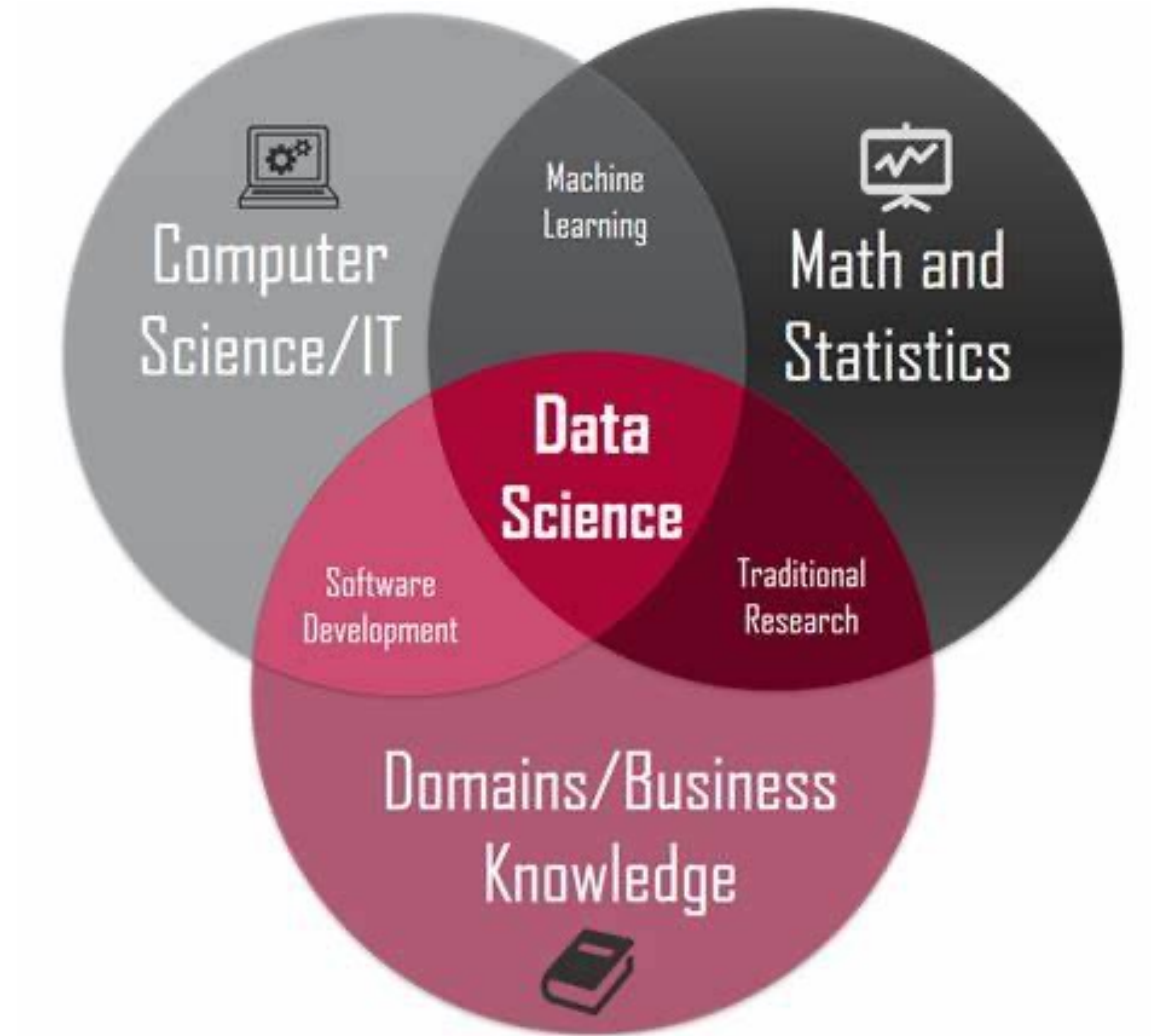


# **STATISTICS REVIEW**

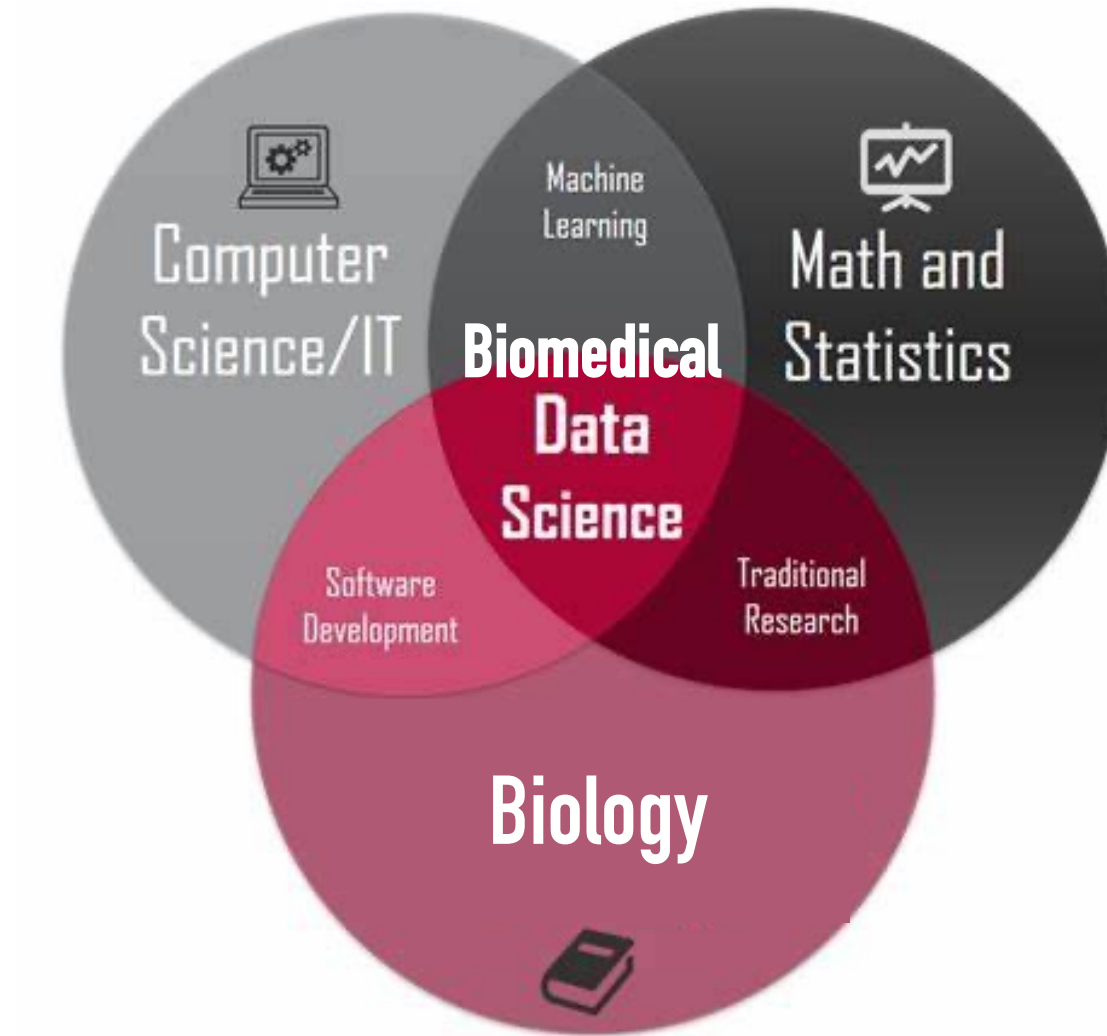
# OUTLINE

- Random variable
- Probability distribution
- Central limit theorem
- Hypothesis testing
- P-value
- Multiple testing correction
- Type I and type II errors
- False discovery rate (FDR)

# Data Science



# Computational Biology



# Why statistics is important

- Statistics is the theoretical foundation of machine learning and data science
- Statistics is the bridge between experiments and theories
- Statistics is the bridge between observations (data) and discoveries (science)

# What is a model

- A model describes the relationship between quantities
  - Quantitative/mathematical models
- Statistical models
  - Relationship between random variables
- “All models are wrong; but some are useful.” (George Box)

# Random variables

- Randomness
- All quantities obtained from experimental measurements can be considered as random variables
- Discrete vs. continuous random variables
- Probability distribution

