# Hypothesis testing

# NHST concepts

- Tell me...

- Tell me...
- **Null hypothesis**: there is no difference between groups.

- Tell me...
- **Null hypothesis**: there is no difference between groups.
- **Alternative hypothesis**: groups are different.

# Are there any differences? A non-sensical question in ecology

*Alejandro Martínez-Abraín*

*IMEDEA (CSIC-UIB), C/Miquel Marquès 21, 07190 Esporles, Majorca, Spain*

### ARTICLE INFO

### ABSTRACT

One of the main questions that ecologists pose in their investigations includes the analysis of differences in some trait between two or more populations. I argue here that asking whether there are differences or not between populations is biologically irrelevant, since no two livings things are ever equal. On the contrary the appropriate question to pose is how large differences are between populations. That is, we urge a shift in interest from statistical significance to biological relevance for proper knowledge accumulation. I empha-

- The probability that the data were produced by random chance alone

https://pollev.com/franciscorod726

- The probability that the data were produced by random chance alone
- The probability of getting results at least as extreme as the ones you observed if H0 was true

https://pollev.com/franciscorod726

- The probability that the data were produced by random chance alone
- The probability of getting results at least as extreme as the ones you observed if H0 was true
- The probability of null hypothesis being true

https://pollev.com/franciscorod726

- Very complicated concept: even statisticians fail to describe it well.

- Very complicated concept: even statisticians fail to describe it well.
- Probability of observing data as or more extreme than these *if H0 was true*.

## P-value

- Very complicated concept: even statisticians fail to describe it well.
- Probability of observing data as or more extreme than these *if H0 was true*.
- Low P-value: data unlikely if H0 was true.

## P-value

- Very complicated concept: even statisticians fail to describe it well.
- Probability of observing data as or more extreme than these *if H0 was true*.
- Low P-value: data unlikely if H0 was true.
- Large P-value: data not unusual if H0 was true.

· the null hypothesis is false, i.e. the alternative hypothesis must be true

https://pollev.com/franciscorod726

- the null hypothesis is false, i.e. the alternative hypothesis must be true
- the alternative hypothesis is false, i.e. the null hypothesis must be true

https://pollev.com/franciscorod726

- the null hypothesis is false, i.e. the alternative hypothesis must be true
- the alternative hypothesis is false, i.e. the null hypothesis must be true
- it's unclear if there are differences between groups

https://pollev.com/franciscorod726

- the null hypothesis is false, i.e. the alternative hypothesis must be true
- the alternative hypothesis is false, i.e. the null hypothesis must be true
- it's unclear if there are differences between groups
- there is no difference between groups

https://pollev.com/franciscorod726

## Are differences *significant*?

- If $p < 0.05$, we **reject** H0.

## Are differences *significant*?

- If $p < 0.05$, we **reject** H0.
- If $p > 0.05$, we **fail to reject** H0

## Are differences *significant*?

- If $p < 0.05$, we **reject** H0.
- If $p > 0.05$, we **fail to reject** H0
- (which is **NOT** the same as 'H0 is true')

## Are differences *significant*?

- If p < 0.05, we **reject** H0.
- If p > 0.05, we **fail to reject** H0
- (which is **NOT** the same as 'H0 is true')
- CAUTION:

## Are differences *significant*?

- If p < 0.05, we **reject** H0.
- If p > 0.05, we **fail to reject** H0
- (which is **NOT** the same as 'H0 is true')
- **CAUTION:**
- P-value is continuous. We must **avoid binary decisions** based on **arbitrary thresholds**.

## Are differences *significant*?

- If $p < 0.05$, we **reject** H0.
- If $p > 0.05$, we **fail to reject** H0
- (which is **NOT** the same as 'H0 is true')
- **CAUTION:**
- P-value is continuous. We must **avoid binary decisions** based on **arbitrary thresholds**.

- If p < 0.05, we **reject** H0.

- If p > 0.05, we **fail to reject** H0

- (which is **NOT** the same as 'H0 is true')

- CAUTION:

- P-value is continuous. We must **avoid binary decisions** based on **arbitrary thresholds**.



https://doi.org/10.1038/d41586-019-00857-9

# Are the heights of local and non-local students different?

```
t.test(h.sevi, h.out)
```

```
        Welch Two Sample t-test

data:  h.sevi and h.out
t = -3.159, df = 11.768, p-value = 0.00842
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -21.986075  -4.013925
sample estimates:
mean of x mean of y
    165.8     178.8
```

https://pollev.com/franciscorod726

| Statistics: Hypothesis Test | Null Hypothesis is True | Null Hypothesis is False |
|---|---|---|
| Reject Null Hypothesis | Type I Error | Correct |
| Fail to Reject Null Hypothesis | Correct | Type II Error |

