

# **STATISTICS REVIEW II**

# OUTLINE

- Sampling Bias
- Simpson's Paradox
- Type I and type II errors
- Frequentist vs. Bayesian
- A case study

# 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.

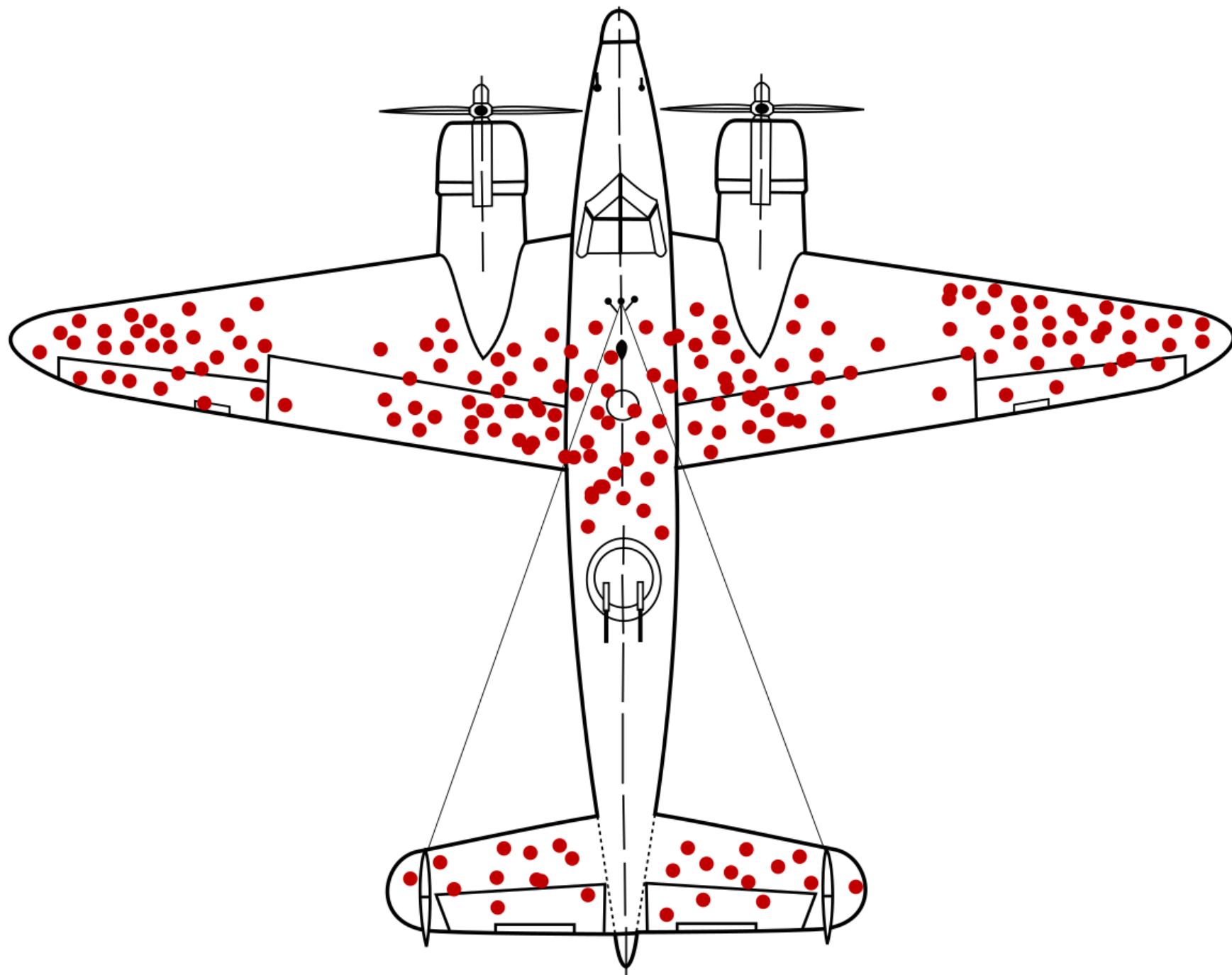
# **ASA statement on statistical significance and p-values**

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**

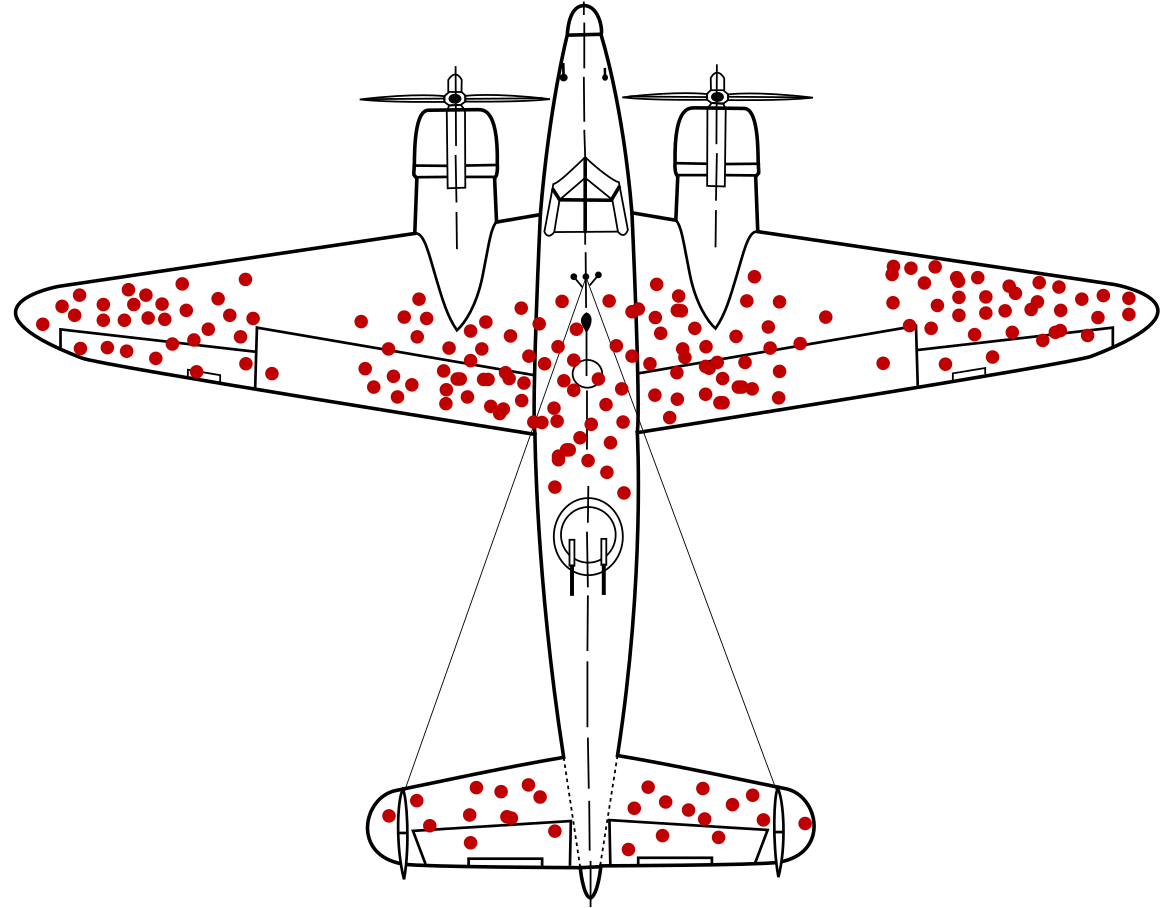
4. Proper inference requires full reporting and transparency.
5. A p-value, or statistical significance, does not measure the size of an effect or the importance of a result.
6. By itself, a p-value does not provide a good measure of evidence regarding a model or hypothesis.

