

# Block 1.4: History & Methods of Psychology

Research Methods & Practices of Psychology

**Oliver Lindemann**

Erasmus University Rotterdam

Research Methods & Practices of Psychology



## Principles of Statistical Hypothesis Testing

Research Methods & Practices of Psychology



### Procedure of Empirical Research in Psychology

Theory → Hypothesis → Experiment → Data → Statistics

Research Methods & Practices of Psychology



### Procedure of Empirical Research in Psychology

Theory → **Hypothesis** → Experiment → **Data** → Statistics

Research Methods & Practices of Psychology



### Classical Statistics



Karl Pearson (1857 – 1936)

Ronald Fisher (1890 – 1962)

Jerzy Neyman (1894 – 1981)

Research Methods & Practices of Psychology



### Classical Statistics

Ronald Fischer



- terms “null-hypothesis” & “significant”
- urged the distinction between sample and population
- degrees of freedom
- suggested  $p < .05$
- random assignment of conditions, random sampling

Neyman and Pearson

- formal decision logic of statistics.
- Power and Type II error
- Effect sizes

Research Methods & Practices of Psychology



### Empirical Behavioural Research

Procedure

Theory → **Hypothesis** → Experiment → **Data** → Statistics

Neyman & Pearson suggested decision rule

- following this rule, in the long run, we will not be often wrong
- error rate ( $\alpha$ ) of the decision process
- e.g.  $p < .05$

Research Methods & Practices of Psychology



### Simple question or not?

What is a  $p$ -value?

Research Methods & Practices of Psychology



## Interpreting $p$

What is a  $p$ -value?

$$p(\text{Data}|H_0)$$

And what do most people want to know from the data?

$$p(H_1|\text{Data})$$

But

$$p(H_1|\text{Data}) \neq p(\text{Data}|H_0)$$

Thus,  $p$  tells us **nothing** about the likelihood of the hypothesis, neither  $H_1$  nor  $H_0$ !

## $p$ -values and strength of evidence

Neyman-Pearson approach:

$p$ -values interpretable as binary decision rule! (effect or not)

- Why can't we use  $p$ -values as measure of evidence?
- Why is a smaller  $p$ -value not more evidence for  $H_1$ ?
- $p$ -values are not consistent measures of evidence
  - It is relative to sample size
  - It is affected by sampling plan and other subjective elements

## Hypothesis Testing with $p$ -values in Practice

## How good is the 5% decision rule?

In psychology, we commonly use for the statistics:

- $\alpha = 0.05$
- $(1 - \beta) = 0.80$  (power of 80%)

If we strictly follow the rules above ...

**How many published research findings are then false?**

## How many research findings are false?

(A)

1000 hypotheses

## How many research findings are false?

(B)

10%  
true

## How many research findings are false?

(C)

10%  
true

5% false-  
positives

## How many research findings are false?

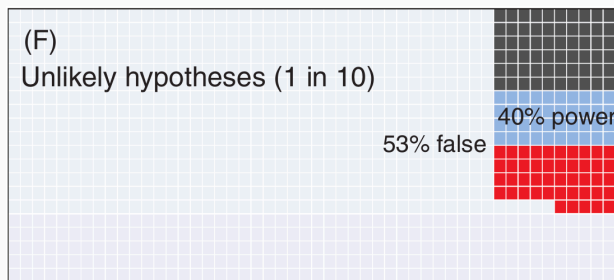
(D)

