# **Introduction to Statistics**

iPQB Bootcamp Statistics Module September 4, 2018 Maureen Pittman





### **Statistics in Practice**

https://en.wikipedia.org/wiki/Statistics •

**Statistics** is a branch of mathematics dealing with the collection, organization, analysis, interpretation and presentation of data.

- How many subjects should I include in my experiment?
- Is there a relationship between the independent and dependent variables?
  - Is the relationship causative or associative?
- Comparing two observed states:
  - Are genes differentially expressed between two groups?
  - Are these genes enriched for some biological process?



https://warwick.ac.uk/fac/sci/moac/people/students/peter\_cock/r/heatmap/

#### **Statistical Goals**

Abstract question  $\rightarrow$  mathematical equation

- Inference: Is a drug effective?  $\rightarrow$  Is p < 0.05?
- Prediction: How effective is a drug?  $\rightarrow$  Central tendency and dispersion of treatment results
- Modeling: How does x process work?  $\rightarrow$  What set of parameters/equations best describe it?
- Estimation: How many people in the US are blond?  $\rightarrow$  What % of a sample are blond?

### Essentials of statistical analysis

Probability

Sampling

Variables

Hypothesis Testing

### **Probability**

• Unconditional probability: the likelihood that some event will occur.

• Conditional probability: the likelihood of event A occurring given that B is true.

 $P(A \mid B)$ 

• Bayes Rule:

$$P(A \mid B) = rac{P(B \mid A) \, P(A)}{P(B)}$$



![](_page_6_Picture_0.jpeg)

Biased sample

![](_page_7_Figure_0.jpeg)

#### Simpson's Paradox

#### • Gender bias at UC Berkeley

	Applicants	Admitted
Men	8442	44%
Women	4321	35%

Significance:  $p \approx 10^{-26}$ 

Department	Men		Women	
	Applicants	Admitted	Applicants	Admitted
Α	825	62%	108	82%
В	560	63%	25	68%
С	325	37%	593	34%
D	417	33%	375	35%
E	191	28%	393	24%
F	373	6%	341	7%

"Small but significant bias in favor of women in most departments"

#### Simpson's Paradox

• Gender bias at UC Berkeley

	Applicants	Admitted
Men	8442	44%
Women	4321	35%

Significance:  $p \approx 10^{-26}$ 

Department	Men		Women	
	Applicants	Admitted	Applicants	Admitted
Α	825	62%	108	82%
В	560	63%	25	68%
С	325	37%	593	34%
D	<mark>4</mark> 17	33%	375	35%
E	191	28%	393	24%
F	373	6%	341	7%

"Small but significant bias in favor of women in most departments"

![](_page_9_Figure_6.jpeg)

Variable types

#### Categorical

• Nominal

![](_page_10_Picture_3.jpeg)

categories with no intrinsic order

• Ordinal

![](_page_10_Picture_6.jpeg)

#### Numerical

• Discrete

![](_page_10_Picture_9.jpeg)

Continuous

![](_page_10_Picture_11.jpeg)

infinite possible values

#### **Data Properties**

#### **Central Tendency**

- Mean
- Median
- Mode

![](_page_11_Figure_5.jpeg)

## SD = 0.5 SD = 1 SD = 2

0

2

#### Dispersion

• Variance

-2

-4

• Standard deviation

#### **Multivariate**

- Covariance
- Correlation

![](_page_11_Figure_13.jpeg)

### **Hypothesis Testing**

- $H_0$ : there is no relationship between the two variables
- *H*<sub>1</sub>: the variables are associated

A neurologist is testing the effect of a drug on response time by injecting 100 rats with a unit dose of the drug, subjecting each to neurological stimulus, and recording its response time. The neurologist knows that the mean response time for rats not injected with the drug is 1.2 seconds. The mean of the 100 injected rats' response times is 1.05 seconds with a sample standard deviation of 0.5 seconds.

https://www.khanacademy.org/math/statistics-probability/significance-tests-one-sample/idea-of-significance-tests/e/

A neurologist is testing the effect of a drug on response time by injecting 100 rats with a unit dose of the drug, subjecting each to neurological stimulus, and recording its response time. The neurologist knows that the mean response time for rats not injected with the drug is 1.2 seconds. The mean of the 100 injected rats' response times is 1.05 seconds with a sample standard deviation of 0.5 seconds.

