



# BIOS14 - Processing and Analysis of Biological Data

---

Per-Erik Isberg, Dept of Statistics



## Introduction

---

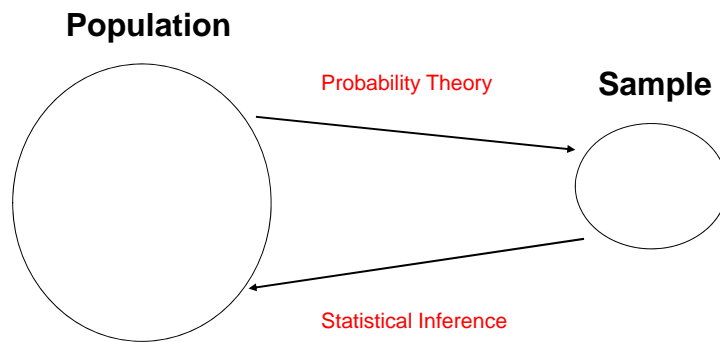
### **Litterature**

Quinn & Keough (2002) – Experimental Design and Data Analysis  
for Biologists, Cambridge University Press (ISBN 0521009766)

$$\Rightarrow y = f(x)$$



## Introduction



## Probability

**Probability  $\Rightarrow$  Probability distributions**

- **Discrete distributions**
  - Binomial
  - Poisson
- **Continuous distributions**
  - Exponential
  - Normal
  - $t$ ,  $\chi^2$ ,  $F$



## Estimation

---

### Example:

We sample  $n=20$  individuals from a population

Calculate the mean in the sample  $\Rightarrow \bar{x}$

Repeat the whole procedure  $\Rightarrow \bar{x}_1, \bar{x}_2, \bar{x}_3, \dots$

**Q:** What distribution?

**A:** Next thursday!



## Estimation

---

### Keywords:

- Confidence intervals
- Standard errors
- Central Limit Theorem

### Methods:

- Maximum Likelihood (ML)
- Ordinary Least Squares (OLS)
- Bootstrap
- Jackknife
- Bayesian



# Hypothesis testing

## Statistical tests

$H_0$ : Null hypothesis (assumed true)

$H_1$ : Alternative hypothesis (one-tailed/two-tailed)

$\alpha$  = significance level of the test (5%)



# Hypothesis testing

## The statistical test

The truth	$H_0$ not rejected	$H_0$ rejected
$H_0$ true	Correct	type I error $\Pr( ) = \alpha = \text{often } 5 \%$
$H_0$ false	type II error $\Pr( ) = \beta$	Correct $\Pr( ) = 1 - \beta = \text{power}$

## Correlation and regression

$$y = f(x_1, x_2, \dots)$$

Numerical variable

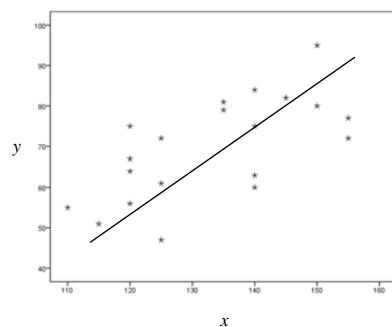
If the x's are numerical  $\Rightarrow$  Regression analysis

If the x's are categorical  $\Rightarrow$  Analysis of Variance (ANOVA)

If we have both types  $\Rightarrow$  Analysis of Covariance (ANCOVA)

## Correlation and regression

Of special interest  $\Rightarrow$  Linear models (OLS)





## Correlation and regression

$$y = a + b \cdot x$$

Diagram illustrating the components of the linear regression equation  $y = a + b \cdot x$ :

- $a$  is labeled as the **intercept**.
- $b$  is labeled as the **slope (regression coefficient)**.
- $y$  is labeled as the **dependent** variable.
- $x$  is labeled as the **independent** variable.

Correlation – measures the strength of the linear relation



## Correlation and regression

### Multiple linear regression

$$y = a + b_1 \cdot x_1 + b_2 \cdot x_2 + \dots + b_p \cdot x_p$$

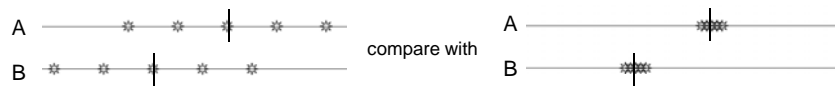
Tests, confidence intervals, predictions, ...