# Discrete probability distributions

- Bernoulli distribution

$$\Pr(X = 1) = p = 1 - \Pr(X = 0) = 1 - q.$$

- Binomial distribution

$$\Pr(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

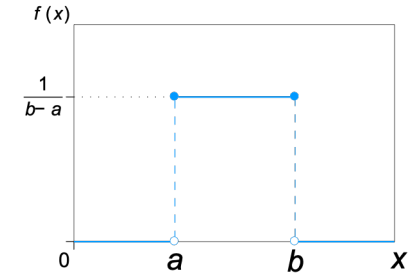
- Negative binomial distribution  $\Pr(X = k) = \binom{k+r-1}{k} (1-p)^k p^r$

- Poisson distribution  $\Pr(X=k) = \frac{\lambda^k e^{-\lambda}}{k!}$



# Continuous probability distributions

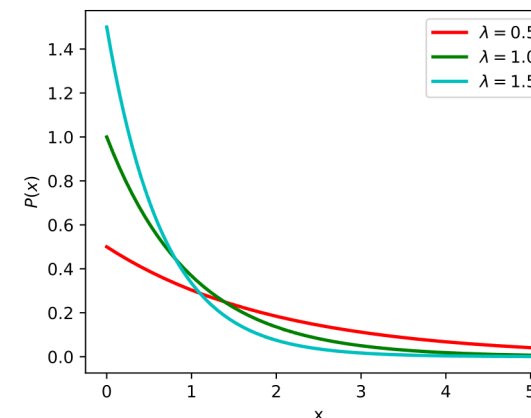
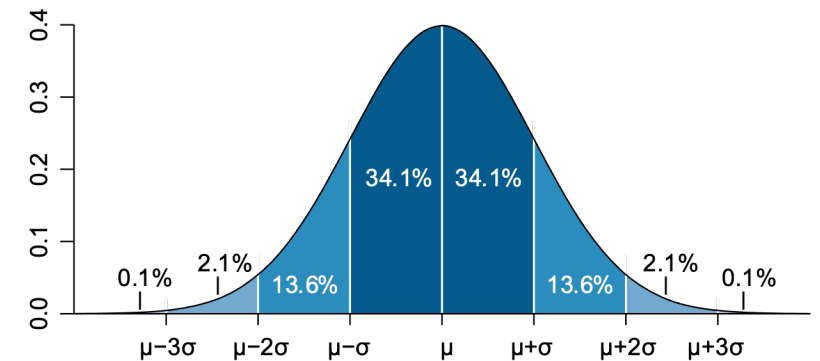
- Uniform distribution
- Normal distribution (Gaussian distribution)



$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

- Exponential distribution

$$f(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0, \\ 0 & x < 0. \end{cases}$$



# Use of probability distributions in computational biology

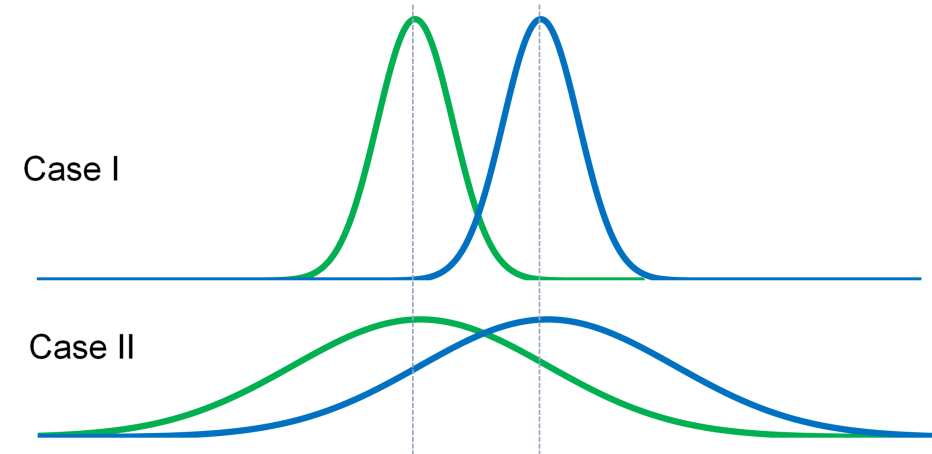
- Null models for hypothesis testing
- Noise or background models for signal detection and differential enrichment analysis
- Error models in regression/machine learning

