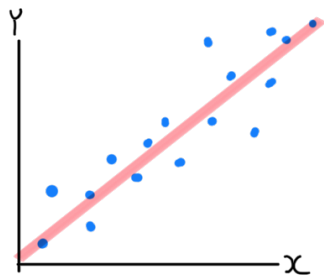


Linear Regression!

What is the main idea?



→ The output variable
Models a **Linear Relationship**
Between input variables



Types of Linear Regression

Simple Linear Regression (SLR)

1 input variable.
Changes in Y is caused by changes
in X

Multiple Linear Regression (MLR)

There are more than 1 input

line is more than 1 variable

★ Terminology ★

Input Variable

aka: predictor variable,
Explanatory variable
Independent variable

Output Variable

aka: Response variable
Explained variable
Dependent variable

Note: Gotta Love Namings in Science Amirite?

Okay, So what does Linear Mean?

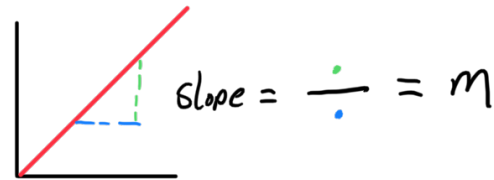


Maths
God

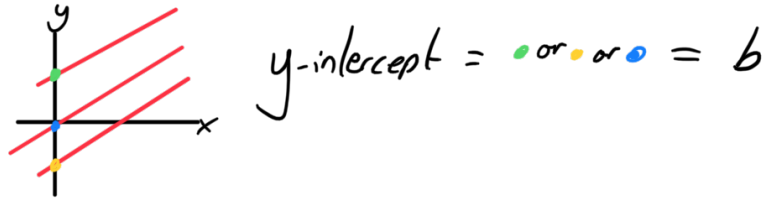
"I hereby declare Linearity to be the property of a mathematical function that can be represented as a straight line"

Nice!

So a straight line has the formula
of $f(x) = mx + b$ where m is
the slope of the line



b is the y-intercept of the line



Now for Linear Regression

$$y = \underset{\substack{\uparrow \\ \text{y-intercept}}}{\beta_0} + \underset{\substack{\uparrow \\ \text{slope}}}{\beta_1} x$$

Note

This is
for SLR

\uparrow we will talk
about MLR later.

$\hat{}$ still the equation of a line!

Linear Regression Goals

Model relationship between
an output variable and input

Variable

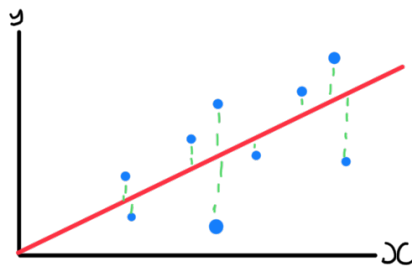
How do we make our Model?

So we know our Model will look like $y = \beta_0 + \beta_1 x$ But what are our β_0 and β_1 ?

In a perfect world we could find β_0 and β_1 exactly but unfortunately for us we can only estimate them.

Our estimations are called $\hat{\beta}_0$ and $\hat{\beta}_1$

Also since we are estimating we will add an error term to our equation to account for the error.



$\vdots = \epsilon = \text{Error / Residual}$

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

Final Model will be.

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

How do we find a model that fits our data the best?

- Some lines fit better than others

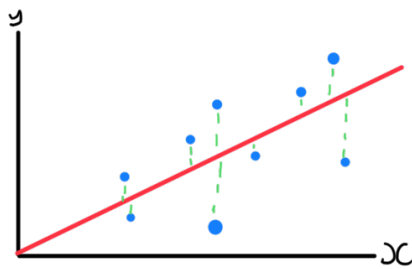
Well one way of checking how well our line fits our data is to check the residuals.

(from our formula)

$$\epsilon_i = y_i - \beta_0 + \beta_1 x_i$$

actual y_i at x_i

line y_i at x_i



So to check how well our line fits we look at the residuals for all our points which we could say is

$$\sum_{i=1}^n \epsilon_i = \sum_{i=1}^n y_i - \beta_0 - \beta_1$$

However!

This doesn't work

↑
this is
wrong

this formula can turn out to be zero
even if we have massive residuals because
Residuals should alternate in sign since they
are above and below the line.
To account for this we square the residual
to get a formula for how well our line
fits the data.

$$\sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

⇓


- If this number is high
our line is not a good
fit.
- if this is close to zero
it is a good fit
- If it equals zero it is
a perfect fit.

Now that we have the formula
for checking how well the line fits
what can we do.

well with a bit of calculus we
can see what happens when

1 11


$\sum (y_i - \beta_0 - \beta_1 x_i)^2$ is Zero and the
Solve for β_0 and β_1 .

 This will give us
the Best fit Parameters

When we do this:

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$
$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{\sum_i^n (y_i - \bar{Y})(x_i - \bar{X})}{\sum_i^n (x_i - \bar{X})^2}$$

But you ask: "Emm why the little hats?"

 Well. These are just Estimates
for β_0 and β_1 and thus the hats

Fitted Simple Linear Regression Model

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

  
iwoooo!

Now, A few things about our model
- we assumed the errors were zero

(it was the best fit)

- we can get 75% confidence intervals for $\hat{\beta}_0$ and $\hat{\beta}_1$
- we can check if β_0 or β_1 is significantly different from 0 with a t-test or ANOVA
- SST = Total sum of Squares
 $\hookrightarrow \sum_i^n (y_i - \bar{y})^2$
- SSR: The Sum Squared deviation from the mean variation in y explained by the Regression Line
 $\hookrightarrow \sum_i^n (\hat{y}_i - \bar{y})^2$
- SSE: Sum of Squared Errors.
variation in y left unexplained
 $\hookrightarrow \sum_i^n (\hat{y}_i - y_i)^2$

Also:

$$\text{SST} = \text{SSR} + \text{SSE}$$

\hookrightarrow The overall variability in y \downarrow Explained by model \downarrow unexplained

Once we are finished we might wish to measure the strength of the linear relationship.

To do this we will use the coefficient of determination
 \Downarrow
 R^2

$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$ which is the fraction of the variation which is explained by the model.

* R^2 does not tell us if changes in x cause changes in y