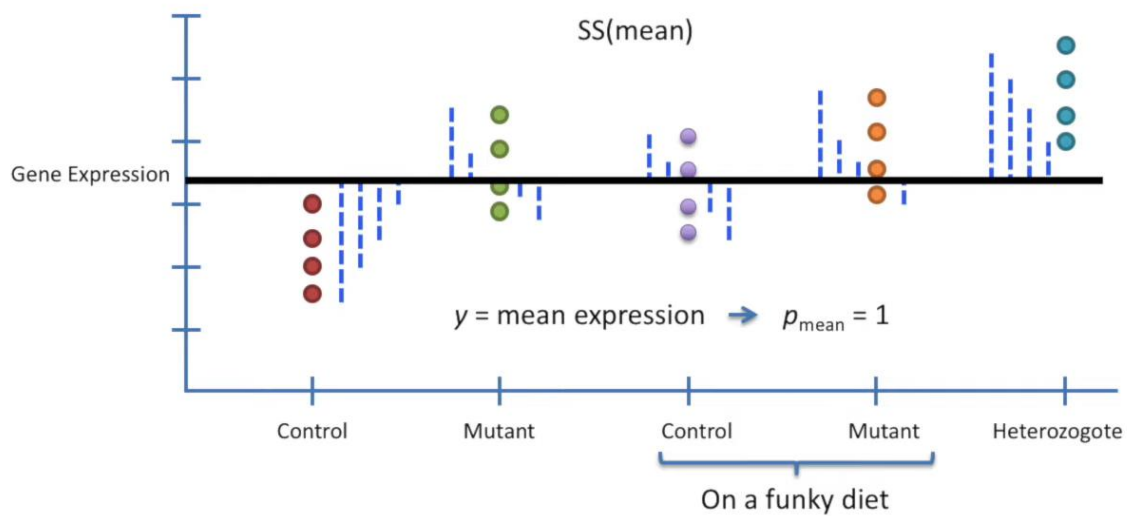
Now let's do an ANOVA!!

ANOVA tests if all 5 categories are the same.

Gene Expression

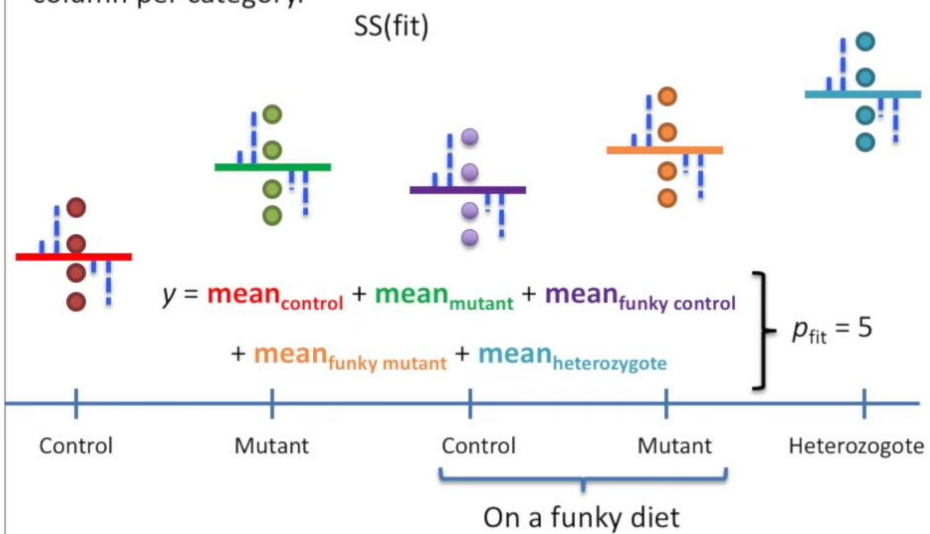






1	0	0	0	0
1	0	0	0	0
1	0	0	0	0
1	0	0	0	0
0	1	0	0	0
0	1	0	0	0
0	1	0	0	0
0	1	0	0	0
0	0	1	0	0
0	0	1	0	0
0	0	1	0	0
0	0	0	1	0
0	0	0	1	0
0	0	0	1	0
0	0	0	0	1
0	0	0	0	1
0	0	0	0	1
0	0	0	0	1

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



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