**What is the control?**  
**What is the null hypothesis?**



# What is the population/baseline?

- Aircrafts that had returned from missions
- All aircrafts that went to missions



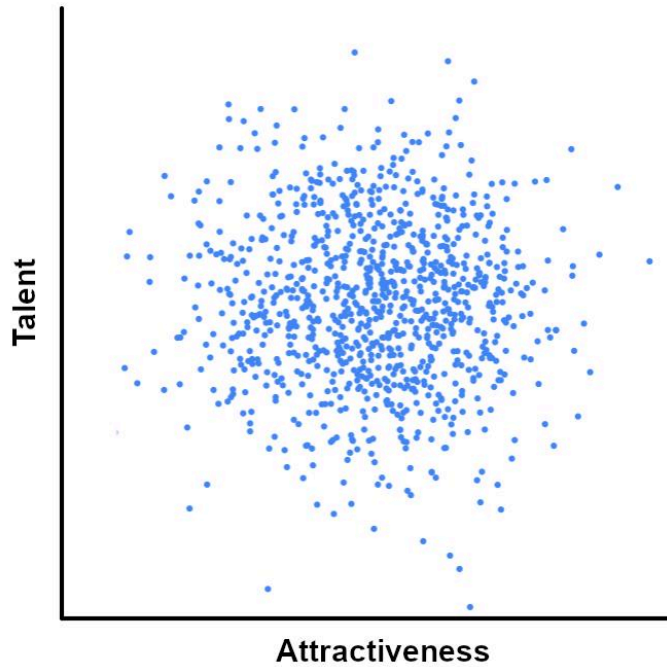


# More examples about sampling bias

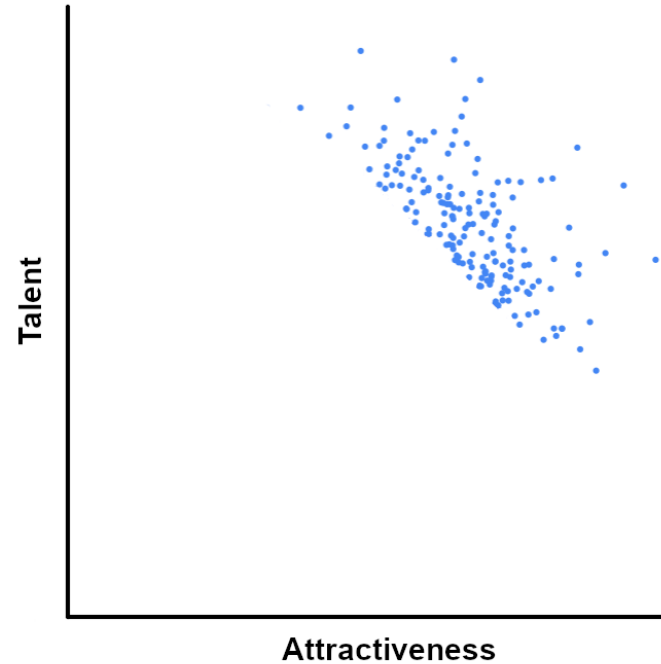
- A survey conducted at a healthcare provider found that 80% of its visitors were diagnosed with a disease.

# More examples about sampling bias

- Are talent and attractiveness negatively correlated?