# Central limit theorem

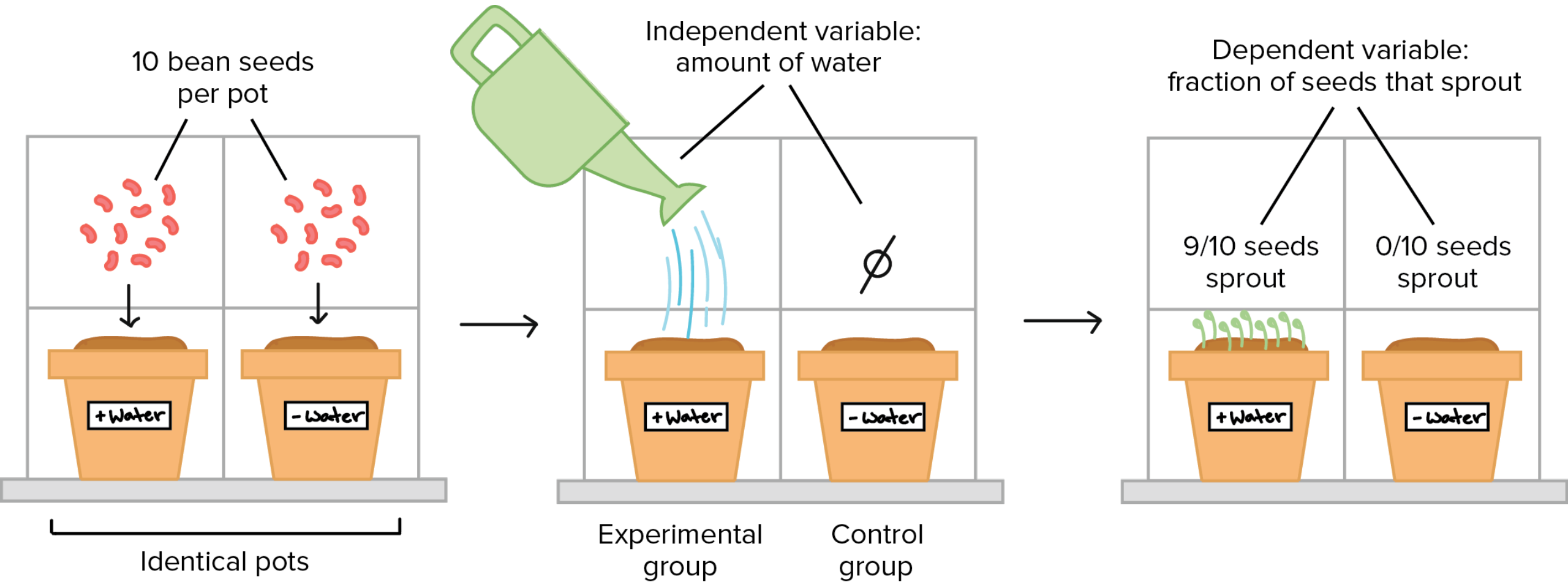
- Regardless of the population distribution, the sample mean will follow a standard normal distribution, if the samples are independent and equal size.
- Theoretical foundation of t-test

# Hypothesis testing

- General thinking:
  - Are they different?
  - Is the difference “statistically significant”?
- Statistical thinking:
  - Null hypothesis
  - Alternative hypothesis



# Negative control in experiments



# Null hypothesis in statistical hypothesis testing

AwesomeFinTech

A null hypothesis is a type of conjecture used in statistics that proposes that there is no difference between certain characteristics of a population or data-generating process.

read more about 📌

**Null Hypothesis : Testing & Examples**

[www.awesomefintech.com/terms/null\\_hypothesis/](https://www.awesomefintech.com/terms/null_hypothesis/)

**The null hypothesis does not depend on a test procedure**

# Null hypothesis in statistical hypothesis testing

AwesomeFinTech

Hypothesis testing provides a method to reject a null hypothesis within a certain confidence level. (Null hypotheses cannot be proven, though.)

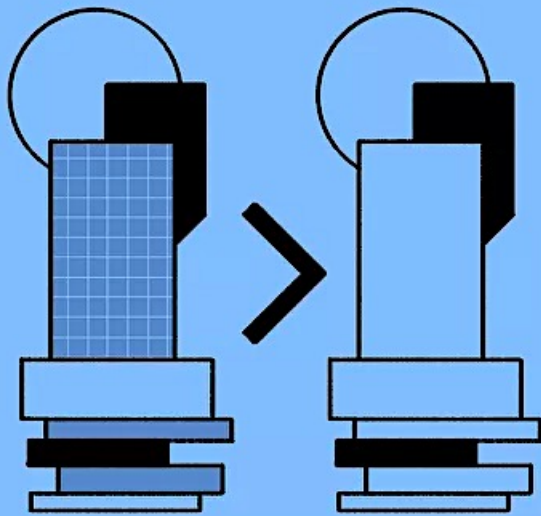
read more about 📌

**Null Hypothesis : Testing & Examples**

[www.awesomefintech.com/terms/null\\_hypothesis/](https://www.awesomefintech.com/terms/null_hypothesis/)

## Statistics cannot tell us everything b/c data is random

# Null hypothesis is often misunderstood



## Null Hypothesis

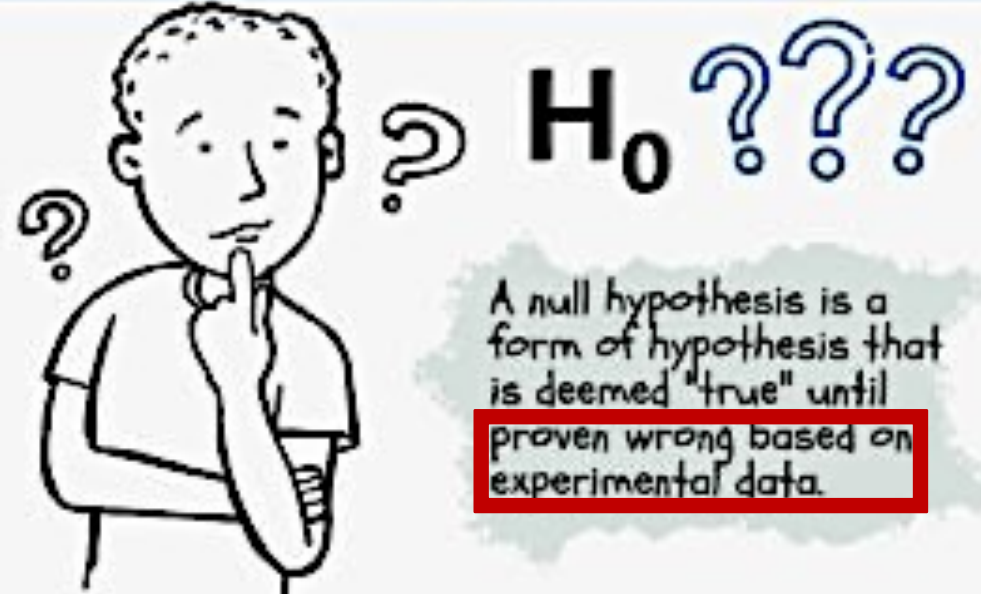
*[ˈnəl hī-ˈpä-thə-səs]*

A hypothesis that proposes that no statistical significance exists in a set of given observations and is used to assess the credibility of a hypothesis by using sample data.

