

THE SIGNIFICANCE LEVEL AS A LONG-RUN ERROR RATE



Significance level



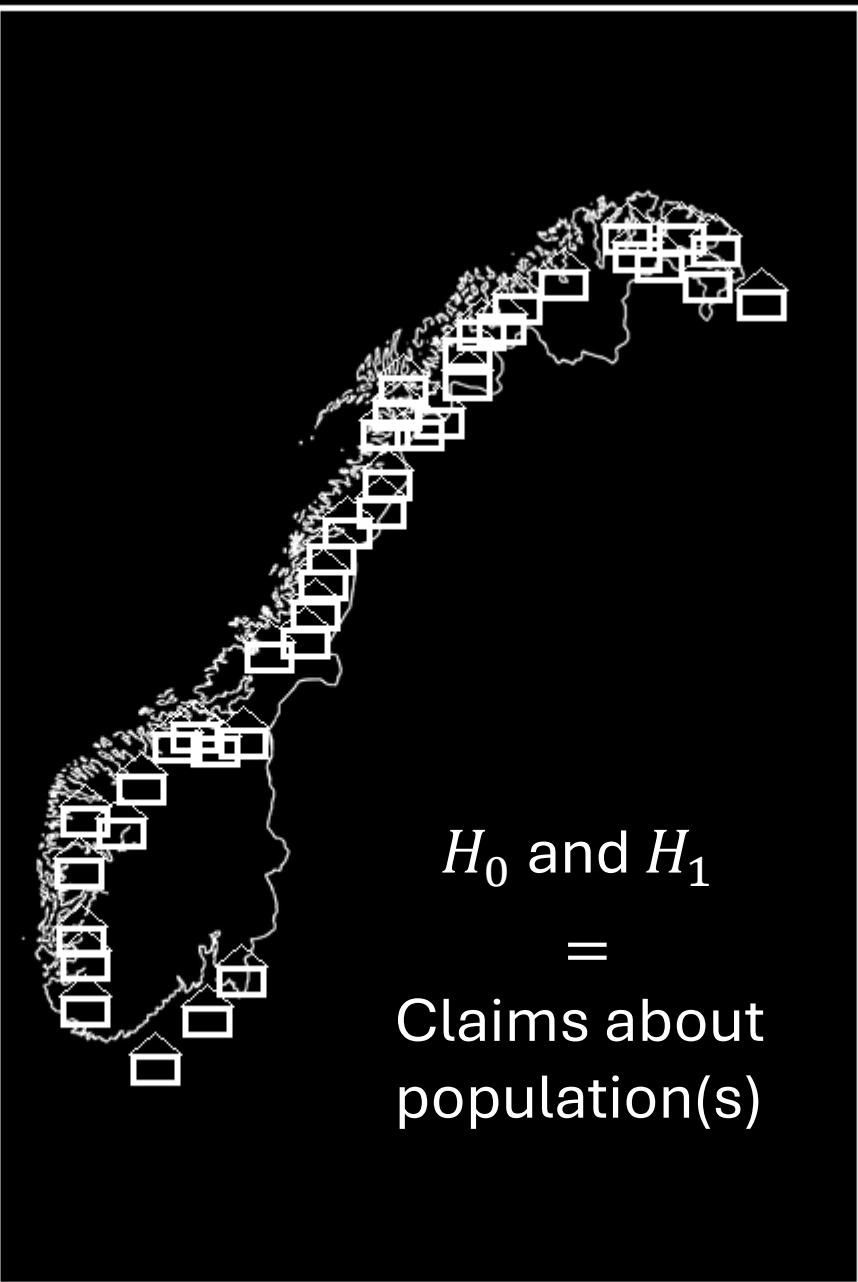
Reject H_0 if $P - value \leq \alpha$

When doing hypothesis testing..

