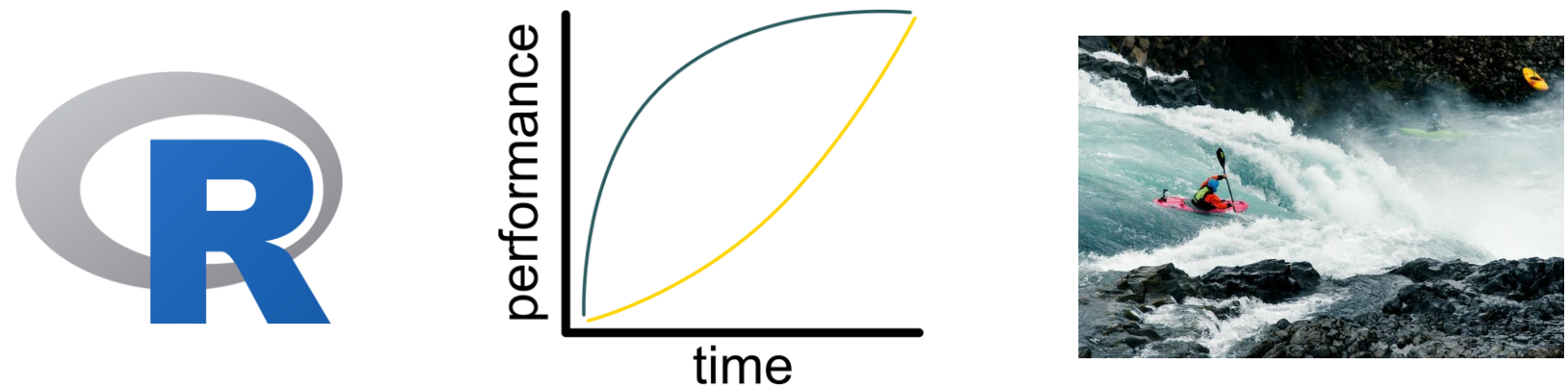
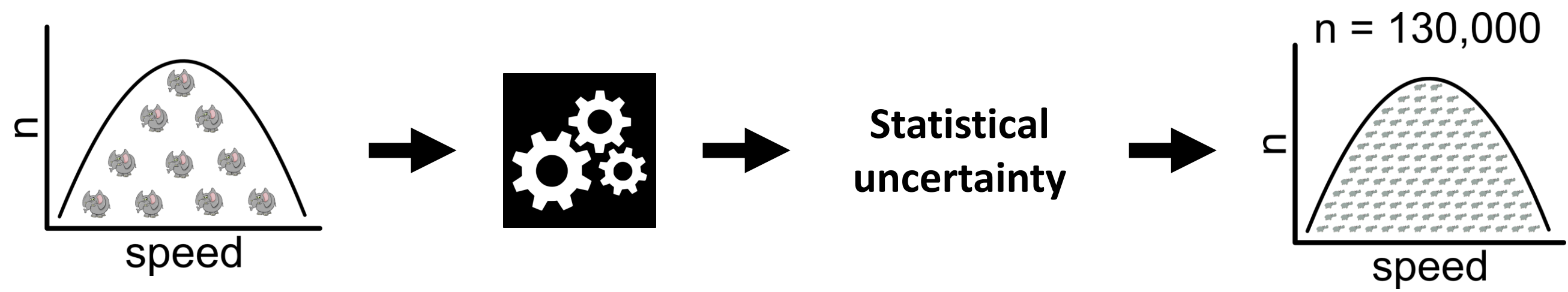
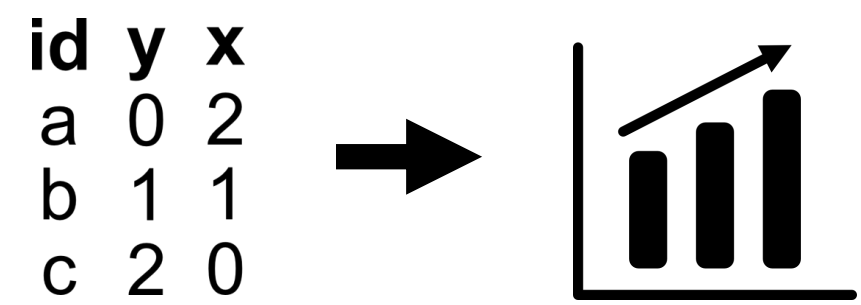


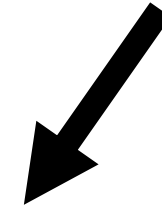
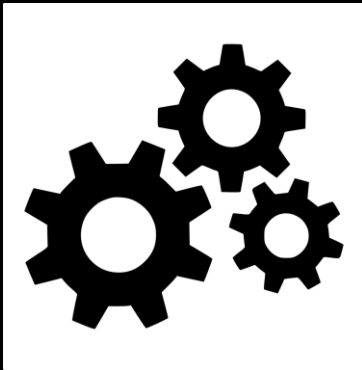
Regression