 Investopedia

Source: [https://www.investopedia.com/terms/n/null\\_hypothesis.asp](https://www.investopedia.com/terms/n/null_hypothesis.asp)

## Null hypothesis



Source: <https://www.biologyonline.com/dictionary/null-hypothesis>



# Common hypothesis tests

- Student's t-test (parametric)
  - Null (2-sample): The sample means are equal.
- Fisher's exact test
  - Null: The two groups are equally likely for an event/feature.
- Wilcoxon (rank-sum) test
  - Null: The two samples  $X$  and  $Y$ ,  $P(X > Y) = P(X < Y) = 0.5$
- Kolmogorov-Smirnov test (K-S test)
  - Null: The two samples have the same cumulative distribution.
- Hypothesis test statistic, p-value

# Which of the following statements about p-values is true?

- A. P-values measure how big the difference is between the datasets compared.
- B. P-value is the probability of observing the data by random chance.
- C. P-value is the least probability of observing the data under the assumption that the null hypothesis is true.

# **ASA statement on statistical significance and p-values**

1. P-values can indicate how incompatible the data are with a specified statistical model.
2. P-values do not measure the probability that the studied hypothesis is true, or the probability that the data were produced by random chance alone.
3. Scientific conclusions and business or policy decisions should not be based only on whether a p-value passes a specific threshold.

# **ASA statement on statistical significance and p-values**

4. Proper inference requires full reporting and transparency.
5. A p-value, or statistical significance, does not measure the size of an effect or the importance of a result.
6. By itself, a p-value does not provide a good measure of evidence regarding a model or hypothesis.

# Multiple testing correction

- High-throughput experiments
  - RNA-seq: 500 Differentially Expressed genes from 20K genes
  - ChIP-seq: 10,000 TF binding sites from the genome

# Bonferroni correction

- Controls family-wise error rate
- $N$  tests
- Adjusted P-value =  $N * P$
- Very strict



Carlo Emilio Bonferroni (1892-1960)

# Benjamini-Hochberg (B-H) Correction

- N observations w/ various p-values
- Rank N
- Adjusted  $P = P * N / R$
- Moderate



Yoav Benjamini (1949-)

and Yosef Hochberg

# Confusion Matrix

		True Class	
		Positive	Negative
Predicted Class	Positive	TP	FP (Type I Error)
	Negative	FN (Type II Error)	TN



# Summary

		CONDITION determined by "Gold Standard"			
TOTAL POPULATION		CONDITION POS	CONDITION NEG	PREVALENCE $\frac{\text{CONDITION POS}}{\text{TOTAL POPULATION}}$	
TEST OUT-COME	TEST POS	True Pos TP	Type I Error False Pos FP	Precision Pos Predictive Value $\text{PPV} = \frac{\text{TP}}{\text{TEST P}}$	False Discovery Rate $\text{FDR} = \frac{\text{FP}}{\text{TEST P}}$
	TEST NEG	Type II Error False Neg FN	True Neg TN	False Omission Rate $\text{FOR} = \frac{\text{FN}}{\text{TEST N}}$	Neg Predictive Value $\text{NPV} = \frac{\text{TN}}{\text{TEST N}}$
ACCURACY ACC $\text{ACC} = \frac{\text{TP} + \text{TN}}{\text{TOT POP}}$		Sensitivity (SN), Recall Total Pos Rate TPR $\text{TPR} = \frac{\text{TP}}{\text{CONDITION POS}}$	Fall-Out False Pos Rate FPR $\text{FPR} = \frac{\text{FP}}{\text{CONDITION NEG}}$	Pos Likelihood Ratio LR + $\text{LR} + = \frac{\text{TPR}}{\text{FPR}}$	Diagnostic Odds Ratio DOR $\text{DOR} = \frac{\text{LR} +}{\text{LR} -}$
		Miss Rate False Neg Rate FNR $\text{FNR} = \frac{\text{FN}}{\text{CONDITION POS}}$	Specificity (SPC) True Neg Rate TNR $\text{TNR} = \frac{\text{TN}}{\text{CONDITION NEG}}$	Neg Likelihood Ratio LR - $\text{LR} - = \frac{\text{TNR}}{\text{FNR}}$	

SHORT REPORT

Open Access



# Exaggerated false positives by popular differential expression methods when analyzing human population samples

Yumei Li<sup>1†</sup>, Xinzhou Ge<sup>2†</sup>, Fanglue Peng<sup>3</sup>, Wei Li<sup>1\*</sup> and Jingyi Jessica Li<sup>2,4,5,6,7\*</sup> 

\*Correspondence:  
wei.li@uci.edu; lijy03@g.  
ucla.edu

<sup>†</sup>Yumei Li and Xinzhou Ge  
contributed equally to this  
work.

