# ECOL 596

Practical and Reproducible Data Science
for Ecology and Evolutionary Biology

# What is "data science?"

a.k.a. biostatistics, biometry, analysis

Input Layer

Multiple hidden Layers

Output Layer

Data science helps us extract the signal from the noise

# Why do EEB grad students need data science?

- Design good studies
- Interpret your data
- Test hypotheses
- Extrapolate/predict
- Have confidence in your results

# The "reproducibility crisis"

# Intentional and unintentional errors



SCIENCEINSIDER | PEOPLE & EVENTS

**Embattled spider biologist resigns university post**

Ecologist Jonathan Pruitt agrees to leave McMaster University, but reasons are still undefined

12 JUL 2022 · 5:30 PM · BY ELIZABETH PENNISI

Retraction Watch Database:
1,179 retractions "error in analyses"

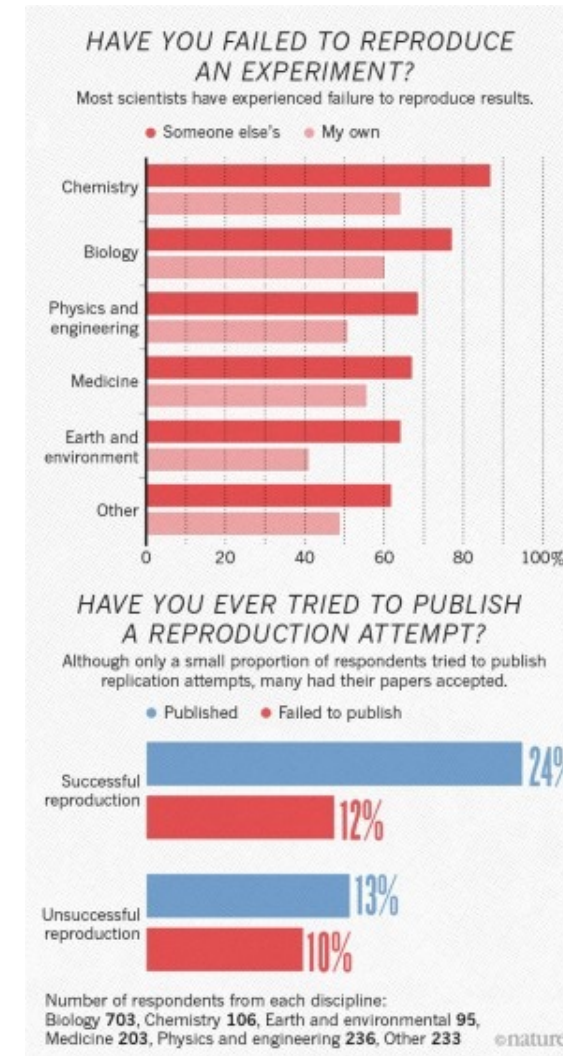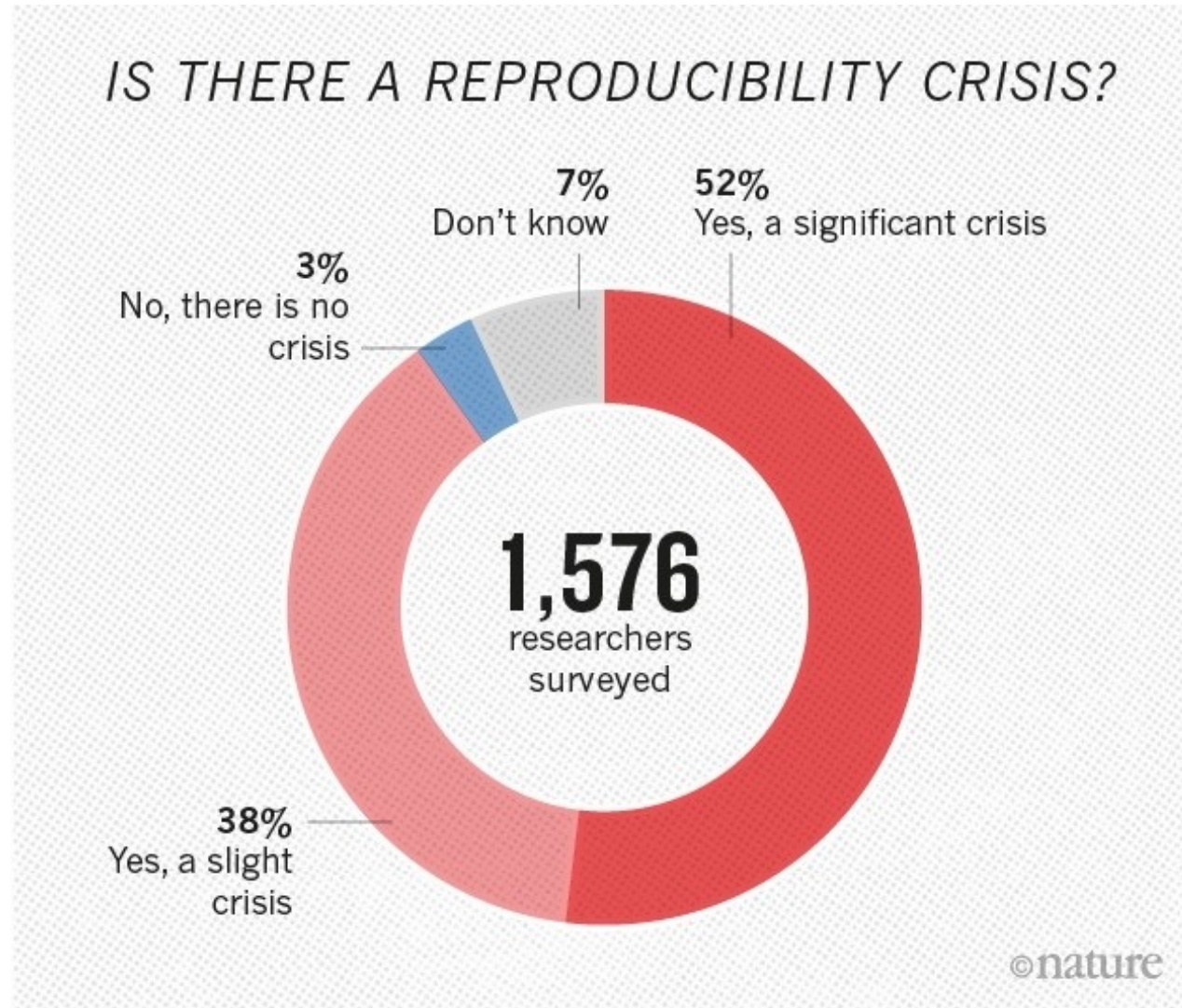The authors have retracted this article because the classification described in the paper was performed on an inaccurate use of a data pre-processing tool, which resulted in an information leak from train to test samples. This has resulted in incorrect classification metrics. However, the statistical differences between

