

5- Introduction to Statistical Inference

Alex Sanchez, Miriam Mota, Ricardo Gonzalo and
Santiago Perez-Hoyos

Statistics and Bioinformatics Unit. Vall d'Hebron Institut de
Recerca

Readme

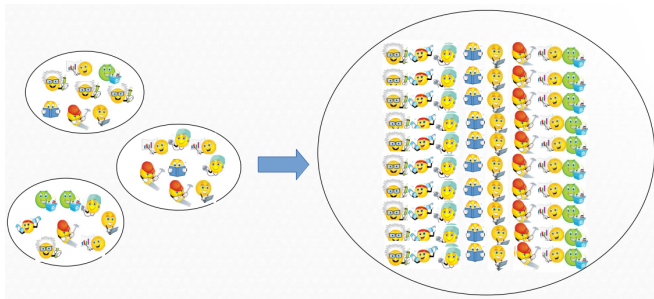
- License: Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License <http://creativecommons.org/licenses/by-nc-sa/4.0/>
- You are free to:
 - **Share** : copy and redistribute the material
 - **Adapt** : rebuild and transform the material
- Under the following conditions:
 - **Attribution** : You must give appropriate credit, provide a link to the license, and indicate if changes were made.
 - **NonCommercial** : You may not use this work for commercial purposes.
 - **Share Alike** : If you remix, transform, or build upon this work, you must distribute your contributions under the same license to this one.

Outline

- The objectives of statistical inference
- Examples
- Point estimation. On incidence and prevalence
- Confidence intervals
- Sample size calculations

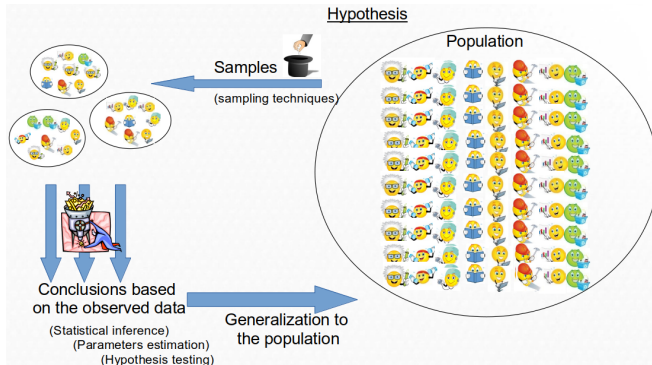
The objectives of Statistical Inference (I)

Taking the observed (measured) values of a group of samples. . .



we aim at determining the properties of the entire population.

The objectives of Statistical Inference (II)



Example

- Consider the data in the “osteoporosis.csv” dataset.
- It can be useful to provide information such as:
 - The percentage of menopausal women with osteoporosis
 - The mean bone density in menopausal or non-menopausal women
 - The existence of significant differences:
 - Observed % of osteoporosis vs “theoretical” population values
 - BUA in menopausal vs non-menopausal
- Answering these questions (and questions like these) is the main goal of Statistical Inference

Two types of statistical inference problems

- ESTIMATION

- When we wish to *learn some characteristics of our population*, such as
 - The percentage of non osteopenic or menopausal women
 - The mean bone density in each of these groups

- HYPOTHESIS TESTING

- When we wish to *check about some statement on some characteristic of the population* or we wish to make some *comparisons*
 - Is it true that the mean bone density is smaller than 75 in menopausal
 - Can we state that non menopausal women have a higher bone density than menopausal?

Estimators: Aproximating the value of population parameters

- Numerical values calculated on a sample that we believe to be a good approximation of a certain real value (parameter) in the population.
- Intuitively, we work with many estimators, such as the mean or a computed percentage of a given sample, that we assume that are somehow characterizing a population.
- It is **not always obvious to decide which is the best estimator for each parameter**
- In order to decide which estimator we use we can rely on the *properties* of the estimators such as **the bias** or the **precision (the variance)** of the estimator.

Example. Computing estimations (1)

- Read the Osteoporosis dataset and turn factors into variables automatically with Rbase function `read.delim`
- Take a sample of size 100 from the original file. Call it 'osteol00' and work with this file from now on.
- Compute the mean value of the variable containing bone density values BUA
- Split the computation between all subgroups from variable `classific` and variable `menop`
- Compute the percentage of menopausal women from variable `menop`

Example. Computing estimations with R

```
library(dplyr)
```

```
##
```

```
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
```

```
##
```

```
##      filter, lag
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
##      intersect, setdiff, setequal, union
```

```
# Read data
```

```
osteoporosis <- read.delim2("~/Dropbox (Nuevo Equipo VHIR10
```

```
# Take subsample
```

```
osteo100 <- sample_n(osteoporosis, 100)
```

```
# mean bone density
```

```
buaMean <- mean(osteo100$bua)
```

```
print(buaMean)
```

Exercise 2

- Read the diabetes dataset. Convert characters into factors before continuing.
- Provide an estimate of
 - The distribution of a numerical variable.
 - a proportion of at least one categorical variable and
 - the mean value of at least one numerical variable.
- Could you have used different estimators?
- How would you decide?

How precise is an estimator?