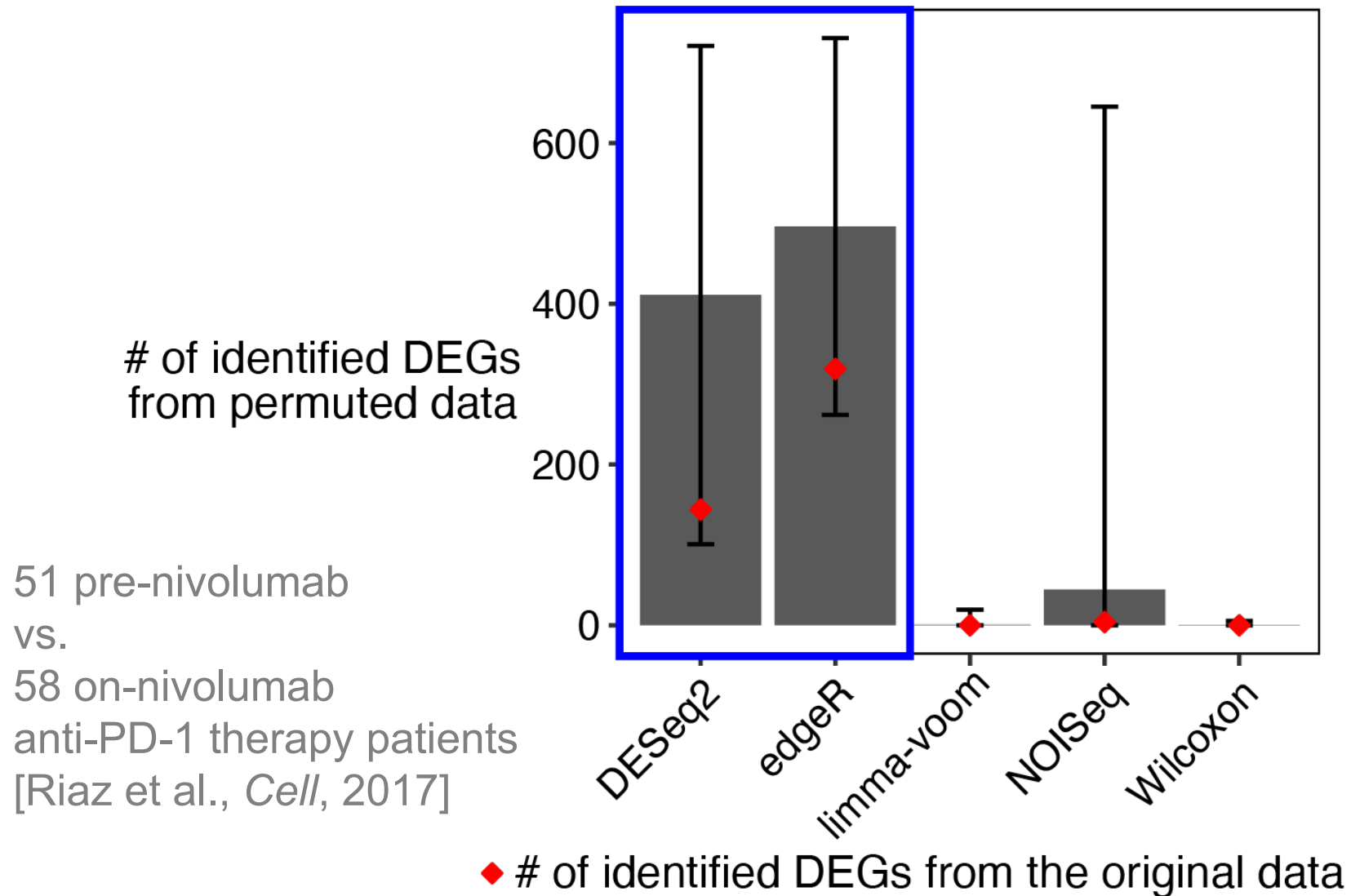
<sup>1</sup> Division of Computational  
Biomedicine, Department  
of Biological Chemistry,  
School of Medicine,  
University of California, Irvine,  
Irvine, CA 92697, USA

<sup>2</sup> Department of Statistics,  
University of California, Los

## Abstract

When identifying differentially expressed genes between two conditions using human population RNA-seq samples, we found a phenomenon by permutation analysis: two popular bioinformatics methods, DESeq2 and edgeR, have unexpectedly high false discovery rates. Expanding the analysis to limma-voom, NOISeq, dearseq, and Wilcoxon rank-sum test, we found that FDR control is often failed except for the Wilcoxon rank-sum test. Particularly, the actual FDRs of DESeq2 and edgeR sometimes exceed 20% when the target FDR is 5%. Based on these results, for population-level RNA-seq studies with large sample sizes, we recommend the Wilcoxon rank-sum test.

# Why there are many DE genes identified from permuted data?



# Important assumption in DESeq2 and edgeR

Both DESeq2 and edgeR assume a **negative binomial (NB)** distribution per gene and condition