**Power**: Probability of detecting true difference (rejecting H0 when it's false).

http://rpsychologist.com/d3/NHST/

```
 [1] 1 0 1 0 1 0 0 1 1 0


    1-sample proportions test without continuity correction

data:  sum(coin) out of ntrials, null probability 0.5
X-squared = 0, df = 1, p-value = 1
alternative hypothesis: true p is not equal to 0.5
95 percent confidence interval:
 0.2365931 0.7634069
sample estimates:
  p
0.5
```

https://pollev.com/franciscorod726

http://rpsychologist.com/d3/correlation/

# Common pitfalls and good practice

CrossMark

ESSAY

**Statistical tests, *P* values, confidence intervals, and power: a guide to misinterpretations**

Sander Greenland[1] · Stephen J. Senn[2] · Kenneth J. Rothman[3] · John B. Carlin[4] ·
Charles Poole[5] · Steven N. Goodman[6] · Douglas G. Altman[7]

https://doi.org/10.1007/s10654-016-0149-3

**esa**

ECOSPHERE

Applied statistics in ecology:
common pitfalls and simple solutions

E. Ashley Steel,[1],† Maureen C. Kennedy,[2] Patrick G. Cunningham,[3] and John S. Stanovick[4]

https://doi.org/10.1890/ES13-00160.1

Also http://www.statisticsdonewrong.com/

Twenty tips for
interpreting
scientific claims

https://doi.org/10.1038/503335a

Visualisation of data and models
is key

- Always

- Always
- Always

- Always
- Always
- Always

https://janhove.github.io/teaching/2016/11/21/what-

All correlations: r(50) = 0

https://janhove.github.io/teaching/2016/11/21/what-
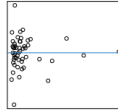
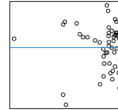All correlations: r(50) = 0.1

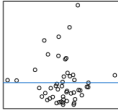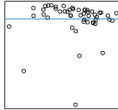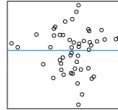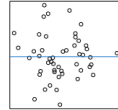(1) Normal x, normal residuals (2) Uniform x, normal residuals (3) +-skewed x, normal residuals (4) --skewed x, normal residuals (5) Normal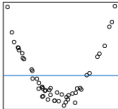 x, +-skewed residuals (6) Nor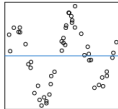mal x, --skewed residuals (7) 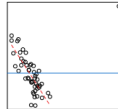Increasing sp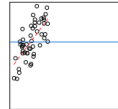read (8) Decreasing spread (9) 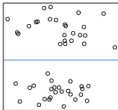Quadratic trend (10) Sinu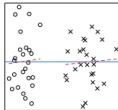soid relationship (11) A sin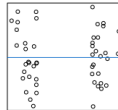gle positive outlier (12) A single negative outlier (13) Bimodal residuals (14) Two groups (15) Sampling at the extremes (16) Coarse data

https://janhove.github.io/teaching/2016/11/21/what-

*Plot. Check models. Plot. Check assumptions. Plot.*

Lavine 2014 *Ecology*

# Inference from observational studies

- In a sample of 10,000 people, it was found that people eating >2 hamburgers a week had 20% higher probability of heart attack.

- In a sample of 10,000 people, it was found that people eating >2 hamburgers a week had 20% higher probability of heart attack.
- Do hamburgers increase heart attacks?

- In a sample of 10,000 people, it was found that people eating >2 hamburgers a week had 20% higher probability of heart attack.
- Do hamburgers increase heart attacks?
- https://pollev.com/franciscorod726

- We found that plants with big flowers produced 30% more seeds...

- We found that plants with big flowers produced 30% more seeds...
- Do big flowers increase reproductive success?

- We found that plants with big flowers produced 30% more seeds...
- Do big flowers increase reproductive success?
- https://pollev.com/franciscorod726

http://tylervigen.com/spurious-correlations

# NHST and p-values

# Are there any differences? A non-sensical question in ecology

*Alejandro Martínez-Abraín*

*IMEDEA (CSIC-UIB), C/Miquel Marquès 21, 07190 Esporles, Majorca, Spain*

ARTICLE INFO

ABSTRACT

One of the main questions that ecologists pose in their investigations includes the analysis of differences in some trait between two or more populations. I argue here that asking whether there are differences or not between populations is biologically irrelevant, since no two livings things are ever equal. On the contrary the appropriate question to pose is how large differences are between populations. That is, we urge a shift in interest from statistical significance to biological relevance for proper knowledge accumulation. I empha-

# Instead of falsifying null model, compare meaningful models

# P-value depends on sample size



Group A: 10 ± 5 (SD)
Group B: 12 ± 5 (SD)

Same real difference is detected as significant or not depending on sample size:



Real difference = 40 g

With big sample size, we can find **highly significant but biologically unimportant** differences.



