

# Computational NeuroEthology

Statistical approaches: From P-Values to Pipelines



Albert Einstein College of Medicine

Class 2

October 6th, 2025

Mikhail Kislin

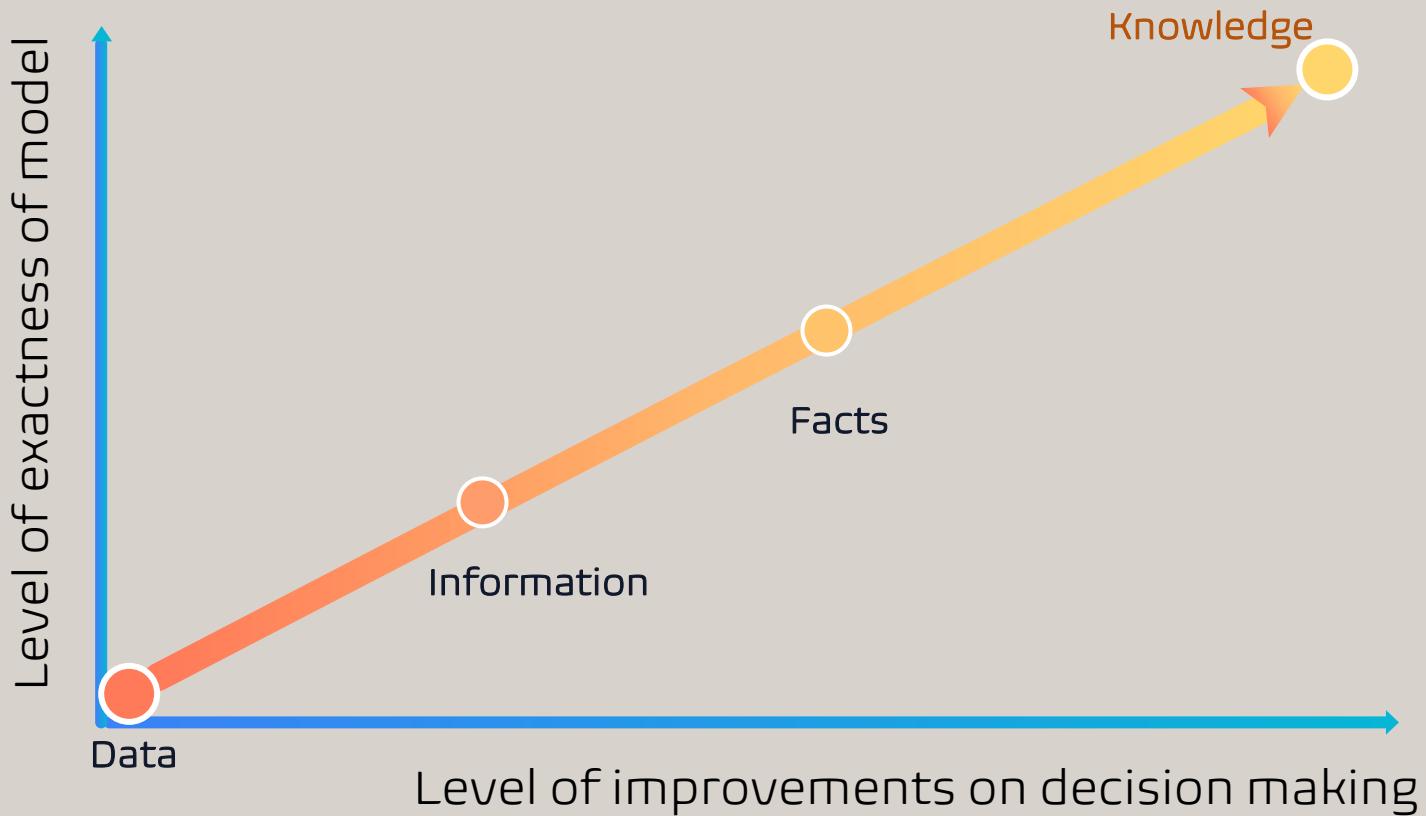
---

# Overview

- Statistical Modeling Frameworks
- Modern Statistical Workflows
- Best Practices and Advanced Topics

# 00

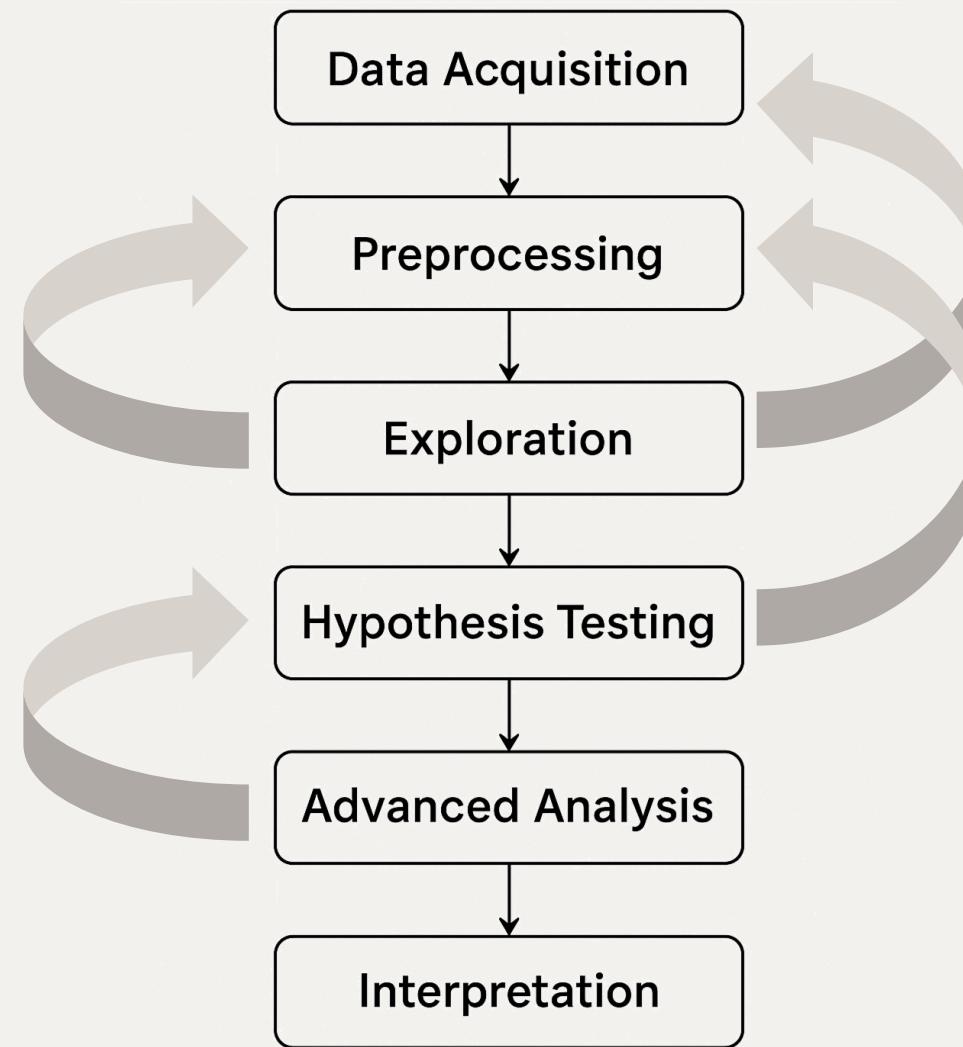
# What can statistics do for us?



1. Learning from data  
stats by aggregation vs ML (all data)
2. Deal with uncertainty and help to make decisions
3. Make inferences about causation
4. Make predictions about new situations based on priors

# The Data Analysis Workflow

micro-meso-macro concept  
of complex systems



# What measurements produce a good data?

A: Reliable and valid



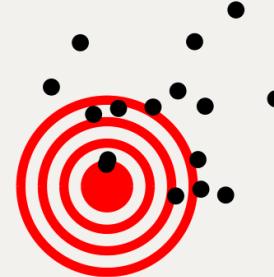
B: Unreliable but valid



C: Reliable but invalid



D: Unreliable and invalid



test-retest reliability

Face validity

Construct validity

Predictive validity

Computational validity: using computation to translate behaviors across species

# Basic probability theory

Probability quantifies uncertainty:

$$P(A) = \frac{\text{Number of favorable outcomes}}{\text{Total outcomes}}$$

Complement Rule:

$$P(\text{not } A) = 1 - P(A)$$

Addition Rule (mutually exclusive):

