

Experimental Design and Data Analysis

Lectures 0 and 1

Eduard Belitser

VU Amsterdam

Lecture Overview

- ① course organization
- ② experimental design
- ③ recap probability theory and basic statistics
- ④ recap: examples in R

Course organisation

- **Prerequisites:** basic statistics course (e.g., Statistical Methods), basic probability, R knowledge.
- The first 1.5 lectures is a recap of what you are supposed to know. **Test your prerequisite knowledge:** exam (quiz) will be available on canvas.
- **All relevant information is on canvas:** schedule, lecture slides, assignments (in due time), R manual(s) and suggestions additional literature.
- **R** is an open software, widely adopted in the academic community, it is a programming language (object oriented), a statistical package.
- **RStudio** is a powerful user interface for R.

Experimental design

What is experimental design?

- Experiments are performed with varied preconditions represented by ind. variables, also referred to as input variables or predictor variables.
 - The change in predictors is hypothesized to result in a change in one or more dep. variables, also referred to as output or response variables. → outcome
 - The experimental design may also identify control variables that must be held constant to prevent external factors from affecting the results.
 - Experimental design involves also planning the experiment under  statistically optimal conditions given the constraints of available resources.
 - Main concerns in experimental design: validity, reliability, replicability, achieving appropriate levels of statistical power and sensitivity.
 - Ronald Fisher: *The Arrangement of Field Experiments* (1926) and *The Design of Experiments* (1935).
→ found. of stats

Experimental design, randomization

- Statistics allows to generalize from data to a true state of nature, but statistical inference requires assumptions and mathematical modeling.
 - The data should be obtained by a carefully designed experiment (or at least it must be possible to think about the data in this way).
 - Any good design involves a chance element: "experimental units" are assigned to "treatments" by chance, or by randomization. The purpose is to exclude other possible explanations of an observed difference.
 - We need probability to quantify the randomization. In practice, randomization is implemented with a random number generator. In R:

```

> x=rep(c("A","B"),each=5); x
[1] "A" "A" "A" "A" "A" "B" "B" "B" "B" "B"
> sample(x) # create a sequence of 5 A's and 5 B's in random order
[1] "A" "B" "A" "B" "B" "A" "B" "A" "A" "B"
> rbinom(10,1,0.5) # toss a fair coin 10 times
[1] 1 0 1 1 1 0 1 0 0 0 generate rand. binom. var.
> rbinom(10,1,0.5) # again toss a fair coin 10 times
[1] 1 0 0 0 0 1 0 1 1 0
> rbinom(5,1,0.8) # toss a biased coin (success probability=0.8) 5 times
[1] 1 1 0 1 1 bigger prob. for 1

```

take person from pop. at rand.
w/o paying att. if it's M/F

to remove the
eff. of gender

Examples, observational studies

choosing plots = control vars

EXAMPLE To compare two fertilisers we prepare 20 plots of land, apply the first fertiliser to 10 **randomly** chosen plots and the second one to the remaining plots. We plant a crop and measure the total yield from each plot.

EXAMPLE To compare two web designs we **randomly** select 50 subjects and measure the time needed to find some information. All 50 subjects perform this task with both designs, but for each subject the order of the two designs is based on **tossing a coin**.

EXAMPLE If an experiment involves subjects, then it could be wrong to assign “task A” to the first 10 subjects who arrive and “task B” to the last 10. (There may be a reason for arriving early.) Instead assign the tasks **at random**. Then an observed difference is due to the task (or chance).

Data obtained by registering an ongoing phenomenon, without randomization or applying other controls, is called **observational**. → you can't change anything, no ctrl.s var

EXAMPLE The incidence of lung cancer among **500 smokers** is observed to be higher than among **500 non-smokers**. Does this finding generalize to the full population? Does this show that smoking causes lung cancer?