After this article [1] was published, the authors identified data analysis errors that led to an overestimation of genomic differentiation among breeds. In light of this issue, the article's results and conclusions are not valid. Therefore, the authors retract this article.

traditional wintering grounds in the Mediterranean?' Due to an analytical error the study only analysed EURING circumstance 20 ring records (see the EURING exchange code) which excludes birds hunted by shot and other ways, but not all hunted individuals. As a result, the effect of illegal hunting on ring
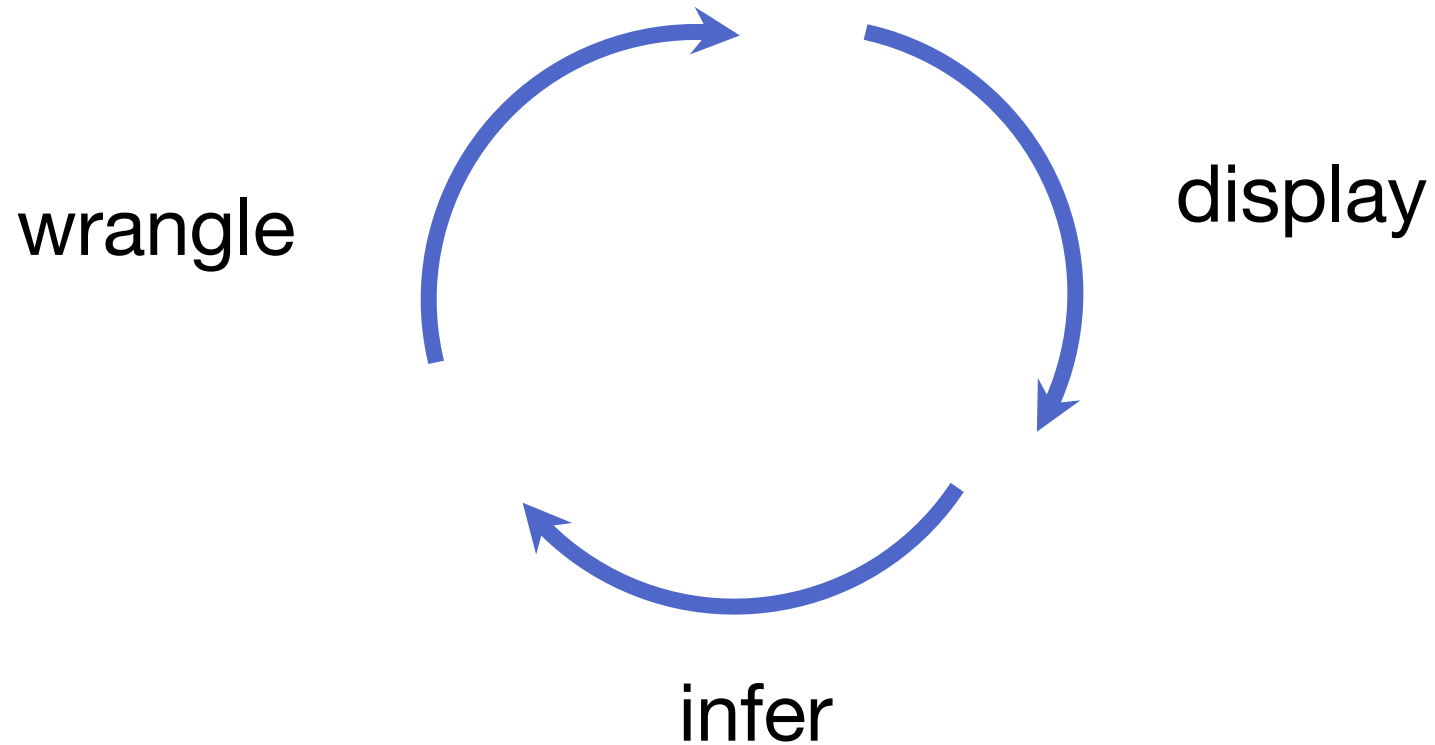
# Scope of this class

✓Good practices for data management

✓Coding in R

✓Common statistical tests

✓**Statistical intuition**

✓**How to teach yourself**

✓**Topics that will help you**

All of the statistical tests ever

Calculus

Other stats languages

# Guiding principle 1: The wheel of analysis

# Guiding principle 2:
# You won't get it until you learn it > 3 times


I GET IT NOW

teach it

reproduce it

follow it