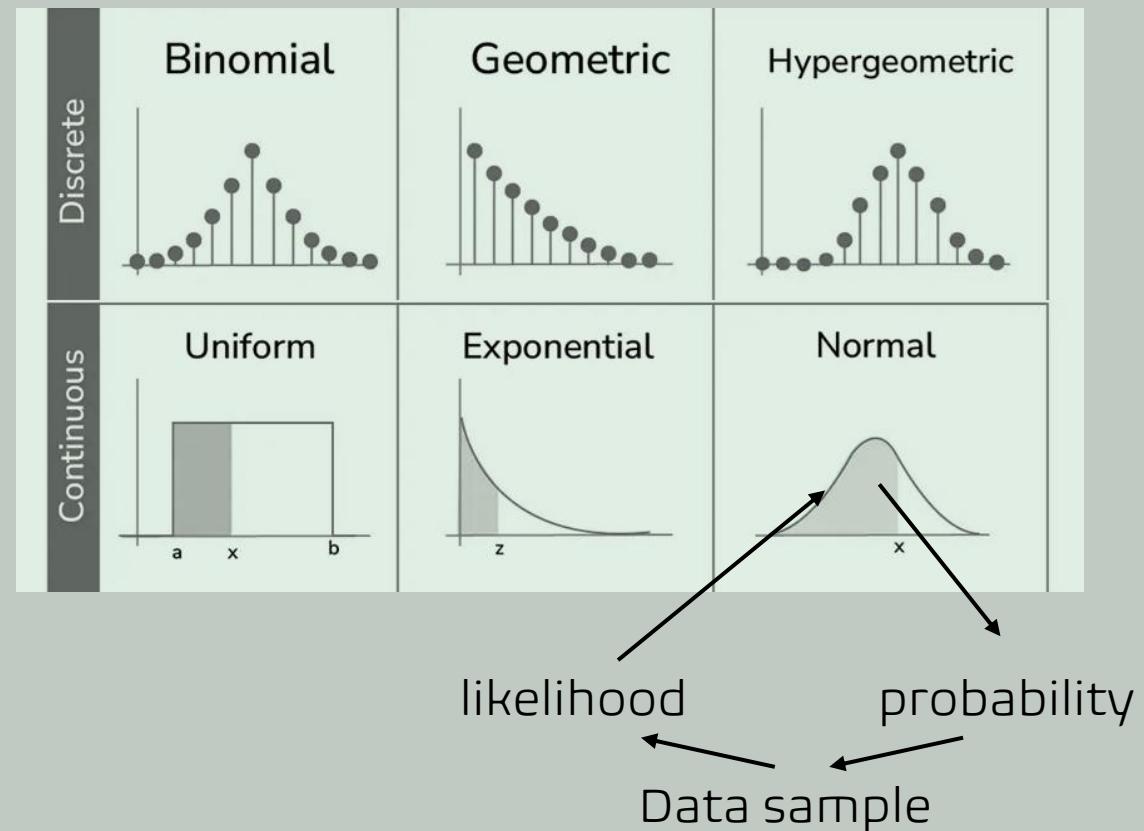
$$P(A \text{ or } B) = P(A) + P(B)$$

Multiplication Rule (independent events):

$$P(A \text{ and } B) = P(A) \times P(B)$$

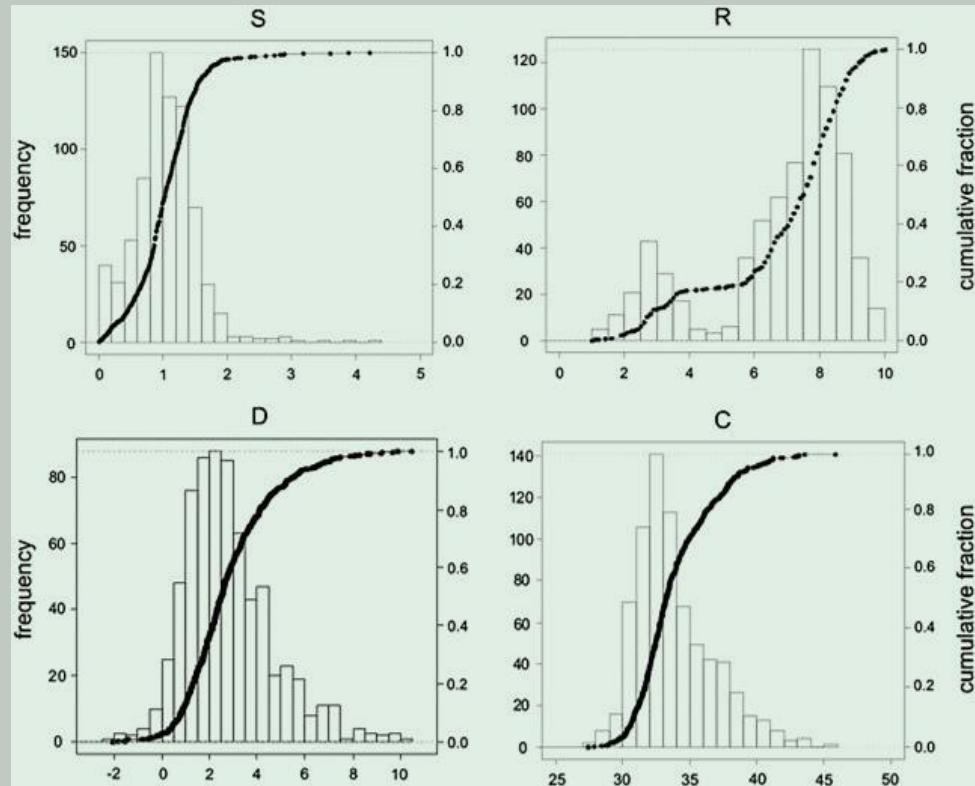
Conditional Probability:

$$P(A|B) = \frac{P(A \text{ and } B)}{P(B)}$$



# Frequentist Inference

Frequentist inference is the process of determining properties of an underlying distribution via the observation of data.



1. Summarize data

Expectation

$$E[X] = \sum_{x \in \mathcal{X}} x P(x)$$

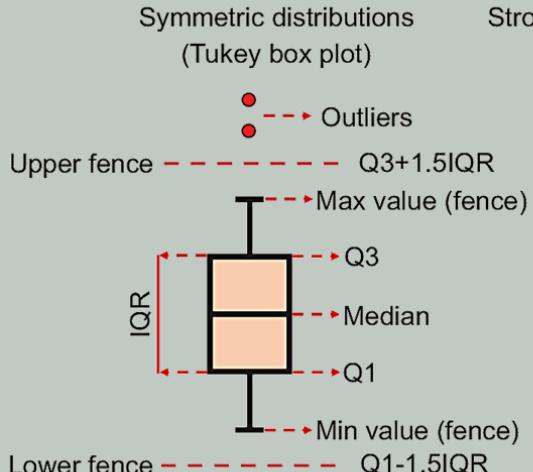
Variance

$$\text{Var}(X) = E[(X - E[X])^2]$$

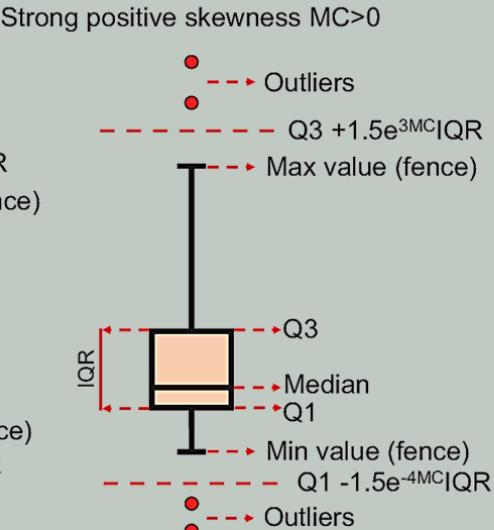
Idealized representations of observed data

# Components and construction of box plots

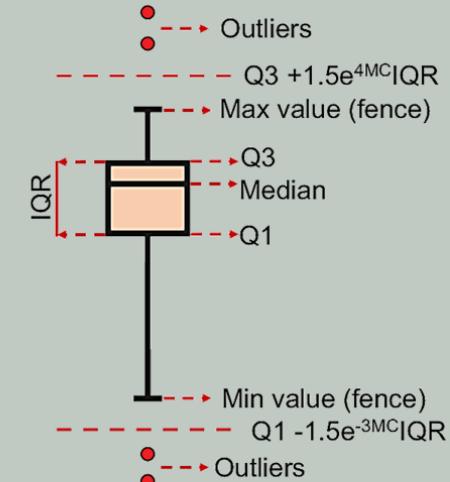
A/



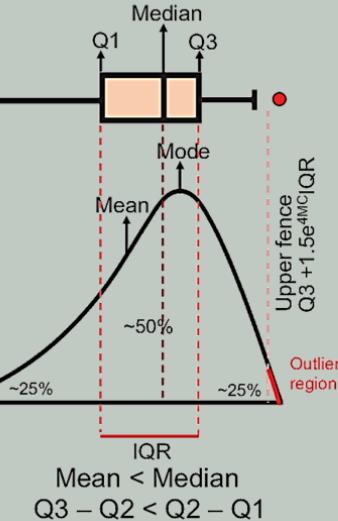
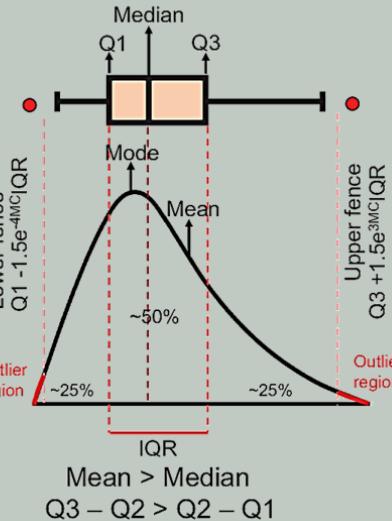
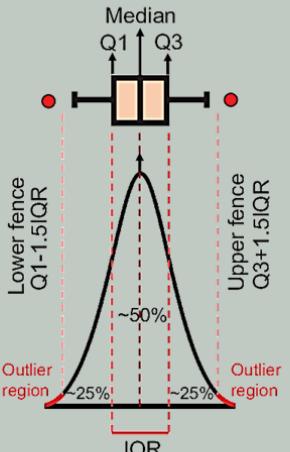
B/