Exp. design
oooo

Recap probab. theory
●oooooooooooooooooooo

Summarizing data
oooooooooooooooooooo

Recap basic stat. concepts
oooooooooooooooooooo

Recap: examples in R
oooooooooooooooooooo

Recap probability theory and basic statistics

(prerequisite for this course, if needed consult *Elementary Statistics*, by Mario Triola)

Probability distributions: continuous, discrete

- descrip. of how your rand. var. leaves , what val.s + prob.s it takes*
- all poss. outcomes + probs of those outcomes*
- A probability distribution P determines the probability of different outcomes of a random variable.
 - Probability distributions for:
 - discrete random variables which have finite or countable sets of possible outcome values (e.g., dice, coins, birthdays);
 - continuous random variables which have infinite sets of possible outcome values (e.g., temperature, length).
 - The corresponding probability distributions: continuous, discrete.

Remark. Actually, there are distributions which neither continuous nor discrete.

Probability density functions

Examples of the probability density p of some continuous distributions (realised also in R with some default parameter values):

- normal distribution `norm` with parameters μ `mean=0` and σ `sd=1`

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

$$\left. \begin{array}{l} 1-f(x) \geq 0 \\ 2-\int f(x)dx = 1 \end{array} \right\} \begin{array}{l} \text{density} \\ \text{func.} \\ x \sim f(x) \end{array}$$

- exponential distribution `exp` with parameter λ (`lambda=1`)

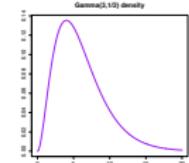
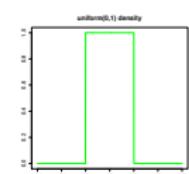
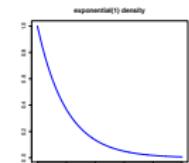
$$p(x) = \lambda e^{-\lambda x}, \quad x > 0.$$

takes only + val.s

- uniform distribution `unif` with parameters minimum (`min=a`) and maximum (`max=b`) of the support interval

$$p(x) = \frac{1}{b-a}, \quad a \leq x \leq b.$$

- Gamma distribution `gamma` with parameters `shape` and rate `rate=1`.



Probabilities of events – continuous distribution

If a random variable X has a distribution with the density $p(x)$, then

prob. of event { $P(X \in I) = \int_I p(x)dx$ for any interval $I \subseteq \mathbb{R}$.
set

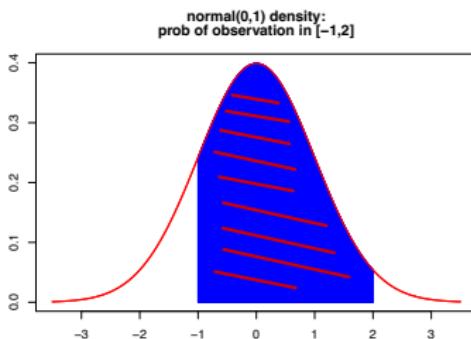
// prob. that In other words, the probability to have an outcome in some interval I is the area under the density function $p(x)$ over that interval.

X belongs to I

Example. For $X \sim N(0, 1)$,

$$\begin{aligned} P(-1 \leq X \leq 2) &= P(X \in [-1, 2]) \\ &= \int_{-1}^2 p(x)dx = \int_{-1}^2 \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx = 0.82. \end{aligned}$$

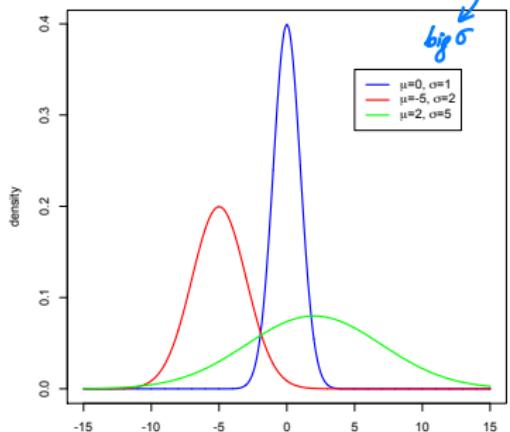
In events for continuous distributions:
 $<$ or \leq ($>$ or \geq) does not matter.



