

Probability and Inferential Statistics

Lecture SS 16

Confidence Interval and Central Limit Theorem

Prof. Dr. Rolf Steyer

Distribution Sample Mean	2
Normal Distribution.....	3
Confidence interval	4
Confidence interval	5
Some t -distributions	6
t -Distribution	7
Confidence interval	8
Confidence interval	9
Confidence interval	10
Central Limit Theorem	11

Distribution of the Sample Mean Assuming a Normal Distribution

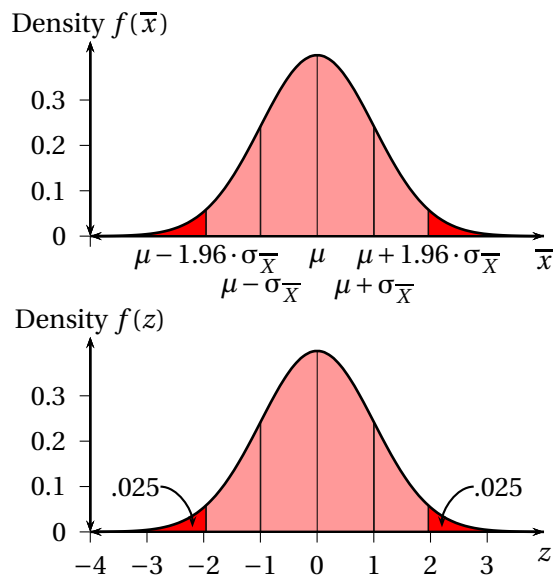
Theorem 1. Let X_1, \dots, X_n be a normally distributed sample of size n with expectation μ and variance σ^2 , that is, a sample with normally distributed random variables $X_i, i = 1, \dots, n$, with expectations $E(X_i) = \mu$ and variance $Var(X_i) = \sigma^2$. Then the sample mean $\bar{X} := (X_1 + \dots + X_n)/n$ is normally distributed as well with expectation μ and variance $\sigma_{\bar{X}}^2 := \sigma^2/n$.

Under the assumptions of this theorem,

$$Z := \frac{(\bar{X} - \mu)}{\sigma_{\bar{X}}} \quad (1)$$

has a *standard normal distribution*. This implies that the 97.5%-quantile ≈ 1.96 of the standard normal distribution can be transformed to 97.5%-quantile $\approx \mu + 1.96 \cdot \sigma_{\bar{X}}$ of the distribution of \bar{X} . Therefore, the probability is $1 - \alpha = .95$ to observe a concrete mean \bar{x} in the interval $\mu \pm 1.96 \cdot \sigma_{\bar{X}}$. Of course, the corresponding implications hold for other quantiles such as the 99.5%-quantile ≈ 2.576 of the standard normal distribution.

Density of a normally distributed sample mean



Confidence interval: Basic idea

Whenever the concrete sample mean \bar{x} is in the interval $\mu \pm 1.96 \cdot \sigma_{\bar{X}}$ (see the light area in the last figure) — and this happens with probability .95 — then the interval

$$\bar{X} \pm 1.96 \cdot \sigma_{\bar{X}} \quad (2)$$

includes μ (see the last figure). This interval is called the 95%-*confidence interval* of the expectation μ .

The lower boundary of this confidence interval is $\bar{X} - 1.96 \cdot \sigma_{\bar{X}}$ and the upper boundary is $\bar{X} + 1.96 \cdot \sigma_{\bar{X}}$. Both boundaries are random variables. Their values depend on the values x_1, \dots, x_n and, therefore, on the concrete mean \bar{x} that realizes if we draw a sample.