Population



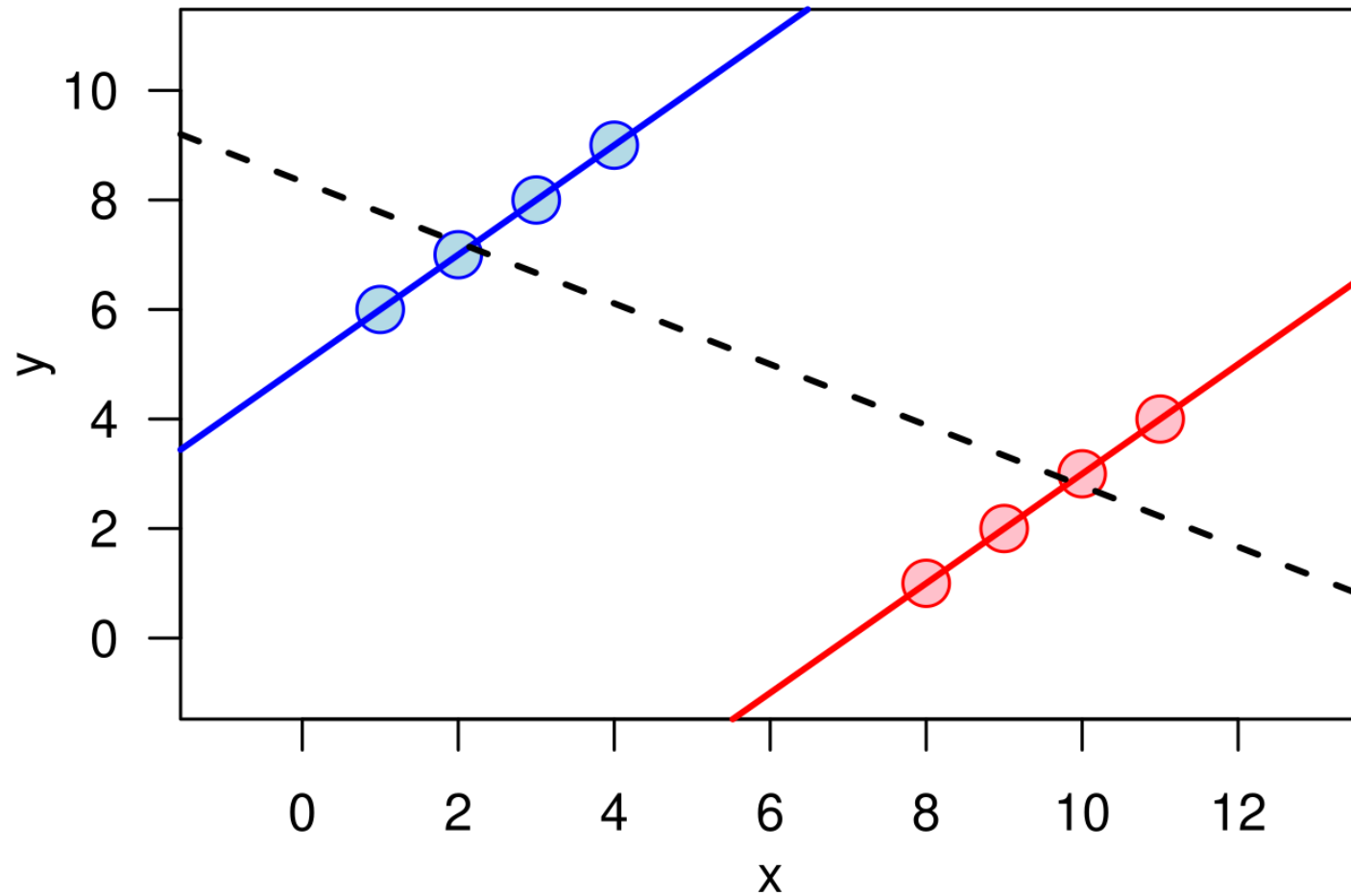
Celebrities

# Simpson's Paradox

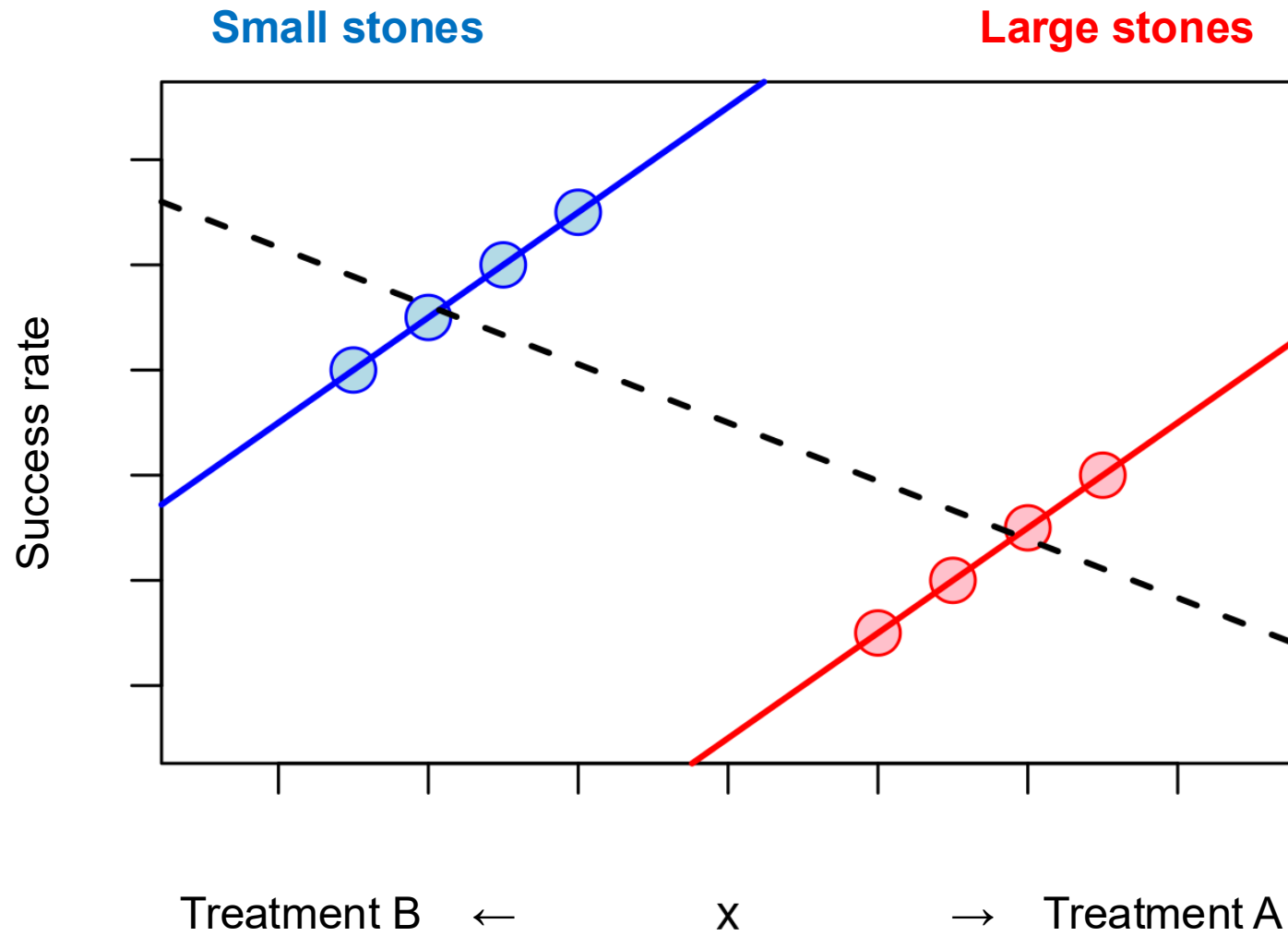
## Kidney stone treatments' success rates

Stone size	Treatment	Treatment A	Treatment B
Small stones		<i>Group 1</i> <b>93% (81/87)</b>	<i>Group 2</i> 87% (234/270)
Large stones		<i>Group 3</i> <b>73% (192/263)</b>	<i>Group 4</i> 69% (55/80)
Both		78% (273/350)	<b>83% (289/350)</b>