Location and scale, normal density

Two important characteristics of a population are location (or mean) μ and scale (or standard deviation) σ .

makes curve more spread/concentr.
normal densities with different location and scale



always bell shaped

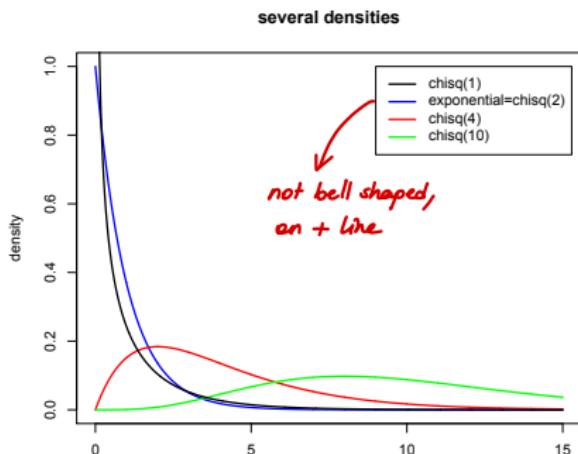
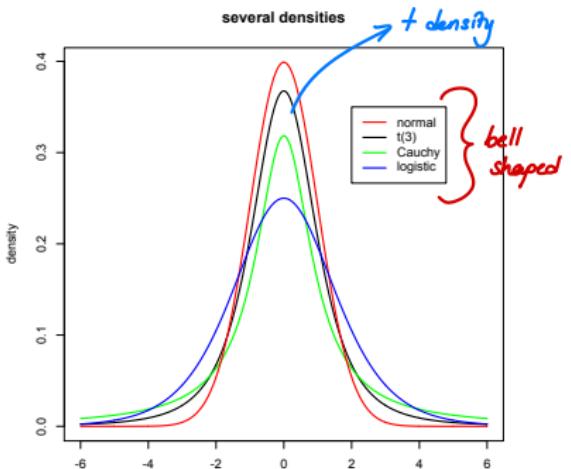
The normal density curve is given by

$$f_{\mu,\sigma}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}(x-\mu)^2/\sigma^2}.$$

The parameters μ and σ are the location and scale. Normal distributions with different μ and σ are still similar in a way.

Remark. The normal curve is very specific! There are many “bell shaped” curves that are not normal.

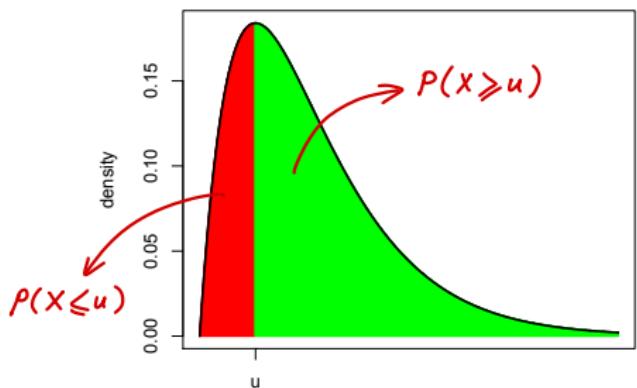
Other symmetric and asymmetric densities



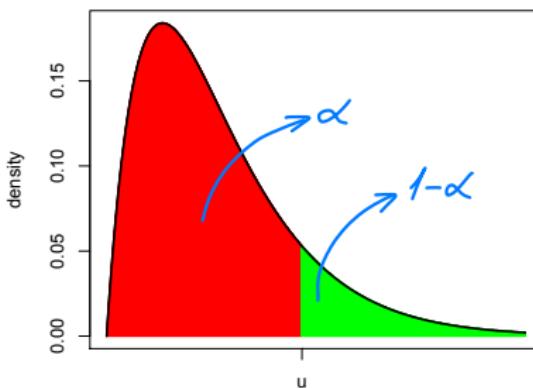
Probabilities and quantiles

If a random variable X is distributed according to a density curve, the probability $P(X \leq u)$ is the (red) area under the density curve **left** of u . Likewise, $P(X \geq u)$ is the (green) area under the density curve **right** of u .

probabilities



probabilities



For distribution P , the quantile of level $\alpha \in (0, 1)$ is the number q_α such that $P(X \leq q_\alpha) = \alpha$, the upper quantile u_α such that $P(X \geq u_\alpha) = \alpha$.

For the standard normal distribution, the quantile and upper quantile are usually denoted by ξ_α and z_α .

