

Contents

| | | |
|-----------|--------------------------------|-----------|
| 1 | Variables and data | 3 |
| 2 | Probability | 6 |
| 3 | Random variable | 14 |
| 4 | Important distributions | 20 |
| 5 | Regression analysis | 25 |
| 6 | Population and Sample | 30 |
| 7 | Estimation | 37 |
| 8 | Testing | 40 |
| 9 | Overview of tests | 46 |
| 10 | Description of tests | 50 |

1 Variables and data

Variable is a quantity that can be measured on a monitored object (speed of a car, severity of an accident etc.)

Random variable is a variable somehow corrupted (e.g. by noise)

Data file is a set of measured values of random variable / more random variables.

$$D = \{x_t\}_{t=1}^N \text{ or } D = \{x_{1;t}, x_{2;t}, \dots, x_{n;t}\}_{t=1}^N$$

Data are

- **discrete** finite number of different values (level of service, color on signal lights)
- **continuous** values from R or frequently R_0^+ (speed of a car)

Ranks of data

| data x_i | ordered data | ranks r_i |
|------------|--------------|-------------|
| 5,2,8,3,6 | 2,3,5,6,8 | 3,1,5,2,4 |

Characteristics of data

- **average**

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i = \frac{1}{N} \sum_X n_i X_i = \sum_X X_i f_i$$

where $f_i = \frac{n_i}{N}$ are relative frequencies

- **variance**

$$s^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 = \sum_X (X_i - \bar{x})^2 f_i$$

- **quantile, critical value** ζ_α, z_α

it is a border separating $\alpha \cdot 100\%$ of the smallest values (quantile) or greatest values (critical values) of a data file.

- **median** $x_{0.5}$ is the $\xi_{0.5}$, i.e. 50% quantile.
- **mode** \hat{x} is the value of dataset which has the higher frequency of repetition.

Graphs

- **time graph:** plots values of x in a discrete time of measurements: $1, 2, 3, \dots$
- **scatter graph** (xy -graph): plots values of y against values of x (used mainly in regression)
- **bar graph:** the values of x in time are plotted as columns.
- **histogram:** is similar to the bar graph but instead of values it plots frequencies of individual values or values from intervals of data.

2 Probability

- **Random experiment** is a trial with defined results which differ (even under the same conditions) - flip of a coin, measurements of car speeds.
- **Sample space** Ω is a set of all possible results
- **Event** E is a subset of the sample space - dice: “even number” $\leftrightarrow \{2, 4, 6\}$
- **Probability** P is a real function defined on the set of events which is
 - non-negative $P(E) \geq 0$
 - normalized $P(\Omega) = 1$
 - additive $P(E_1 \cup E_2) = P(E_1) + P(E_2)$, for $E_1 \cap E_2 = \emptyset$
... must hold even for all countable sets.Is a normalized measure on sets.

Definitions of probability

1. Classical definition

- (a) only for finite number of equally probable results
- (b) based on number of possible and all options (possibilities)

$$P = \frac{m}{n}$$

2. Statistical definition

- (a) based on number of positive and all results of performed experiments

$$P = \frac{M}{N}$$

Example

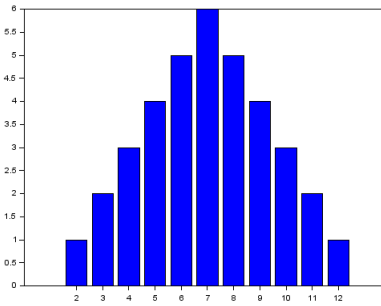
Roll the dice. Request - even number. Classical and statistical definition.

Example

What are results of an experiment: sum on two dice?

| | | | | | | | | | | |
|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|
| 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| $\frac{1}{36}$ | $\frac{2}{36}$ | $\frac{3}{36}$ | $\frac{4}{36}$ | $\frac{5}{36}$ | $\frac{6}{36}$ | $\frac{5}{36}$ | $\frac{4}{36}$ | $\frac{3}{36}$ | $\frac{2}{36}$ | $\frac{1}{36}$ |

Graph



The result is an ordered couple. The shape is because e.g. 3 can be $[1, 2]$ or $[2, 1]$ etc.

Conditional probability

$$P(E_1|E_2) = \frac{P(E_1 \cap E_2)}{P(E_2)}$$

is a probability of event E_1 when we know that the sample space is restricted by the event E_2 .

Explanation

$$P(E_1|E_2) = \frac{P(E_1 \cap E_2)}{P(E_2)} = \frac{\frac{n(E_1 \cap E_2)}{n(\Omega)}}{\frac{n(E_2)}{n(\Omega)}} = \frac{n(E_1 \cap E_2)}{n(E_2)}$$

where $n(\cdot)$ is number of results and Ω is sample space.

Thus, the conditional probability is given by the ration of positive results of E_1 within E_2 and number of results in E_2 .

Remark: E_2 is the new sample space.

Example

Statistical conditional probability for flipping two coins.

We measured

| | | | | | | | | |
|--------|---|---|---|---|---|---|---|---|
| Coin 1 | H | H | T | T | T | H | T | H |
| Coin 2 | T | H | T | H | T | T | T | H |
| C2=H | | x | | x | | | | x |
| C1,2=H | | x | | | | | | x |

1. What is the **conditional** probability $P(\text{Coin 1} = H | \text{Coin 2} = H)$

$$n(\text{Coin 2} = H) = 3, \quad n(\text{Coin 1} = H \wedge \text{Coin 2} = H) = 2$$

$$P = 2/3$$

2. What is the **joint** probability $P(\text{Coin 1} = H \wedge \text{Coin 2} = H)$

$$P = 2/8 = 1/4$$

Example

Roll a dice. $E_1 = \{2, 4, 6\}$ (“even number”), $E_2 = \{3, 4, 5, 6\}$ (“greater than 2”).

Determine $P(E_1|E_2)$

$$E_1 \cap E_2 = \{4, 6\}$$

$$P(E_1|E_2) = \frac{\frac{2}{6}}{\frac{4}{6}} = \frac{2}{4} = \frac{1}{2}$$

| Ω | 1 | 2 | 3 | 4 | 5 | 6 |
|----------|---|---|---|---|---|---|
| E_1 | | x | | x | | x |
| E_2 | | | x | x | x | x |

within E_2 there are two positive results of E_1 .