# Simpson's Paradox



# Simpson's Paradox



# Simpson's Paradox: An scRNA-seq example

Proportion	Pre-treatment	Post-treatment
Subpopulation A	0.04	0.80
Subpopulation B	0.16	0.16
Subpopulation C	0.80	0.04
Total	1.00	1.00

Gene X expression	Pre-treatment	Post-treatment	Log2 Fold Change
Subpopulation A	0.10	0.30	+1.58
Subpopulation B	1.50	1.80	+0.26
Subpopulation C	3.00	3.50	+0.22
Population Average	2.64	0.67	-1.98

# 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 - $\text{LR} - = \frac{\text{TNR}}{\text{FNR}}$	



# Example: Rare disease screening

Suppose:

- Disease prevalence: **1 in 10,000** (0.01%)
- Test sensitivity: **99%** (correctly detects 99% of cases)
- Test specificity: **99%** (correctly rules out 99% of non-cases)

Now test **1,000,000** people:

- **True cases** = 100
  - True positives = 99 (99% of 100)
  - False negatives = 1
- **Non-cases** = 999,900
  - False positives = 9,999 (1% of 999,900)
  - True negatives = 989,901
- So total positives reported by the test = **99 + 9,999 = 10,098**.  
Only **99** of those are real.
- The **Positive Predictive Value (PPV)** =  $99 / 10,098 \approx 0.98\%$ .  
That means **99.02% of the “positive” results are false alarms**.

# Why screening does not work well for rare diseases with imperfect tests?

- **Key issue:** Even if a test has "good" accuracy (say, 99% sensitivity and 99% specificity), when the disease is rare, most positive results will actually be **false positives** rather than **true positives**.
- This is because the **prevalence** (base rate) of the disease is very low, so the number of healthy individuals vastly outnumbers the true cases.
- **Example: COVID antibody tests** in the early stage of the COVID pandemic: When prevalence was  $<5\%$  in most populations, even tests with 95% specificity yielded more false positives than true positives.

# The Statistics Behind It

- The key relationship is given by **Bayes' theorem**:

$$PPV = \frac{(\text{sensitivity}) \times (\text{prevalence})}{(\text{sensitivity} \times \text{prevalence}) + (1 - \text{specificity}) \times (1 - \text{prevalence})}$$

- When prevalence is very small, the denominator is dominated by false positives (the  $(1 - \text{specificity}) \times (1 - \text{prevalence})$  term).
- This drives PPV close to zero, even for high-quality tests.

# Summary

- Population-wide screening for rare diseases with tests of “ordinary” accuracy does not work because the **false positives overwhelm the true positives**.
- Instead, targeted screening of **higher-risk subgroups** (increasing effective prevalence) may dramatically improve predictive value.

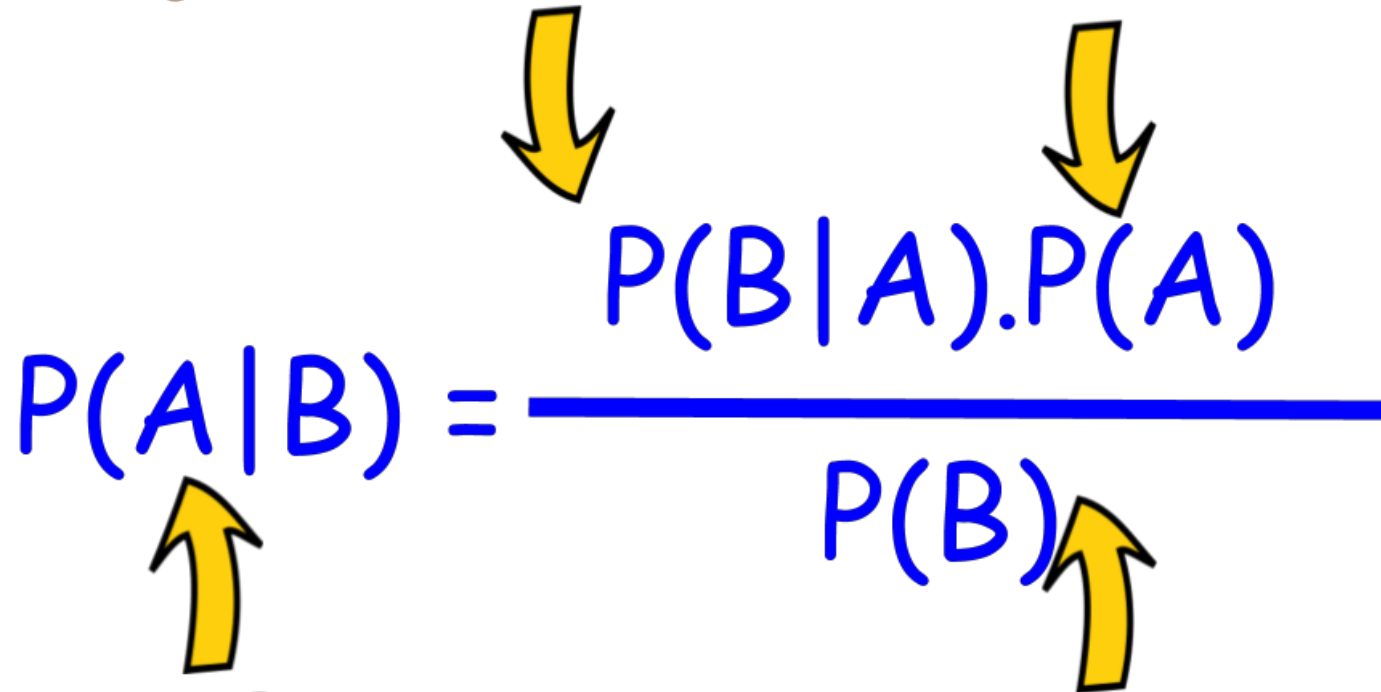
# Bayes' Theorem

## LIKELIHOOD

The probability of "B" being True, given "A" is True

## PRIOR

The probability "A" being True. This is the knowledge.



The diagram shows the Bayes' Theorem formula with four yellow arrows pointing to its components: one from 'LIKELIHOOD' to  $P(B|A)$ , one from 'PRIOR' to  $P(A)$ , one from 'POSTERIOR' to  $P(A|B)$ , and one from 'MARGINALIZATION' to  $P(B)$ .

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

## POSTERIOR

The probability of "A" being True, given "B" is True

## MARGINALIZATION

The probability "B" being True.