Probability of events – discrete distribution

For discrete distributions we have a probability mass function p

$$p(x) = P(X = x).$$

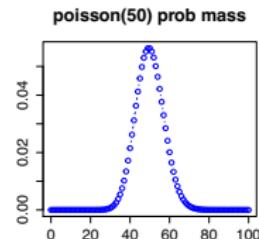
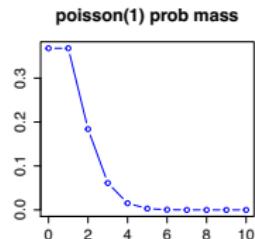
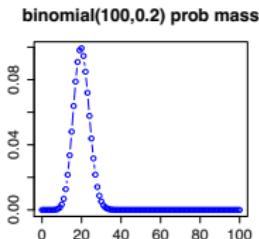
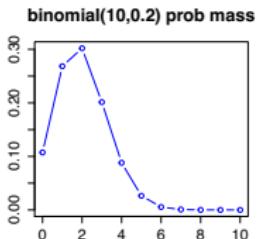
computes prob. in
finately many points

The probability to have an outcome in some set A is the sum

$$P(X \in A) = \sum_{x \in A} p(x).$$

sum is pos. +
can be at most 1.

Examples of discrete distributions are binomial and Poisson.



Probability mass functions for some discrete distributions

Discrete distributions (realised also in R):

- Binomial distribution `binom` with parameters n **size** and p **prob**

$$p(x) = \frac{n!}{x!(n-x)!} p^x (1-p)^{n-x}.$$

*x takes val.s
b.w. (0, n).*

of trials
prob. of success.
trials

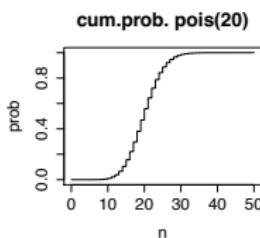
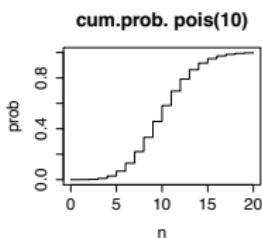
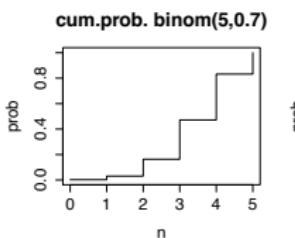
- Poisson distribution `pois` with parameter λ `lambda`

$$p(x) = \frac{\lambda^x}{x!} e^{-\lambda}.$$

*x takes oo many pos.
int. val.s*

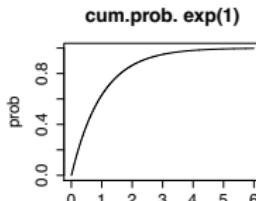
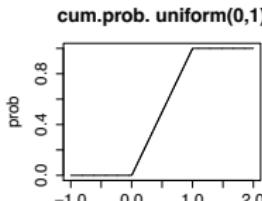
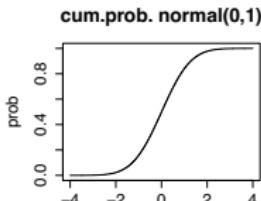
Cumulative distribution/probability function

- The cumulative distribution function (CDF) (sometimes also called cumulative probability function) of a random variable X is $F(u) = P(X \leq u) = \text{pdist}(u, \text{par})$ (continuous and discrete) bigger u = bigger prob. \rightarrow non-decr. func. b.u 0 and 1
- Continuous distr.: $F(u) = \int_{-\infty}^u p(x)dx$; discrete: $F(u) = \sum_{x \leq u} p(x)$.
- Any other probability can be computed via $F(u)$, e.g., for any $a \leq b$, $P(a < X \leq b) = P(X \leq b) - P(X \leq a) = F(b) - F(a)$.



① discrete dist. =
step func.

② cont. dist. =
cont. func.