80%  
power

36% false

## How many research findings are false?

19



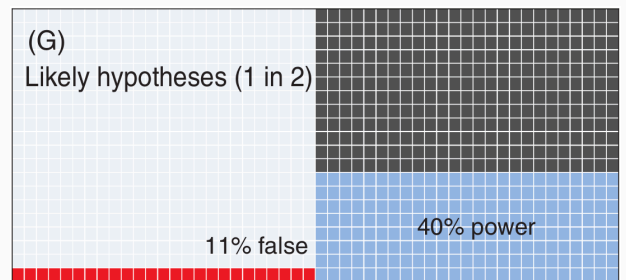
Research Methods & Practices of Psychology

from Forstmeier et al., 2017



## How many research findings are false?

20



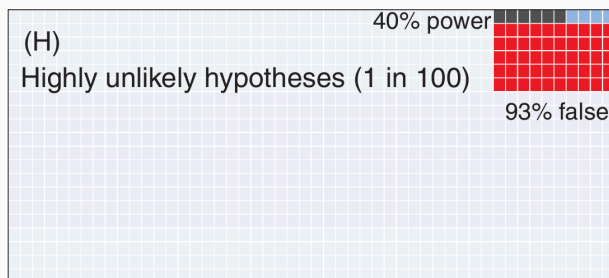
Research Methods & Practices of Psychology

from Forstmeier et al., 2017



## How many research findings are false?

21



Research Methods & Practices of Psychology

from Forstmeier et al., 2017



## Scientific Publication Practice

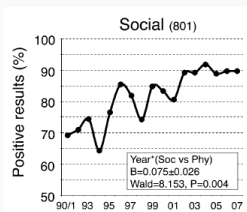
Research Methods & Practices of Psychology



## How many Published Effects are Significant?

24

Study by Fanelli (2012)



- Negative results are disappearing from the literature
- this happens from most disciplines and countries

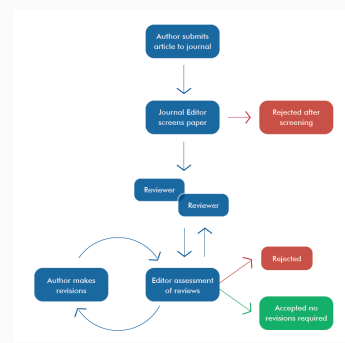
Why are most published effects significant?

Research Methods & Practices of Psychology



## Publishing Research in Journals: The Review Process

25



Research Methods & Practices of Psychology



## Which Studies and Results will be Published?

26

- Journals want to publish the most **exciting** and surprising findings
- Unfortunately, the review process and selection of the editor is affected by many non-scientific factors

### Publication bias

- occurs when the publication of research results depends **not just on the quality** of the research but also on the **hypothesis tested**, and the **significance** and direction of effects detected.
- is usually a bias towards reporting significant results

Research Methods & Practices of Psychology



## “File Drawer Problem”

- Studies showing significant effects and supporting the hypotheses will be published.
- Studies with no effect and no support for the hypothesis end up in the researchers’ file drawer.



As a result, the amount of significant results in most studies is overestimated.

## Psychological Reasons for Publication Biases

Why tend researchers to focus on significant effects so strongly?

- Confirmation bias
- Problem of Incentives in Science
- Researchers evaluation and career depend mainly on the amount of publications
- This is a general problem that ultimately affects the quality of research.

*“There is no cost in to getting things wrong,  
the cost is not getting them published”*

Brain Nosek

## Replicability

## Skepticism and Replications

Demarcation criterion between science and non-science

### Replicability

- the amount of consistency in results when scientific studies are repeated
- a basic element of critical scrutiny of claims
- an engine to the advancement of **self-correcting** science

### Advantages

- confirms scientific findings
- specifies the conditions under which the effect is registered
- more accurate estimates of the strength of the effect (Brandt et al, 2013)

## What does it look like in real science?

Meta-science study by Mackel, Pluncker & Hegarty (2012)

- Analysis of ALL articles in top 10 psychological journals from 1900
- The term “*replication*” occurred how many articles?