Pärt Prommik, PhD

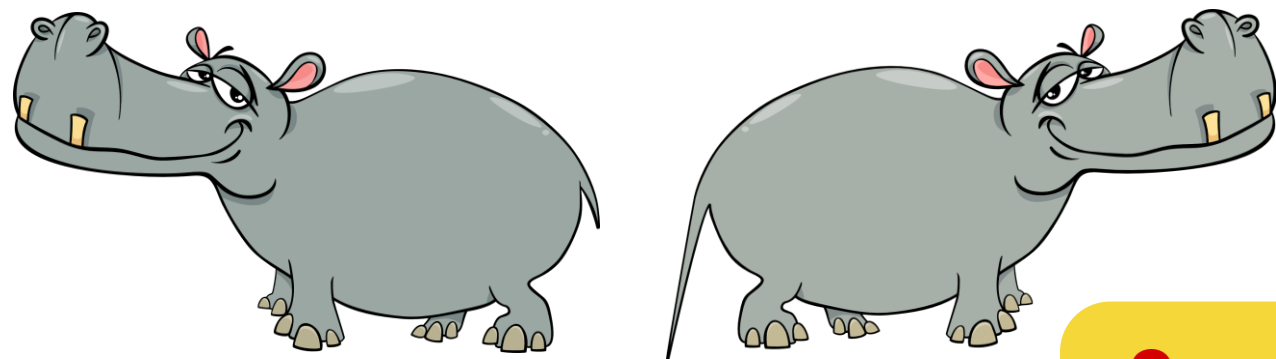
Ülo Maiväli, PhD

Reminder





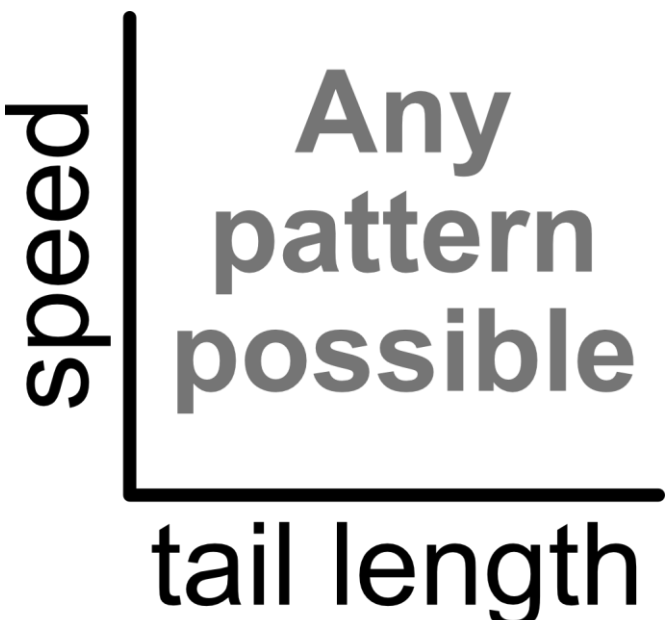
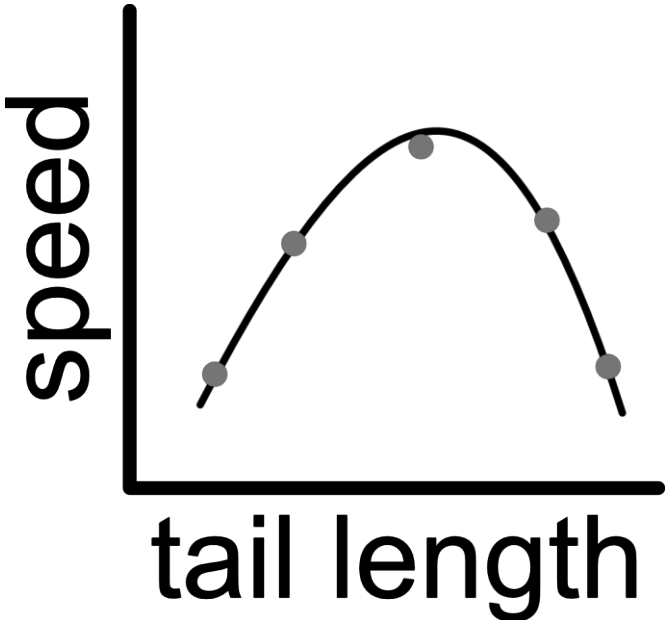
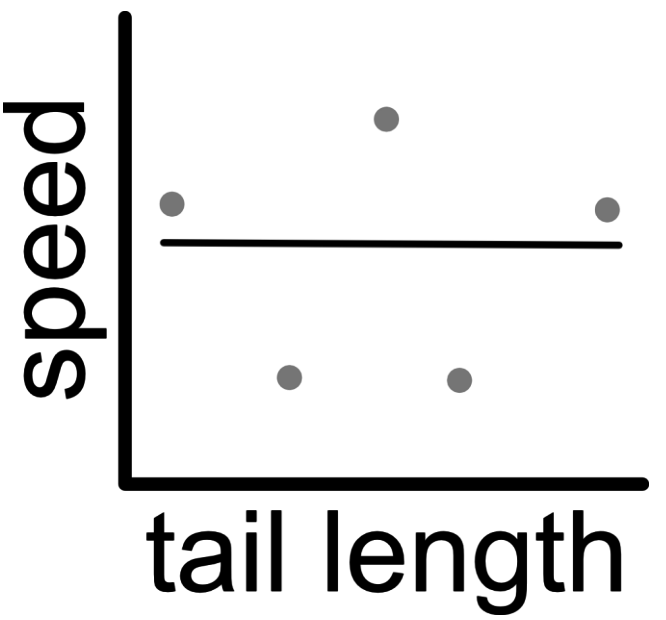
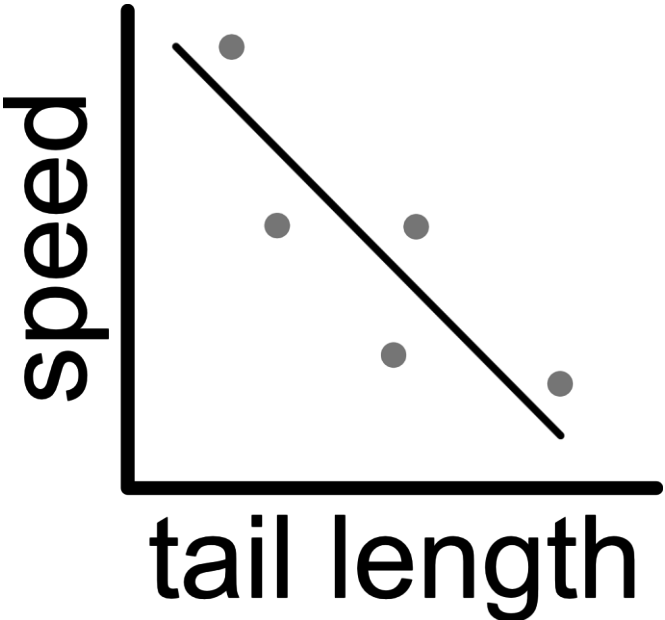
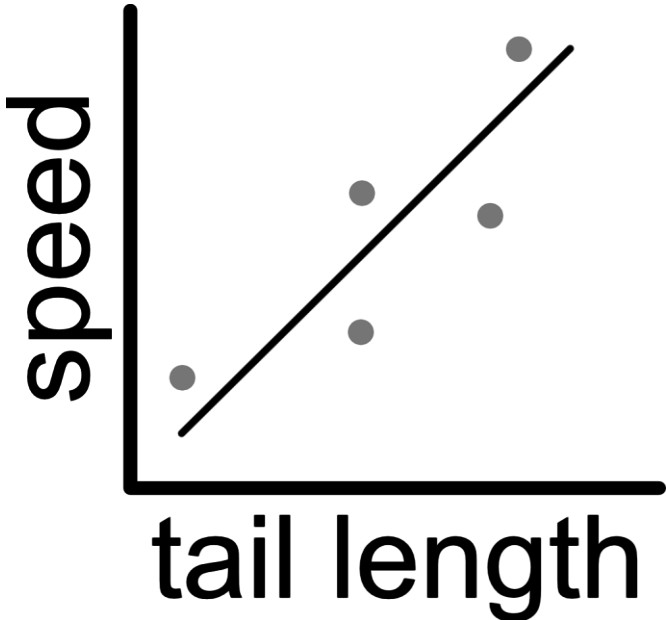
Analysing relationships