Confidence interval if the standard deviation is unknown

If we want to use $\bar{X} \pm 1.96 \cdot \sigma_{\bar{X}}$, then we have to know the true variance $\sigma_{\bar{X}}$ of the sample mean. If this variance is unknown, then we can neither use the normal distribution with expectation μ and standard deviation $\sigma_{\bar{X}}$ nor the standard normal distribution. However, we can estimate the standard error of the sample mean via $\hat{\sigma}_{\bar{X}}$ and consider the random variable

$$T := \frac{\bar{X} - \mu}{\hat{\sigma}_{\bar{X}}}.$$

This random variable has a central t -distribution with $df = n - 1$ degrees of freedom. Therefore, we can use the 97.5%-quantile $t_{n-1}(.975)$ of the distribution of T , and specify the interval

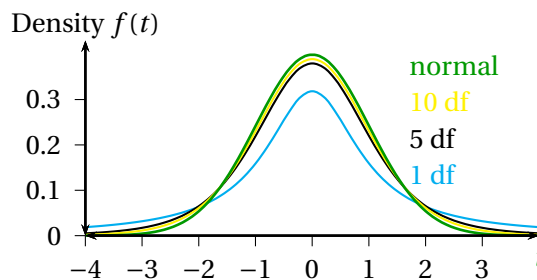
$$\bar{X} \pm t_{df}(.975) \cdot \hat{\sigma}_{\bar{X}}, \quad (3)$$

which includes the expectation μ with probability .95. Note that the distribution of \bar{X} is the normal distribution. Because we do not know $\sigma_{\bar{X}}$, we have to use the t -distribution, which yields a larger confidence interval.

Densities of some t -distribution with $df = 1, 5, 10$ and the standard normal distribution

The t -distribution depends on the degrees of freedom $df = n - 1$, which in turn depends on the sample size. Considering the random variable $T := (\bar{X} - \mu) / \hat{\sigma}_{\bar{X}}$ the degrees of freedom is $n - 1$, because $n - 1$ values of $X_i - \bar{X}$ can realize, if we consider the estimated standard error of the sample mean $\hat{\sigma}_{\bar{X}} = \hat{\sigma} / \sqrt{n}$, where

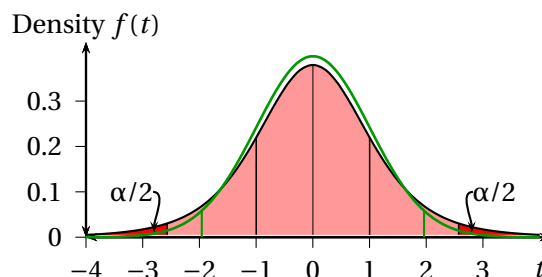
$$\hat{\sigma} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2}.$$



t -Distribution with five degrees of freedom

As mentioned before, the random variable $T := (\bar{X} - \mu) / \hat{\sigma}_{\bar{X}}$ has a t -distribution with $n - 1$ degrees of freedom. Therefore, we observe with probability $1 - \alpha = .95$ a concrete mean \bar{x} in the interval $\mu \pm t_{df}(.975) \cdot \hat{\sigma}_{\bar{X}}$, where $t_{df}(.975) \approx 2.57$ is the 97.5%-quantile of the t -distribution.

If we have a normally distributed sample of size $n = 6$, then we observe a concrete sample mean \bar{x} in the interval $\mu \pm 2.57 \cdot \hat{\sigma}_{\bar{X}}$ (see the light area in the figure below) with probability .95. Therefore, the interval $\bar{X} \pm 2.57 \cdot \hat{\sigma}_{\bar{X}}$ includes the expectation μ with probability .95 as well. If we have a normally distributed sample with unknown variance of size $n = 6$, then this interval is called the 95%-confidence interval of the expectation μ . For sample size n , the confidence interval is: $\bar{X} \pm t_{df}(.975) \cdot \hat{\sigma}_{\bar{X}}$.



Confidence interval: Difference variable

Equation (3) for the confidence interval can also be applied to a difference variable $Y_1 - Y_2$. For example, $Y_1 - Y_2$ can represent the difference between a pretest and a post-test. We might then be interested in the true mean change, for example, after a psychological intervention. If we assume a normally distributed sample $D_1 := Y_{11} - Y_{12}, \dots, D_n := Y_{n1} - Y_{n2}$, then the interval

$$\bar{D} \pm t_{df}(1 - \alpha/2) \cdot \hat{\sigma}_{\bar{D}} \quad (4)$$

includes the expectation $\mu_D := E(D_i) = E(Y_{i1} - Y_{i2})$ with probability $1 - \alpha$. Under the assumptions (normally distributed sample), the interval (4) is the $100 \cdot (1 - \alpha)\%$ -confidence interval of the expectation μ_D .

