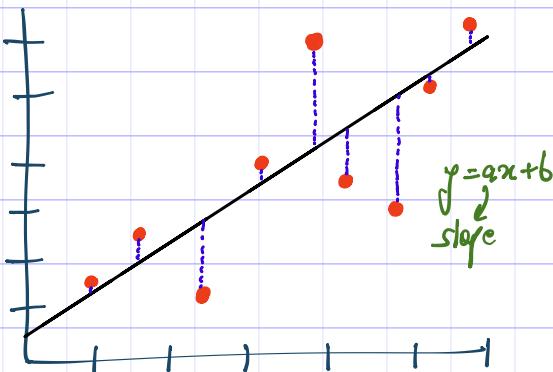


Linear Regression

- fitting a line to data
- Least squares
- Linear Regression

Concepts:

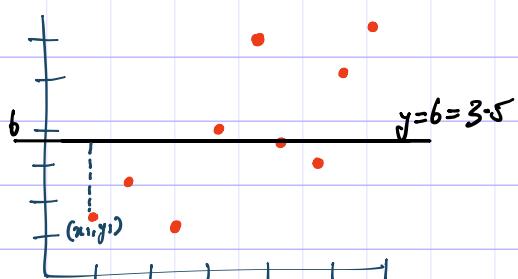
- #1 We want to minimize the distance between the observed values and the line



- #2 we do this by taking the derivative and finding where it is equal to 0

- #3 The final line minimizes the sum of squares (It gives the "least squares") between it & the real data.

We can measure how well any line fits the data by seeing how close it is to the data points.



The distance between the line & the first data point is $b - y_1$, the second data point is $b - y_2$

$$\text{Total distance} = (b - y_1) + (b - y_2) + (b - y_3) + (b - y_4) + \dots$$

* we can take the absolute values $|b - y_i|$ but it will make the math tricky.

so they ended up squaring each term. Squaring ensures that each term is positive.

$$(b - y_1)^2 + (b - y_2)^2 + (b - y_3)^2 + \dots \\ = \Sigma (b - y_i)^2$$

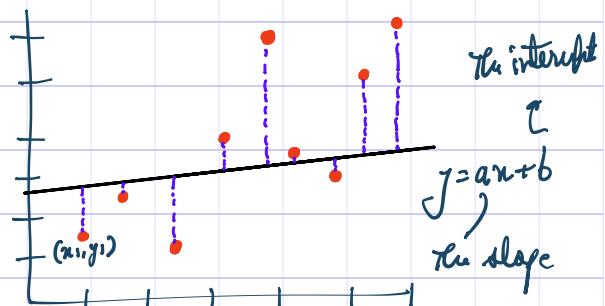
* This is our measure of how well this line fits the data.

* It's called the sum of squared residuals

↳ difference between the real data & the line

- # if we rotate the line a little bit, the sum of squared residuals = 18.72

(better than before)



- # we want to find the optimal values of \textcircled{a} & \textcircled{b} so that we minimize the sum of squared residuals.

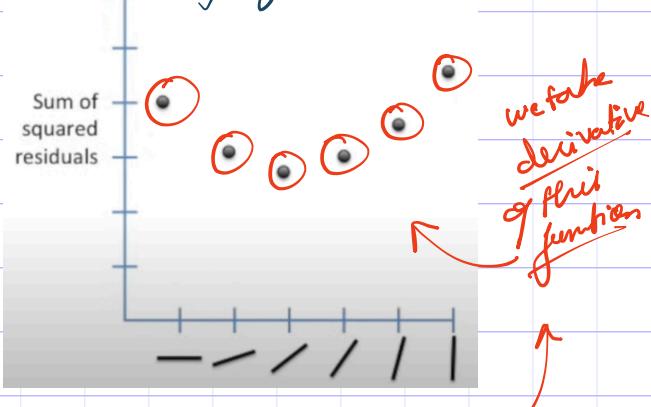
sum of squared residuals = $(ax_1 + b) - y_1^2 + (ax_2 + b) - y_2^2 + \dots$

↳ value of the line at x_1 ↳ observed value at x_1

Calculating the distance between the line & the observed value

Least Squares: we want the line that will give us the smallest sum of squares, this method of finding the best values of \textcircled{a} & \textcircled{b} is called "Least Squares"

If we plotted the sum of squared residuals vs. each rotation we'd get something just like this:

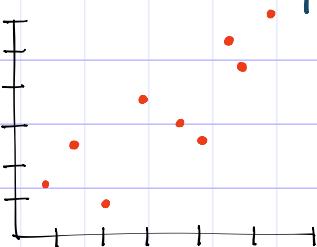


To find the optimal rotation for the line the derivative tells us the slope of the function at every point.