speed	tail length
20	10
31	32
50	17
28	23
44	39



Lines are good for
expressing relationships



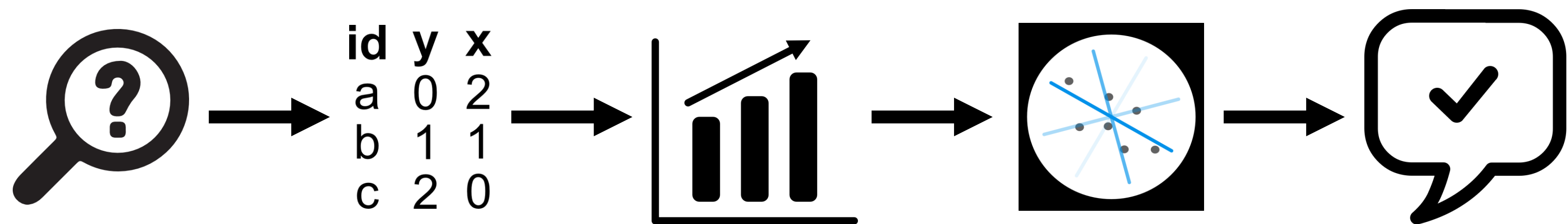
Line drawing task

I'll draw few data points.

You come and draw a **STRAIGHT LINE** that fits the data points:
it must be **AS CLOSE AS POSSIBLE** to all points.

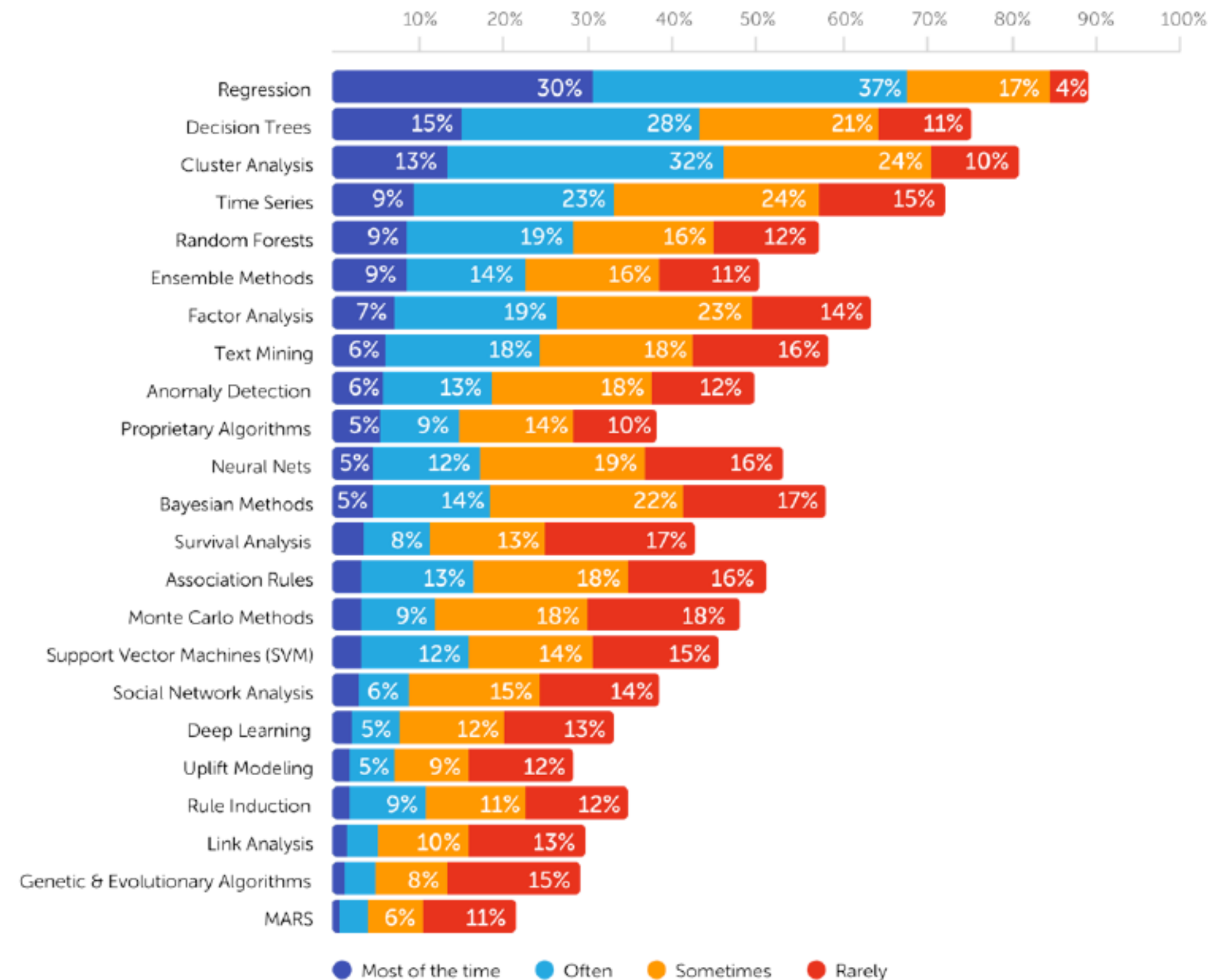
You just did a **regression**

Regression is a statistical method that estimates relationship between two or more variables by finding a line of best fit.