# Frequentist vs. Bayesian

## Frequentist

- P-value
- Confidence
- Maximum Likelihood Estimation (MLE)

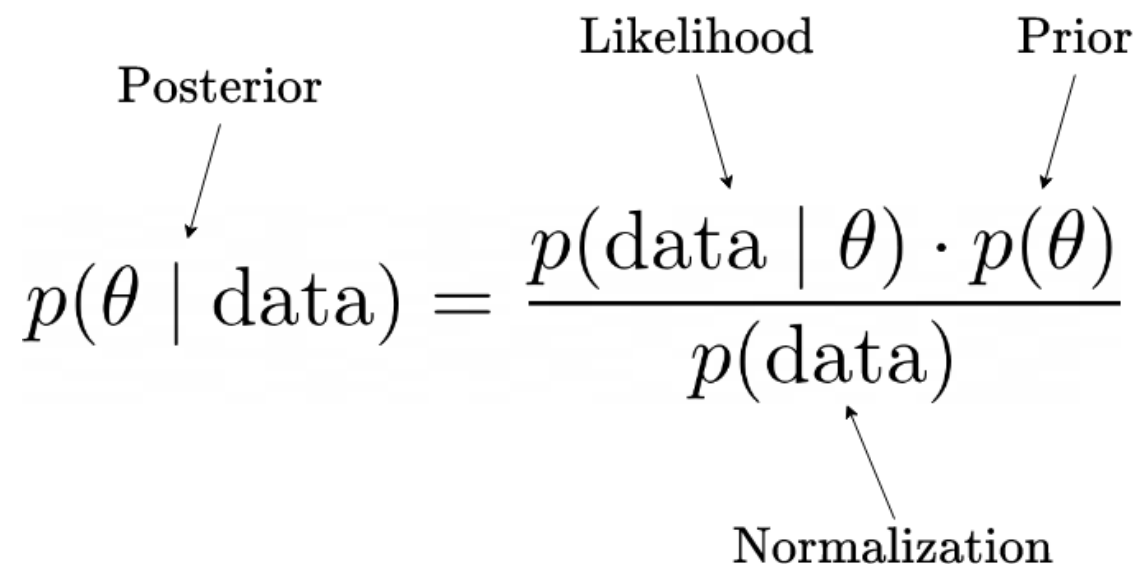
$$\mathcal{L}_n(\theta) = \mathcal{L}_n(\theta; \mathbf{y}) = f_n(\mathbf{y}; \theta)$$

$$\hat{\theta} = \arg \max_{\theta \in \Theta} \mathcal{L}_n(\theta; \mathbf{y})$$

$$p(\text{data}|\theta)$$

## Bayesian

- Bayes' Theorem



The diagram illustrates Bayes' Theorem with arrows pointing from labels to parts of the equation. 'Posterior' points to  $p(\theta | \text{data})$ , 'Likelihood' points to  $p(\text{data} | \theta)$ , 'Prior' points to  $p(\theta)$ , and 'Normalization' points to  $p(\text{data})$ .

$$p(\theta | \text{data}) = \frac{p(\text{data} | \theta) \cdot p(\theta)}{p(\text{data})}$$

# Frequentist vs. Bayesian

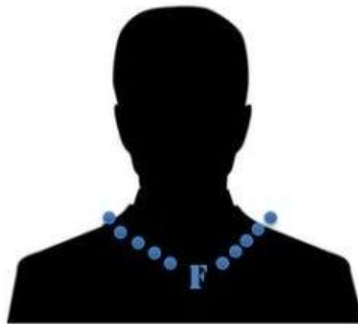
	Frequentist	Bayesian
Hypothesis test	$p$ value (null hypothesis significance test)	Bayes factor
Estimation with uncertainty	maximum likelihood estimate with confidence interval (The "New Statistics")	posterior distribution with highest density interval

# Frequentist vs. Bayesian

%

Probability of the  
**events observed**  
given a theory

**FREQUENTIST**  
**STATISTICS**



%

Probability of the  
**multiple theories**  
given the observed events

**BAYESIAN**  
**STATISTICS**





# A case study

## Article

# Spatial transcriptomics reveal neuron–astrocyte synergy in long-term memory

<https://doi.org/10.1038/s41586-023-07011-6>

Received: 16 March 2023

Accepted: 21 December 2023

Published online: 7 February 2024

Open access



Check for updates

Wenfei Sun<sup>1,2,6</sup>, Zihui Liu<sup>2,3,6</sup>, Xian Jiang<sup>2</sup>, Michelle B. Chen<sup>1</sup>, Hua Dong<sup>4</sup>, Jonathan Liu<sup>5</sup>, Thomas C. Südhof<sup>2,3</sup>✉ & Stephen R. Quake<sup>1,5</sup>✉

Memory encodes past experiences, thereby enabling future plans. The basolateral amygdala is a centre of salience networks that underlie emotional experiences and thus has a key role in long-term fear memory formation<sup>1</sup>. Here we used spatial and single-cell transcriptomics to illuminate the cellular and molecular architecture of the role of the basolateral amygdala in long-term memory. We identified transcriptional signatures in subpopulations of neurons and astrocytes that were memory-specific and persisted for weeks. These transcriptional signatures implicate neuropeptide