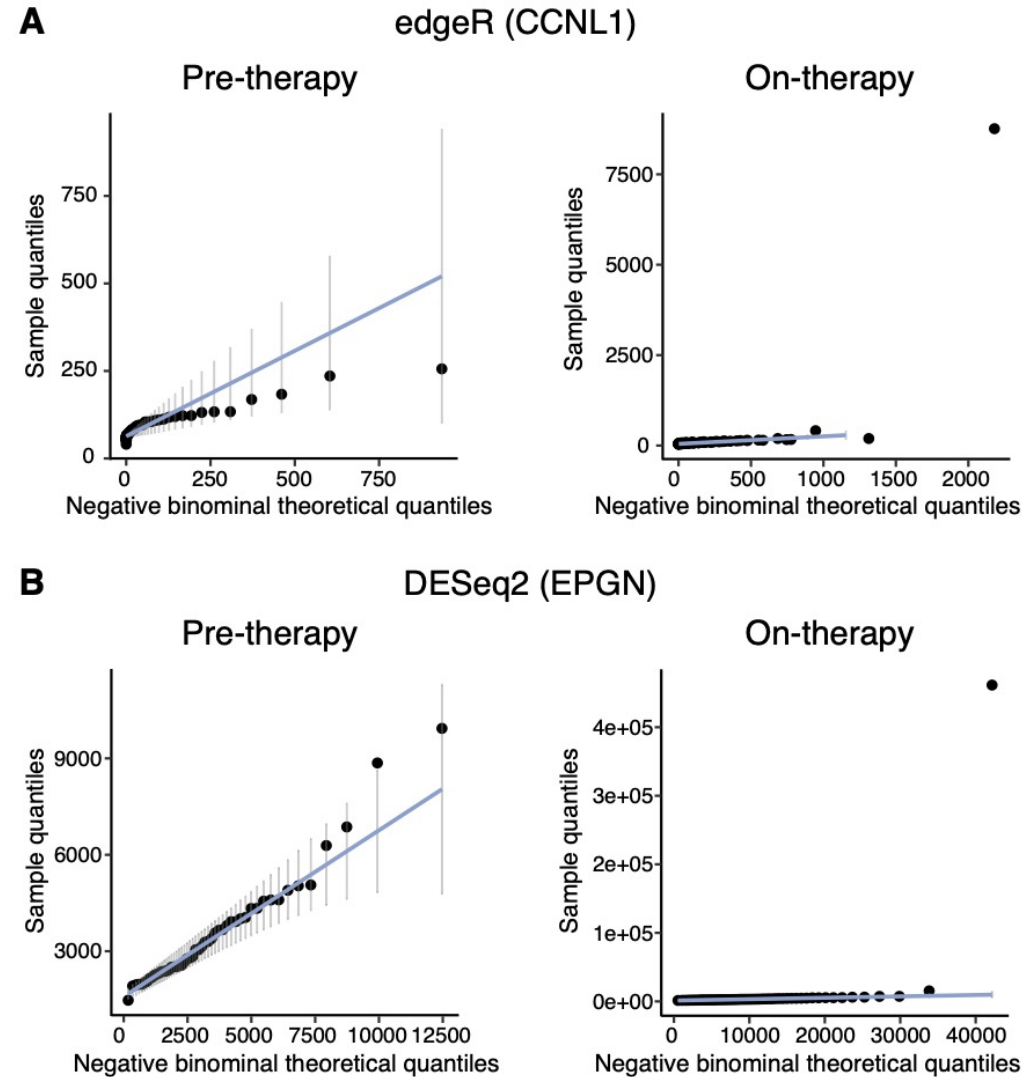
For each gene,

- Condition 1:  $X_i \stackrel{\text{ind}}{\sim} \text{NB}(\mu_1 s_i, \sigma_1), i = 1, \dots, n$
- Condition 2:  $Y_j \stackrel{\text{ind}}{\sim} \text{NB}(\mu_2 s_j, \sigma_2), j = 1, \dots, m$

**Null hypothesis**  $H_0 : \mu_1 = \mu_2$

appropriate ***only if*** the NB assumption is reasonable

# Gene expression can deviate from NB distribution



# Why does Wilcoxon test work in this scenario?

For each gene, the normalized counts

Condition 1:  $\tilde{X}_i, i = 1, \dots, n$

Condition 2:  $\tilde{Y}_j, j = 1, \dots, m$

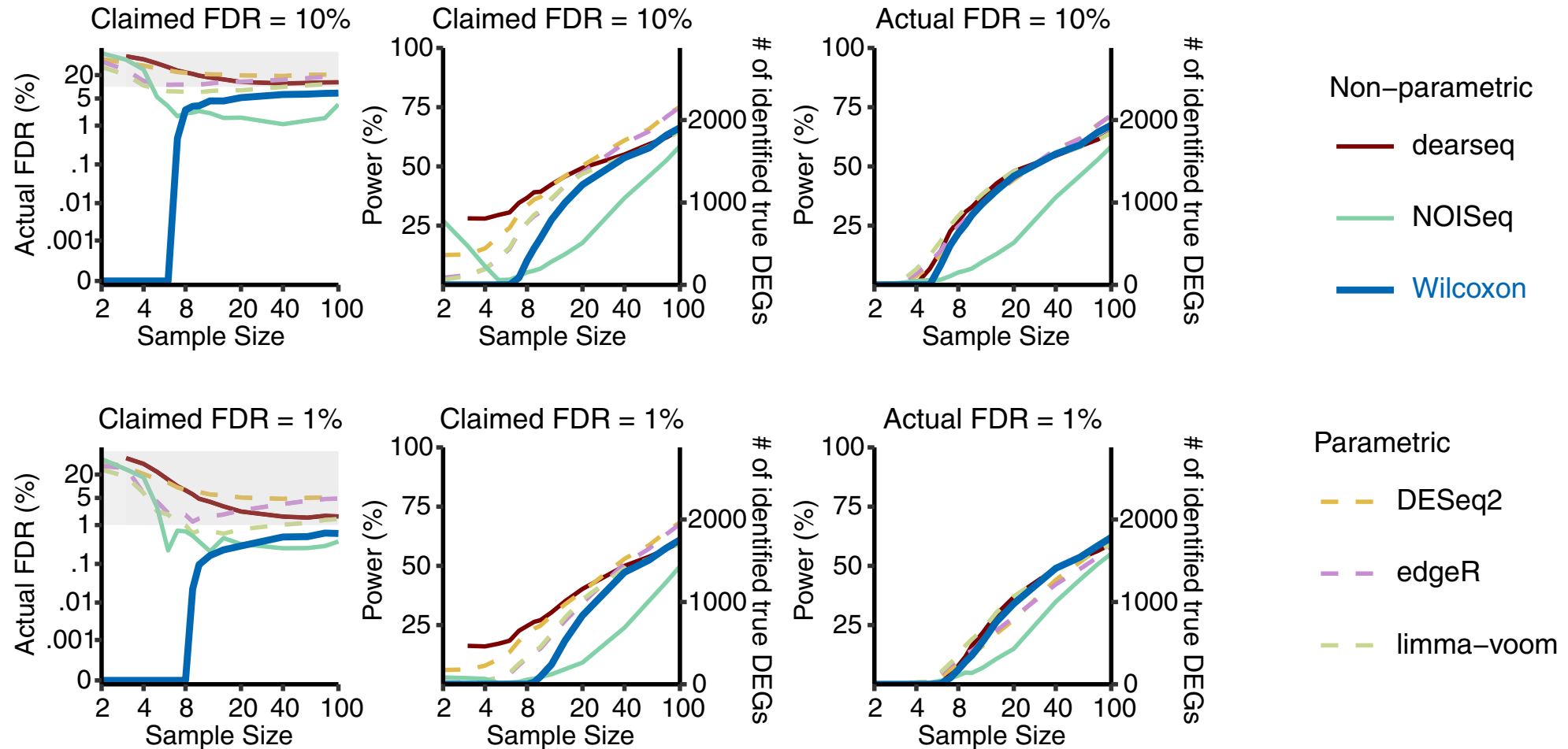
**Null hypothesis (approximate, ignoring ties):**

