

Dividing your data set

- Do I need to divide my data set?
- How to choose the size of the fractions
- Randomisation

Do I need to divide my data set

YES

Data: $\{x_1, x_2, \dots, x_n\}$

n : # observations (# women)

$x_i := \begin{bmatrix} \text{ge1} \\ \text{ge2} \\ \vdots \\ \text{ge7018} \\ \text{coffee1} \\ \text{coffee2} \\ \text{coffee3} \end{bmatrix}$

Research question:

Which genes are related to coffee consumption?

Research question



Working hypotheses



Final hypotheses

Before accessing the data

After accessing the data

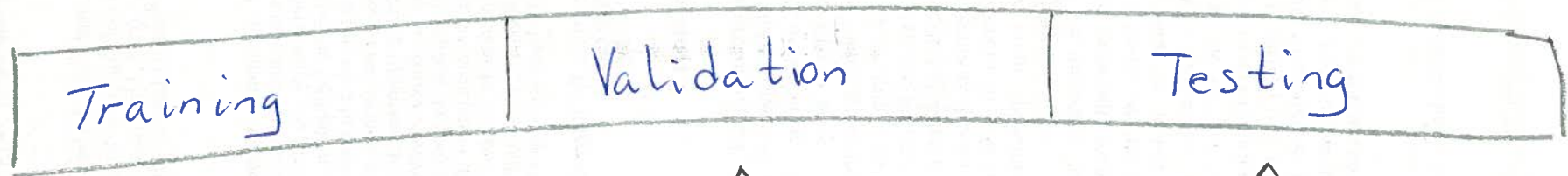
ANY adjustment of your method after accessing the data gives a HIGH risk of wrong conclusions due to overfitting.

ANY adjustment:

outlier removal, pre-processing, exclusion criteria,
p-value versus false discovery rate, grouping,
mean versus median,

HIGH risk:

The probability of making a wrong conclusion is
inaccessible.