- We all are familiar with “forks” associated with voting results.
 - They usually start “wide” and tend to disappear as more votes are counted.
- Imagine you are given an estimate of 18% for the incidence of a certain disease.
- Is it a good estimate?
- Hard to know without more information
 - 18 ± 2 is probably useful
 - 18 ± 12 is probably too wide to be considered useful
- So given an estimator and a n estimation (a value) **how can we provide a measure of how precise this estimation is?**

The *Standard Error* of an estimator

- An obvious question when we choose an estimator is *how precise it is to approximate the value of the population parameter*.
- This can be answered using the **standard error of the estimator**
- The standard error is a great quantity :
 - It informs about the *precision* of our estimates
 - Helps build another type of estimators: *confidence intervals*
 - Helps find formulae to compute *sample size* for estimation

Some standard errors

- Standard error of the sample mean

$$SEM = \frac{\hat{s}}{\sqrt{n}}$$

- Standard error of the sample proportion

$$SEP = \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

Computing the standard error with R

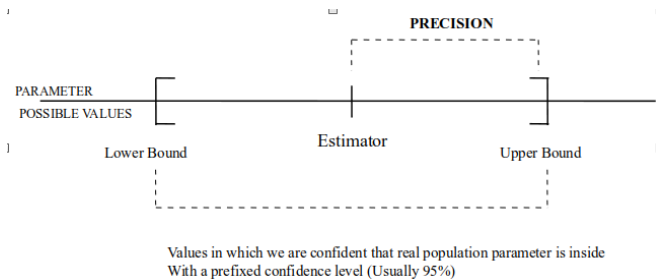
- R does not usually include a function for standard errors, although it can be easily programmed.

```
SEM <- function (x){sd(x)/sqrt(length(x))}
```

```
SEP <- function (x){  
  ssize <- length(x)  
  p <- sum(x)/ssize  
  return(sqrt(p*(1-p)/ssize))  
}
```

Confidence intervals

- Confidence intervals are based on standard errors



Formulae for confidence intervals

- Confidence interval for the mean

$$\underbrace{\bar{X} - t_{\epsilon/2} \frac{\hat{s}}{\sqrt{n}}}_{\text{Precision}} \leq \mu \leq \bar{X} + t_{\epsilon/2} \frac{\hat{s}}{\sqrt{n}} = \bar{\mathbf{X}} \pm \mathbf{t}_{\epsilon/2} \cdot \text{SEM}$$

- Confidence interval for the proportion

$$\underbrace{\hat{p} - z_{\epsilon/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}}_{\text{Precision}} \leq \mu \leq \hat{p} + z_{\epsilon/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}} = \hat{\mathbf{p}} \pm \mathbf{z}_{\epsilon/2} \cdot \text{SEM}$$

Example. Computing Confidence Intervals with R

- In general R does not compute (has no functions) for the direct calculation of confidence intervals
- This can be done by calling the corresponding tests functions such as `t.test` or `prop.test`
- Some R commander plugins such as EZR allow this computations directly

Example. Computing Confidence Intervals with R (2)

```
t.test(osteo100[["bua"]])  
##  
## One Sample t-test  
##  
## data:  osteo100[["bua"]]  
## t = 42.137, df = 99, p-value < 2.2e-16  
## alternative hypothesis: true mean is not equal to 0  
## 95 percent confidence interval:  
##  68.77151 75.56849  
## sample estimates:  
## mean of x  
##      72.17
```

Example. Computing Confidence Intervals with R (3)

```
cntMenop <- table(osteo100[["menop"]])["SI"]
ssize <- length(osteo100[["menop"]])
prop.test (x=cntMenop, n=ssize)

##
## 1-sample proportions test with continuity correction
##
## data:  cntMenop out of ssize, null probability 0.5
## X-squared = 13.69, df = 1, p-value = 0.0002156
## alternative hypothesis: true p is not equal to 0.5
## 95 percent confidence interval:
##  0.5885509 0.7766330
## sample estimates:
##      p
## 0.69
```

Sample Size for estimation (1)

- The standard error informs of how precise an estimation is **if one knows the variability and the sample size**

$$SE = \frac{\hat{\sigma}}{\sqrt{n}}$$

- We can proceed in the opposite sense: assuming we know:
 - ① the variability (e.g. from a pilot study) and
 - ② the highest precision we wish to attain (“arm length” of a confidence interval:

$$\Delta = z_{\epsilon_2} \cdot SE = z_{\epsilon_2} \cdot \frac{\hat{\sigma}}{\sqrt{n}}$$

Sample Size for estimation (2)

- The sample size needed to attain this precision can be isolated from the previous equation:

$$n = \frac{z_{\epsilon_2}^2 \hat{\sigma}^2}{\Delta^2}$$

Sample size formulae for estimating a mean or a proportion

The previous formula becomes, for specific questions:

$$n = \frac{t_{n-1, \epsilon_2}^2 \hat{s}^2}{\Delta^2} \quad (1), \quad n = \frac{z_{\epsilon_2}^2 \hat{p}(1 - \hat{p})}{\Delta^2} \quad (2), \quad n = \frac{z_{\epsilon_2}^2}{4 \Delta^2} \quad (3)$$

- 1 Mean of a normal population with a given precision Δ .
- 2 Proportion p , with a given precision Δ and with an estimate, \hat{p} available, from a pilot study.
- 3 Proportion p , with a given precision Δ and assuming the *worst case* $p = q = 0.5$.

Sample size calculations with R