

1 Statistics

Linear combinations of random variables

Continuous random variables

A continuous random variable X has a pdf f such that:

1. $f(x) \geq 0 \forall x$
2. $\int_{-\infty}^{\infty} f(x) dx = 1$

$$\Pr(X \leq c) = \int_{-\infty}^c f(x) dx$$

Expectation theorems

For some non-linear function g , the expected value $E(g(X))$ is not equal to $g(E(X))$.

$$E(X^n) = \sum x^n \cdot p(x) \quad (\text{non-linear function})$$

$$\neq [E(X)]^n$$

$$E(aX \pm b) = aE(X) \pm b \quad (\text{linear function})$$

$$E(b) = b \quad (\text{for constant } b \in \mathbb{R})$$

$$E(X + Y) = E(X) + E(Y) \quad (\text{for two random variables})$$

Variance theorems

$$\text{Var}(aX \pm bY \pm c) = a^2 \text{Var}(X) + b^2 \text{Var}(Y)$$

Linear functions $X \rightarrow aX + b$

$$\begin{aligned} \Pr(Y \leq y) &= \Pr(aX + b \leq y) \\ &= \Pr\left(X \leq \frac{y-b}{a}\right) \\ &= \int_{-\infty}^{\frac{y-b}{a}} f(x) dx \end{aligned}$$

$$\text{Mean:} \quad E(aX + b) = aE(X) + b$$

$$\text{Variance:} \quad \text{Var}(aX + b) = a^2 \text{Var}(X)$$

Linear combination of two random variables

$$\text{Mean:} \quad E(aX + bY) = aE(X) + bE(Y)$$

$$\text{Variance:} \quad \text{Var}(aX + bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) \quad (\text{if } X \text{ and } Y \text{ are independent})$$

Sample mean

Approximation of the **population mean** determined experimentally.

$$\bar{x} = \frac{\Sigma x}{n}$$

where

n is the size of the sample (number of sample points)

x is the value of a sample point

On CAS

1. Spreadsheet
2. In cell A1: `mean(randNorm(sd, mean, sample size))`
3. Edit → Fill → Fill Range
4. Input range as A1:An where n is the number of samples
5. Graph → Histogram

Sample size of n

$$\bar{X} = \sum_{i=1}^n \frac{x_i}{n} = \frac{\Sigma x}{n}$$

Sample mean is distributed with mean μ and sd $\frac{\sigma}{\sqrt{n}}$ (approaches these values for increasing sample size n).

For a new distribution with mean of n trials, $E(X') = E(X)$, $sd(X') = \frac{sd(X)}{\sqrt{n}}$

On CAS

- Spreadsheet → Catalog → `randNorm(sd, mean, n)` where n is the number of samples. Show histogram with Histogram key in top left
- To calculate parameters of a dataset: Calc → One-variable

Normal distributions

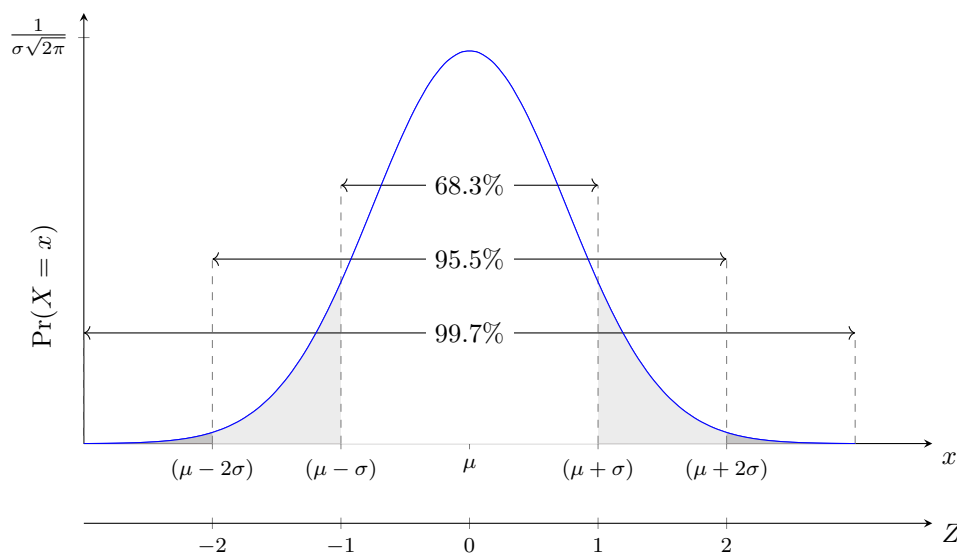
$$Z = \frac{X - \mu}{\sigma}$$

Normal distributions must have area (total prob.) of 1 $\implies \int_{-\infty}^{\infty} f(x) dx = 1$

mean = mode = median

Central limit theorem

If X is randomly distributed with mean μ and sd σ , then with an adequate sample size n the distribution of the sample mean \bar{X} is approximately normal with mean $E(\bar{X})$ and $sd(\bar{X}) = \frac{\sigma}{\sqrt{n}}$.



Confidence intervals

- **Point estimate:** single-valued estimate of the population mean from the value of the sample mean \bar{x}
- **Interval estimate:** confidence interval for population mean μ
- $C\%$ confidence interval $\implies C\%$ of samples will contain population mean μ

95% confidence interval

For 95% c.i. of population mean μ :

$$x \in \left(\bar{x} \pm 1.96 \frac{\sigma}{\sqrt{n}} \right)$$

where:

\bar{x} is the sample mean

σ is the population sd

n is the sample size from which \bar{x} was calculated

Always express z as +ve. Express confidence interval as ordered pair.

On CAS

Menu \rightarrow Stats \rightarrow Calc \rightarrow Interval

Set *Type* = *One-Sample Z Int*

and select *Variable*

Margin of error

For 95% confidence interval of μ :

$$M = 1.96 \times \frac{\sigma}{\sqrt{n}}$$

$$\implies n = \left(\frac{1.96\sigma}{M} \right)^2$$

General case

For $C\%$ c.i. of population mean μ :

$$x \in \left(\bar{x} \pm k \frac{\sigma}{\sqrt{n}} \right)$$

where k is such that $\Pr(-k < Z < k) = \frac{C}{100}$

Confidence interval for multiple trials

For a set of n confidence intervals (samples), there is 0.95^n chance that all n intervals contain the population mean μ .

2 Hypothesis testing

Note hypotheses are always expressed in terms of population parameters

Null hypothesis H_0

Sample drawn from population has same mean as control population, and any difference can be explained by sample variations.

Alternative hypothesis H_1

Amount of variation from control is significant, despite standard sample variations.

p -value

Probability of observing a value of the sample statistic as significant as the one observed, assuming null hypothesis is true.

Distribution of sample mean

If $X \sim N(\mu, \sigma)$, then distribution of sample mean \bar{X} is also normal with $\bar{X} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$.

Statistical significance

Significance level is denoted by α .

If $p < \alpha$, null hypothesis is **rejected**

If $p > \alpha$, null hypothesis is **accepted**

z -test

Hypothesis test for a mean of a sample drawn from a normally distributed population with a known standard deviation.

On CAS

Menu → Statistics → Calc → Test.

Select *One-Sample Z-Test* and *Variable*, then input:

- μ cond: same operator as H_1
- μ_0 : expected sample mean (null hypothesis)
- σ : standard deviation (null hypothesis)
- \bar{x} : sample mean
- n : sample size