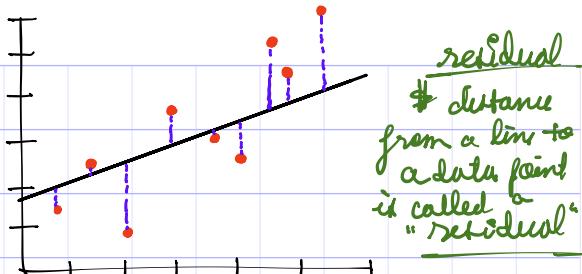
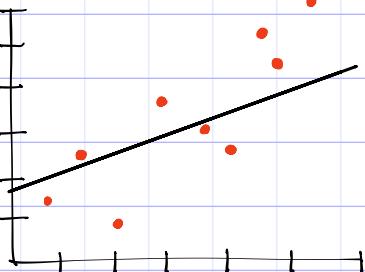
- # The slope at the best point is zero
- # different rotations are just different values of (① slope) & (② the intercept).

Main Ideas behind linear Regression

- ① use least squares to fit a line to the data
- ② calculate R^2
- ③ calculate a p-value for R^2



* draw a line through the data



measure the distance of the data from the line, square each distance & add them up.

$$y = 0 \cdot 1 + 0 \cdot 78x$$

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

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

- ① Calculate the average mouse size
- ② Now, sum the squared residuals
- ③ find $SS(\text{mean})$ "sum of squares around the mean"

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

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

sample size

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

average sum of squares per mouse

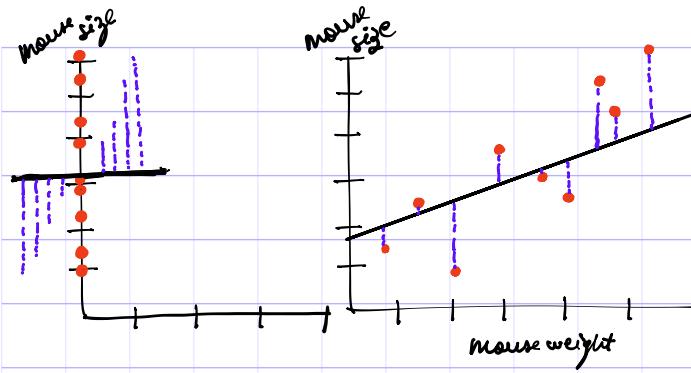
$SS(\text{fit}) \Rightarrow$ sum of squares around the least squares fit

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

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

distance between the data point & the line squared & added

Variance(Something) = $\frac{\text{sum of squared}}{\text{the no. of those things}}$



- # R^2 tells us how much of the variation in mouse size can be explained by taking mouse weight into account.
- heavier mice are bigger
 - lighter mice are smaller

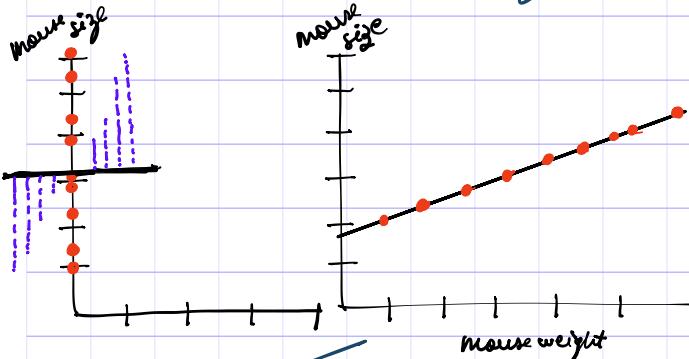
$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{fit})}{\text{Var}(\text{mean})}$$

$$R^2 = \frac{11.1 - 4.4}{11.1} = 0.6 = 60\%$$

* 60% reduction in variance when we take the mouse weight into account

OR

mouse weight explains 60% of the variance in mouse size



In this case knowing mouse weight means you can make a perfect prediction of mouse size.