$P(E_1) = P(E_1|E_2)$ - knowledge of E_2 brings nothing. E_1 does not depend on E_2 .

Continuation of the example

New event $E_3 = \{4, 5, 6\}$ (greater than 3). Determine $P(E_1|E_3)$

| Ω | 1 | 2 | 3 | 4 | 5 | 6 |
|----------|---|---|---|---|---|---|
| E_1 | | x | | x | | x |
| E_3 | | | | x | x | x |

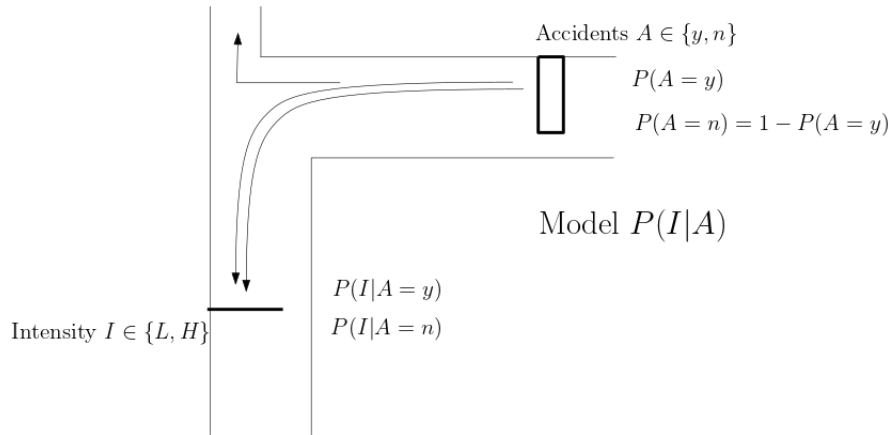
Now the new sample space is E_3 and in it there are two favorable outcomes (4,6). So the probability is

$$P(E_1|E_3) = \frac{2}{3} \neq 0.5.$$

Result: E_1 and E_3 are dependent as knowledge of E_3 leads to the change of probability of E_1 .

Independence

$$P(E_1|E_2) = P(E_1) \quad \text{or} \quad P(E_2|E_1) = P(E_2) \quad \text{or} \quad P(E_1 \cap E_2) = P(E_1)P(E_2)$$



1. Total probability

$$P(I) = P(I|A = y)P(A = y) + P(I|A = n)P(A = n)$$

2. Bayes rule

$$P(A = y|I) = \frac{P(I|A = y)P(A = y)}{P(I)}$$

3 Random variable

- **random variable (rv)** corresponds to random experiment.

Its results are always numbers.

It can be **discrete** or **continuous**.

$X = [2, 1, 3, 3, 2, 3, 1, \dots]$, $X = [5.24, 8.13, 6.78, \dots]$

random vector - corresponds to more measured variables.

- **distribution function** (for both discrete and continuous rv)

Probability of past values; $X \in (-\infty, x)$

$$F_X(x) = P(X \leq x)$$

- **probability function** (for discrete rv)

$$f_X(x) = P(X = x) \quad \longleftrightarrow \quad F(x) = \sum_{t \leq x} f(t)$$

- **probability density function** (for continuous rv)

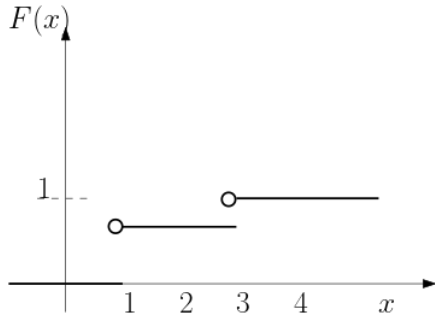
$$f_X(x) = \frac{dF_X}{dx} \iff F(x) = \int_{-\infty}^x f(t) dt \quad \rightarrow \quad P(X \in (a, b)) = \int_a^b f(x) dx$$

Proof

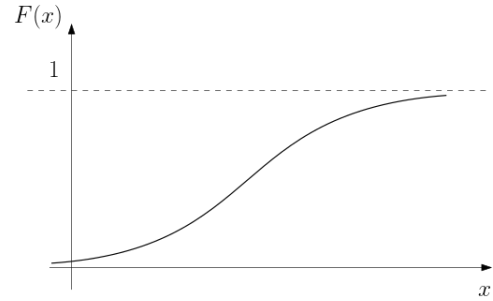
$$\begin{array}{ccc|ccc|}
 P(X < b) & - & - & - & - & - & - & - \\
 P(X < a) & - & - & - & & & & \\
 \hline
 & & & & a & & & b
 \end{array}$$

$$P(X \in (a, b)) = P(X < b) - P(X < a) = F(b) - F(a)$$

Distribution function



Discrete distribution function

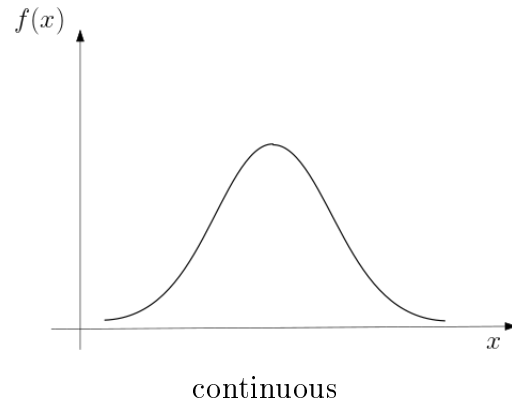
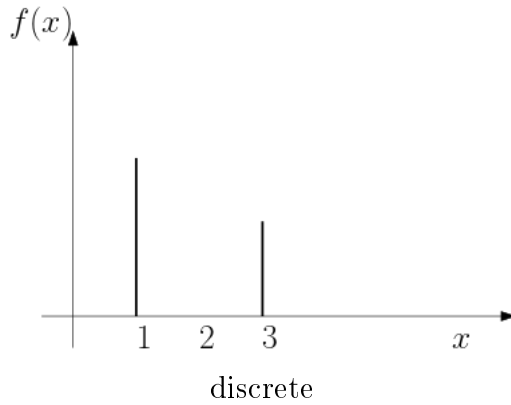


Continuous distribution function

Properties

- (i) $F(-\infty) = 0$, (ii) $F(\infty) = 1$, (iii) non decreasing (constant or grows)