⇒ Non-linear regression

# Analysis of Variance

$$y = f(x)$$

Numerical      Categorical



**Keywords:** between group variation – within group variation

# Analysis of Variance

⇒ ANOVA table

	Sum of Squares	df	Mean Square
Between Groups			
Within Groups			
Total			

⇒ One-way analysis of variance

⇒ if just two groups = t-test (for independent samples)

## Analysis of Variance

Sometimes the groups consists of the same individuals measured at several occasions

⇒ Repeated measurements

⇒ if measured twice = t-test (paired)



**Keywords:** within individual variation

## Analysis of Variance

### Model I - Fix effects

Conclusions are only valid for the observed levels of the factor ⇒ Treatment A, B, C

### Model II - Random effects

The levels used are a sample from a population of levels ⇒ Operator 1, 2, 3, 4

**Mixed models** include both fix and random effects





## Analysis of Variance

---

### **Generalizations**

- Factorial designs
- Nested designs
- Repeated measurements
- General Linear Models (GLM)



## Analysis of Variance

---

### **Assumptions**

- Independence
- Equal variances
- Normality

### **⇒ Non-parametric methods**

- Mann-Whitney
- Kruskal-Wallis
- Kolmogorov-Smirnov



## Analysis of frequencies

### Goodness of fit

- $\chi^2$ -test
- G-test
- Fisher's exact test

### Contingency tables (two-way)

Count

		Disease		Total
		Yes	No	
Exposed to riskfactor	Yes			
	No			
Total				



## Analysis of frequencies

Higher order tables (Three-way and more)

⇒ Log-linear models

**Keywords:** Dependence - Independence

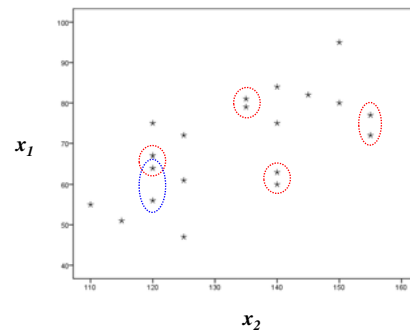
⇒ Logistic regression

$$y = f(x)$$

↑  
Probability (0/1)

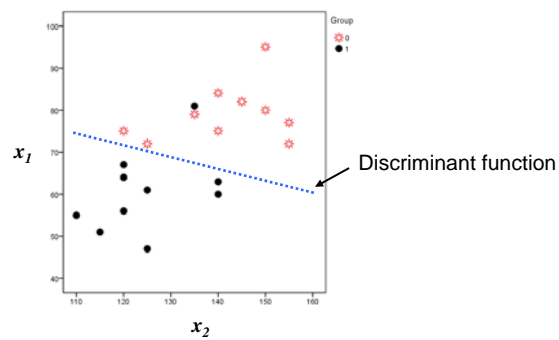
## Multivariate analysis

**Cluster analysis** – Similarity among individuals



## Multivariate analysis

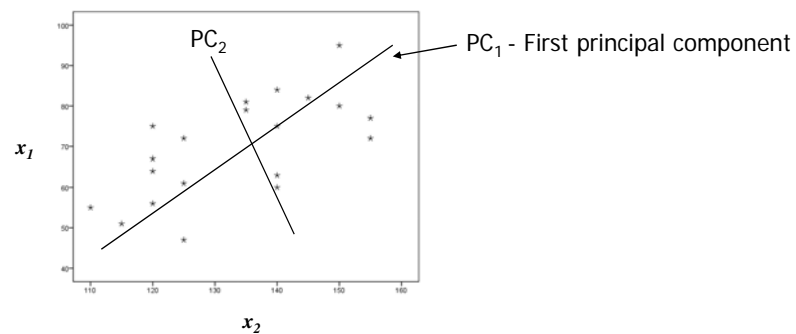
**Discriminant analysis** – Prediction of group membership





## Multivariate analysis

**Principal component** – Summarization of information



## Survival analysis



We have a group of individuals that are "alive" and want to study **when** they move into state "dead".

**Of special interest:**

$T$  = Time for death,  $T \geq 0$



## Survival analysis

---

**Special feature:**

Censored observations – We can't observe the true death time. The individual is still alive at the end of the study.

**Example:** The time for completing a PhD!