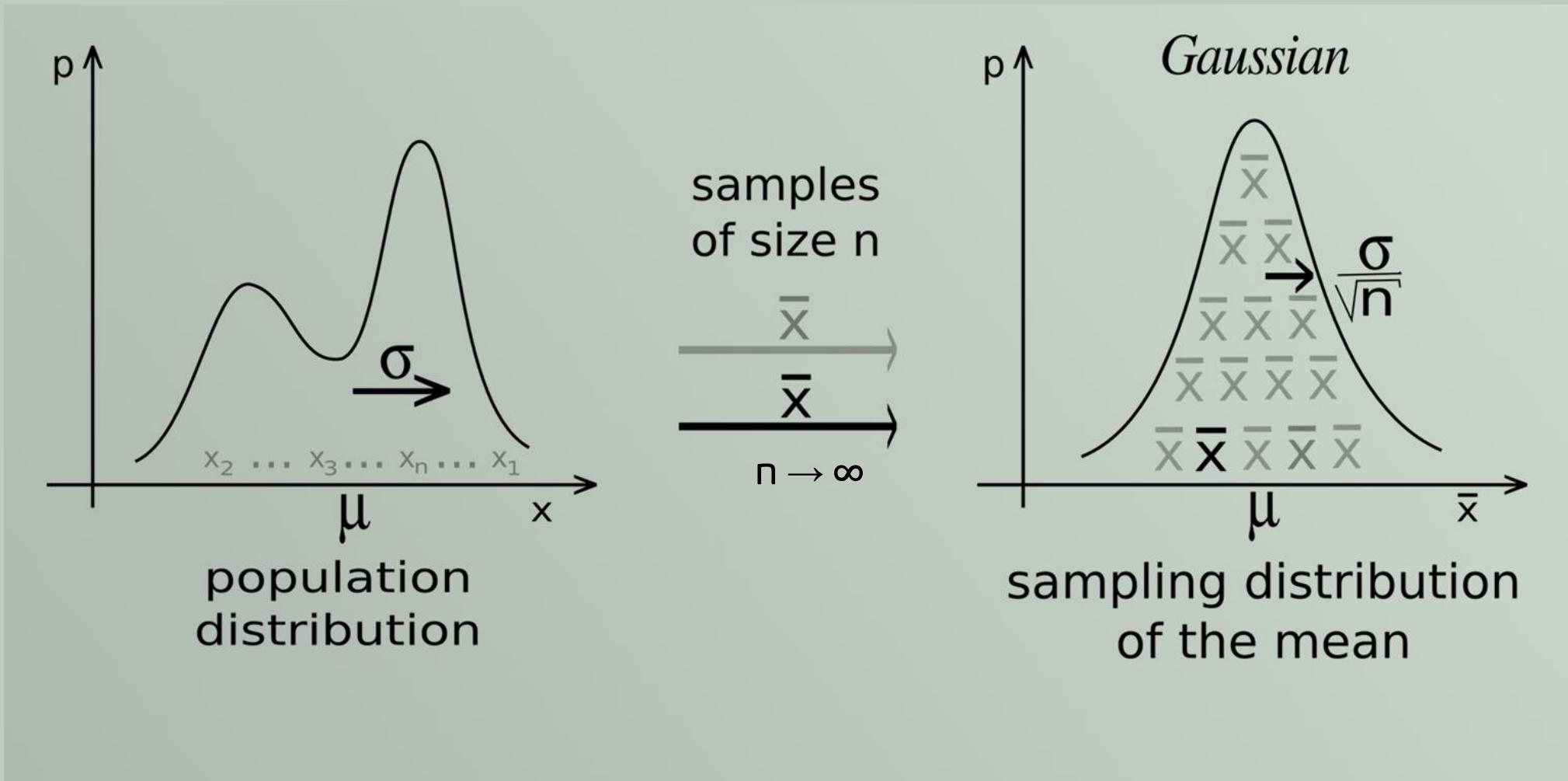
Strong negative skewness  $MC < 0$



C/

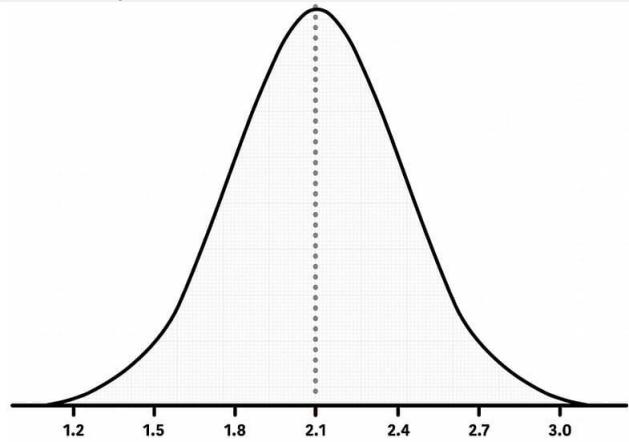


# Central Limit Theorem



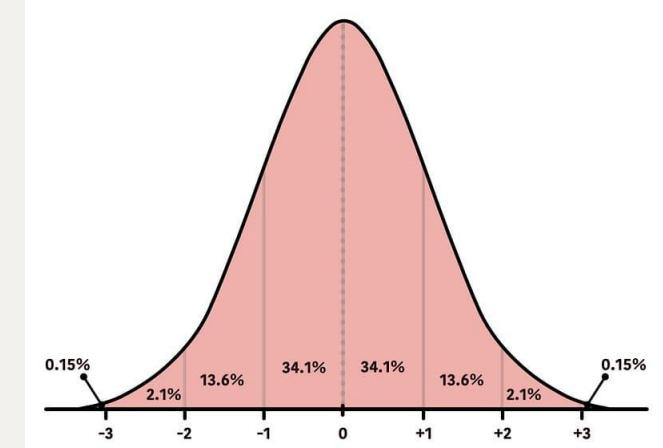
# Hypothesis testing

Population distribution

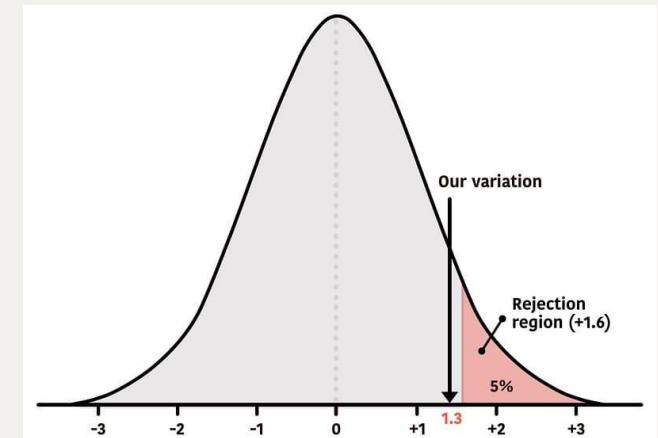
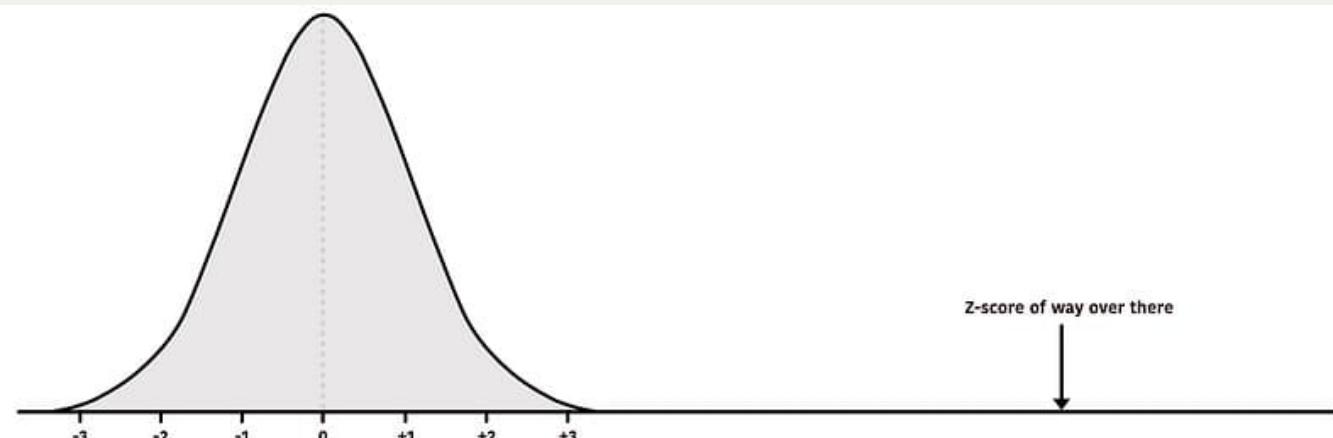


$$z = \frac{x - \mu}{\sigma}$$

The equation shows the formula for calculating a z-score. A red arrow points from the text "raw score" to the variable  $x$ . Another red arrow points from the text "mean" to the symbol  $\mu$ . A third red arrow points from the text "standard deviation" to the symbol  $\sigma$ .