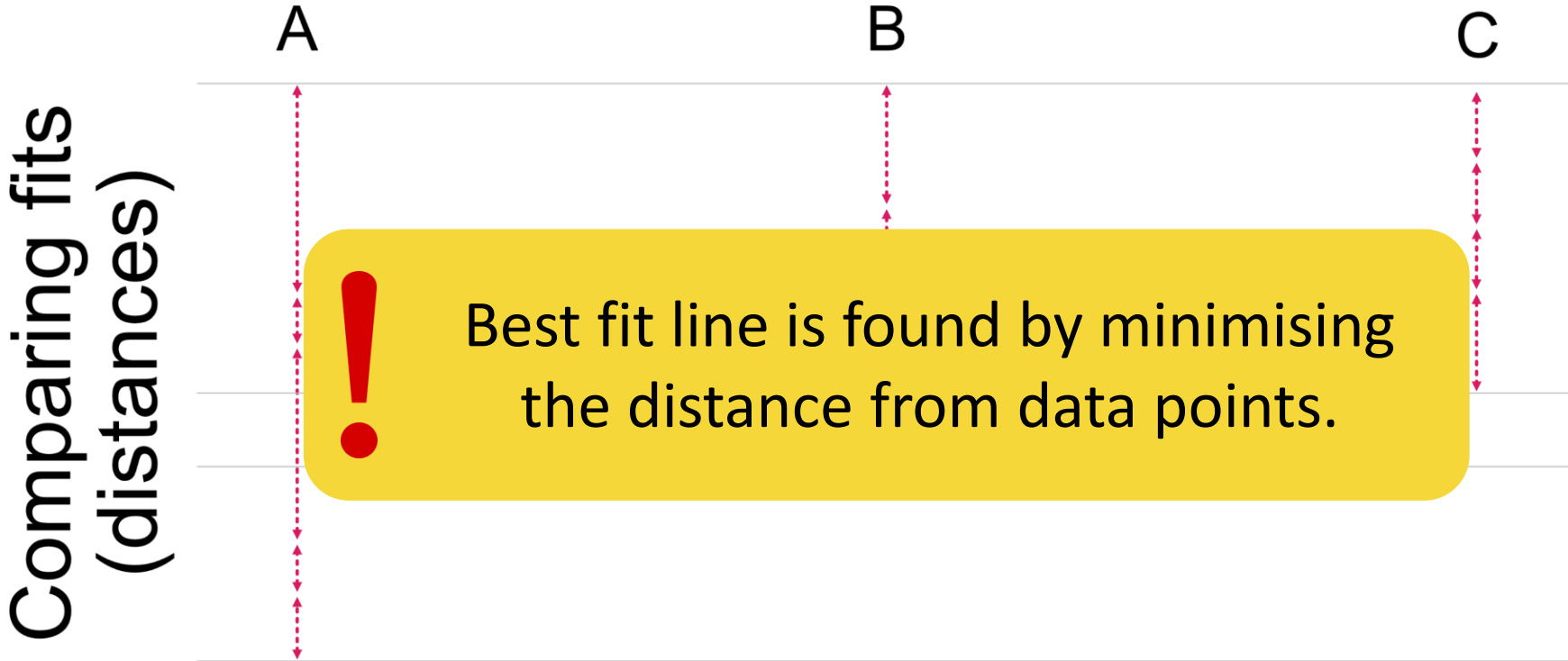
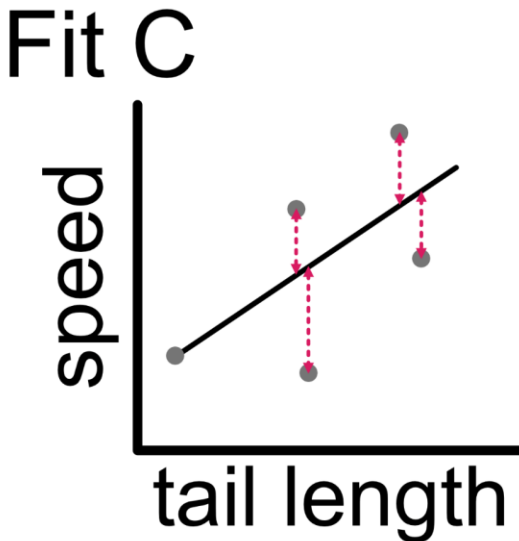
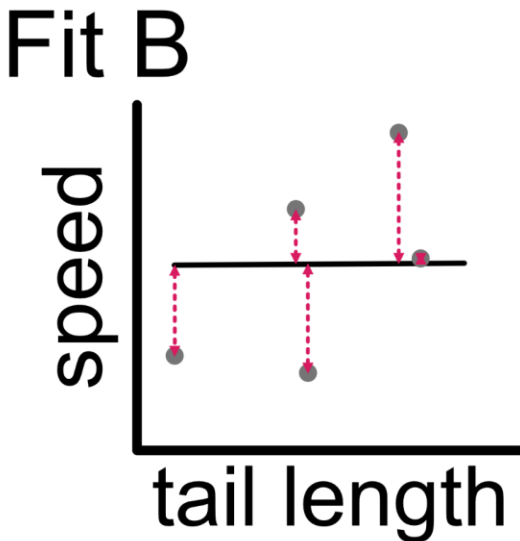
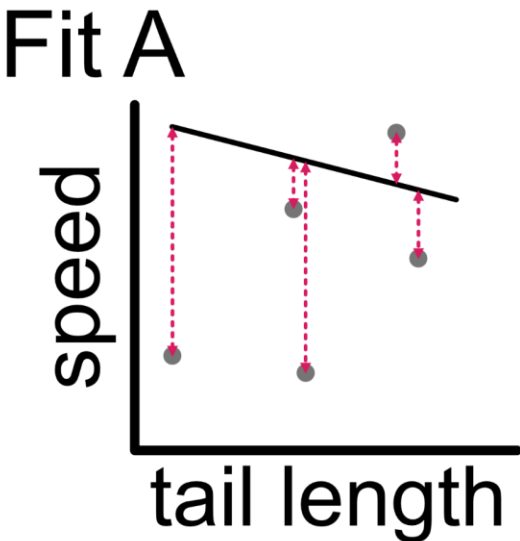
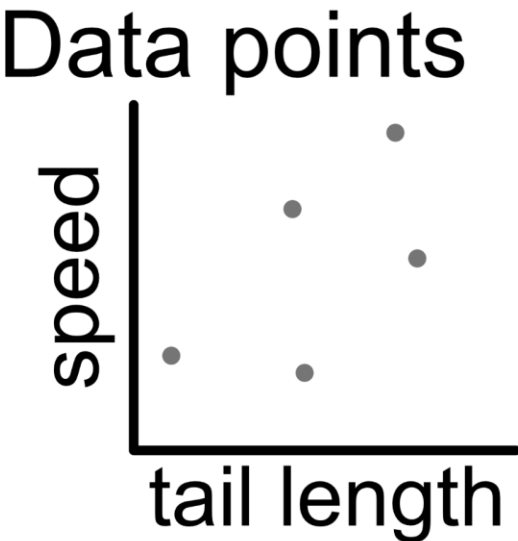
Why regression?