Probability and density function



Properties

- (i) $f(x) \geq 0$, (ii) $\sum_i x_i = 1$ or $\int f(x) dx = 1$

Random vector

is a vector of random variables

$$X = [X_1, X_2, \dots, X_n]$$

Joint distribution (for $n = 2$)

$$f(x_1, x_2)$$

Marginal distribution

$$f(x_1) = \sum_{x_2} f(x_1, x_2) \text{ or } \int_{-\infty}^{\infty} f(x_1, x_2) dx_2$$

Conditional distribution

$$f(x_1|x_2) = \frac{f(x_1, x_2)}{f(x_2)}$$

It holds (**chain rule**)

$$f(x_1, x_2) = f(x_2|x_1) f(x_1)$$

generally $f(x_1, x_2, x_3 \dots x_n) = f(x_n|x_{n-1} \dots x_1) f(x_{n-1}|x_{n-2} \dots x_1) \dots f(x_2|x_1) f(x_1)$

Characteristics

Expectation

| | | |
|--|------------------------|--|
| $\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$ | $E[X] = \sum_X x f(x)$ | $E[X] = \int_{-\infty}^{\infty} x f(x) dx$ |
|--|------------------------|--|

Variance

| | | |
|--|-----------------------------------|---|
| $s^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$ | $D[X] = \sum_X (x - E[X])^2 f(x)$ | $D[X] = \int_{-\infty}^{\infty} (x - E[X])^2 f(x) dx$ |
|--|-----------------------------------|---|

Covariance

| | | |
|---|--|---|
| $c = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$ | $C[X, Y] = \sum_{XY} (x - E[X]) \times (y - E[Y]) f(x, y)$ | $C[X] = \int_{-\infty}^{\infty} (x - E[X]) \times (y - E[Y]) f(x) dx$ |
|---|--|---|

Quantile, critical value

| | | |
|---|---|--|
| border of $\alpha \cdot 100\%$ smallest , greatest | $\alpha = \sum_{x \leq \zeta_\alpha} f(x),$ $\alpha = \sum_{x \geq z_\alpha} f(x)$ | $\alpha = \int_{-\infty}^{\zeta_\alpha} f(x) dx,$ $\alpha = \int_{z_\alpha}^{\infty} f(x) dx$ |
|---|---|--|

4 Important distributions

Discrete random variable

Bernoulli distribution

EXAMPLE: a car turns to left or right.

Probability function

$$P(x; \pi) = \pi^x (1 - \pi)^{1-x}, \quad x = 0, 1 \quad (1)$$

$$E = \pi, \quad D = \pi(1 - \pi).$$

Binomial distribution

EXAMPLE: n times toss a coin. $x = 3$ means we demand so that head comes three times.

Probability function

$$P(x; n, \pi) = \binom{n}{x} \pi^x (1 - \pi)^{n-x}, \quad x = 0, 1, \dots, n$$

$$E = n\pi, D = n\pi(1 - \pi).$$

Poisson distribution

Models the number of times an event appears in a given interval of time or space.

Probability function

$$P(x; \lambda) = e^{-\lambda} \frac{\lambda^x}{x!}, \quad x = 0, 1, 2, \dots$$

$$E = D = \lambda.$$

Geometric distribution

EXAMPLE: A shooter shots at a small target. The probability of hitting is $p = 0.2$. What is the probability that the first hit will occur at the x -th shot.

Probability function

$$P(x, p) = p(1 - p)^x, \quad x = 0, 1, 2, \dots$$

$$E = \frac{1-p}{p}, \quad D = \frac{1-p}{p^2}$$

General discrete (categorical) distribution

Probability function

$$P(x; p_1, p_2, \dots, p_n) = p_x, \quad x \in \{x_1, x_2, \dots, x_n\} \quad \text{or} \quad \begin{array}{c|cccc} x & 1 & 2 & \cdots & n \\ \hline f(x) & p_1 & p_2 & \cdots & p_n \end{array}$$

$$E = \sum_{x=1}^n x p_x, \quad D = \sum_{x=1}^n (x - E)^2 p_x$$

Continuous random variable

Uniform distribution

Probability density function

$$f(x; a, b) = \frac{1}{b - a}, \quad x \in (a, b), \quad \text{and zero otherwise}$$

$$E = \frac{a+b}{2}, \quad D = \frac{(b-a)^2}{12}$$

Normal distribution

Probability density function

$$f(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}, \quad x \in R$$

$$E = \mu, \quad D = \sigma^2$$

Exponential distribution

Probability density function

$$f(x; \delta) = \frac{1}{\delta} e^{-\frac{x}{\delta}}, \quad x > 0$$

$$E = \delta, \quad D = \delta^2$$

Sample distributions

χ^2 -distribution

$$\chi^2(n) = \sum_{i=1}^n (N_i(0, 1))^2$$

Student t -distribution

$$St(n) = \frac{N(0, 1)}{\sqrt{\chi^2(n)/n}}$$

F -distribution

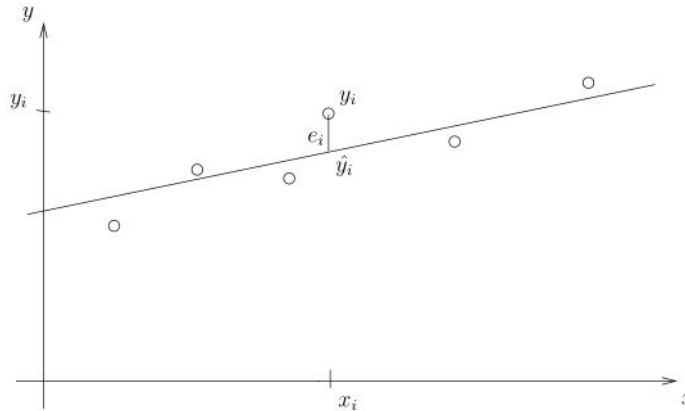
$$F(m, n) = \frac{\chi_1^2(m)/m}{\chi_2^2(n)/n}$$

5 Regression analysis

Two variables x and y .

Measurement $[x_i, y_i]$ - point in a plane.

Can the points $i = 1, 2, \dots, N$ be interpolated by a line?