How unusual your observed data is?

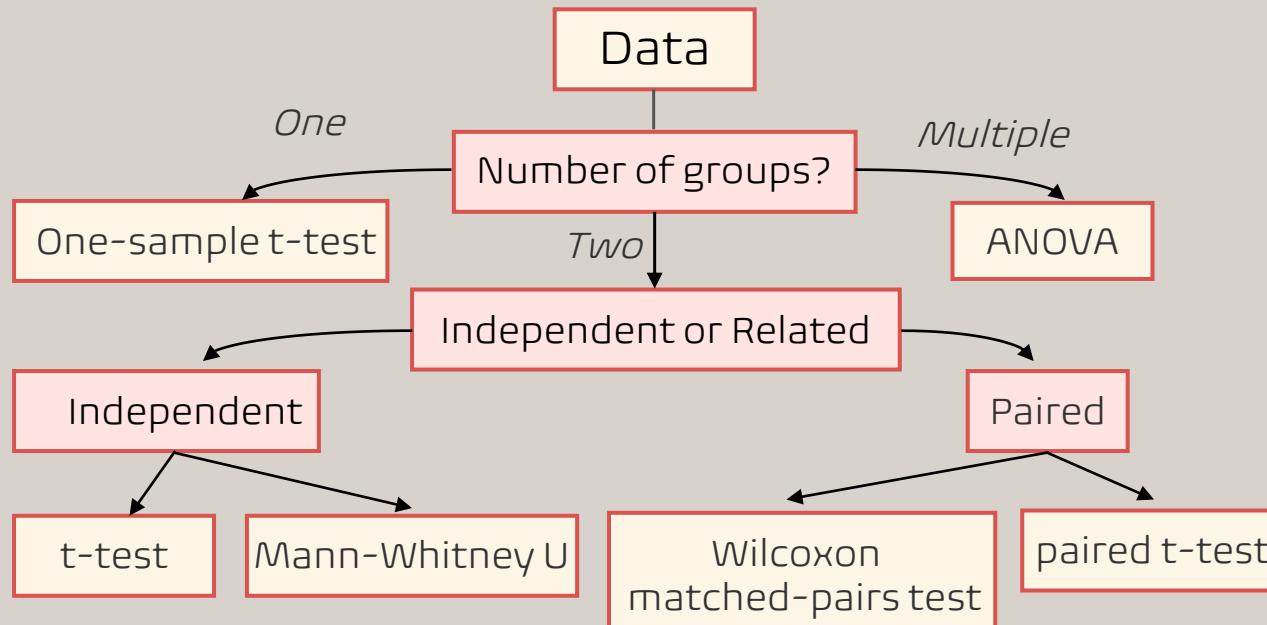


# What is the problem with the "Cookbook" approach to statistics?

## The "Cookbook" Approach

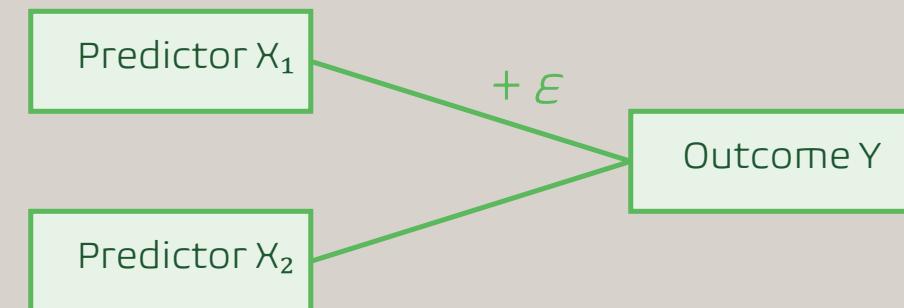
Focus:

*Which test should I run?*



## The Principled (Modeling) Approach

*What process generated my data?*

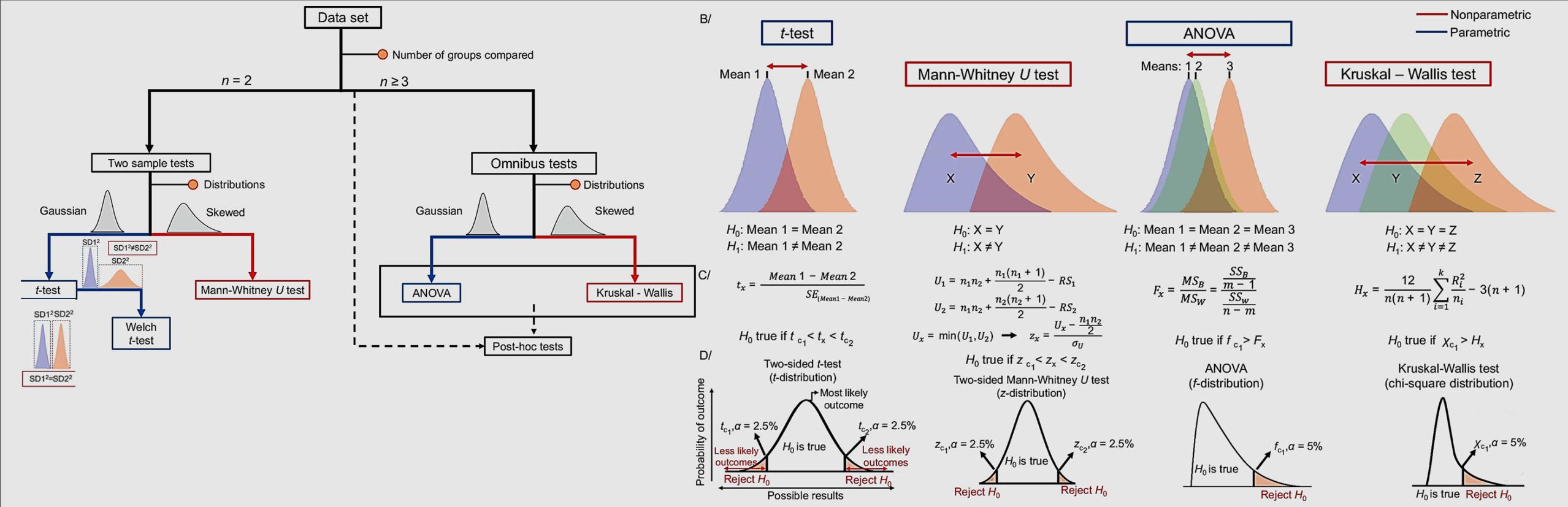


Consequence:

Shaky conclusions, risk of p-hacking, and a disconnect from the scientific question

Flexible insights, robust understanding of uncertainty, direct link to the scientific hypothesis.

# Commonly used statistical tests



Idkowiak, J., et al. Nat Commun 2025

## 02

# A Unified Framework: The General Linear Model

The GLM explains variation in Y as a linear combination of predictors

Vector form

$$y_i = \theta^\top x_i + \eta \quad \forall i = 1, \dots, N$$

neural response  
 $\downarrow$   
 $y_i$   
 $\nearrow$   
 $\theta = [\theta_0, \theta_1, \theta_2, \dots]^\top$   
 linear weights  
 $x_i = [1, x_{i,1}, x_{i,2}, \dots]^\top$   
 multiple stimulus features  
 (e.g., orientation, contrast,  
 etc.)  
 Gaussian noise  
 $\eta \sim \mathcal{N}(0, \sigma^2)$   
 number of data points