- The most used technique
- Powerful
- Flexible

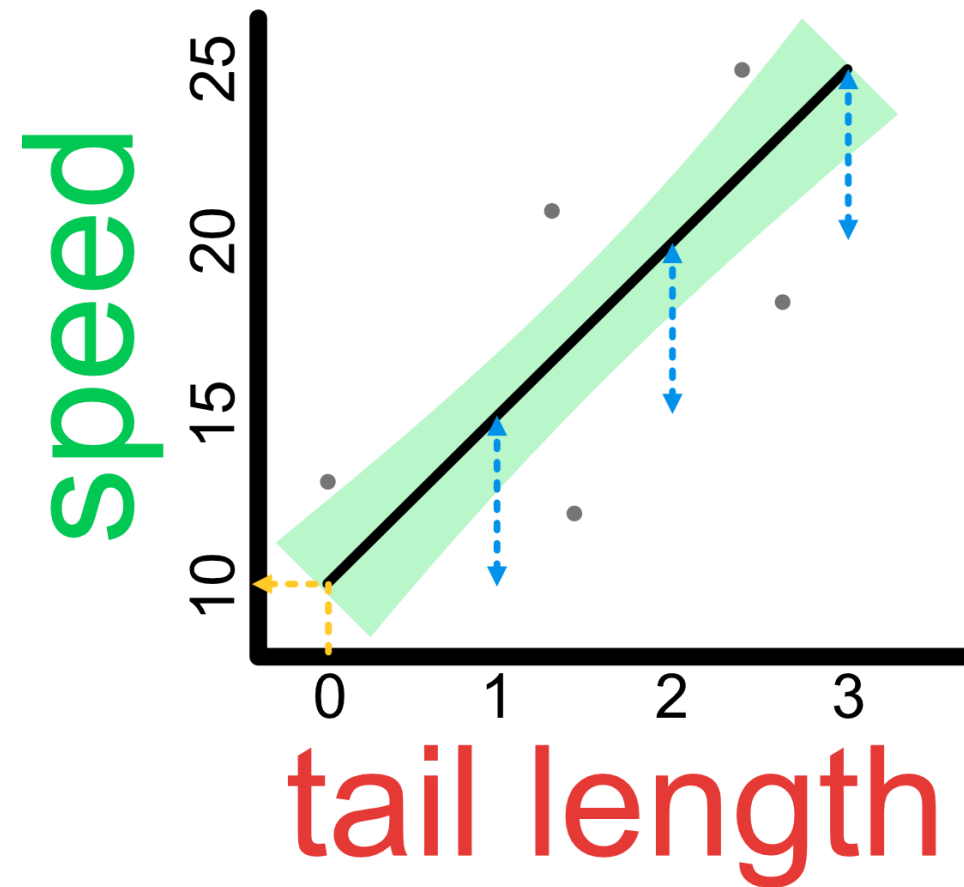


How the best fit line is found?

! Residual: the difference between the observed value and the estimated value.



A line can be described mathematically as an **equation**



Y
Predicted output
Outcome variable
Response variable
Dependent variable

speed

Intercept
Constant

Coefficients

Slope

X
Explanatory variable
Predictor variable
Independent variable

Residual standard deviation
Residual standard error
Error
Disturbance

$\text{speed} = a + b \times \text{tail_length} + e$

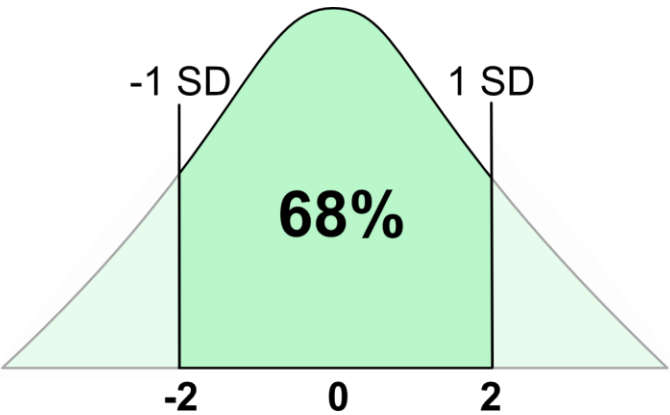
- 1. **Intercept:** starting point
- 2. **Slope:** inclination
- 3. **Residual SD:** accuracy of relationship