$$H_0 : \mathbb{P}(\tilde{X}_i > \tilde{Y}_j) = 0.5, \text{ for all } i, j$$

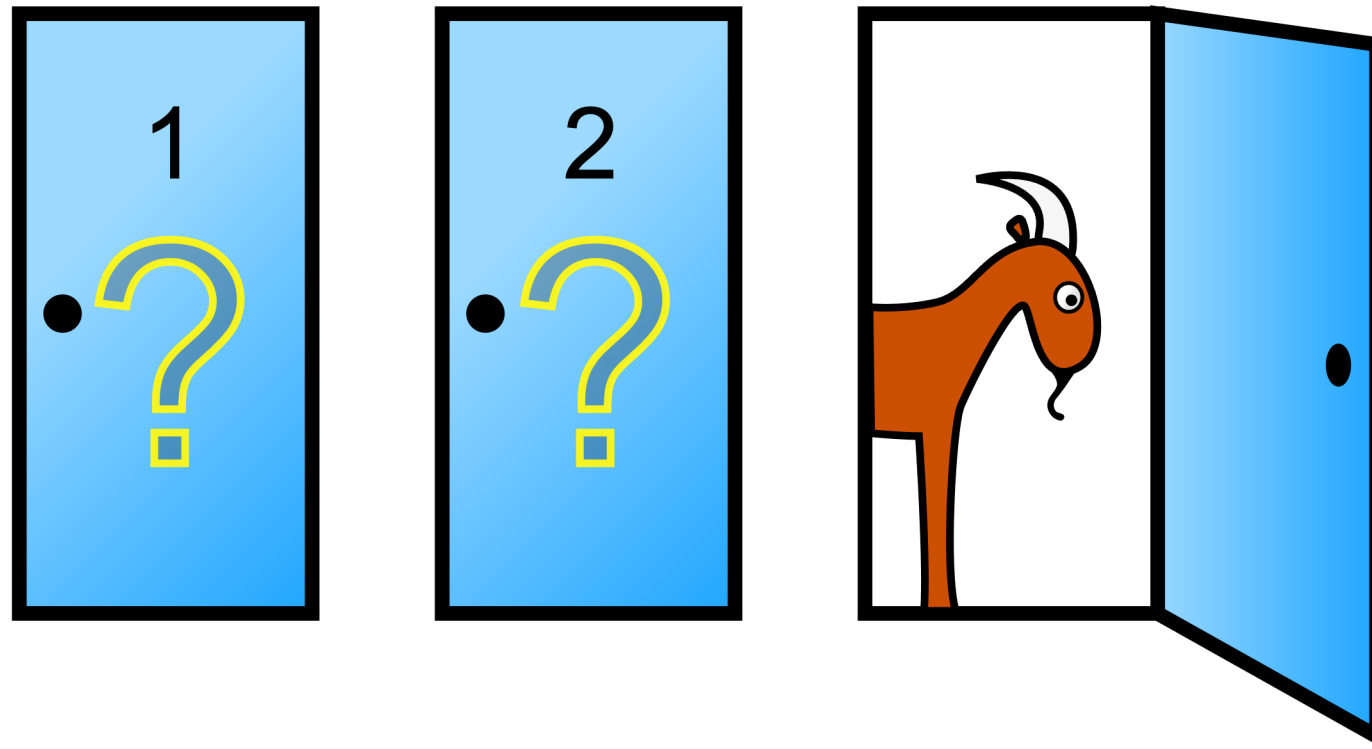
which does NOT have the NB assumption



# Wilcoxon test is better when sample size > 8



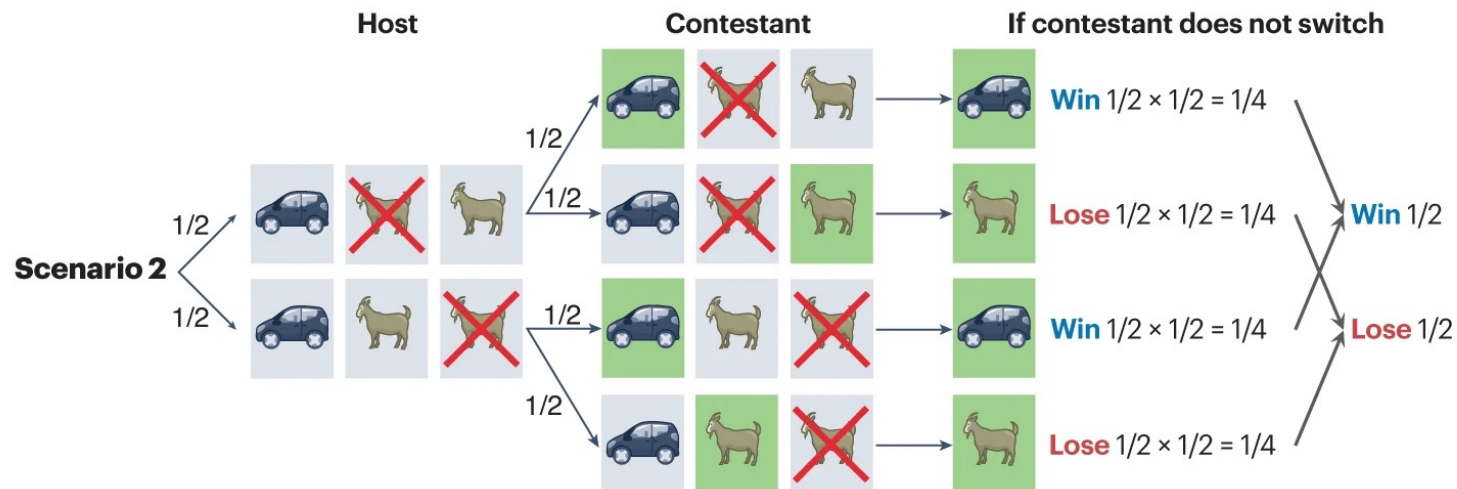
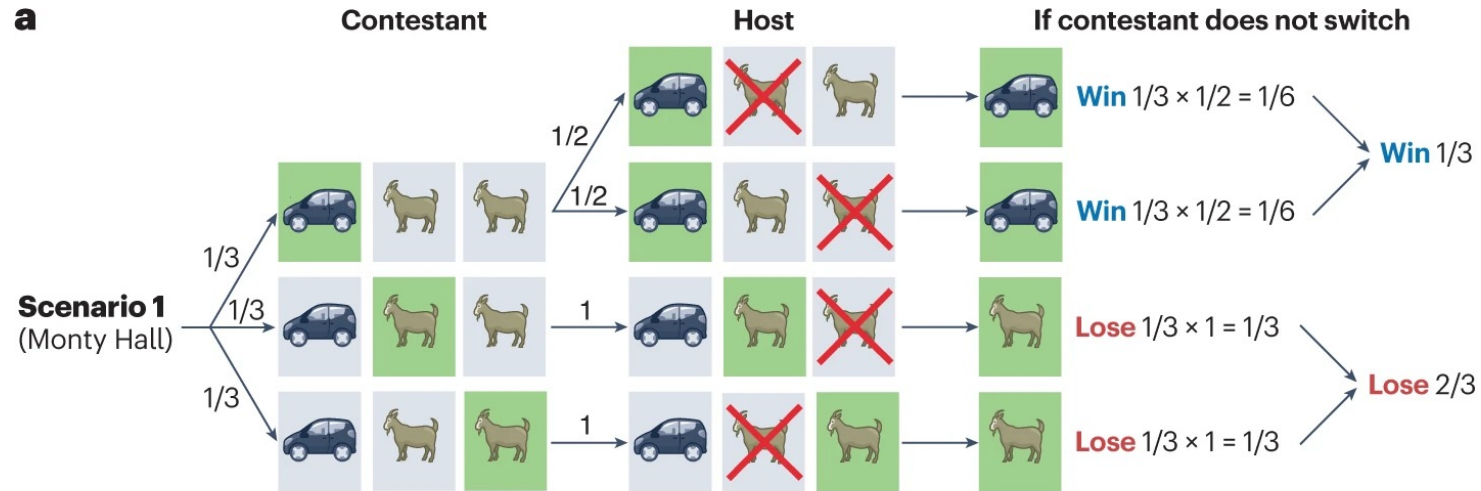
# Monty Hall Problem