\hat{y}_i - prediction (lies at the line)

e_i - residuum $e_i = y_i - \hat{y}_i$ (error in approximation)

Optimal approximation - criterion

$$\sum_{i=1}^N e_i^2 \rightarrow \min$$

Regression line

$$\hat{y} = b_0 + b_1 x$$

where

$$b_1 = \frac{S_{xy}}{S_x}, \quad b_0 = \bar{y} - b_1 \bar{x}$$

$$S_x = \sum (x_i - \bar{x})^2, \quad S_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y})$$

Indication of derivation

$$y_i = b_0 + b_1x_i + e_i \rightarrow e_i = y_i - b_0 - b_1x_i$$

$$\begin{aligned} \sum_{i=1}^N e_i^2 &= \sum_{i=1}^N (y_i - b_0 - b_1x_i)^2 = \\ &= \sum_{i=1}^N (y_i^2 + b_0^2 + b_1^2x_i^2 - 2b_0y_i - 2b_1x_iy_i + 2b_0b_1x_i) = \\ &= \underbrace{\sum_{i=1}^N y_i^2}_A + Nb_0^2 + b_1^2 \underbrace{\sum_{i=1}^N x_i^2}_B - 2b_0 \underbrace{\sum_{i=1}^N y_i}_C - 2b_1 \underbrace{\sum_{i=1}^N x_iy_i}_D + 2b_0b_1 \underbrace{\sum_{i=1}^N x_i}_E = \\ &= A + Nb_0^2 + Bb_1^2 - 2Cb_0 - 2Db_1 + 2Eb_0b_1 \end{aligned}$$

... and we look for minimum in variables b_0 and b_1 .

Multivariate regression

Regression line

$$\hat{y} = b_0 + b_1x_1 + b_2x_2 + \cdots + b_mx_m$$

Construction

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \cdots \\ y_N \end{bmatrix}, \quad X = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1m} \\ 1 & x_{21} & x_{22} & \cdots & x_{2m} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ 1 & x_{N1} & x_{N2} & \cdots & x_{Nm} \end{bmatrix}$$

$$\theta = (X'X)^{-1} X'Y$$

where $\theta = [b_0, b_1, b_2 \cdots b_m]'$.

Nonlinear regression

– polynomial (solved as multivariate)

$$\hat{y} = b_0 + b_1x + b_2x^2 + \cdots + b_px^p$$

– exponential (solve by linearization)

$$\hat{y} = b_0 \exp(b_1x)$$

linearization: taking logarithm $\underbrace{\log(\hat{y})}_{\hat{Y}} = \underbrace{\log(b_0)}_{B_0} + b_1x \rightarrow \hat{Y} = B_0 + b_1x$

6 Population and Sample

Population (variable) is the whole set of values from which we take realizations of random variable, and the probabilistic rules of generation). It is the **generator**.

Sample is a set of sampled values from the population. The set of **generated data**.

Examples:

- Estimating average age of people living in Prague.
Population: a set of all people living in Prague. *Sample:* a set of people we ask about their age.
- Estimation of average speed of passing cars.
Population: all possible speeds. *Sample:* speeds of cars that we have measured.
- Throwing the dice
Population: probability function. *Sample:* results of a hundred rolls of the dice.

Sample

The population is clear, with no secrets. The sample is more complex.

Sample - is a vector of independent and equally distributed random variables

$$\mathbf{X} = [X_1, X_2, \dots, X_N]$$

1. independent - random choice (representativeness),
2. equally distributed - come from the same experiment.

it takes into account all possible samples

Sample realization - is a vector of numbers; realizations of the sample.

As the sample is random, the sample realization differs at each selection.

Characteristics

- **Population** is random variable. Its characteristic are probability or density function. Usually it depends on some parameter which is unknown and is to be estimated.
- **Sample realization** is an ordinary set of numbers. For its description we use the characteristics from the descriptive statistics.

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

- **Sample** similarly as population is random (vector). Its description is given by characteristic of random variable.

$$E[X] = \int_{-\infty}^{\infty} x f(x) dx$$

For the sample $X=[X_1, X_2 \cdots X_n]$ we define:

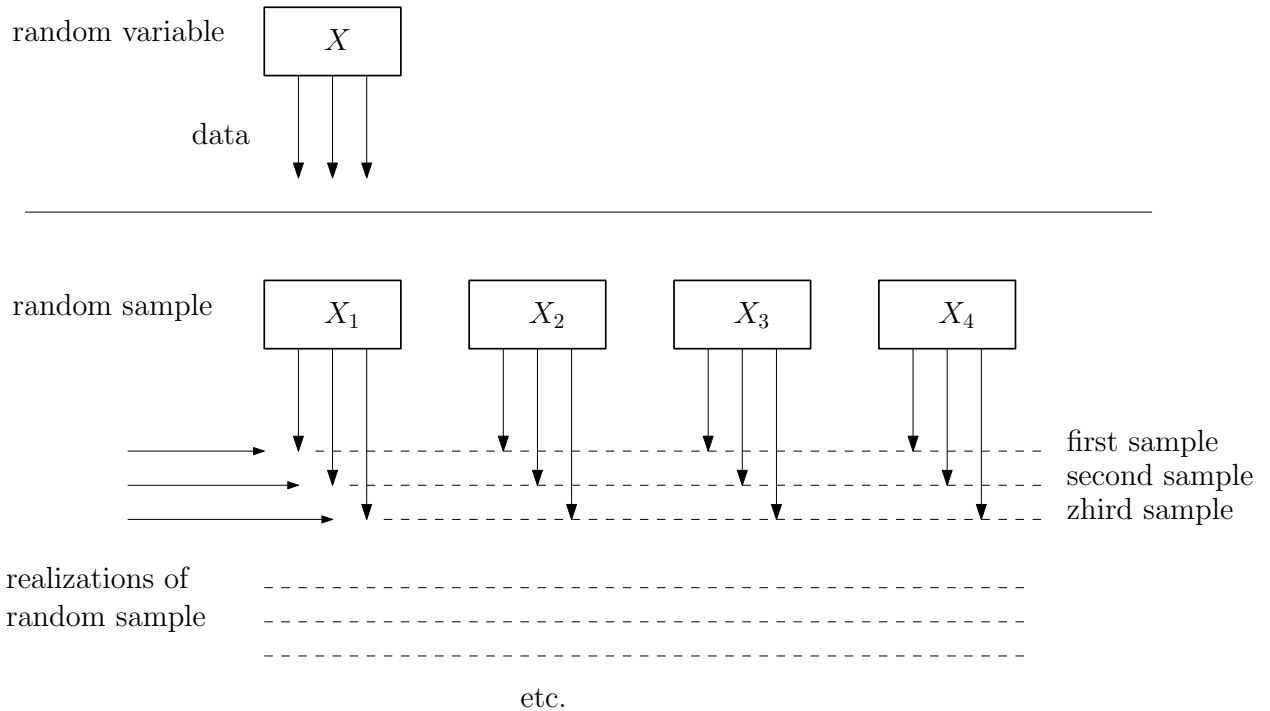
Sample average

$$\bar{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n X_i$$

Sample variance

$$\mathbf{S}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{\mathbf{X}})^2$$