R-commands for distributions

- in R, there is a number of continuous and discrete distributions `dist` with parameters `par`.
 - Let $p(x)$ denote the density for continuous distribution and the probability mass function for discrete distribution.
 - `ddist(x,par)` computes $p(x)$ (i.e., either density or mass function),
 - `pdist(u,par)` computes the CDF $F(u) = P(X \leq u)$,
 - `qdist(a,par)` ($a \in [0, 1]$) computes the value q such that $pdist(q,par)=a$, this is the a -quantile. The α -quantile q_α is such number that $P(X \leq q_\alpha) = \alpha$. \rightarrow reverse of CDF
 - `rdist(size,par)` yields a random sample from `dist` with parameter `par` of size `size`.
- x = point you want to compute
par = pars of dist.
dist = name of dist.*
- sample size*

Examples in R

*norm.
dist.*

```
> pnorm(2,mean=0,sd=1)-pnorm(-1,mean=0,sd=1) # $P(-1 < X < 2) = P(X < 2) - P(X < -1)$ 
[1] 0.8185946
> pnorm(2)-pnorm(-1) # no need to set the default mean=0, sd=1
[1] 0.8185946
> rnorm(4) # generate 4 standard normals
[1] 0.5592590 -0.3570060 -0.7276720 0.8368255
> dbinom(1,size=5,prob=0.2) # this is  $P(X=1)$  → prob. mass func.
[1] 0.4096
> pbinom(1,size=5,prob=0.2) # this is  $P(X \leq 1)$  → here  $X$  takes 0 and 1.
[1] 0.73728
> dbinom(0,5,0.2)+dbinom(1,5,0.2) # indeed,  $P(X \leq 1) = P(X=0) + P(X=1)$ 
[1] 0.73728
> rpois(3,lambda=5)
[1] 6 7 2
```

Expectation

- The expectation or mean $E(X)$ of a random variable X with probability distribution P is a location parameter of distribution P .
 - For discrete random variable: $E(X) = \sum_x xp(x)$.
 - For continuous random variable: $E(X) = \int xp(x)dx$.
- not a rand. val.
but char. of dist.*

Examples

Throwing a dice: $E(X) = \sum_x xp(x) = \underbrace{1 \times \frac{1}{6}} + \dots + \underbrace{6 \times \frac{1}{6}} = 3\frac{1}{2}$.

Normal distribution: $E(X) = \int xp(x)dx = \int x \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}} dx = \dots = \mu$.

does not have to be 1 of val.s e.g. 1, 2, ..., 6

mean = expect.

Variance and standard deviation

char. of dist.

- The variance of a probability distribution is a scale (or spread) parameter.
- For discrete random variable: $\text{Var}(X) = \sum_x (x - E(X))^2 p(x)$.
- For continuous random variable: $\text{Var}(X) = \int (x - E(X))^2 p(x) dx$.
- Definition: the standard deviation σ is the square root of the variance $\sigma = \sqrt{\text{Var}(X)}$.

Examples

Throwing dice:

$$\text{Var}(X) = \sum_x (x - 3\frac{1}{2})^2 p(x) = (\underbrace{1 - 3\frac{1}{2}}_{})^2 \times \frac{1}{6} + \dots + (\underbrace{6 - 3\frac{1}{2}}_{})^2 \times \frac{1}{6} = 2.92.$$

Normal distribution:

$$\text{Var}(X) = \int (x - \mu)^2 p(x) dx = \int (x - \mu)^2 \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}} dx = \dots = \circled{\sigma^2}.$$

Expectation and variance for some distributions

	Expectation	Variance	
Uniform(a, b)	$(a + b)/2$	$(b - a)^2/12$	<i>for unif. std. dist: $\text{var}(x) = 1/12$</i>
Normal(μ, σ^2)	μ	σ^2	
Exponential(λ)	$1/\lambda$	$1/\lambda^2$	★
Binomial(n, p)	np	$np(1 - p)$	
Poisson(λ)	λ	λ	

Exp. design
oooo

Recap probab. theory
oooooooooooooooooooo

Summarizing data
●oooooooooooooooooooo

Recap basic stat. concepts
oooooooooooooooooooo

Recap: examples in R
oooooooooooooooooooo

Summarizing data and exploring distributions

Population and sample

set of cert. obj.

- A population can be an actual population, e.g., the heights of all men in the Netherlands.
- It can also be the (imaginary) infinite number of outcomes obtained by repeating an experiment over and over, e.g., throwing a dice many times.
- A sample is a set of values (randomly) selected from a population.
- The population has a certain distribution, called the population distribution.
- From the sample we want to gain/extract information about this unknown population distribution.
- This is the main problem of statistics/data analysis.

e.g. for norm.
dist → pop. is
whole real line

Types of data summaries

A good summary of a data set shows the **relevant information** in a data set.