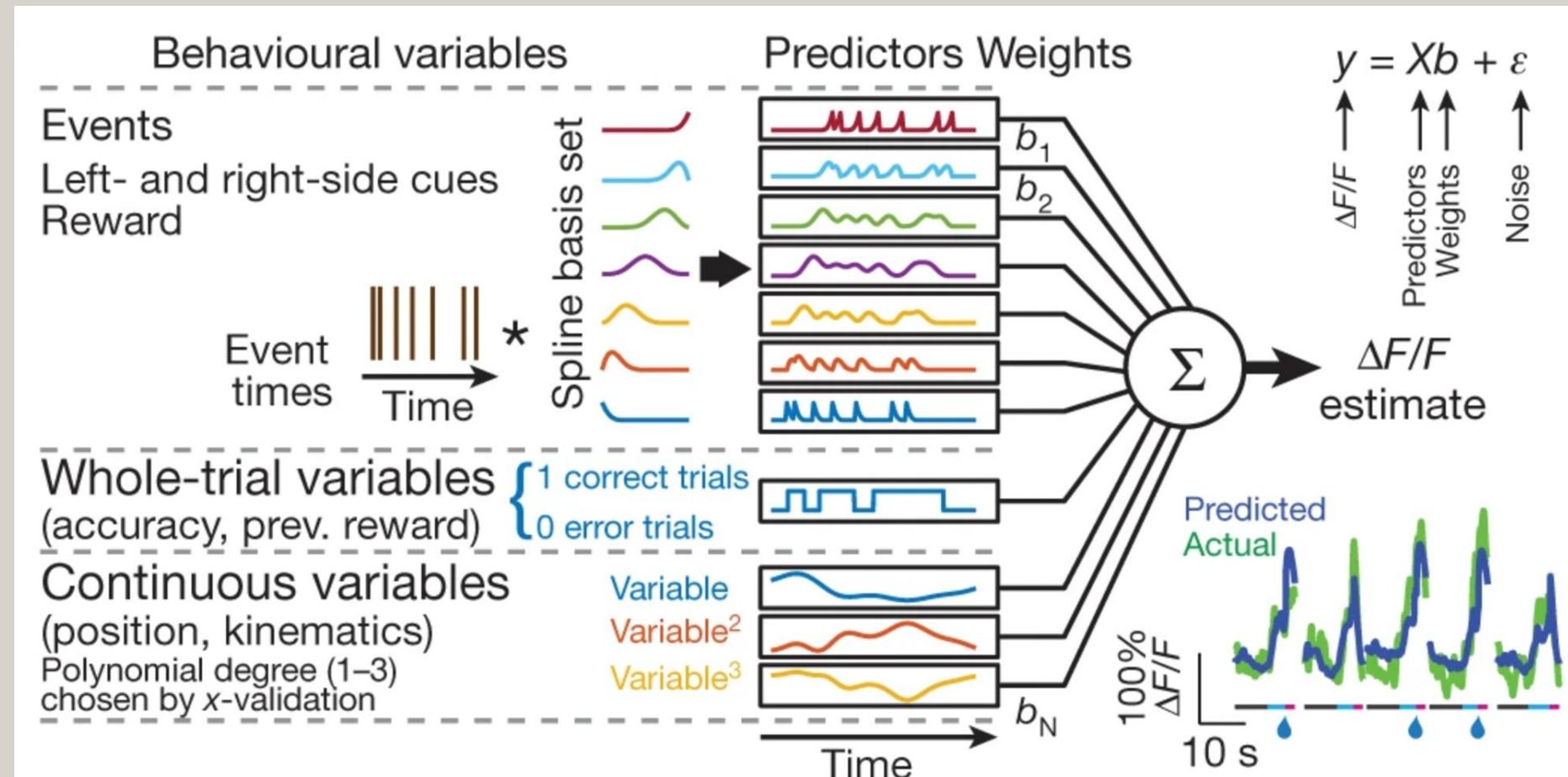
Matrix form

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \boldsymbol{\eta}$$

Index  $i$   
 $\downarrow$   
 $\begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ \vdots \end{bmatrix} = \begin{bmatrix} \cdots & \mathbf{x}_0 & \cdots \\ \cdots & \mathbf{x}_1 & \cdots \\ \cdots & \mathbf{x}_2 & \cdots \\ \vdots & \vdots & \vdots \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \end{bmatrix} + \begin{bmatrix} \eta_0 \\ \eta_1 \\ \eta_2 \\ \vdots \end{bmatrix}$   
 $\underbrace{\qquad\qquad\qquad}_{\text{design matrix}}$   
 $\boldsymbol{\eta} \sim \mathcal{N}(0, \sigma^2)$

# GLM for quantification responses to specific behavioral variables

encoding model to quantify the relationship between behavioral variables and the activity of neurons



# Common statistical tests are linear models

t-tests, ANOVA, regression – are all special cases of one model:  $Y=X\beta+\varepsilon$

Common name	Function in <code>scipy.stats</code>	Equivalent linear model in <code>smf.ols</code>	Exact?	The linear model in words	Icon
<b>y is independent of x</b> P: One-sample t-test N: Wilcoxon signed-rank	<code>scipy.stats.ttest_1samp(y)</code> <code>scipy.stats.wilcoxon(y)</code>	<code>smf.ols("y ~ 1", data)</code> <code>smf.ols("y ~ 1", signed_rank(data))</code>	✓ <a href="#">for N &gt; 14</a>	One number (intercept, i.e., the mean) predicts <b>y</b> . - (Same, but it predicts the <i>signed rank</i> of <b>y</b> .)	
P: Paired-sample t-test N: Wilcoxon matched pairs	<code>scipy.stats.ttest_rel(y1, y2)</code> <code>scipy.stats.wilcoxon(y1, y2)</code>	<code>smf.ols("y2_sub_y1 ~ 1", data)</code> <code>smf.ols("y2_sub_y1 ~ 1", signed_rank(data))</code>	✓ <a href="#">for N &gt; 14</a>	One intercept predicts the pairwise <b>y<sub>2</sub>-y<sub>1</sub></b> differences. - (Same, but it predicts the <i>signed rank</i> of <b>y<sub>2</sub>-y<sub>1</sub></b> .)	
<b>y ~ continuous x</b> P: Pearson correlation N: Spearman correlation	<code>scipy.stats.pearsonr(x, y)</code> <code>scipy.stats.spearmanr(x, y)</code>	<code>smf.ols("y ~ 1 + x", data)</code> <code>smf.ols("y ~ 1 + x", rank(data))</code>	✓ <a href="#">for N &gt; 10</a>	One intercept plus <b>x</b> multiplied by a number (slope) predicts <b>y</b> . - (Same, but with <i>ranked x</i> and <b>y</b> )	
<b>y ~ discrete x</b> P: Two-sample t-test P: Welch's t-test N: Mann-Whitney U	<code>scipy.stats.ttest_ind(y1, y2)</code> N/A in Python, but <a href="#">see R version</a> <code>scipy.stats.mannwhitneyu(y1, y2)</code>	<code>smf.ols("y ~ 1 + group", data)<sup>A</sup></code> N/A in Python, but <a href="#">see R version</a> <code>smf.ols("y ~ 1 + group", signed_rank(data))<sup>A</sup></code>	✓ <a href="#">for N &gt; 11</a>	An intercept for <b>group 1</b> (plus a difference if <b>group 2</b> ) predicts <b>y</b> . - (Same, but with one variance <i>per group</i> instead of one common.) - (Same, but it predicts the <i>signed rank</i> of <b>y</b> .)	
P: One-way ANOVA N: Kruskal-Wallis	<code>scipy.stats.f_oneway(a, b, c)</code> <code>scipy.stats.kruskal(a, b, c)</code>	<code>smf.ols(y ~ 1 + G<sub>2</sub> + G<sub>3</sub> + ... + G<sub>N</sub>)<sup>A</sup></code> <code>smf.ols(rank(y) ~ 1 + G<sub>2</sub> + G<sub>3</sub> + ... + G<sub>N</sub>)<sup>A</sup></code>	✓ <a href="#">for N &gt; 11</a>	An intercept for <b>group 1</b> (plus a difference if $group \neq 1$ ) predicts <b>y</b> . - (Same, but it predicts the <i>rank</i> of <b>y</b> .)	
P: One-way ANCOVA	N/A in Python, but <a href="#">see R version</a>	<code>smf.ols("y ~ 1 + G<sub>2</sub> + G<sub>3</sub> + ... + G<sub>N</sub> + x", data)<sup>A</sup></code>	✓	- (Same, but plus a slope on <b>x</b> ). Note: this is <i>discrete AND continuous</i> . ANCOVAs are ANOVAs with a continuous <b>x</b> .	

# Different GLMs can solve different problems

## Output type

real values

## Likelihood

Gaussian

$$P(y) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp -\frac{(y-\mu)^2}{2\sigma^2}$$

## Nonlinearity

identity

$$\mu = \theta^\top x$$



Linear regression

discrete counts  
0,1,2,3...