Random variable and random sample



Characteristics of sample average

For random variable X with expectation μ and variance σ^2 and sample \mathbf{X} of the length n , it holds

$$E[\bar{\mathbf{X}}] = \mu, \quad D[\bar{\mathbf{X}}] = \frac{\sigma^2}{n}$$

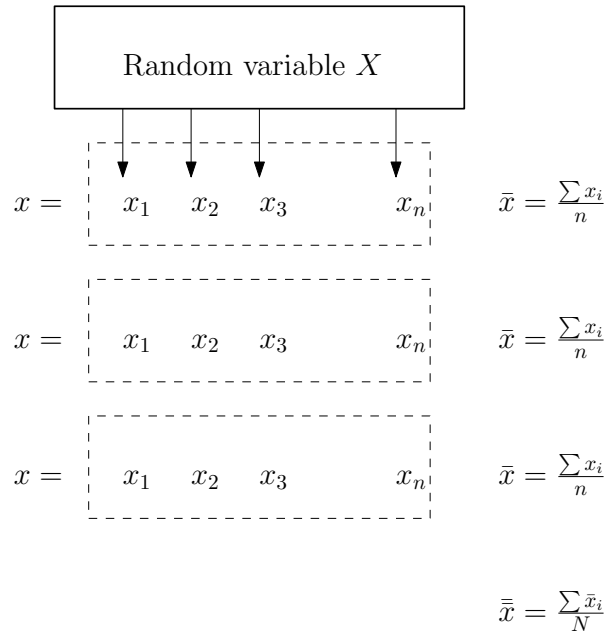
Proof

$$E\left[\frac{1}{n} \sum_i X_i\right] = \frac{1}{n} \sum_i E[x_i] = \frac{1}{n} n\mu = \mu$$

... how does it work??? For repetitive sampling.

$$D\left[\frac{1}{n} \sum_i X_i\right] = \frac{1}{n^2} \underbrace{\sum_i D[X_i]}_{x_i \text{ independent}} = \frac{1}{n^2} n\sigma^2 = \frac{\sigma^2}{n}$$

... variance falls down with the sample length.



Limit theorems

Law of large numbers

For sample length going to infinity, the values of sample characteristics converge to corresponding population ones.

E.g.

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n x_i \rightarrow \mu$$

Central limit theorem

For sample length going to infinity the sum characteristics have approximately normal distribution, no matter what distribution has the random variable.

E.g. For $x \sim \text{binomial}$, $f(\bar{x}) \doteq N_{\bar{x}}\left(\mu, \frac{\sigma^2}{n}\right)$ where μ and σ^2 are expectation and variance of x and n is length of the sample ($n \rightarrow \infty$).

7 Estimation

We have random variable X with the description $f(x, \theta)$
and the sample realization $x = [x_1, x_2, \dots, x_N]$.

Statistics T for estimation of the parameter θ is a function of random sample - mostly a sample characteristic.

Expectation - average, variance - sample variance, independence - covariance, etc.

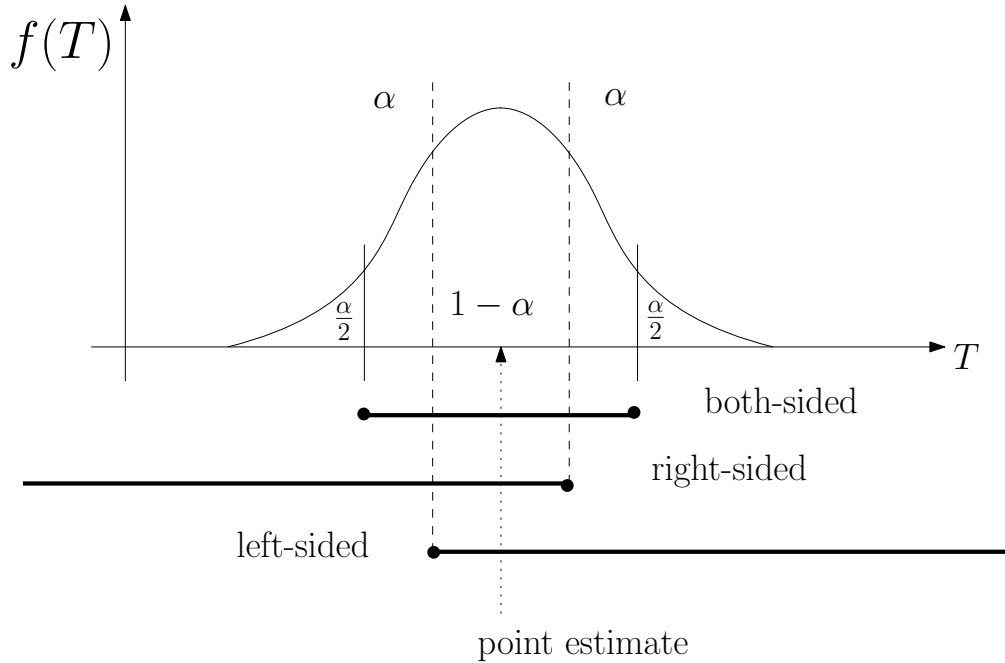
Point estimate $\hat{\theta}$ of a parameter θ is the value of the statistic with the sample realization inserted.