# Spatial transcriptomics reveal neuron–astrocyte synergy in long-term memory

<https://doi.org/10.1038/s41586-023-07011-6>

Received: 16 March 2023

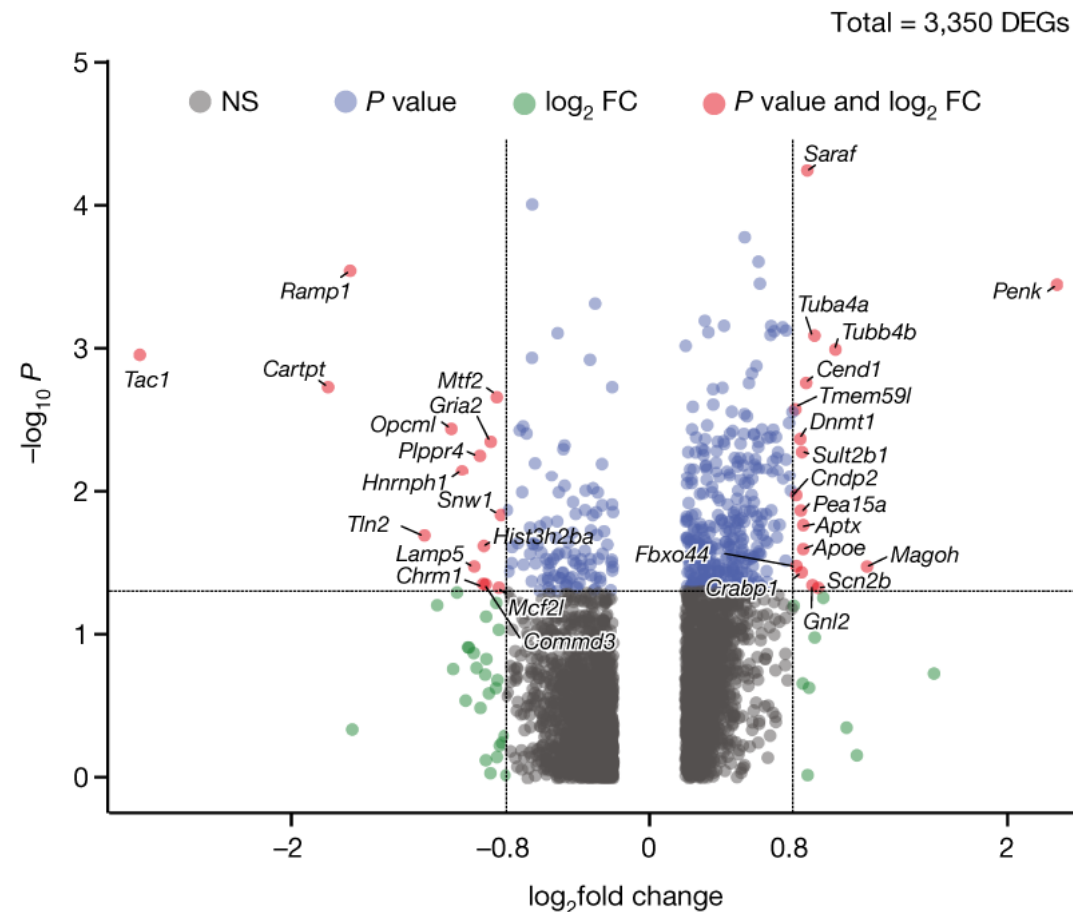
Accepted: 21 December 2023

Published online: 7 February 2024

Open access

 Check for updates

Wenfei Sun<sup>1,2,6</sup>, Zhihui Liu<sup>2,3,6</sup>, Xian Jiang<sup>2</sup>, Michelle B. Chen<sup>1</sup>, Hua Dong<sup>4</sup>, Jonathan Liu<sup>5</sup>,  
Thomas C. Südhof<sup>2,3,✉</sup> & Stephen R. Quake<sup>1,5,✉</sup>



# False positives in study of memory-related gene expression

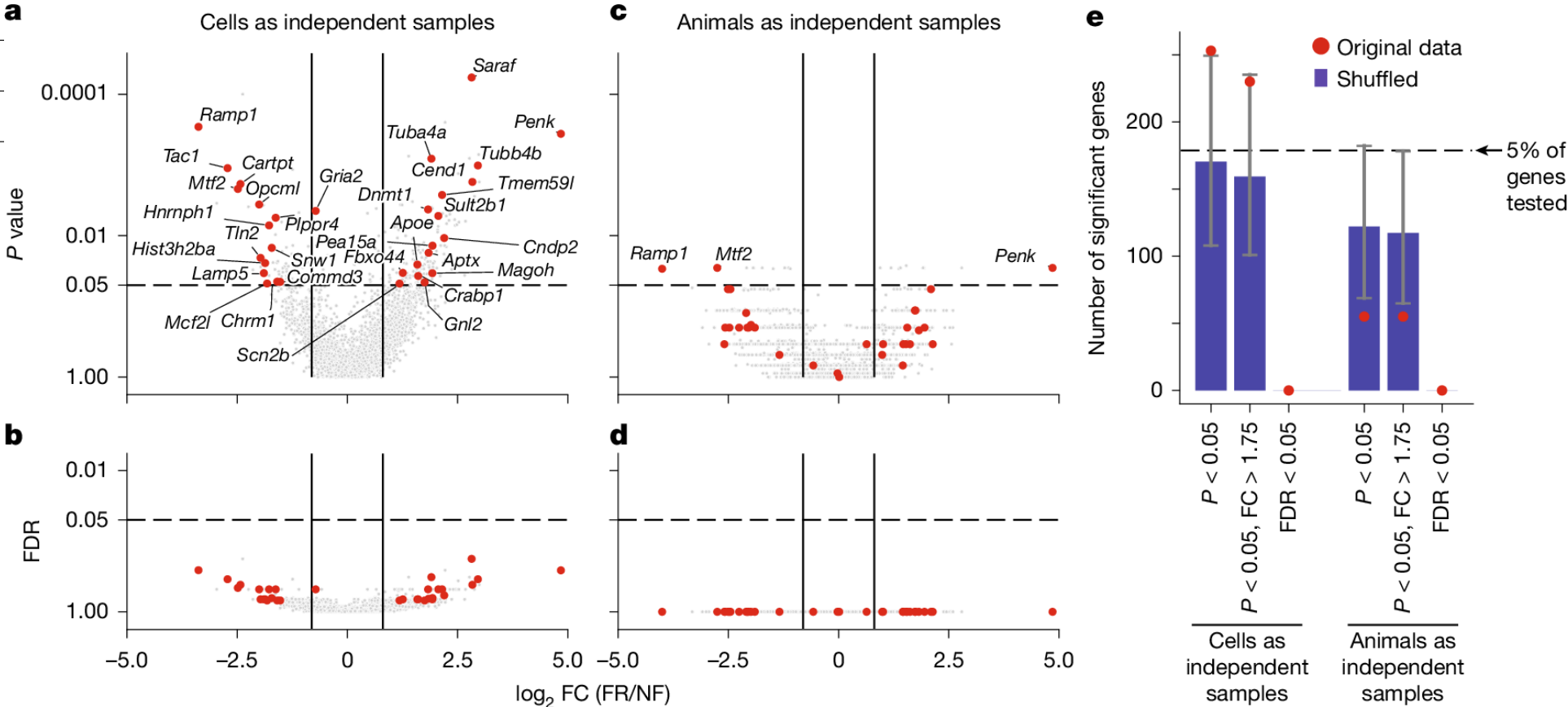
<https://doi.org/10.1038/s41586-025-08988-y>