Poisson

$$P(y) = \frac{\lambda^y \exp(-\lambda)}{y!}$$

exponential

$$\lambda = \exp(\theta^\top x)$$



Poisson GLM

binary 0,1

Bernoulli

$$P(y) = p^y (1-p)^{1-y}$$

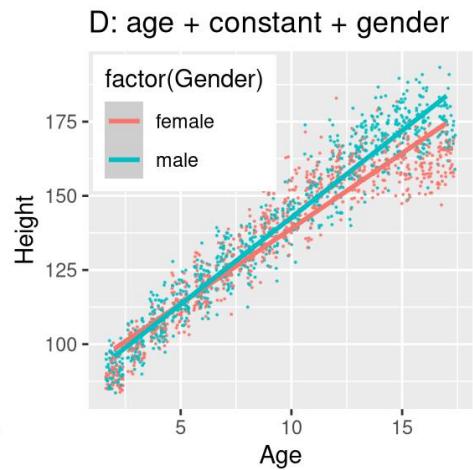
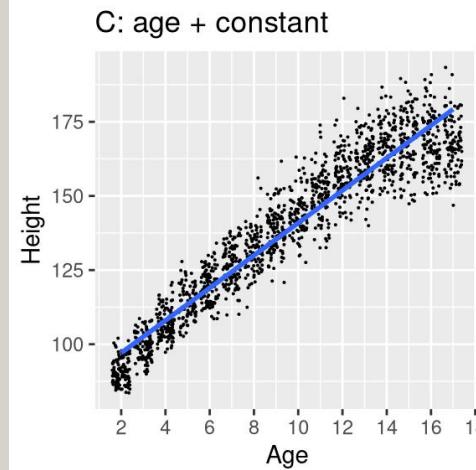
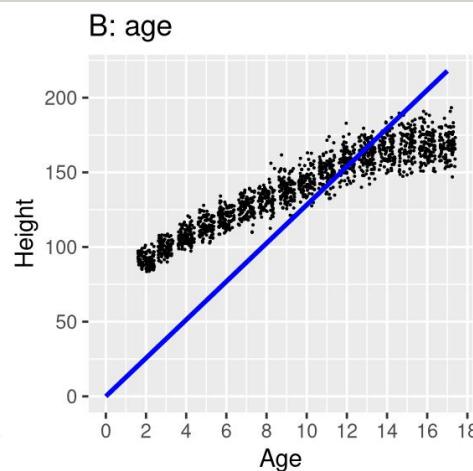
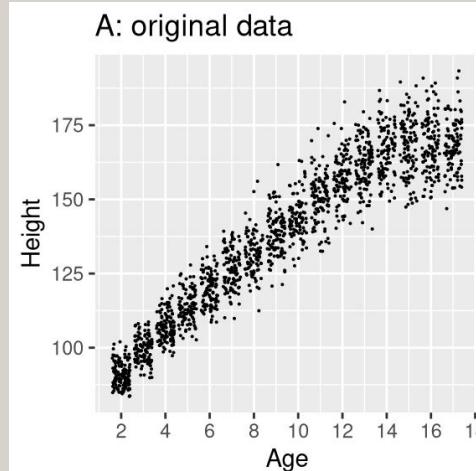
logistic

$$p = \sigma(\theta^\top x)$$

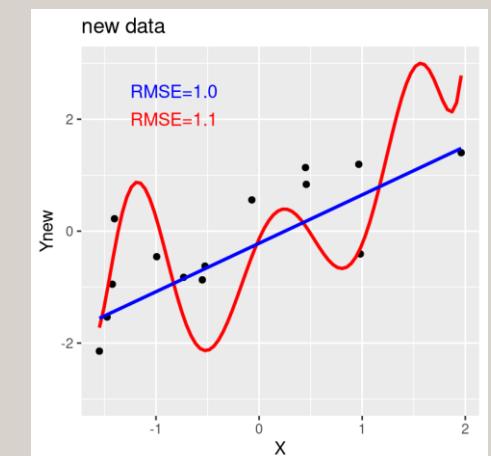
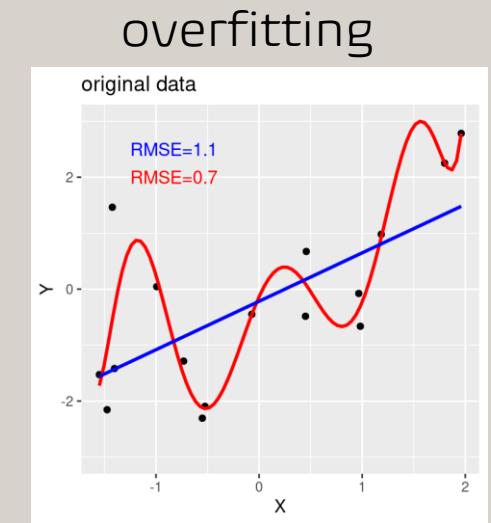
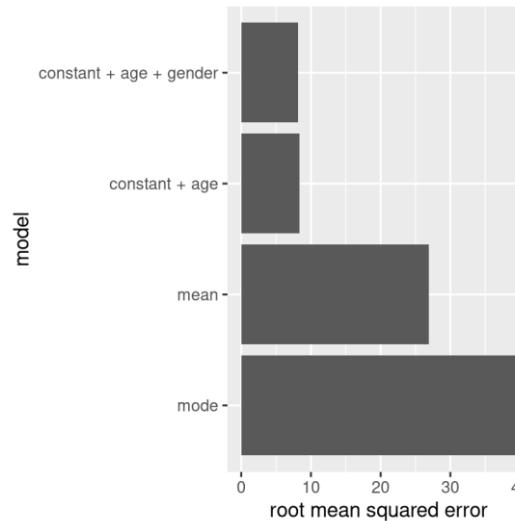


Logistic regression  
(Bernoulli GLM)

# Improving GLM



$$RMSE = \sqrt{\sum_{i=1}^n \frac{(\hat{y}_i - y_i)^2}{n}}$$



Model Comparison Metrics AIC/BIC, adjusted R<sup>2</sup>, likelihood ratio tests, cross-validation

---

# Diagnostic tools to confirm you used the right statistical test

Test Type	Key Assumptions to Check	Diagnostic Tools
<i>t-test / ANOVA</i>	Normality, equal variances	<ul style="list-style-type: none"><li>• Q-Q plot</li><li>• Shapiro-Wilk test</li><li>• Levene's test</li></ul>
<i>Regression / GLM</i>	Linearity, independence, homoscedasticity, normal residuals	<ul style="list-style-type: none"><li>• Residual plots</li><li>• Durbin-Watson test</li><li>• Variance Inflation Factor (VIF)</li><li>• Breusch-Pagan test</li></ul>
<i>Logistic regression</i>	Linearity of logit, absence of separation	<ul style="list-style-type: none"><li>• Box-Tidwell test</li><li>• ROC curve</li><li>• Hosmer-Lemeshow test</li></ul>
<i>Chi-square</i>	Expected counts $\geq 5$	<ul style="list-style-type: none"><li>• Inspect contingency table</li></ul>

Use bootstrapping or permutation tests to verify that your inference is not dependent on assumptions

# 03 The Challenge of High-Dimensional Data

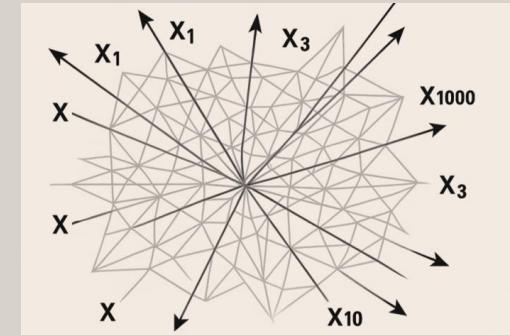
Datasets often have many more features ( $f$ ) than observations ( $n$ ).

Traditional statistics break down when  $f \gg n$ .

## Curse of Dimensionality

Data become sparse in high-dimensional space.

Distances and similarities lose meaning.



## Overfitting

Models capture noise, not signal. Perfect fit  $\neq$  generalizable insight.

## Multicollinearity

Predictors overlap  $\rightarrow$  unstable estimates of effects.