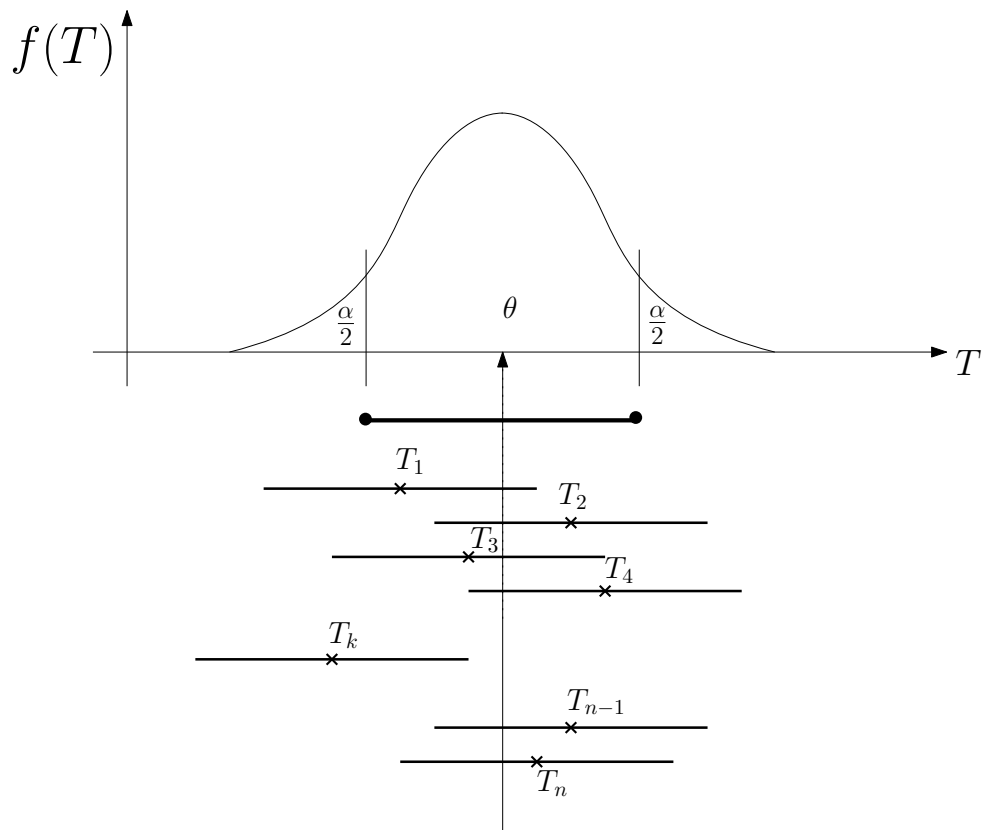
Interval estimate CI is an interval to which the true value of the estimated parameter θ belongs with a defined probability $1 - \alpha$

- both-sided interval $CI = (L, U)$
- left-sided interval $CI = (L, \infty)$
- right-sided interval $CI = (-\infty, U)$

Construction of CI

$f(x) \rightarrow$ data; $f(T) \rightarrow$ parameter (T is estimator)





8 Testing

Test about the value of some parameter.

H_0 : **zero hypothesis** - stand up for the current state (nothing has changed)

$$\mu = \mu_0$$

H_A : **alternative hypothesis** - denies H_0

both-sided test $\mu \neq \mu_0$

left-sided test $\mu < \mu_0$

right-sided test $\mu > \mu_0$

Principle of testing

It is based on CI .

1. Perform CI for given parameter **according** H_0 (both/left/right)
2. If the value of realized statistics lies
 - (a) within CI the H_0 is not rejected.
 - (b) outside of CI - the H_0 is rejected.

Example

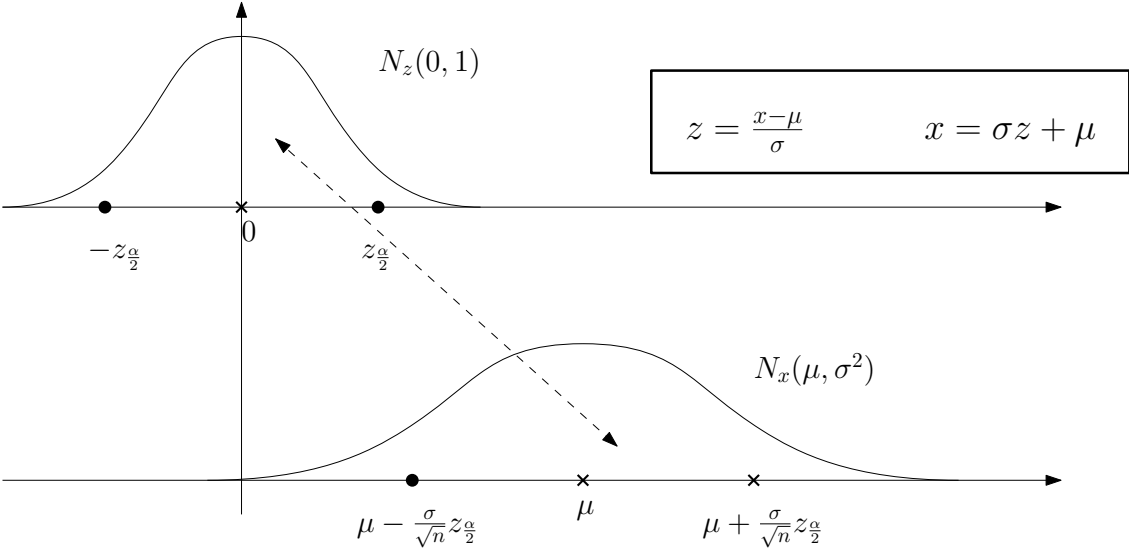
Test H_0 hypothesis $\mu_0 = 1$ against $H_A \mu \neq \mu_0$ if CI constructed from the sample realization of the length 35 is $(-0.7, 3.2)$ and the sample estimate $\hat{\theta} = T_t = 1.7$.

Solution: H_0 is not rejected as $T_r \in CI$. I.e. we do not reject that the true expectation can be 1.