- **numerical summaries** (of what it estimates/investigates)

- sample mean (**population mean**)
- sample median (**population median**)
- sample standard deviation (**population standard deviation**)
- sample variance (**population variance**)
- sample correlation(s) (**population correlation(s)**)
- ...

O sample:

you draw indep. copies of rand. var. w/ cert. dist = #'s that could be thought as realis. of cert. dist

- **graphical summaries**

- histogram (estimates probability density or probability mass)
- boxplot (assess symmetry, range, outliers)
- scatter plot(s) (assess relations between variables)
- normal QQ-plot (checks normality)
- empirical distribution function (cumulative prob. function)
- ...

est. of dist. func.

Data summaries and some useful R -commands

- Densities, probabilities and quantiles of many distributions can be computed in R. Commands in R: `dnorm(u,par)`, `pnorm(q,par)`, `qnorm(a,par)`, `rnorm(size,par)`, etc.
- Numerical summaries: sample mean, sample variance, sample median, sample standard deviation, sample α -quantile, etc. Commands in R: `mean(x)`, `var(x)`, `med(x)`, `sd(x)`, `quantile(x,a)`, `summary(x)`, `range(x)`, etc.
- Graphical summaries: histogram, boxplot, (normal) QQ-plot, scatter plot(s), empirical distribution function (cumulative histogram), etc. Commands in R: `hist(x)`, `boxplot(x)`, `qqnorm(x)`, `plot(x,y)`, `plot(ecdf(x))`, etc.

Study Assignment 0.

The **boxplot** of a sample is a box with whiskers and (possibly) extremes, from which you can see the **scale** of the data, its **symmetry**, whether there are **extreems** (outliers).

Complement graphical summaries with numerical summaries and vice versa.

Some numerical summaries: reminder

sample size	n
location	$\bar{x} = n^{-1} \sum_{i=1}^n x_i$
	$\text{med}(x) = \begin{cases} x_{((n+1)/2)}, & \text{if } n \text{ odd} \\ (x_{(n/2)} + x_{(n/2+1)})/2, & \text{if } n \text{ even} \end{cases}$
scale	$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$
	$s = \sqrt{s^2}$

Here $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$ is the ordered sample.

Interpretation of location measures:

- mean – average value
- median – middle value in sorted values

dep.s on sample size
being even or odd

Interpretation of scale measures:

- variance – average squared deviation from mean
- standard deviation – square root of variance

Histogram

The histogram of a sample x_1, x_2, \dots, x_n is a barplot composed of cells, where the height of the bar of cell C is either the count $\#\{1 \leq i \leq n : x_i \in C\}$ of observations in that cell C , or its fraction normalized by the cell size:

$$\frac{\text{number of observations in cell } C}{\text{sample size} * \text{cell size}} = \frac{\#\{1 \leq i \leq n : x_i \in C\}}{n\delta_C}.$$

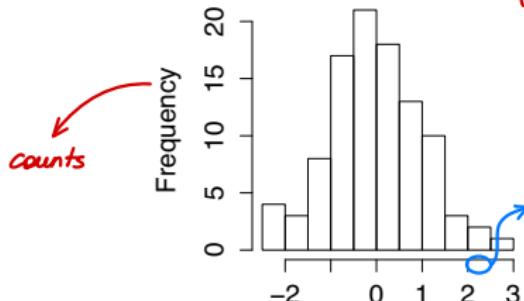
size of cell
can be chosen
or R handles
 δ_C

```
> x=rnorm(100); par(mfrow=c(1,2)) # two plots next to each other
> hist(x)} # frequencies on y-axis
> hist(x,prob=T) # probabilities on y-axis
```

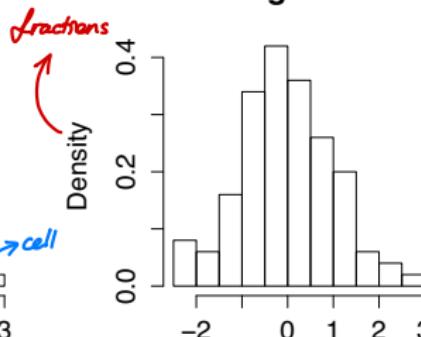
Why are the Y-axes different? Because $\delta_C = 1/2$.

