

Introduction to statistics

Computational biology week - London, April 2022

Jesús Urtasun Elizari, PhD



Imperial College
London

Outline

① Probability in a nutshell

- Discrete and continuous probability
- Mean and variance of a distribution
- Exercises in R

② Linear models

- What is a linear model
- Fitting a linear model
- Prediction vs inference
- Exercises in R

③ Hypothesis testing

- Probability distributions
- Statistical tests
- Real genetics problem
- Exercises in R

Chapter I

Probability in a nutshell

Probability in a nutshell

Discrete and continuous probability

- Discrete output \rightarrow probability well defined, compute by counting
- Continuous output \rightarrow need for probability distributions
- Definition of probability in both scenarios
- Computing probabilities in both scenarios

Probability in a nutshell

Discrete probability

Compute probability of getting 4 heads in a row

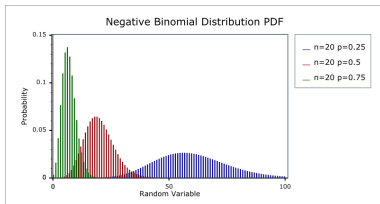
Probability in a nutshell

Discrete probability

Compute probability of getting a specific genotype

Probability in a nutshell

Binomial distribution

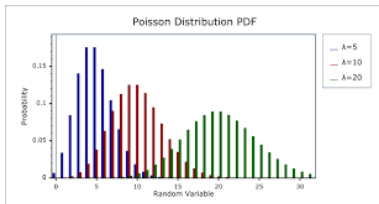


- Binomial distribution: probability of getting k events in n trials

$$B(x = k) = \binom{n}{x} p^k (1 - p)^{n-k}$$

Probability in a nutshell

Poisson distribution



- Poisson distribution: probability of counting k events in a given interval

$$P(x = k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

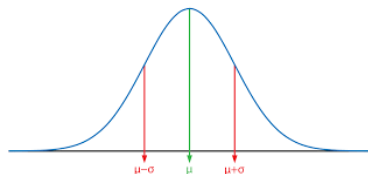
Probability in a nutshell

Continuous probability

Compute probability of measuring a particular height
Used when comparing means

Probability in a nutshell

Gaussian distribution

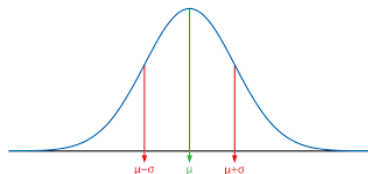


- Gaussian distribution: continuous distribution for a real-valued random variable

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

Probability in a nutshell

Mean and variance



- Momenta of a distribution
- Mean, variance, skewness, kurtosis

Probability in a nutshell

Standard Deviation vs Standard Error

Performing multiple sets of measurements and compare means

- Standard deviation (σ) quantifies variation over one set of measurements (second momentum of the probability distribution)
- Standard error (SE) quantifies **variation of means from multiple sets!**

$$SE = \frac{\sigma}{\sqrt{N}}$$

- Source of confusion: SE can be estimated from single set of measurements, even though it describes variation of means

Probability in a nutshell

Exercises in R

https://lmsbioinformatics.github.io/LMS_StatisticsInR/course/CBW_StatisticsInR_course

Chapter II

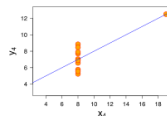
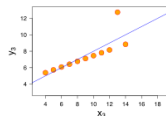
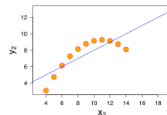
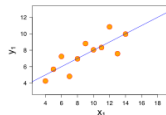
Linear models

Linear models

What is a linear model

- 1 Find a function that describes a set of observations

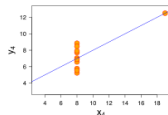
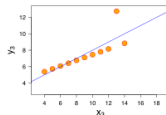
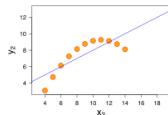
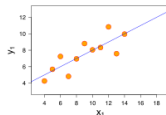
$$y(x) = b_0 \cdot x + b_1$$



Linear models

Fitting a linear model

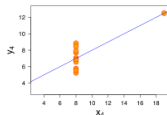
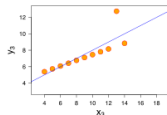
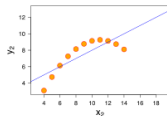
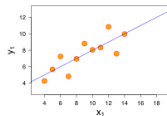
- 1 Get the overall mean
- 2 Compute sum of residuals (SS) over the mean
- 3 Find time with the smallest SS \rightarrow fit
- 4 Evaluate fit - correlation coefficient R^2



Linear models

Generalized linear models

- 1 Fit count data (non-linearly distributed)
- 2 Poisson regression (assume variance = mean)
- 3 Negative binomial

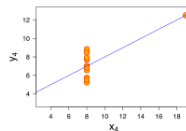
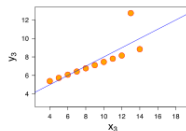
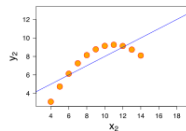
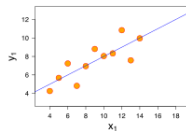


Linear models

Advanced topics

1 ANOVA

2 (...)



Linear models

Exercises in R

https://lmsbioinformatics.github.io/LMS_StatisticsInR/course/CBW_StatisticsInR_course

Chapter III

Hypothesis testing

Hypothesis testing

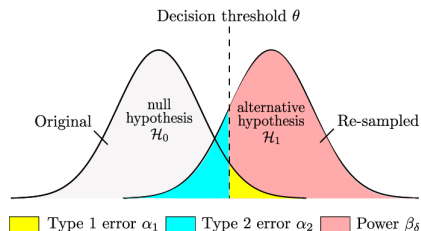
Introduction

- Null hypothesis and alternative hypothesis
- Statistic tests and p-values
- χ^2 -test, t -test, Wald test
- Exercises in R

Hypothesis testing

Compare different models

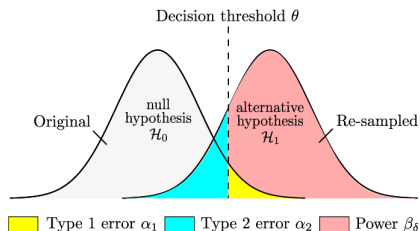
1 Null and alternative hypothesis



Hypothesis testing

Statistic test

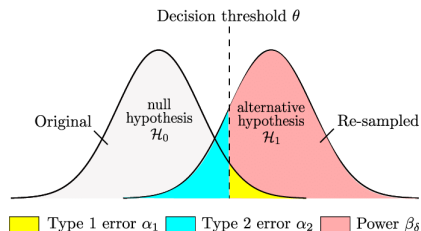
- 1 Quantify the significance of an observation
- 2 Certainty when accepting / rejecting a hypothesis



Hypothesis testing

p-values

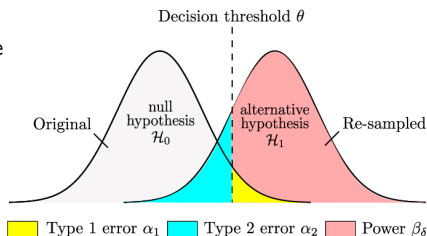
1 p-values



Hypothesis testing

Power calculations

- 1 Do not just add samples until getting a good p-value → increases chance of reporting a false positive
- 2 "Power": probability that a test will correctly give a small p-value
- 3 4 factors affect power (effect size, variation in data, sample size, test used)



Hypothesis testing

Exercises in R

https://lmsbioinformatics.github.io/LMS_StatisticsInR/course/CBW_StatisticsInR_course