Eran A. Mukamel<sup>1</sup>✉ & Zhaoxia Yu<sup>2</sup>

Received: 18 July 2024

Accepted: 7 April 2025

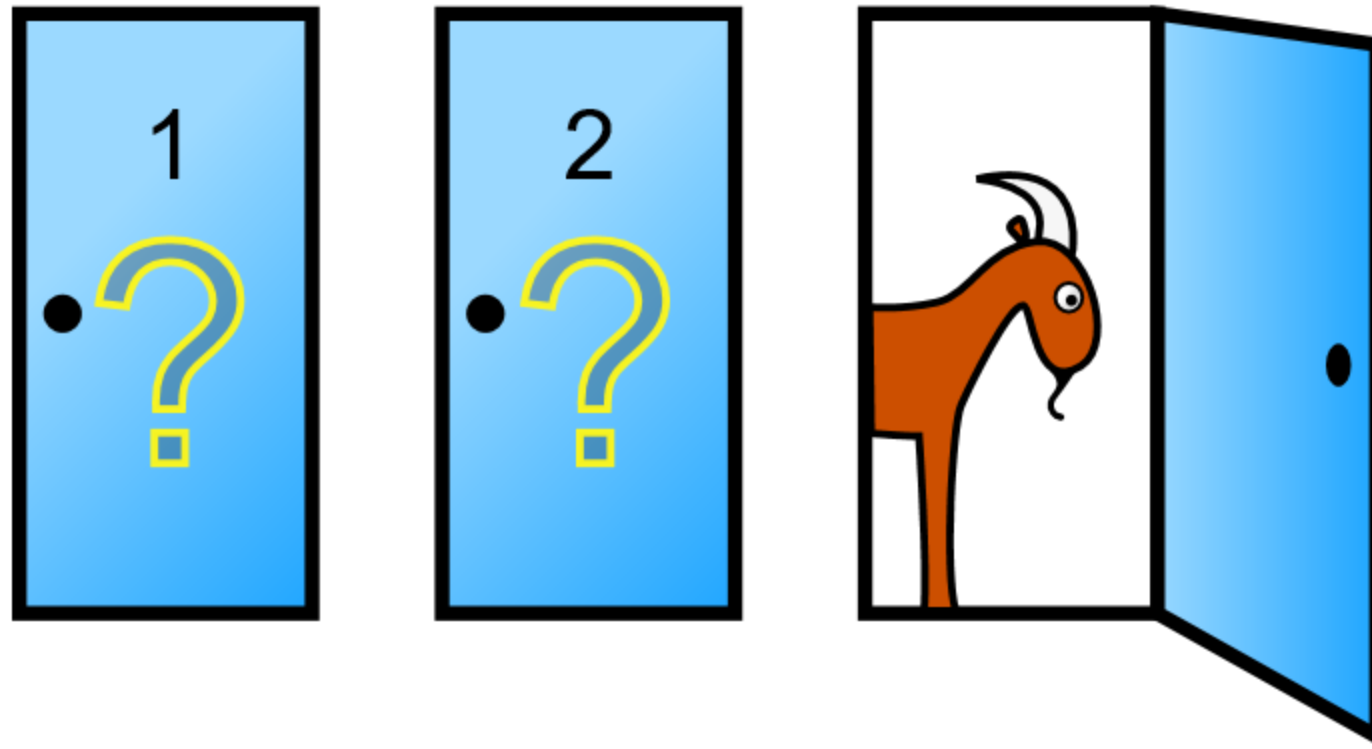
Published online: 4 June 2025



# False positives in study of memory-related gene expression

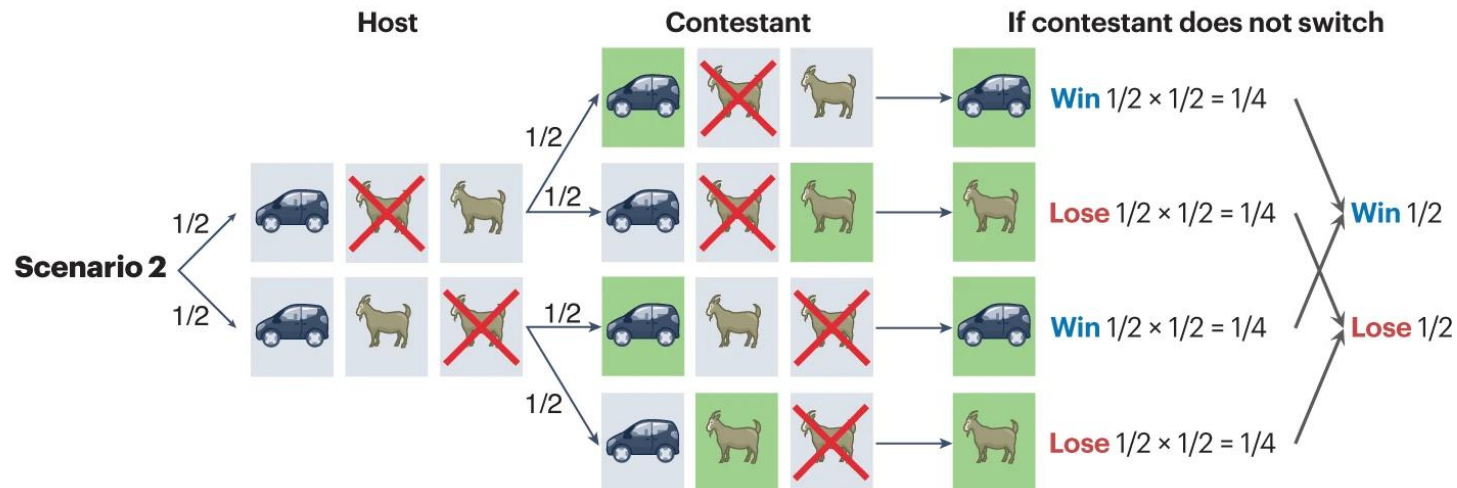
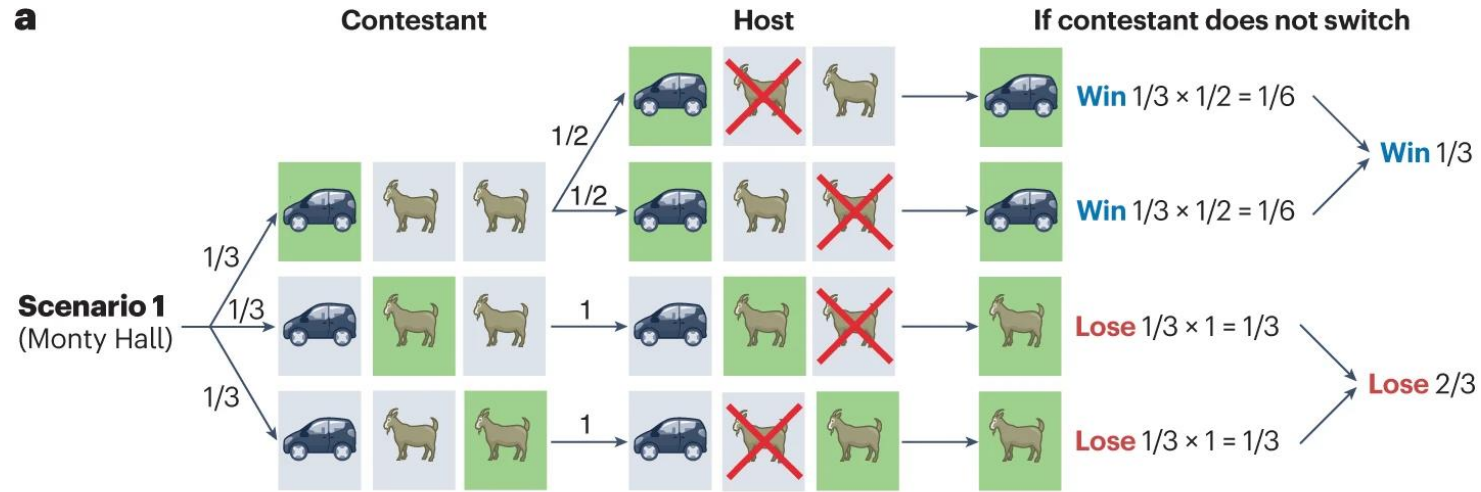
- Multiple testing correction (FDR) was not applied correctly.
  - Sun et al. “used a series of criteria to pre-select 56 candidate genes of interest, thus reducing the burden of multiple hypothesis testing.”
  - This is double-dipping!
- Fail to use animal as sample.
  - “Treatment of individual cells as independent samples.”
  - Cells correlate from the same animal/sample.