Equations are good for answering different study questions

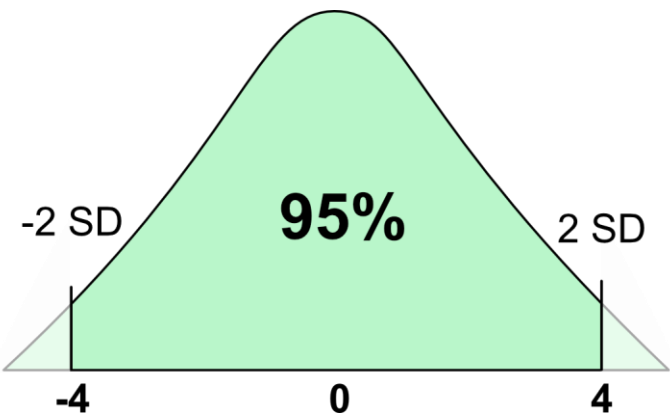
Descriptive

tail length

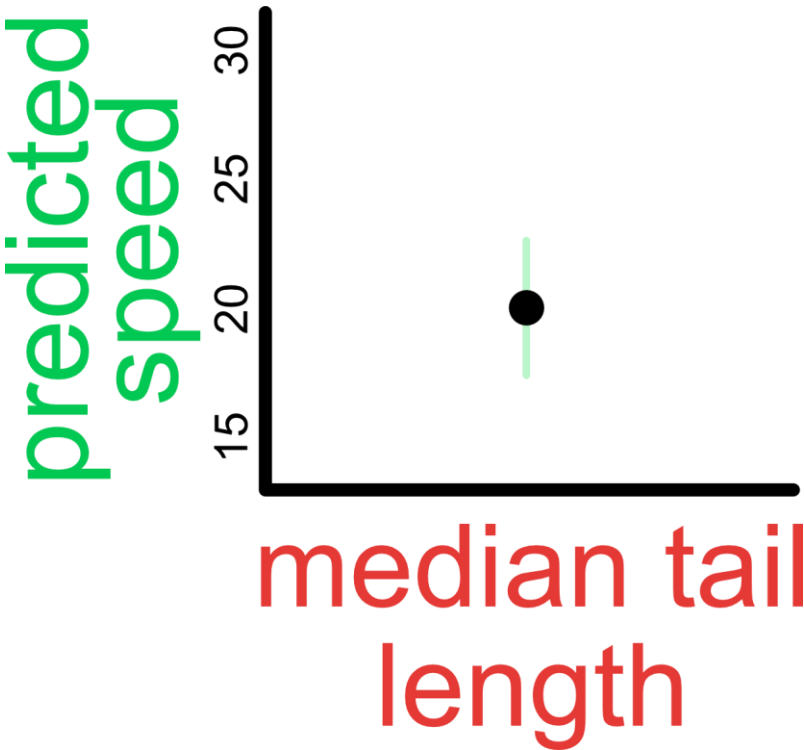
$$20 = 10 + 5 \times \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 5 \end{matrix} + 2$$



Lower confidence limit: $20 - 2 = 18$
Upper confidence limit: $20 + 2 = 22$



Lower confidence limit: $20 - 4 = 16$
Upper confidence limit: $20 + 4 = 24$

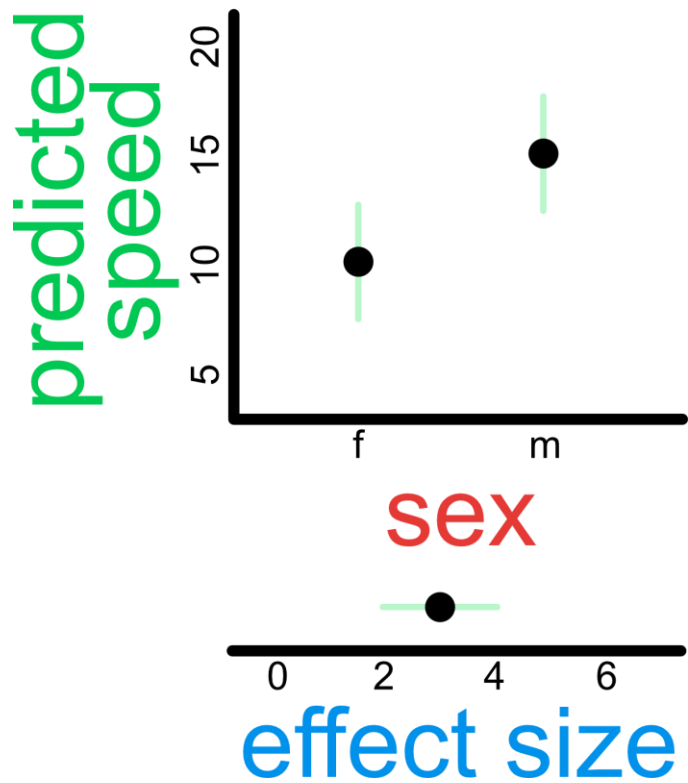


Equations are good for answering different study questions

Comparative

$$10 = 10 + 5 \times \begin{matrix} \text{sex} \\ f \\ m \end{matrix} + 2$$

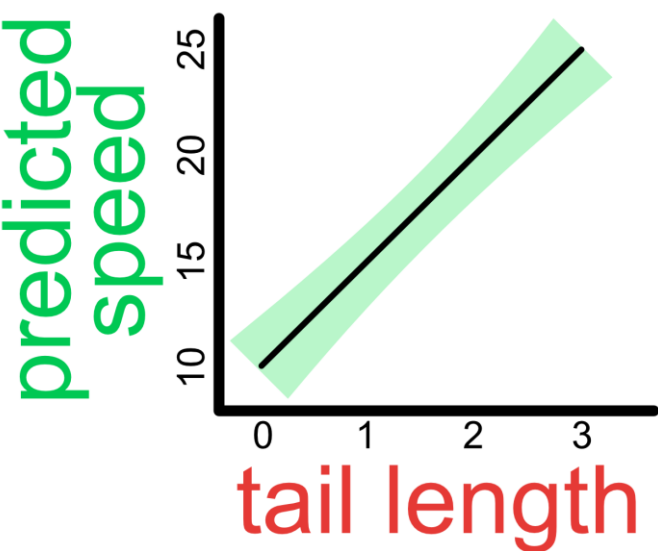
$$15 = 10 + 5 \times \begin{matrix} \text{sex} \\ f \\ m \end{matrix} + 2$$