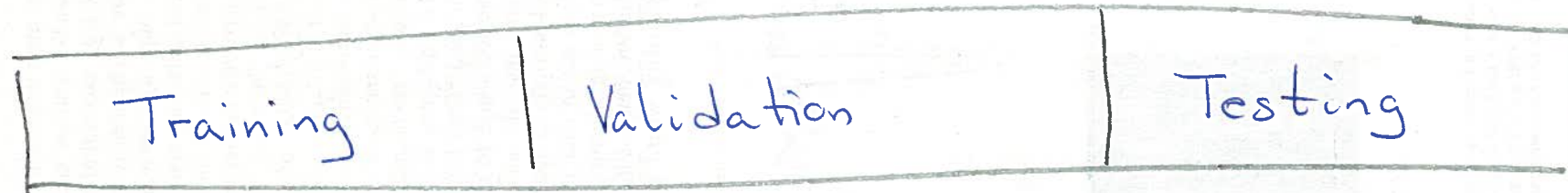


↑
Exploration
Choice of method

↑
Parameter setting
Grouping

↑
Calculation of
sensitivity/specificity,
p-values

Exploratory data analysis

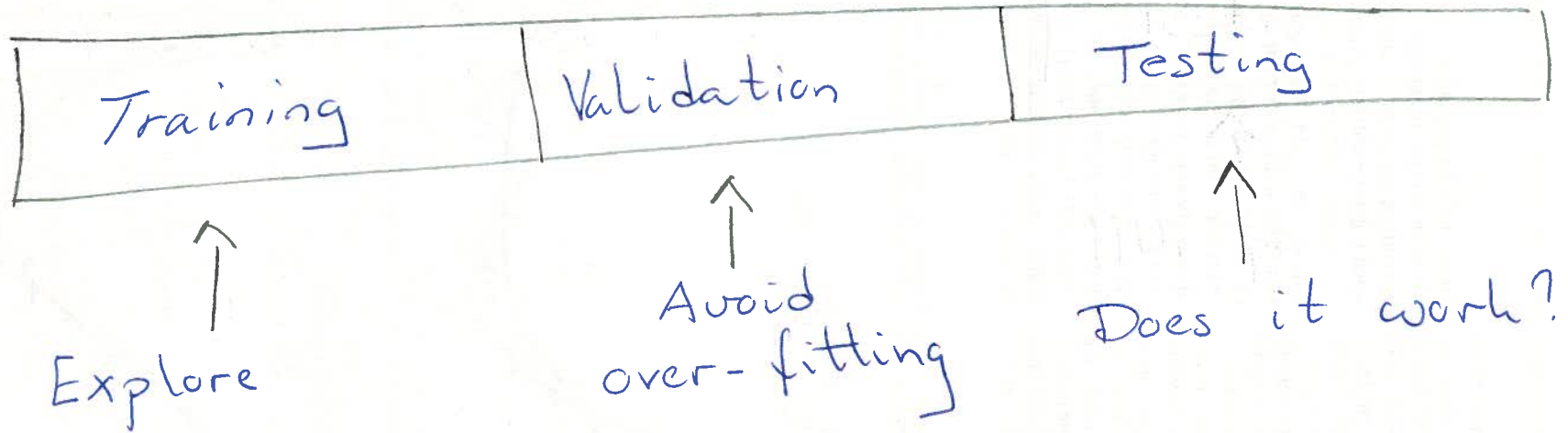


↑
1st paper

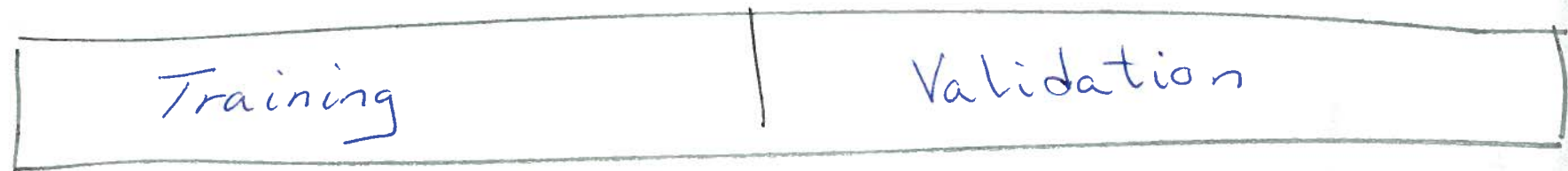
↑
2nd paper

↑
11th paper

Exploratory data analysis



Exploratory data analysis



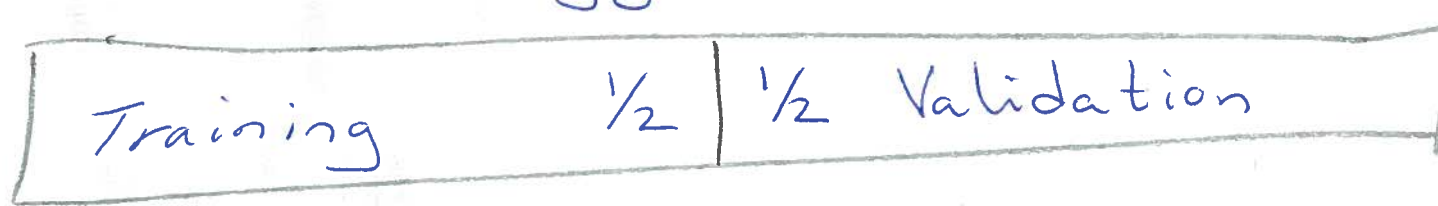
$n_{\text{train}} \rightarrow n$: better model

$n_{\text{validate}} \rightarrow n$: more confidence in the result

$n_{\text{train}} \rightarrow n$: model cannot be validated

$n_{\text{validate}} \rightarrow n$: confident that the model is wrong

Suggestion



Go explore

Set aside

Cross-validation: Estimate variance

POSITIVELY BIASED estimate of the performance of your models.

Increase training set?



Training	$\frac{2}{3}$ / $\frac{1}{3}$ Validation
----------	--

Limited by the number of observations
in the smallest class

Training $\frac{1}{1}$

Report methods, not results

↑
all

Cross-validation

"Wrong use of cross-validation is the single most costly error in medicine, biology, chemistry, physics, and related fields of research."

Correct use of cross-validation:

Any adjustment must be done independently for each repetition.

Exploratory research: Erase the memory of the researcher.

How to use Cross-Validation and Bootstrapping in exploratory research?

- Overly optimistic estimate of how your method works on an independent data set.
- Choose method, estimate variation in outcome

Randomisation

- Random not consecutive
- Stratified or not

Stratify if smallest class is very small.
Stratification gives positive bias.

Summary

- Change of plans (or no plan)

→ Validation

- Too small data set

→ Methods, not results