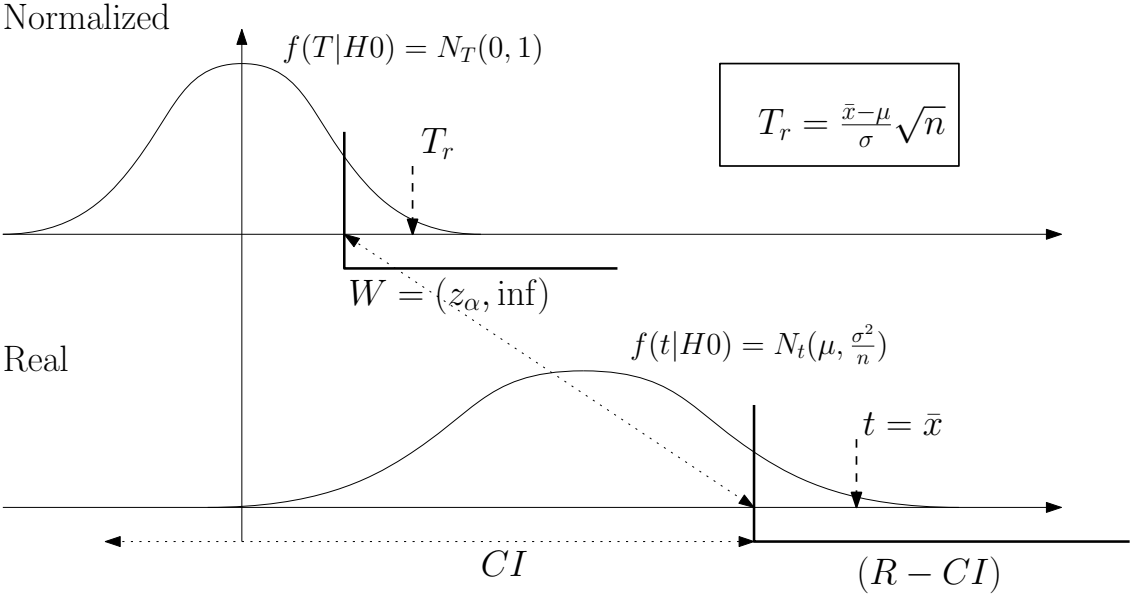
Normalization

For $\bar{x} \rightarrow \mu$

Normalization



Performing tests



W is critical region. If $T_r \in W$: reject H_0 .

For $x \sim N_x(\mu, \sigma^2)$ we have $\hat{\mu} = \bar{x}$. $E[\bar{x}] = \mu$, $D[\bar{x}] = \frac{\sigma^2}{n}$.

In real units

$$\bar{x} \in (R - CI_\alpha) \text{ - supplement to CI}$$

In normalization (tests)

$$T_r = \frac{\bar{x} - \mu_0}{\sigma} \sqrt{n} \in (z_\alpha, \infty) = W$$

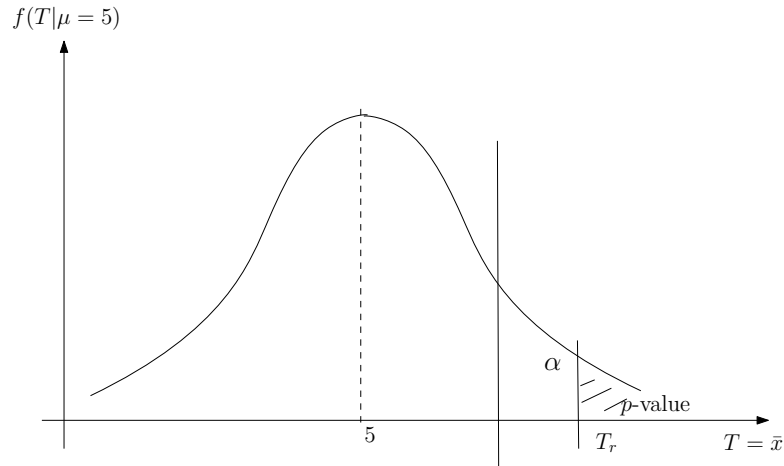
Evaluation of a test

Critical region W is a complement to confidence interval CI (according to H_0), T_r is the value of the statistics with the sample realization inserted.

Test: $T_r \in W$ reject H_0 otherwise do not reject.

p-value is the probability that repeated experiments will give results that are similar or worse with respect to not rejecting H_0 .

Test: if p -value is small (smaller with respect to α) then reject H_0 otherwise do not reject.



9 Overview of tests

How to determine a proper test in Statext

- parametric or nonparametric
- number of variables: 1, 2, more
- tested parameter (property)
- important: H_0 and p-value

Tests with one sample

Parametric tests (normality required)

- expectation (known \times unknown variance)
- proportion
- variance

Nonparametric tests (normality is not required)

- Wilcoxon test: tests median

Tests of distribution type

- w/s test of normality
- Kolmogorov-Smirnov test: tests given distribution.
- Chi-square test of homogeneity: test of distribution type.

Tests with two samples

Parametric

- two expectations: independent \times paired
- two proportions
- two variances

Nonparametric

- Mann-Whitney test: equality of two medians (independent samples)
- Wilcoxon test: two medians (paired samples)
- McNemar test: improvement after some action. Binary data.

Tests with more samples

Parametric

- Analysis of variance: equality of several expectations
- Anova with two factors: equality in columns and rows
- Bartlett - equality of more variances
- Scheffé - detects different expectations

Nonparametric

- Kruskal-Wallis: nonparametric anova.
- Friedman - block test of equality of medians

Tests of independence

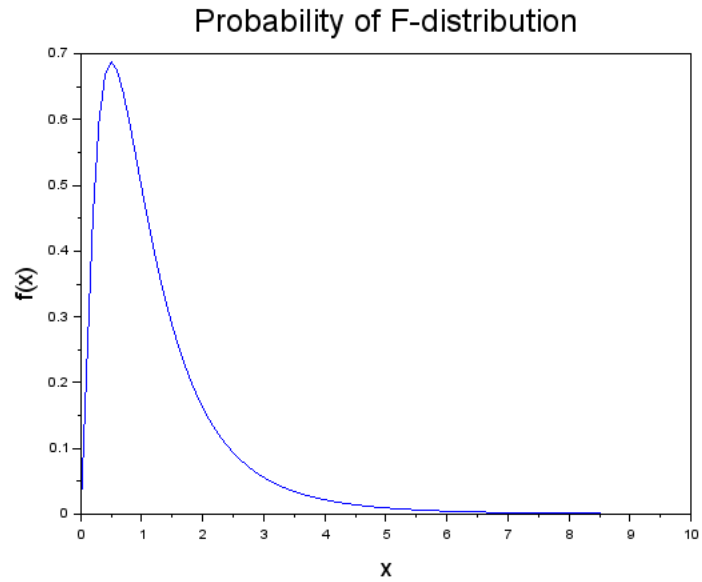
- Lambda coefficient: association of two discrete rvs.
- Pearson test: independence of two rvs (parametric).
- Spearman test: nonparametric Pearson.
- Chi-square independence: test of independence of two rvs.