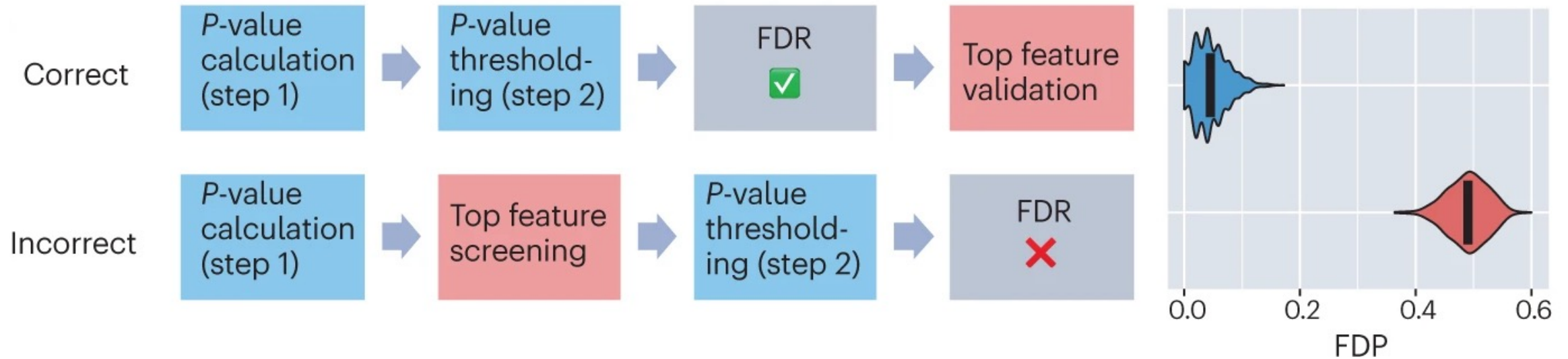


# Monty Hall Problem

**a**



# The order of action matters



EDITORIAL

# Ten Simple Rules for Effective Statistical Practice

**Robert E. Kass<sup>1</sup>, Brian S. Caffo<sup>2</sup>, Marie Davidian<sup>3</sup>, Xiao-Li Meng<sup>4</sup>, Bin Yu<sup>5</sup>, Nancy Reid<sup>6\*</sup>**

**1** Department of Statistics, Machine Learning Department, and Center for the Neural Basis of Cognition, Carnegie Mellon University, Pittsburgh, Pennsylvania, United States of America, **2** Department of Biostatistics, Bloomberg School of Public Health, Johns Hopkins University, Baltimore, Maryland, United States of America, **3** Department of Statistics, North Carolina State University, Raleigh, North Carolina, United States of America, **4** Department of Statistics, Harvard University, Cambridge, Massachusetts, United States of America, **5** Department of Statistics and Department of Electrical Engineering and Computer Science, University of California Berkeley, Berkeley, California, United States of America, **6** Department of Statistical Sciences, University of Toronto, Toronto, Ontario, Canada

\* [reid@utstat.utoronto.ca](mailto:reid@utstat.utoronto.ca)

# Ten Simple Rules for Effective Statistical Practice

1. Statistical methods should enable data to answer scientific questions.
2. Signals always come with noise.
3. Plan ahead, really ahead.
4. Worry about data quality.
5. Statistical analysis is more than a set of computations.

# Ten Simple Rules for Effective Statistical Practice

6. Keep it simple.
7. Provide assessments of variability.
8. Check your assumptions.
9. When possible, replicate!
10. Make your analysis reproducible.

# HOW TO: DRAW A HORSE

BY VAN OKTOP

---



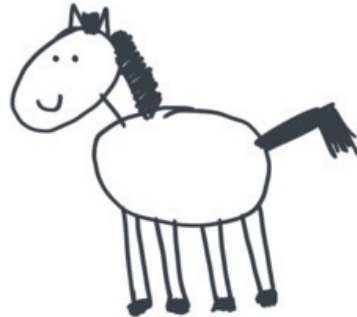
① DRAW 2 CIRCLES



② DRAW THE LEGS

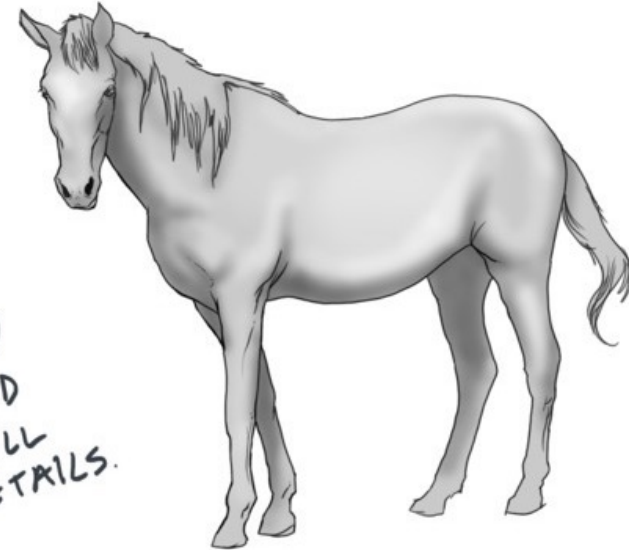


③ DRAW THE FACE



④ DRAW THE HAIR

*Record procedure details!*



⑤  
ADD  
SMALL  
DETAILS.

# About Assignment 1 (Section B)

- Record all code and results
- Submit in any format (RMD preferred)
- Due Feb 19, 2024.
- Assignment 1 can be emailed to [zang@virginia.edu](mailto:zang@virginia.edu)