







Given that a sample is randomly selected from a source population, conclusions drawn from a sample can always be wrong.

Statistical inference is done with "awareness of inherent limitations":

1) One tries to estimate the probability to be wrong.

2) One tries to limit the probability to be wrong.

"As you are a human being, never state what tomorrow will bring", Semonides of Amorgos.

"State contente umane genti al quia, che se possuto aveste veder tutto, mestier non era parturir Maria", Dante Alighieri Current scientific approach: We cannot find the truth, but we can approach it.

To make statistical inference, we need to recall previous information:

1) The normal distribution, in particular which interval comprises 95% of the distribution.

2) The sampling distribution of the mean.







Point estimate consists in a single value. However:

- 1) This value nearly never coincides with the true value (parameter) of the population;
- 2) Different samples yield different point estimates.

Interval estimate consists in an interval:

- 1) This interval has a given probability (usually 95%) to comprise the true value (parameter) of the population;
- 2) Intervals, obtained from different samples, are at least partly superimposed.

DEFINITION of CONFIDENCE INTERVAL

The confidence interval of the population parameter Θ is an interval which has a given probability (1- α) to comprise the true population parameter:

$$p(L_{low} < \Theta < L_{up}) = 1 - \alpha$$

where:

 $L_{low} = lower limit; \quad L_{up} = upper limit$

1- α = level of confidence; α = error probability





















Confidence interval

 $\theta = \mu$

confidence level = 0.95

$$\overline{x} - 1.96 * \sigma / \sqrt{n} < \mu < \overline{x} + 1.96 * \sigma / \sqrt{n}$$

for whatever confidence level = $1 - \alpha$

$$\overline{x}$$
 - $Z_{\alpha/2}$ * $\sigma / \sqrt{n} < \mu < \overline{x} + Z_{\alpha/2}$ * σ / \sqrt{n}

for whatever parameter θ

$$\stackrel{\mathsf{A}}{\theta} - Z_{\alpha/2} * S.E.(\theta) < \theta \ < \theta + Z_{\alpha/2} \ * S.E.(\theta)$$





4th problem: Using the Confidence Interval to plan sample size.

A study aims at computing the prevalence (probability) of asthma in a population. Preliminary data from the current literature suggest that asthma prevalence could be about 6%. Which is the sample size necessary to estimate asthma prevalence with a width of the 95% confidence interval not greater than 2%?

Data: p = 0.06 $1-\alpha = 95\%$ $z_{\alpha/2} = 1.96$ CI width $\le 2\%$ n = ? $(p + z_{\alpha/2} \cdot \sqrt{p(1-p)/n}) - (p - z_{\alpha/2} \cdot \sqrt{p(1-p)/n}) \le \delta$ $2 \cdot z_{\alpha/2} \cdot \sqrt{p(1-p)/n} \le \delta$ dividing by $2 \cdot z_{\alpha/2}$ $\sqrt{p(1-p)/n} \le \delta / (2 \cdot z_{\alpha/2})$ squaring the equation $p(1-p)/n \le \delta^2 / (2 \cdot z_{\alpha/2})^2$ multiplying by n and dividing by the 2nd member $p(1-p) (2 \cdot z_{\alpha/2})^2 / \delta^2 \le n$ $n \ge 0.06^* 0.94^* (2^* 1.96)^2 / 0.02^2$ $n \ge 0.0564^* (3.92)^2 / 0.0004$ $n \ge 0.0564^* 15.36 / 0.0004$ $n \ge 2166.58$ $n \ge 2167$