Real difference = 1 g

N = 2000, p = 0.0008

- Statistically significant = unlikely to be zero

- Statistically significant = unlikely to be zero
- Good read: *significantly misleading*

- Statistically significant = unlikely to be zero
- Good read: *significantly misleading*
- My suggestion: avoid significant/not significant (and maybe p-values too)

- Statistically significant = unlikely to be zero
- Good read: *significantly misleading*
- My suggestion: avoid significant/not significant (and maybe p-values too)
- Beyond significance, look at *effect sizes.*

- Absence of evidence != Evidence of absence

P >> 0.05

P << 0.05

- "We were **unable to find evidence** against the hypothesis that A = B **with the current sample size**" (Harrell)

- "We were **unable to find evidence** against the hypothesis that A = B **with the current sample size**" (Harrell)
- "Differences between groups were **not statistically clear**" (Dushoff et al)

- Right turn not allowed: 308 accidents

https://www.statisticsdonewrong.com/power.html#the-wrong-turn-on-red

## Is it safe to allow right turn with red lights?

- Right turn not allowed: 308 accidents
- Right turn allowed: 337 accidents

https://www.statisticsdonewrong.com/power.html#the-wrong-turn-on-red

- Right turn not allowed: 308 accidents
- Right turn allowed: 337 accidents
- No *significant* difference, hence safe

https://www.statisticsdonewrong.com/power.html#the-wrong-turn-on-red

- Right turn not allowed: 308 accidents
- Right turn allowed: 337 accidents
- No *significant* difference, hence safe
- Misinterpretation of underpowered study cost lives

https://www.statisticsdonewrong.com/power.html#the-wrong-turn-on-red

**The Difference Between "Significant" and "Not Significant" is not Itself Statistically Significant**

Andrew GELMAN and Hal STERN

http://dx.doi.org/10.1198/000313006X152649

http://dx.doi.org/10.1002/prp2.93

1. Test multiple variables, then report the ones that are significant.

1. Test multiple variables, then report the ones that are significant.
2. Artificially choose when to end your experiment.

1. Test multiple variables, then report the ones that are significant.

2. Artificially choose when to end your experiment.

3. Add covariates until effects are significant.

1. Test multiple variables, then report the ones that are significant.
2. Artificially choose when to end your experiment.
3. Add covariates until effects are significant.
4. Test different conditions (e.g. different levels of a factor) and report the ones you like.

1. Test multiple variables, then report the ones that are significant.

2. Artificially choose when to end your experiment.

3. Add covariates until effects are significant.

4. Test different conditions (e.g. different levels of a factor) and report the ones you like.

- To read more: Simmons et al 2011

https://www.youtube.com/watch?v=ZaNtz76dNSI

- P-values do not measure the **probability of hypothesis** being true, or the probability that the data were produced by **random chance** alone.

https://doi.org/10.1080/00031305.2016.1154108

- P-values do not measure the **probability of hypothesis** being true, or the probability that the data were produced by **random chance** alone.
- Scientific conclusions or policy decisions should NOT be based only on **whether a p-value passes a specific threshold**.

https://doi.org/10.1080/00031305.2016.1154108

- P-values do not measure the **probability of hypothesis** being true, or the probability that the data were produced by **random chance** alone.
- Scientific conclusions or policy decisions should NOT be based only on **whether a p-value passes a specific threshold**.
- P-value, or statistical significance, does not measure the **size of an effect** or the **importance** of a result.

https://doi.org/10.1080/00031305.2016.1154108

- P-values do not measure the **probability of hypothesis** being true, or the probability that the data were produced by **random chance** alone.
- Scientific conclusions or policy decisions should NOT be based only on **whether a p-value passes a specific threshold**.
- P-value, or statistical significance, does not measure the **size of an effect** or the **importance** of a result.
- By itself, a p-value does NOT provide a good **measure of evidence** regarding a model or hypothesis.

https://doi.org/10.1080/00031305.2016.1154108

Aim for estimation of effects and their uncertainty (SE, CI...)



*General Article*

# The New Statistics: Why and How

**Geoff Cumming**
La Trobe University

http://dx.doi.org/10.1177/0956797613504966

- **Type I**: False positive (incorrect rejection of null hypothesis).

- **Type I**: False positive (incorrect rejection of null hypothesis).
- **Type II**: False negative (failure to reject false null hypothesis).

# How many types of errors?

- **Type I**: False positive (incorrect rejection of null hypothesis).
- **Type II**: False negative (failure to reject false null hypothesis).
- **Type S (Sign)**: estimating effect in opposite direction.

- **Type I**: False positive (incorrect rejection of null hypothesis).
- **Type II**: False negative (failure to reject false null hypothesis).
- **Type S (Sign)**: estimating effect in opposite direction.
- **Type M (Magnitude)**: Misestimating magnitude of the effect (under or overestimating).

- **Type I**: False positive (incorrect rejection of null hypothesis).
- **Type II**: False negative (failure to reject false null hypothesis).
- **Type S (Sign)**: estimating effect in opposite direction.
- **Type M (Magnitude)**: Misestimating magnitude of the effect (under or overestimating).
- Beyond Power Calculations: Assessing Type S (Sign) and Type M (Magnitude) Errors