1.6%

- Analysis of 500 randomly chosen articles from this 1.6%:
  - 68% of articles using the term replication are designed to replicate

## Replication Rate in Psychology?

## Open Science Collaboration 2015

- 100 direct replications of experimental and correlational studies
- Direct replication = recreate conditions that are thought to suffice to obtain original effect
- Close to original studies (consultation of authors, use of original materials and internal review)
- Studies were matched with interests and expertise of replication team

## Replication project 2015: Method Details

- Quasi-random sample
  - 2008 articles from 3 journals:  
Journal of Personality and Social Psychology (JPSP), Psychological Science (PSCI),  
Journal of Experimental Psychology: Learning, Memory and Cognition (JEP:LMC)
  - From chosen articles, one study was selected
  - from this study only 1 statistical result was tested
- No standard exists to assess replication success
- Used
  - Significance and P-values
  - Effect sizes (transformed into r)
  - Subjective assessments of replication success
  - Meta analysis of effect sizes

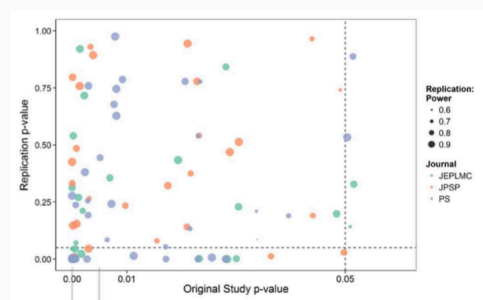
## Results: Significance and $p$ -Values

Replication effect tested against Null-hypothesis of no effect

- 97 studies originally significant
- Expected: 89 positive results

**Only 35 studies could be replicated**  
→ a replication rate of 36%

## Results: Significance and $p$ -Values



## Differences between Subdisciplines

Replication success rate for

1. Social psychology: 25%
2. Cognitive psychology: 50%

Possible explanation

- Weaker original effects for social psychology
- Higher power of test in cognitive psychology (e.g., within-subject designs)

## Reasons for Low Replication Rates?



## Fraud & Fabricating Data

*Example:* The case of Diederik Stapel (Tilburg University)

- fabricated data for at least 30 publications.
- young researchers as the whistleblowers
- suspended from his duties as Professor and returned his Ph.D.

Several other cases are reported ...

**But fraud, as a general problem in society, is a very isolated problem.**

**Fraud does not explain the low replication rate**

## Questionable Research Practices (QRP)

### “p-hacking”

Analyzing your data multiple ways and selectively reporting only those that result in  $p < .05$ .

1. outcome-dependent analysis

- researchers degrees of freedom (see next slides)

- special case: *Optional stopping*

- Peek into the data frequently and stop analysing if result is significant
- Collecting more data until results become significant

## Researchers Degrees of Freedom while Data Analysis

Researchers flexibility in

- Selecting dependent variables
- Selecting the participants
- Choosing covariates
- Analysis only subsets
- Exclude outliers selectively
  - Choosing to conduct analyses with different outlier criteria

Demo: [Hack Your Way To Scientific Glory](#)

## Further Types of $p$ -Hacking

### 2. Selective reporting

- Selectively reporting treatment groups and covariates
- Reporting only significant variables
- only reporting studies that show an effect (File drawer problem, problem 7)

### 3. **HARKing**: Hypothesizing After Results are Known

John Oliver on P-Hacking ([YouTube](#), 1:44–7:55)

## What shall we do?

## Guidelines for Researchers (... and students writing theses)

- Be clear: Exploratory or confirmatory analysis
- Confirmatory research → specify hypothesis **in advance**
- Report data collection practices
- Determine sample size in advance
  - Include at least 20 participants
- List all variables, experimental conditions and covariates
- Specify analysis procedure beforehand

→ **Study pre-registration!**

## Preregistration

- Prevents most types of p-hacking
  - e.g., outcome switching, garden of forking paths, adaptive outlier dropping, exclusion of conditions
- Clear distinction between confirmatory & exploratory research
  - No HARKing
  - Make p-values
- Minimizes publication bias
  - Even if pre-registered studies are not published at the end, the registry can be searched

How to do it? → use websites such as [aspredicted.org](#) or [osf.io](#)

## Questions?

Thank you very much