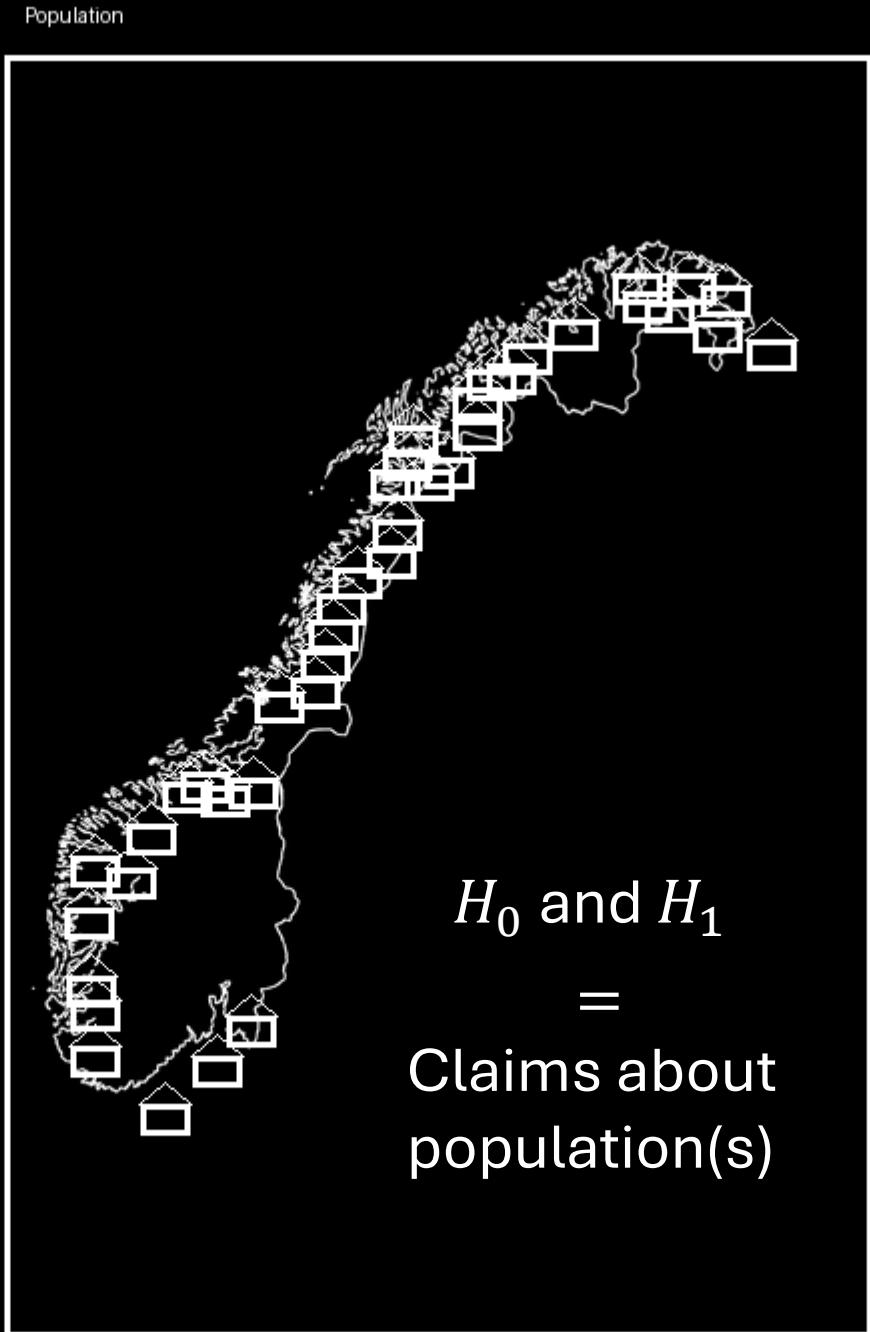
MISTAKES WILL SOMETIMES BE MADE



Population



Our conclusions about H_0 and H_1
are based on:



J. Neyman and Pearson 1933:

“...Without hoping to know whether each separate hypothesis is true or false, we may search for rules to govern our behaviour with regard to them, in following which we insure that, in the long run of experience, we shall not often be wrong”



States of reality:

A 2x2 matrix diagram illustrating hypothesis testing decisions. The columns represent the state of reality: H_0 is true or H_0 is false. The rows represent possible decisions: Reject H_0 or Retain H_0 . The cell where H_0 is true and the decision is to Reject H_0 is labeled "Type I Error". The cell where H_0 is false and the decision is to Retain H_0 is labeled "Type II Error". All other cells are labeled "Correct decision". A bracket on the left indicates the rows represent "Possible decisions", and a bracket at the top indicates the columns represent "States of reality".

		States of reality:	
		H_0 is true	H_0 is false
Possible decisions	Reject H_0	Type I Error	Correct decision
	Retain H_0	Correct decision	Type II Error

Decision rule:

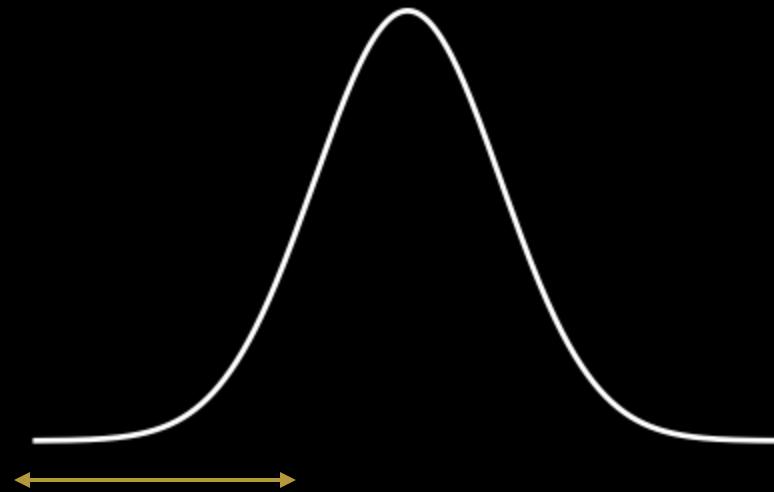
Reject H_0 if $P - value \leq \alpha$



Decision rule:

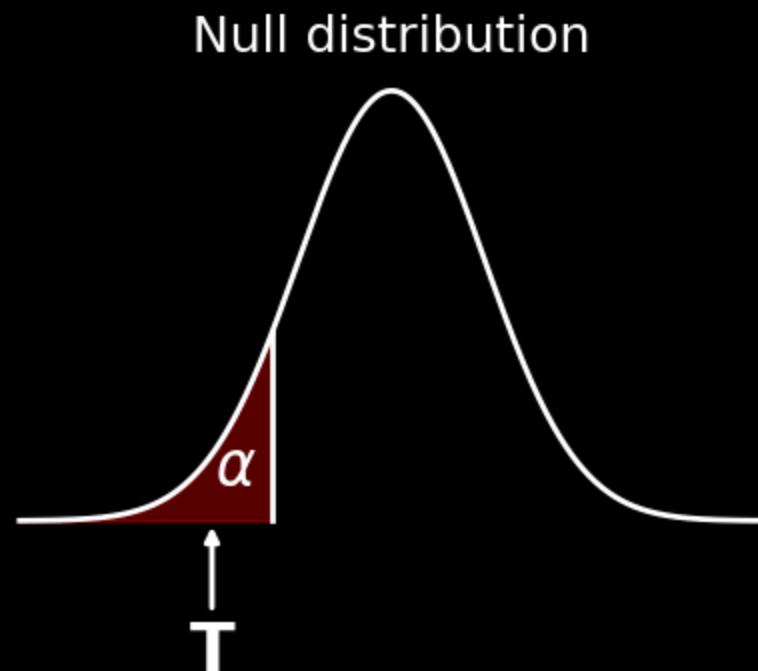
Reject H_0 if $P - value \leq \alpha$

Null distribution

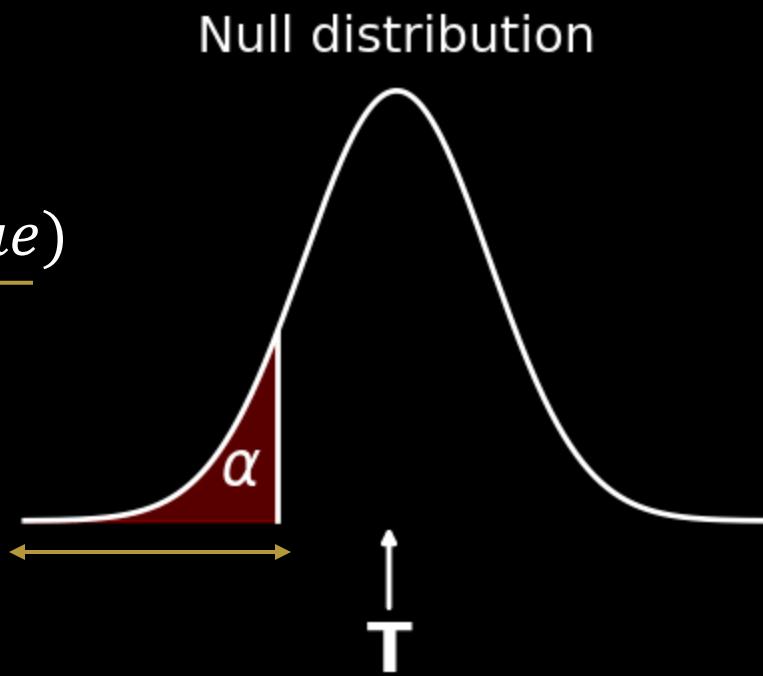


Decision rule:

Reject H_0 if $P - value \leq \alpha$



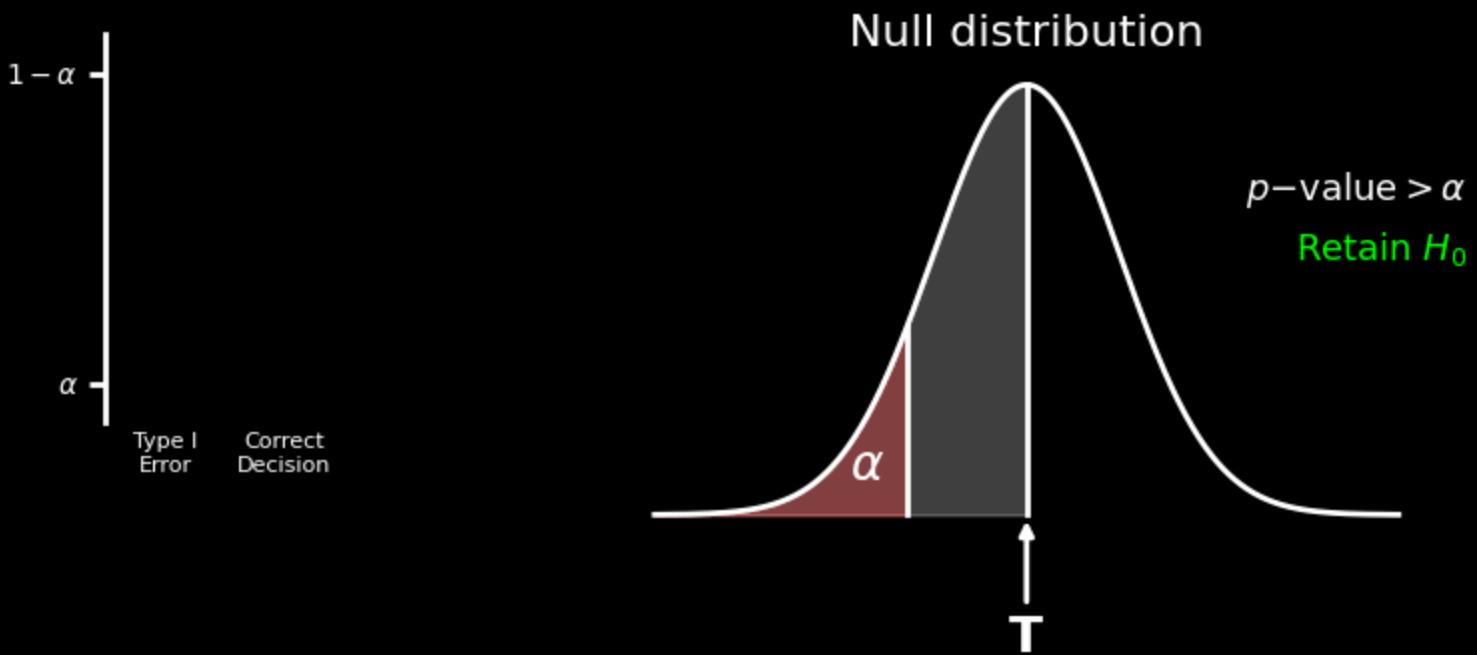
Type I Error = $P(\text{Reject } H_0 | H_0 \text{ true})$
Reject H_0 if $P\text{-value} \leq \alpha$



If H_0 is true:

1. Reject H_0 = Type I Error
2. T follows the Null distribution

$$\Rightarrow P(\text{Type I Error}) = \alpha$$



$$\Rightarrow P(\text{Type I Error}) = \alpha$$

Hypotheses:

H_0 : The batch of airbags is **defective**.

H_A : The batch of airbags is **safe**.

Action: If we reject H_0 the batch is approved for installation.

Type I error: Approving a batch of defective airbags → **Risk of injury or death**.

Very low α → Ensures defective airbags rarely pass quality checks.



Hypotheses:

H_0 : The batch of airbags is **defective**.

H_A : The batch of airbags is **safe**.

Common misconception:

- “A proportion α of all approved batches are defective.”

The significance level is a conditional error rate:

- α is the proportion of defective batches that we mistakenly approve in the long run.

1000 A/B tests with hypotheses:

$$H_0 : p_{new} = p_{current}$$

$$H_A : p_{new} > p_{current}$$

$1000 \times 0.05 = 50$ Significant results



«Multiple testing problem»

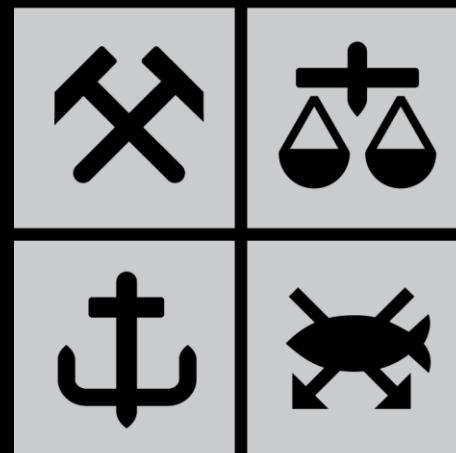
Why not set $\alpha = 0$, so that $P(\text{Type I Error}) = 0$?

	H_0 is true	H_0 is false
Reject H_0	Type I Error	Correct decision
Retain H_0	Correct decision	Type II Error

Neyman–Pearson framework:

- *We cannot avoid mistakes.*
- *We can aim to control how often we make them in the long run.*

NHH TECH3



Sondre Hølleland
Geir Drage Berentsen