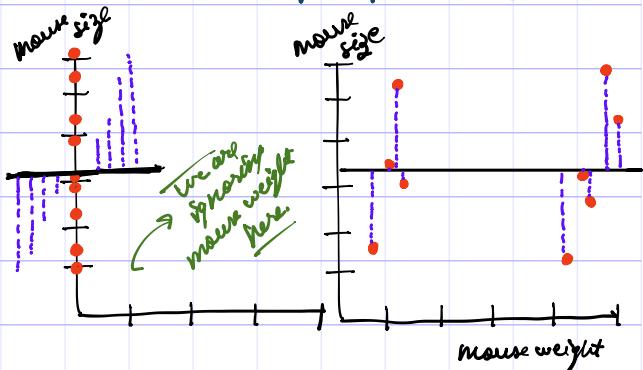
$$\text{Var}(\text{mean}) = 11.1 \quad \text{Var}(\text{fit}) = 0$$

$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{fit})}{\text{Var}(\text{mean})} = 100\%$$

mouse weight "explains" 100% of the variation in mouse size.

Additionally, knowing mouse weight doesn't help us predict mouse size.

- heavy mice could be large or small
- light mice could be large or small (with equal probability)



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

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

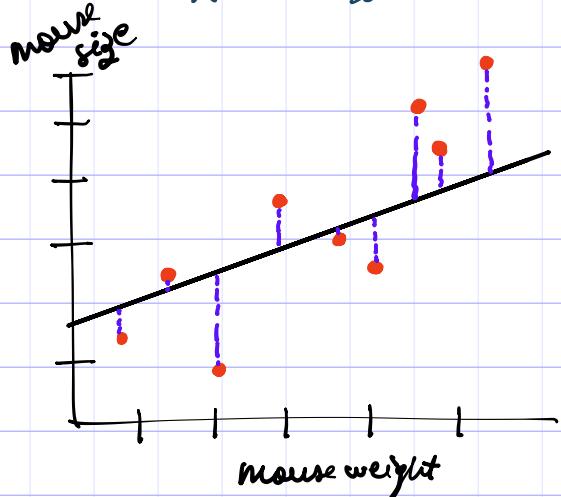
$$R^2 = \frac{\text{var}(\text{mean}) - \text{var}(\text{fit})}{\text{var}(\text{mean})} \Rightarrow 0\%$$

* mouse weight doesn't "explain" any of the variation around the mean.

we applied R^2 to a simple equation of line

$$y = 0.1 + 0.78x$$

$$R^2 = 60\%$$



meaning # 60% of the variation in mouse size could be explained by mouse weight.

measure, square & sum the distance from the data to the mean

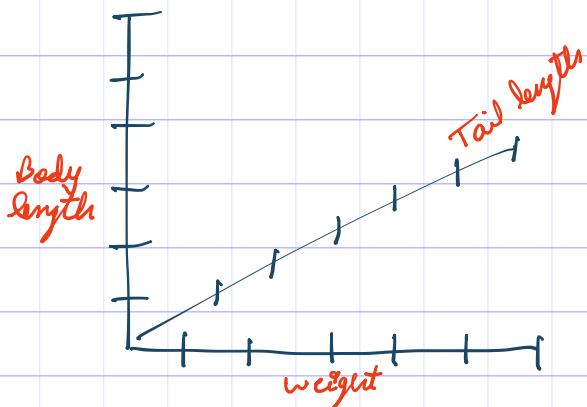
$$R^2 = \frac{\text{SS}(\text{mean}) - \text{SS}(\text{fit})}{\text{SS}(\text{mean})}$$

measure, square & sum the distance from the data to the complicated equation

more complicated example 😊

we wanted to know if mouse weight & tail length did a good job predicting the length of the mouse's body.

3-D graph

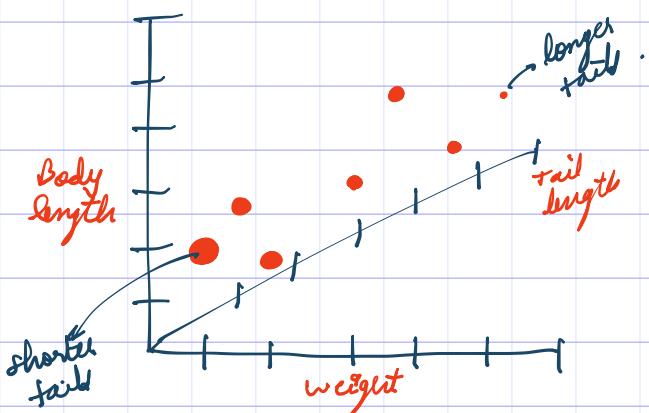


multi variable model

we measure a bunch of mice

	Mouse weight	Tail length	Body length
3.5	2.9	3.1	
4.3	2.1	2.8	
5.9	4.1	6.1	
4.8	3.2	3.8	
...	