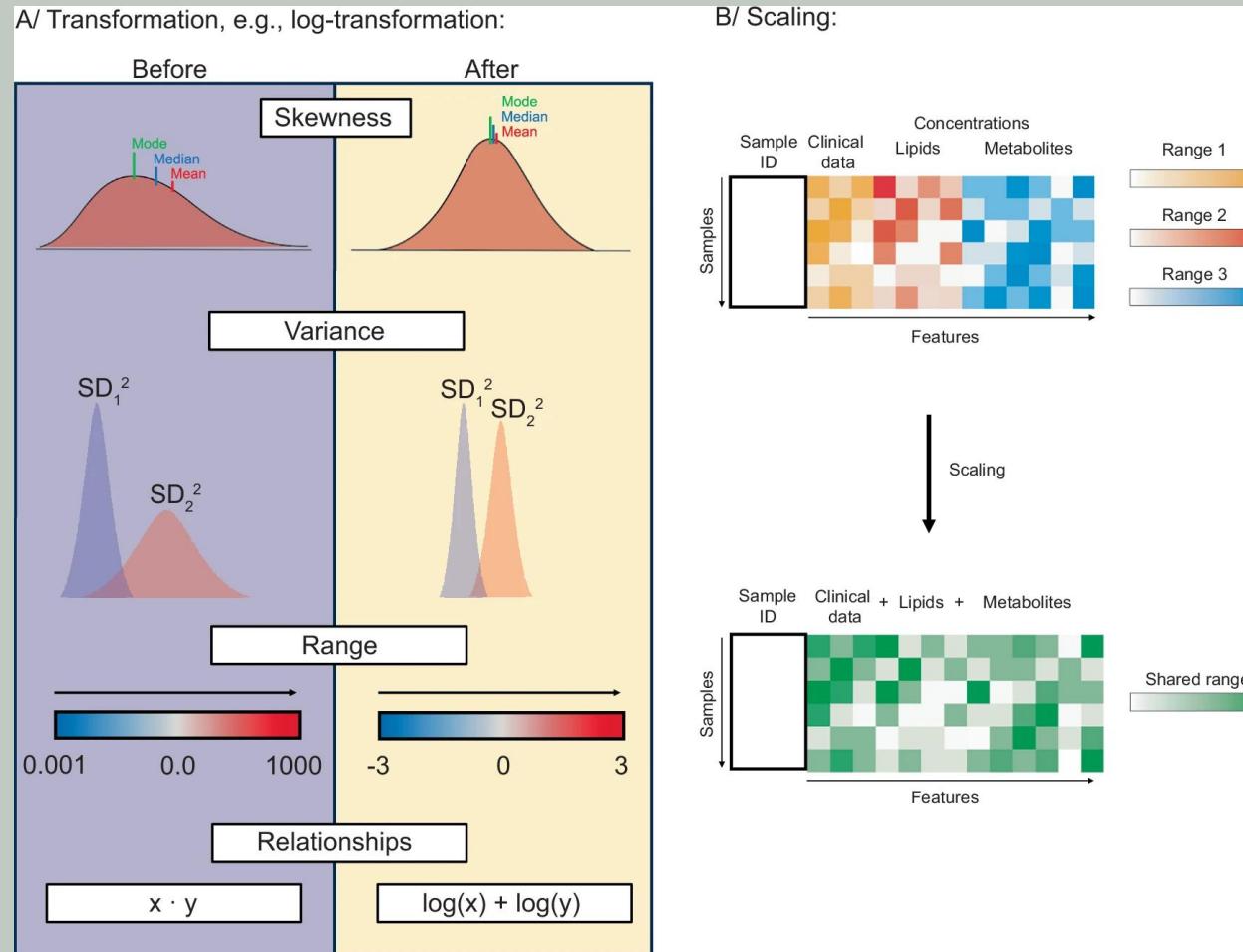
## Interpretability vs. Complexity Trade-off

Simple models = understandable, limited power

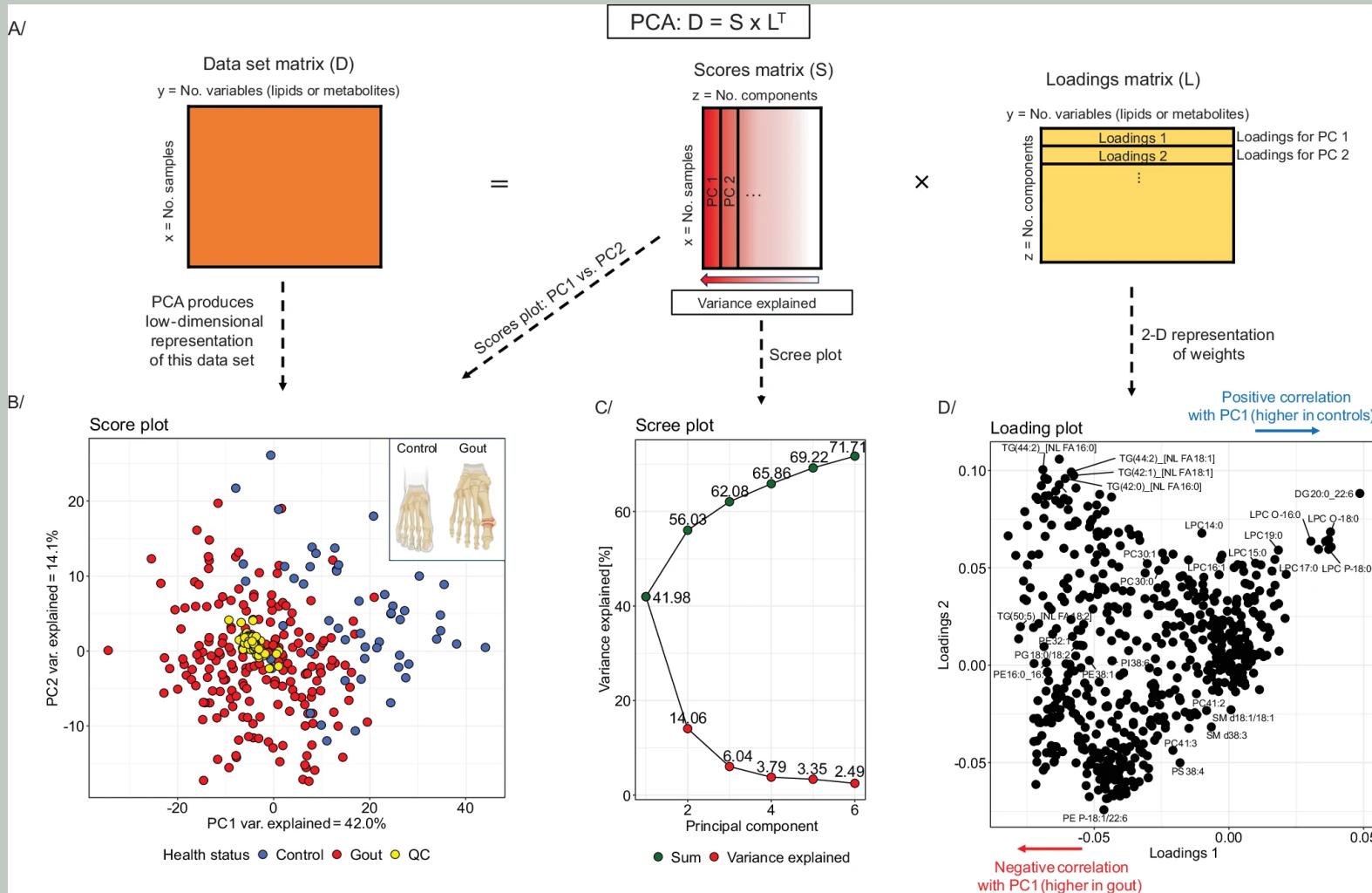
Complex models = powerful, hard to interpret