10 Description of tests

F-test

Compares explained variance V_E and unexplained variance V_U . The statistics is

$$F = \frac{V_E}{V_U} \sim F \text{ Fisher distribution}$$



If $F = 0$, p-value=1 - $V_E = 0$ nothing is explained.

If $F \rightarrow \infty$, p-value $\rightarrow 0$ - $V_U \rightarrow 0$ all is explained.

ANOVA I

We have data from several sources (populations). We test, if the expectations of the populations are equal.

The data are X_i in the following table.

| X_1 | X_2 | X_3 |
|-------------|-------------|-------------|
| x | x | x |
| x | x | x |
| x | x | x |
| \bar{x}_1 | \bar{x}_2 | \bar{x}_3 |
| v_1 | v_2 | v_3 |

We compute averages \bar{x}_i and variances v_i . Then:

- Average of variances v_i corresponds to **unexplained variance** V_U - it describes the overall variance in the data.
- Variance of the averages \bar{x}_i corresponds to **explained variance** V_E - it expresses the variance between classes.

If the explained variance V_E is sufficiently larger with respect to the unexplained one V_U then we conclude that the classes are not equal.

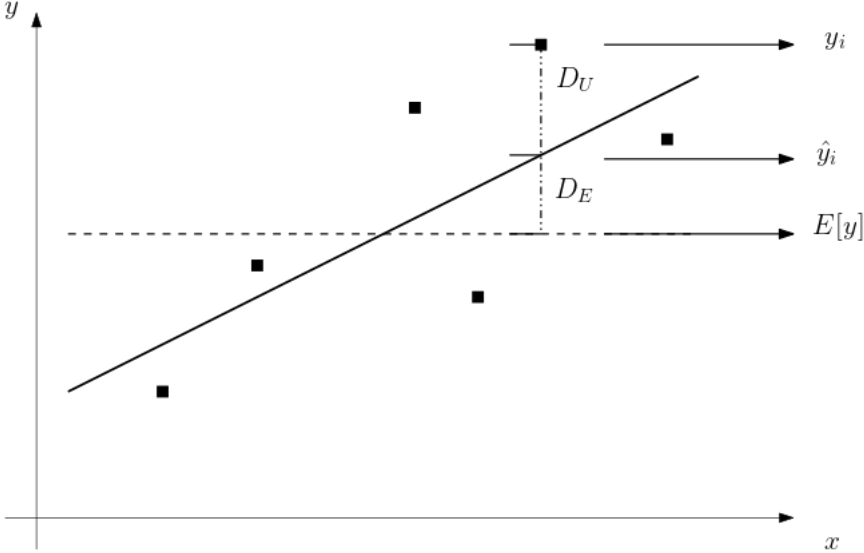
Statistics:

$$F = \frac{V_E}{V_U} \sim F \text{ distribution (right-sided test)}$$

ANOVA II

| | X_1 | X_2 | X_3 | | |
|-------|-------------|-------------|-------------|-------------|----------|
| Y_1 | x | x | x | \bar{y}_1 | v_{y1} |
| Y_2 | x | x | x | \bar{y}_2 | v_{y2} |
| Y_3 | x | x | x | \bar{y}_3 | v_{y3} |
| Y_4 | x | x | x | \bar{y}_4 | v_{y4} |
| | \bar{x}_1 | \bar{x}_2 | \bar{x}_3 | | |
| | v_{x1} | v_{x2} | v_{x3} | | |

Regression



Friedman test

Example

A quality of five shops is checked by three evaluators. Each evaluator rates each shop. The results are in the table

| Evaluator/shop | 1 | 2 | 3 | 4 | 5 |
|----------------|---|---|---|---|---|
| 1 | x | x | x | x | x |
| 2 | x | x | x | x | x |
| 3 | x | x | x | x | x |

We are interested in evaluation of the shops (not evaluators).

Remark: The shops are called treatment, the evaluators are subjects. In the Statext, we choose “Each data set is for subject” which means, the data in rows come from individual subjects.

Chi-square test

Works for discrete data or continuous ones discretized on intervals.

I based on

- O observed absolute frequencies - from measured data,
- E expected absolute frequencies - constructed so that:
 - H_0 is precisely fulfilled,
 - number of data (sum of frequencies) is the same as for measured ones.

Statistics with χ^2 distribution is

$$\chi^2 = \sum_i \frac{(O_i - E_i)^2}{E_i}$$

EXAMPLE

Test if the number of traffic accidents is uniformly distributed during the week, if we measured the following number of accidents

| Weekdays | Saturday | Sunday |
|----------|----------|--------|
| 587 | 98 | 103 |
| O_1 | O_2 | O_3 |

Sum $\sum O_i = 788$. Time axis has 7 intervals (days). Weekdays has 5, Sat and Sun have 1. So we need to divide 788 into groups with proportion 5/1/1.

$$E_1 = \frac{788}{7}5 = 562.86, \quad E_2 = \frac{788}{7}1 = 112.57, \quad E_3 = \frac{788}{7}1 = 112.57$$

$$\chi^2 = 3.73; \quad pv = 0.155$$

As $pv > 0.05$ the H_0 is not rejected at the confidence level $\alpha = 0.05$.

Pearson and Spearman tests

Pearson test tests the correlation coefficient

$$R = \frac{C[X, Y]}{\sqrt{D[X] D[Y]}}$$

If $R = 0$ the random variables X and Y are uncorrelated. If $R \rightarrow -1$ or 1 , the variables are strongly correlated. The test is both sided.

Its statistics is the sample correlation coefficient

$$r = \frac{\text{cov}(x, y)}{\sqrt{\text{var}(x) \text{var}(y)}} \sim t \text{ Student distribution}$$

H_0 : The variables are independent.

For p-value small, the independence is rejected - variables are dependent.

Spearman test is a nonparametric variant of Pearson. Instead of data it uses their ranks. Can be used for discrete variables.

McNemar test

We measure binary (yes/not) data before and after some action. We test if the action caused any change.

Example

We ask 20 people if they have a cold. The answers are yes/not. Then we give them some medicine and ask again. The question is whether there was any change after the application of the drug (either positive or negative).

H0: no change.

→ if p -value is small, a change has been detected.

11 Validation in regression

- Graph xy
- Test of independence x and y (Pearson, Spearman).
- F test of prediction.
- Independence and autocorrelation of residuals ($e_t = ae_{t-1} + \epsilon_t$)
- Prediction error $RPE = \text{var}(y - y_p) / \text{var}(y)$.