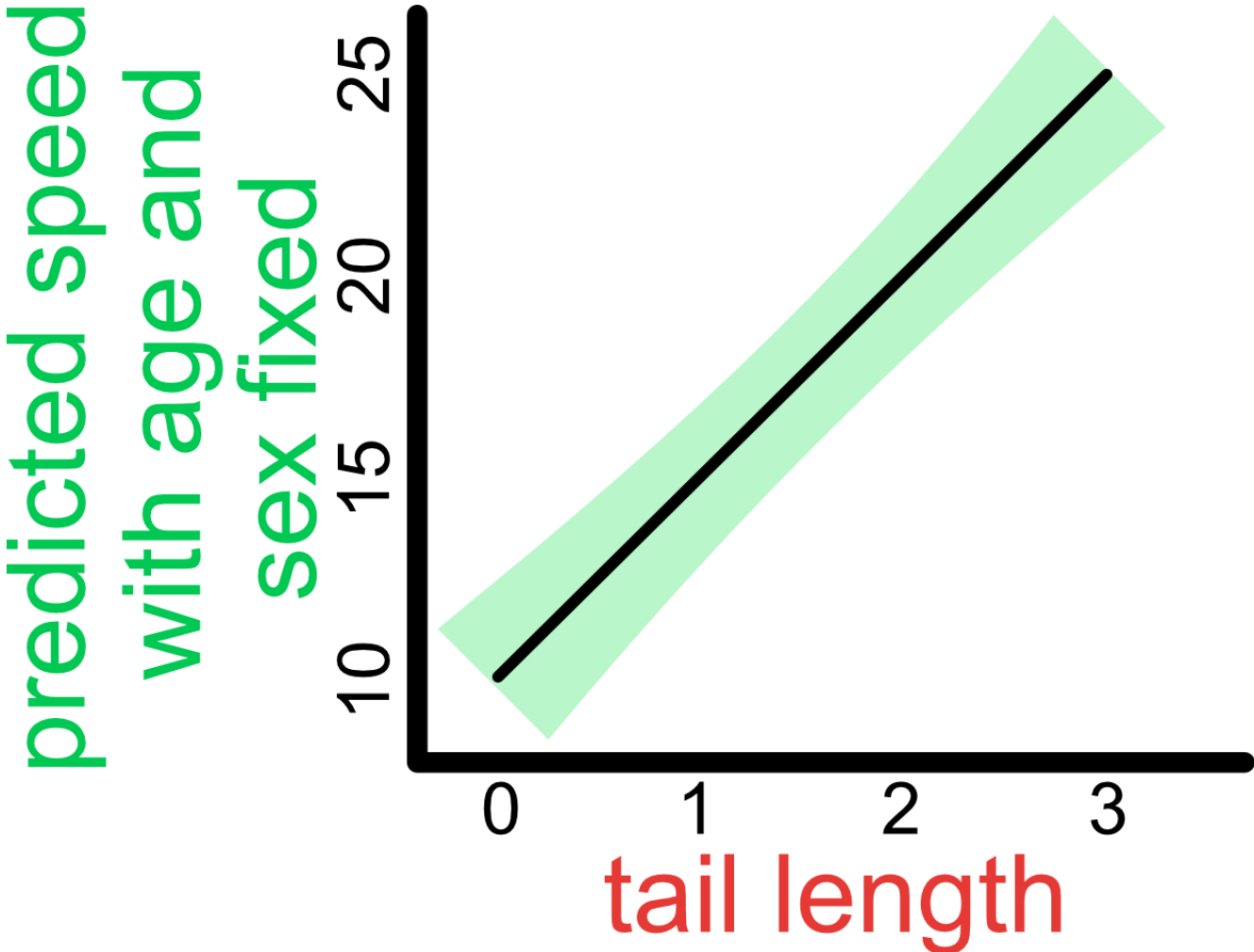
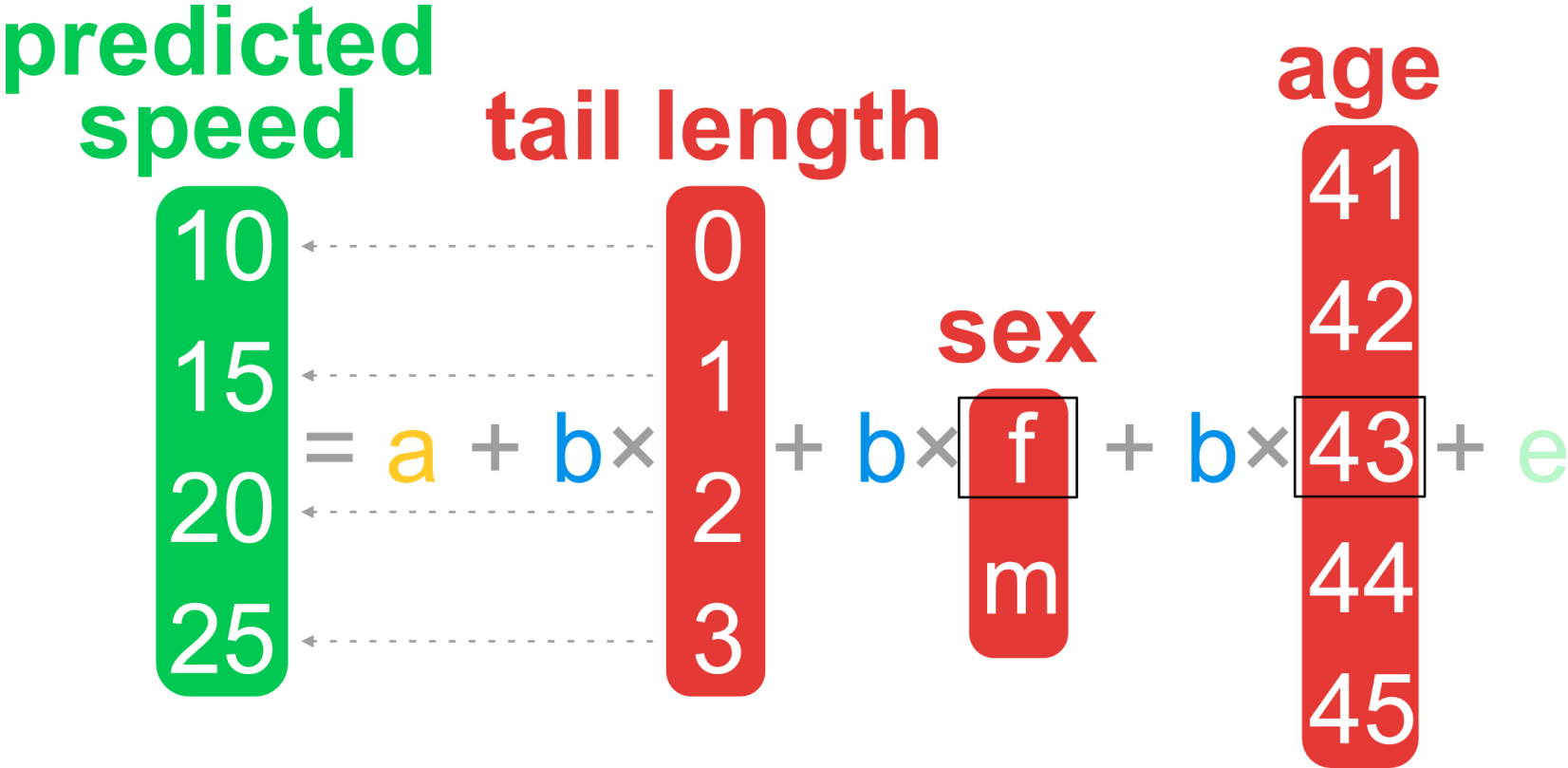
Relationship