# Data Preprocessing and Exploration

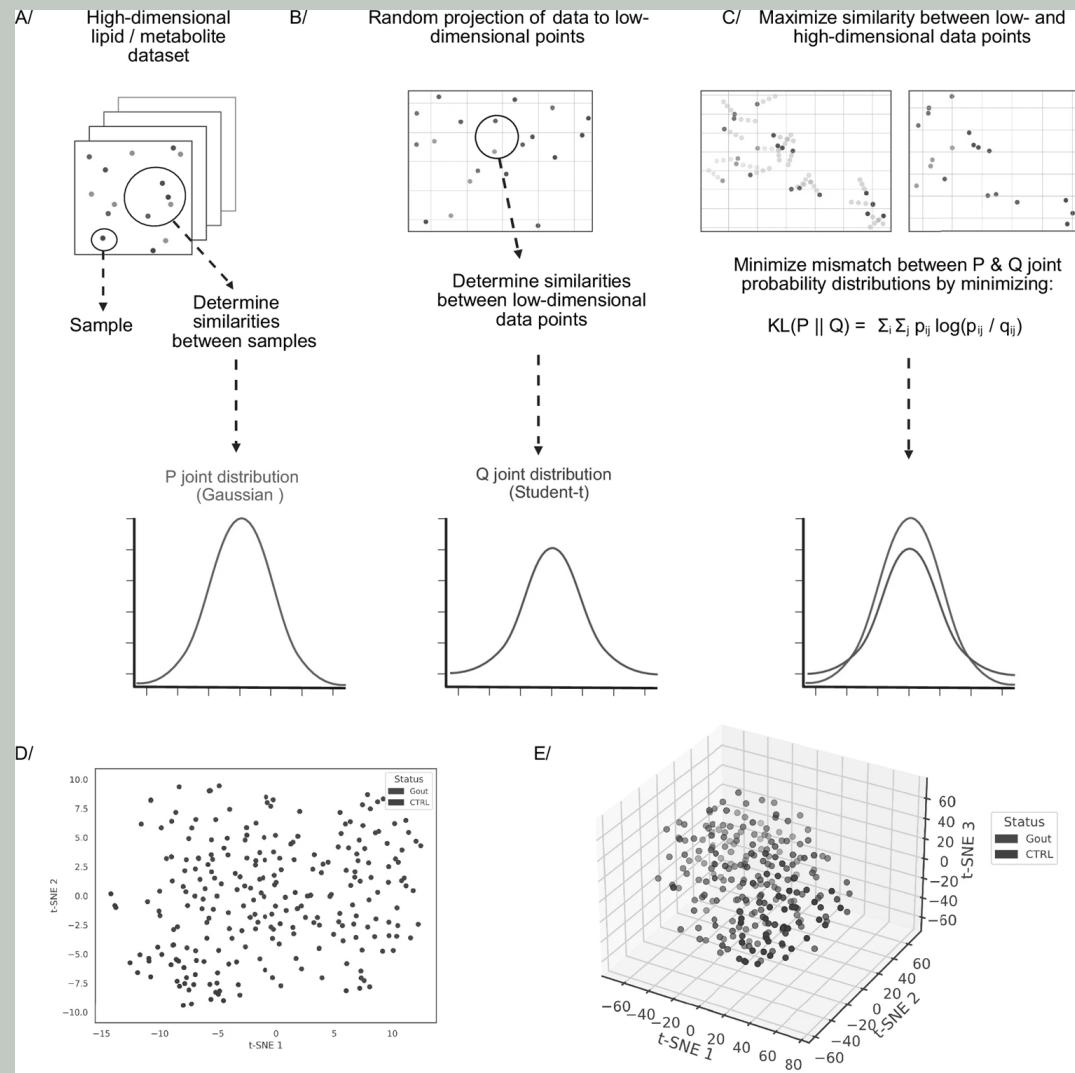
## Data transformation and scaling



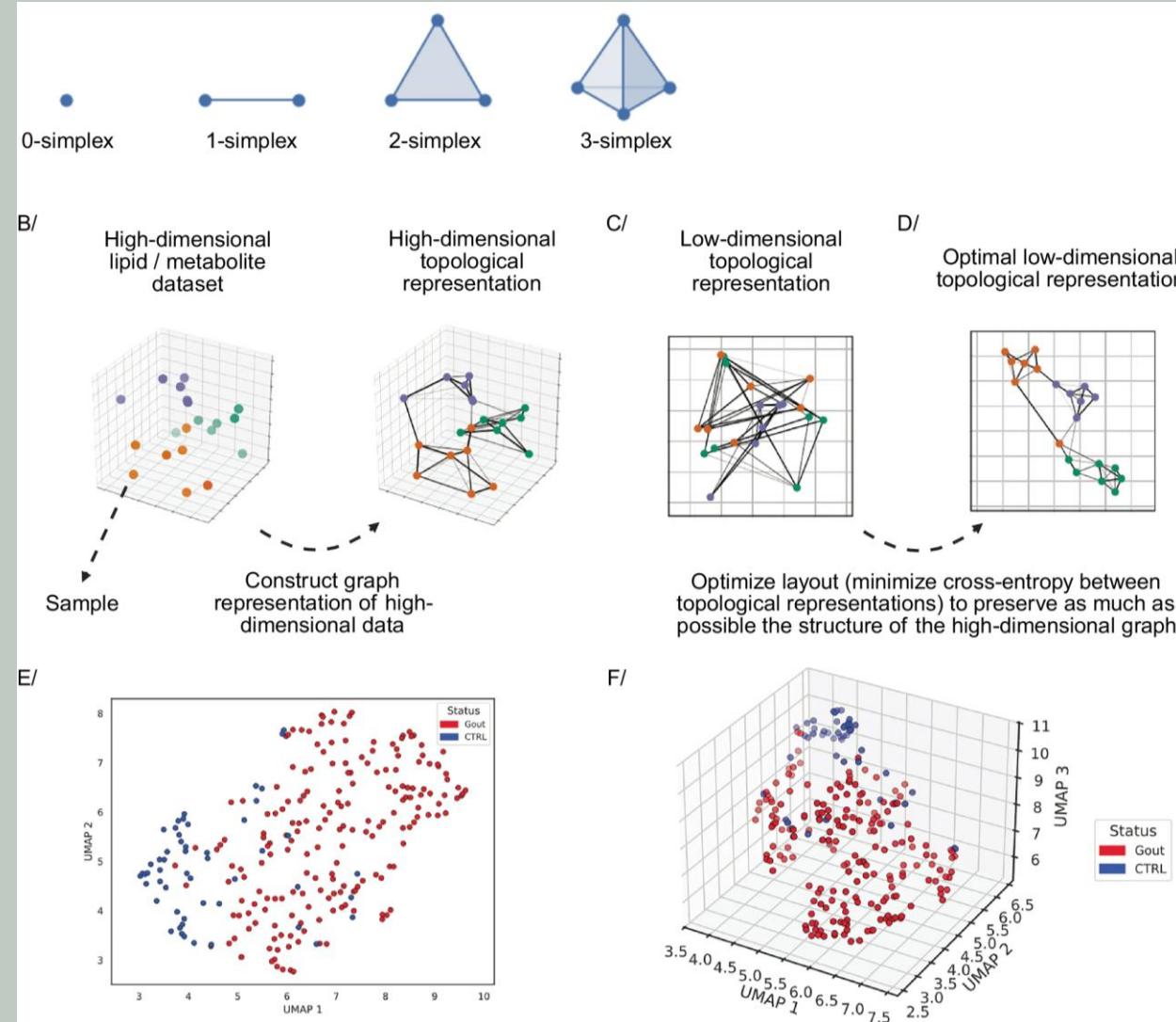
# Overview of principal component analysis



# t-Distributed Stochastic Neighbor Embedding

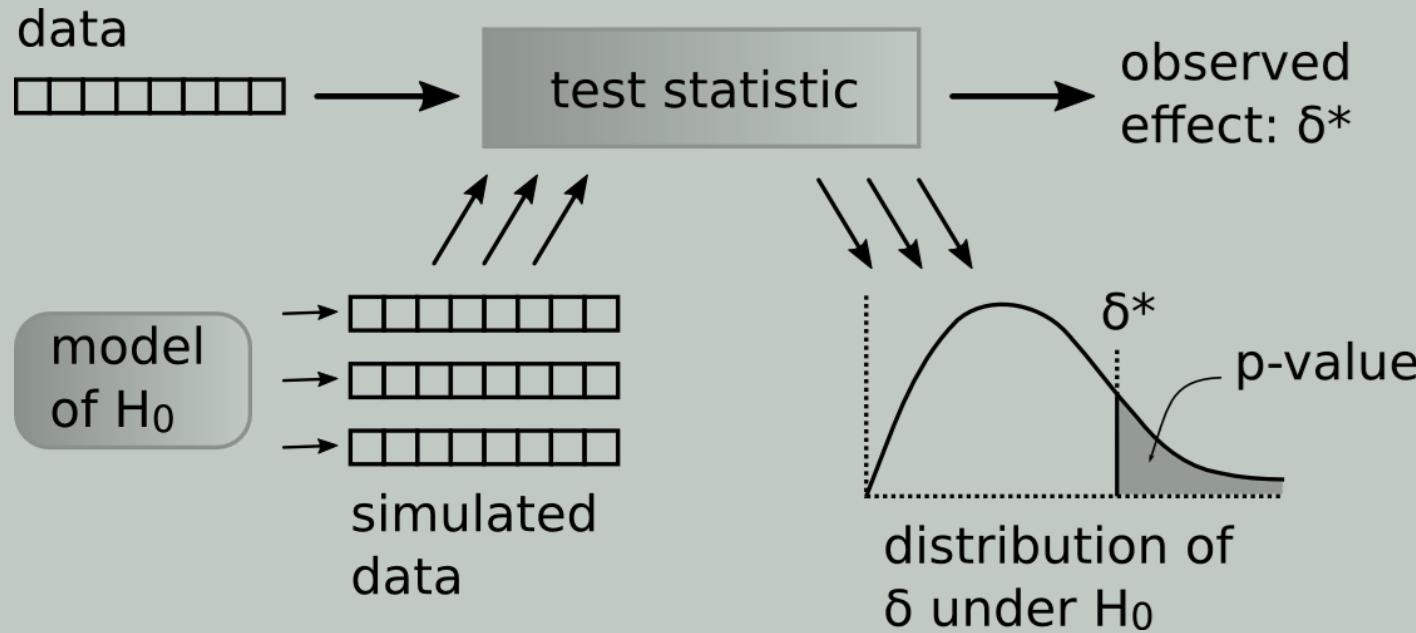
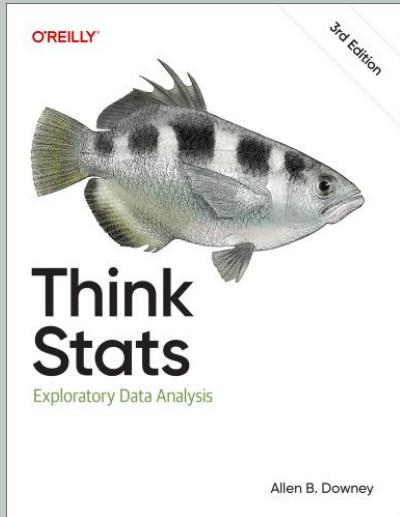


# Uniform manifold approximation and projection



# There is only one test!

Allen Downey



---

# Libraries for statistical analysis



SciPy



**pingouin**

---

Time to practice  
Open the Google Colaboratory (Colab)