we want to know, how well weight and tail length ... predict body length



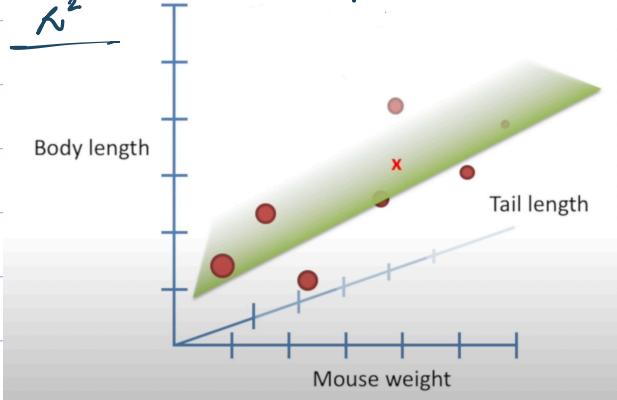
since we have 3-D graph we will fit plane.

$$\hat{y} = 0.1 + 0.7x + 0.5z$$

Y intercept (Tail length, mouse weight equal to 0)

weight Tail length
body length

Just like before, we can measure the residuals, square them & then add them up to calculate R^2



Consider equation $\rightarrow \hat{y} = 0.1 + 0.7x + 0.5z$

→ If the tail length is useless? doesn't make $SS(fit)$ smaller then least-squares will ignore it by making that parameter = 0

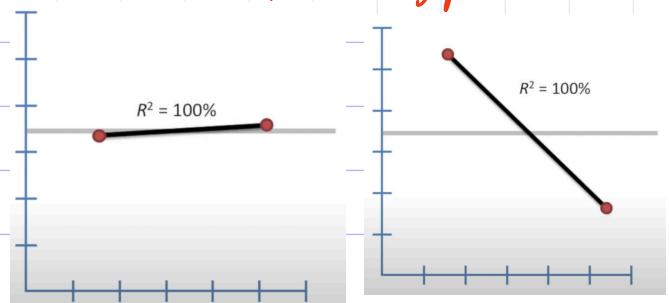
→ plugging the tail length into the equation would have no effect on predicting the mouse size.

This means, equations with more parameters will never make $SS(fit)$ worse than equations with fewer parameters.

mouse size = 0.3 → mouse weight + flip a coin + fav color + astrological sign.

The more parameters we add to the equations, the more opportunities we have for random events to reduce $SS(fit)$ & result in a better R^2

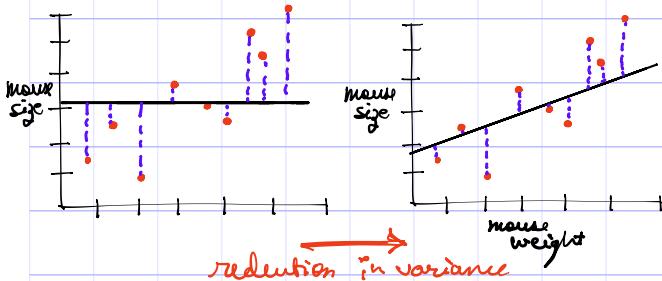
"adjusted R^2 " → scales R^2 by the number of parameters



100% is a great number, but any two random points will give us the same thing.

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

we need a p-value



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

$$= \frac{\text{Var}(\text{mouse size}) - \text{Var}(\text{after taking weight into account})}{\text{Var}(\text{mouse size})}$$

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

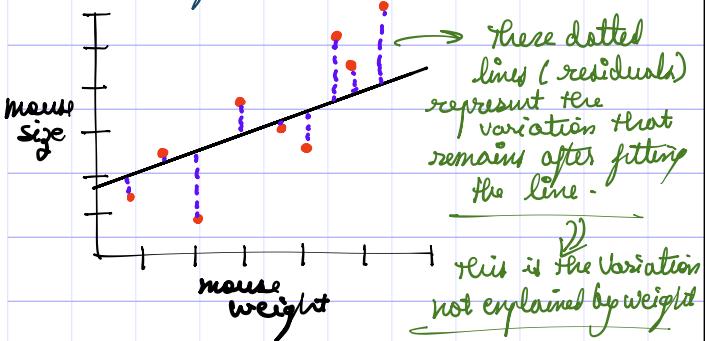
The p-value of R^2 comes from something called "F"

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

Numerators for R^2 & F are same \rightarrow

variation in mouse size explained by weight

which is the reduction in variance when we take weight into account



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

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

Numerators are same

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

or

$$F = R^2 * \text{degree of freedom}$$

$$\boxed{DOF = \frac{n - p_{\text{fit}}}{p_{\text{fit}} - p_{\text{mean}}}}$$

$p_{\text{fit}} \Rightarrow$ the no. of parameters in the fit line

fit line $\Rightarrow y = a_1x + b$.