# Monty Hall Problem

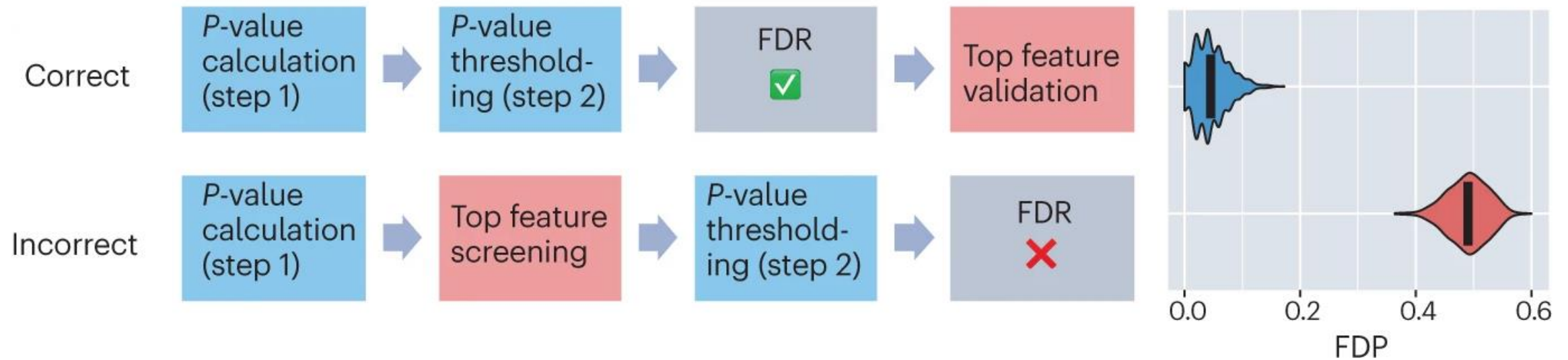


# Monty Hall Problem

**a**



# The order of action matters



# SUMMARY

1. Avoid sampling bias.
2. Carefully plan the study design.
3. Beware Simpson's paradox.
4. Think before you analyze.
5. Statistical analysis is more than a set of computations.



# Ordinary



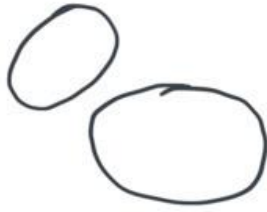
# James-Stein



# HOW TO: DRAW A HORSE

BY VAN OKTOP

---



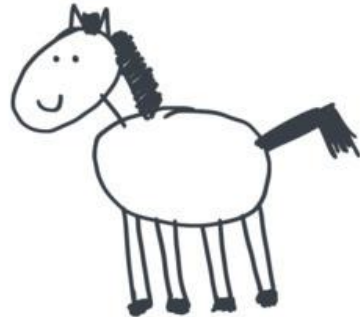
① DRAW 2 CIRCLES



② DRAW THE LEGS

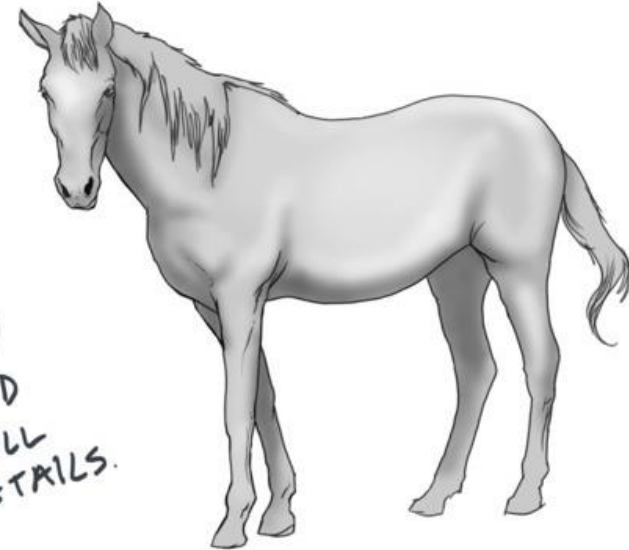


③ DRAW THE FACE



④ DRAW THE HAIR

*Record procedure details!*



⑤  
ADD  
SMALL  
DETAILS.