![](_page_13_Figure_1.jpeg)

## **Error Types**

- Type I: the null hypothesis was true, but you rejected it
- Type II: the null hypothesis was false, but you fail to reject it

Table of error types		Null hypothesis ( <i>H</i> <sub>0</sub> ) is		
		True	False	
Decision About Null Hypothesis ( <i>H</i> <sub>0</sub> )	Fail to reject	Correct inference (True Negatives)	Type II error (False Negative)	
	Reject	Type I error (False Positive)	Correct inference (True Positives)	

### Intuitive definitions

- P-value: the probability of seeing a result as extreme (or more extreme) as the one observed, assuming the null hypothesis is true.
- P-value: the probability of making a Type I error (rejecting the null when it is true).
- P-value cutoff (also called α, often set to 0.05): the level of uncertainty you're willing to accept in order to reject the null.

#### **Other kinds of distributions**

![](_page_16_Figure_1.jpeg)

Y = In(X) has a normal distribution

Example: Normalized RNA-seq counts

![](_page_16_Figure_4.jpeg)

Discrete probability distribution of a series of true/false questions

Example: Inheritance of an allele to the next generation

negative binomial distribution, n=10, p=.5

![](_page_16_Figure_8.jpeg)

Binomial distribution of the number of successes before a specified number of failures

Example: overdispersed data (data with high variance)

### **Other kinds of distributions**

Multinomial distribution

![](_page_17_Picture_2.jpeg)

Generalization of the binomial distribution

Example: inheritance of an allele at a triallelic site

![](_page_17_Figure_5.jpeg)

Hypergeometric distribution, n=400, p=.75, k=100

Discrete probability of sampling k times from a group of n, without replacement

Example: modeling the number of activated synapses in ganglion

#### Poisson distribution, lambda=3

![](_page_17_Figure_9.jpeg)

Models how many times an event will occur in an interval given a rate (lambda)

Example: rare events with predictable variance (de novo mutation arising)

http://zoonek2.free.fr/UNIX/48\_R/07.html

### Crimes against data

Confounding variables

Multiple hypotheses

**Publication bias** 

#### **Multiple Hypotheses**

![](_page_19_Figure_1.jpeg)

![](_page_20_Picture_0.jpeg)

### **Correction for Multiple Hypotheses**

- Bonferroni correction: divide α by the number of hypotheses, so that the new cutoff is α/m
- Sidak correction:  $\alpha_{SID} = 1 (1 \alpha)^{\frac{1}{m}}$
- Family-Wise Error Rate (FWER)

![](_page_21_Figure_4.jpeg)

#### **Publication Bias**

- Studies with positive results are more likely to:
  - Be published
  - Be cited by others
  - Produce multiple publications
- Replication crisis

#### **Retraction Watch**

Tracking retractions as a window into the scientific process

![](_page_22_Figure_8.jpeg)

### Summary

- Statistics is math for data analysis
- Populations sampling and bias
- Variables
  - Roles: Explanatory, Response, Extraneous, Confounding
  - Types: Categorical (nominal & ordinal); Numerical (discrete & continuous)
- Data properties
  - Central tendency
  - Dispersion
  - Correlation/association

- Hypothesis testing
  - What does a p-value actually mean?
- Types of distributions
  - Normal, binomial, negative binomial
  - Multivariate, hypergeometric, Poisson
- Pitfalls
  - Confounding variables
  - Multiple hypotheses
  - Publication bias

### Acknowledgements

![](_page_24_Picture_1.jpeg)

![](_page_24_Picture_2.jpeg)

![](_page_24_Picture_3.jpeg)

### Summary

- Statistics is math for data analysis
- Populations sampling and bias
- Variables
  - Roles: Explanatory, Response, Extraneous, Confounding
  - Types: Categorical (nominal & ordinal); Numerical (discrete & continuous)
- Data properties
  - Central tendency
  - Dispersion
  - Correlation/association

- Hypothesis testing
  - What does a p-value actually mean?
- Types of distributions
  - Normal, binomial, negative binomial
  - Multivariate, hypergeometric, Poisson
- Pitfalls
  - Confounding variables
  - Multiple hypotheses
  - Publication bias

#### Questions?