2 parameters

$p_{\text{mean}} \Rightarrow$ the no. of parameters in the mean line

$y = b \Rightarrow 1 \text{ parameter}$

$$p_{\text{fit}} - p_{\text{mean}} = 2 - 1 = 1$$

$$\boxed{\text{ss}(\text{mean}) - \text{ss}(\text{fit}) / (p_{\text{fit}} - p_{\text{mean}})}$$

Variance explained by extra parameters

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

more parameters you have in equation, the more data you need to estimate them

Ex: you need two parameters to estimate a line, but you need 3 parameters to estimate a plane

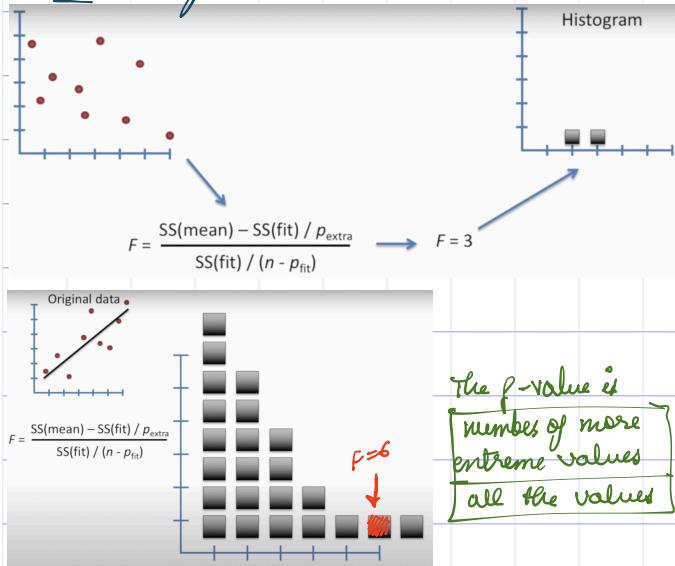
If the fit is good:

$$F = \frac{\text{variation explained by extra parameters in the "fit"}}{\text{variation not explained by extra parameters in the "fit"}}$$

variation not explained by extra parameters in the "fit" \rightarrow large number / small number

$F = \text{really large number}$

find the values of F for different set of random data & create a histogram.



Summary (linear regression)

given some data that you think are related ...

- 1) Quantifies the relationship in the data (this is R^2)
→ This needs to be large

- 2) Determine how reliable that relationship is (this is the p-value that we calculate with F)
→ This needs to be small

You need both to have an interesting result

linear regression in R

```
mouse.data <- data.frame(
  weight = c(...), size = c(...))
```

mouse.data

```
plot(mouse.data$weight, mouse.data
  $size)
```

mouse.regression ←

```
lm(size ~ weight, data = mouse.data)
```

y-values

x-values

$$y = ax + b$$

↓ slope

$$\text{size} = y\text{-intercept} + \text{slope} \times \text{weight}$$

```
> mouse.regression <- lm(size ~ weight, data = mouse.data)
> summary(mouse.regression)
```

Call:

```
lm(formula = size ~ weight, data = mouse.data)
```

Residuals:

Min	1Q	Median	3Q	Max
-1.5482	-0.8037	0.1186	0.6186	1.8852

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.5813	0.9647	0.603	0.5658
weight	0.7778	0.2334	3.332	0.0126 *

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.19 on 7 degrees of freedom
Multiple R-squared: 0.6133, Adjusted R-squared: 0.558

F-statistic: 11.1 on 1 and 7 DF p-value: 0.01256

original call to the lm function

summary of the residuals (closer to zero)

weight gives us the reliable estimate of size

→ least squares estimates for the fitted line

$$\text{size} = y\text{-intercept} + \text{slope} \times \text{weight}$$

0.5813 value of the intercept

$$\text{size} = 0.5813 + \text{slope} \times \text{weight}$$

0.7778 value for the slope

$$\text{size} = 0.5813 + 0.7778 \times \text{weight}$$

we want the p-value to be < 0.05 for statistically significant