Confidence interval: Mean difference of two independent samples

Theorem 2. Let X_{11}, \dots, X_{n_1} and X_{12}, \dots, X_{n_2} be two *independent* normally distributed samples of size n_1 and n_2 with expectations μ_1 and μ_2 , respectively, and variance $\sigma_1^2 = \sigma_2^2 := \sigma^2$. According to this assumption, *all* $n_1 + n_2$ random variables $X_{11}, \dots, X_{n_1}, X_{12}, \dots, X_{n_2}$ are independent. Then

$$T := \frac{\bar{X}_1 - \bar{X}_2 - (\mu_1 - \mu_2)}{\hat{\sigma}_{\bar{X}_1 - \bar{X}_2}} \quad (5)$$

has a t -distribution with $df = n_1 + n_2 - 2$ degrees of freedom. In this equation

$$\hat{\sigma}_{\bar{X}_1 - \bar{X}_2} := \sqrt{\left(\frac{1}{n_1} + \frac{1}{n_2}\right) \cdot \frac{n_1 S_1^2 + n_2 S_2^2}{n_1 + n_2 - 2}}, \quad (6)$$

where

$$S_1^2 = \frac{1}{n_1} \sum_{i=1}^{n_1} (X_{i1} - \bar{X}_1)^2 \quad \text{and} \quad S_2^2 = \frac{1}{n_2} \sum_{j=1}^{n_2} (X_{j2} - \bar{X}_2)^2 \quad (7)$$

are the two sample variances. Therefore,

$$(\bar{X}_1 - \bar{X}_2) \pm t_{df}(1 - \alpha/2) \cdot \hat{\sigma}_{\bar{X}_1 - \bar{X}_2} \quad (8)$$

is the $100 \cdot (1 - \alpha)\%$ -confidence interval of the difference $\mu_1 - \mu_2$ of expectations, where $df = n_1 + n_2 - 2$.

General Definition of a Confidence Interval

Definition 1. Let X_1, \dots, X_n be sample and θ a parameter of the distributions of the random variables X_i . Furthermore, let $\mathbf{X} := (X_1, \dots, X_n)$ be the vector of the random variables X_1, \dots, X_n and $f_1(\mathbf{X}), f_2(\mathbf{X})$ numerical functions of \mathbf{X} . If

$$P(f_1(\mathbf{X}) \leq \theta \leq f_2(\mathbf{X})) = 1 - \alpha, \quad 0 < \alpha < 1,$$

the interval $[f_1(\mathbf{X}), f_2(\mathbf{X})]$ is called *the* $100 \cdot (1 - \alpha)\%$ -confidence interval of θ .

Under the standard normal distribution and the t -distribution the confidence intervals are always symmetric. Under other distributions, this is not necessarily the case. An example is the binomial distribution. There we may ask for the interval including (e.g., with probability .95) the unknown probability p . In this case, the confidence intervals are not symmetric any more around the relative frequency.

Central Limit Theorem for the Sample Mean

Theorem 3. Let X_1, X_2, \dots be a sequence of independent identically distributed (i. i. d.) random variables with finite expectation $E(X_i) := \mu$ and finite variance $Var(X_i) := \sigma^2$. Furthermore, let $\bar{X}_n := (X_1 + \dots + X_n)/n$ denote the mean of the first n random variables X_i . Then the distribution of the standardized sample mean $(\bar{X}_n - \mu)/\sigma_{\bar{X}_n}$ converges to the standard normal distribution for $n \rightarrow \infty$. Hence, if $\Phi(z)$ denotes the distribution function of standard normal distribution, then

$$\lim_{n \rightarrow \infty} P\left(\frac{\bar{X}_n - \mu}{\sigma_{\bar{X}_n}} \leq z\right) = \Phi(z).$$

In Theorem 3 we do not assume a distribution of a specific kind. Hence, the random variables X_i , $i = 1, \dots, n$ can have a binomial distribution, a uniform distribution, a normal distribution, or any other kind of distribution. The theorem also applies even if the sample has a bimodal distribution. The only assumption is that X_1, \dots, X_n is a sample, that is, the random variables X_i are i. i. d. However, for a bimodal distribution, the convergence of the distribution of the standardized sample mean is slower than its convergence if the X_i have a normal distribution.