predicted speed		tail length
10	←	0
15	←	1
20	←	2
25	←	3

$$= 10 + 5 \times + 2$$



Multiple variables are not a problem

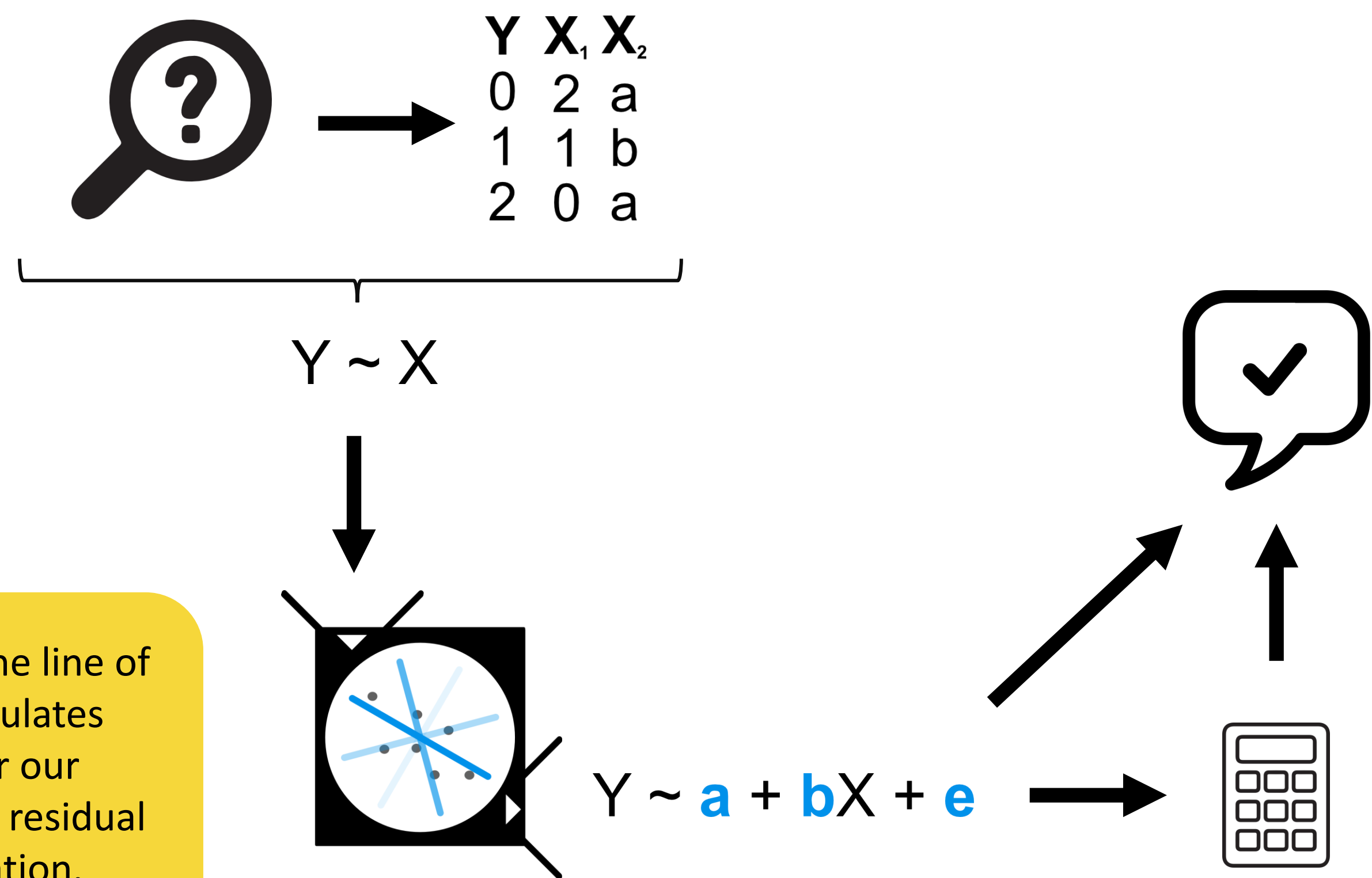


! Solves many kind of study questions!

Regression equations allow analysing the world



Line fitting machine's input and output



! Regression finds the line of best fit and calculates coefficients for our input formula and residual standard deviation.

Regression in R language

	speed	tail_length
1	20	15
2	25	12
3	32	18
4	37	14
5	34	24
6	39	22

Use variable names and tilde for determining formula

`speed ~ tail_length`

Use the formula inside a regression function

`function_name(speed ~ tail_length, data = yourdata)`