○ histogram = estimates
pdf or prob. mass func.

Histogram of x

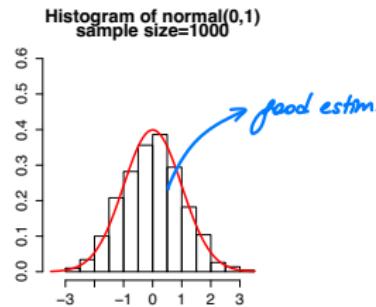
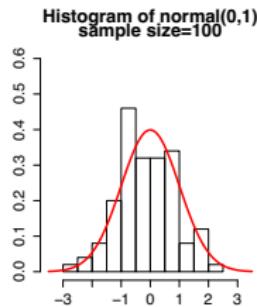
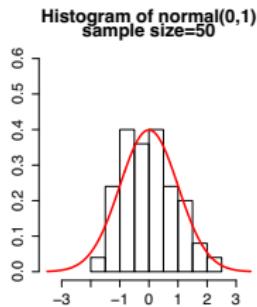
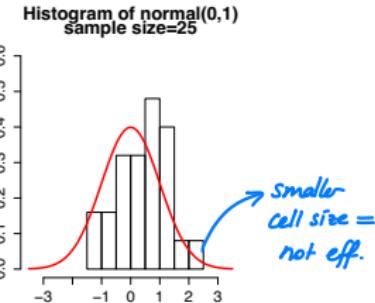
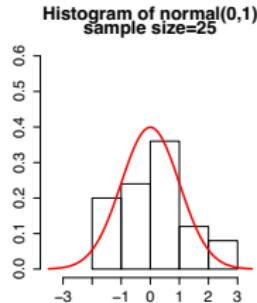
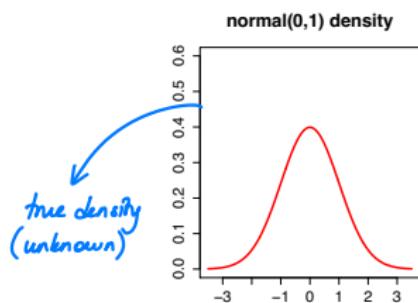


Histogram of x



Histogram versus density (1)

The histogram of a sample (from the true density p) varies around p . The smaller the sample, the bigger this variation.

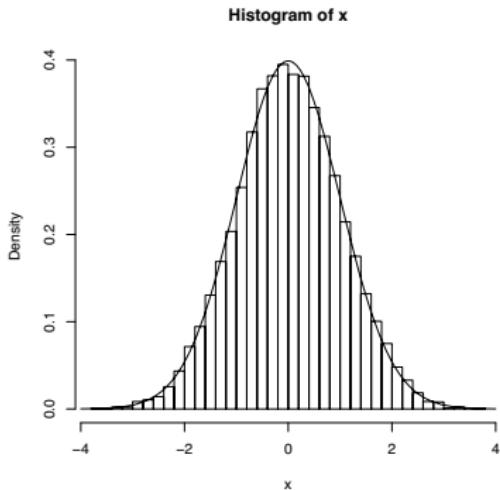


Histogram versus density (2)

For continuous distributions, the true population density can be seen as the smoothed (or limiting as sample size $\rightarrow \infty$) histogram of the population values.

The resemblance between the true $\text{normal}(0,1)$ density and the histogram of a sample of size 10000.

You can think of the population here as consisting of infinitely many values.



Covariance, correlation, sample correlation

can be also comp. from joint dist.

- The covariance between two random variables X and Y is
 $\text{Cov}(X, Y) = E[(X - EX)(Y - EY)] = EXY - EX \cdot EY$

always b.w. -1 and 1!

- The correlation between two variables X and Y quantifies the linear relation between them:

$$\rho = \text{Cor}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} = \frac{E[(X - EX)(Y - EY)]}{\sqrt{\text{Var}(X)\text{Var}(Y)}}.$$

can be also comp. from joint dist.

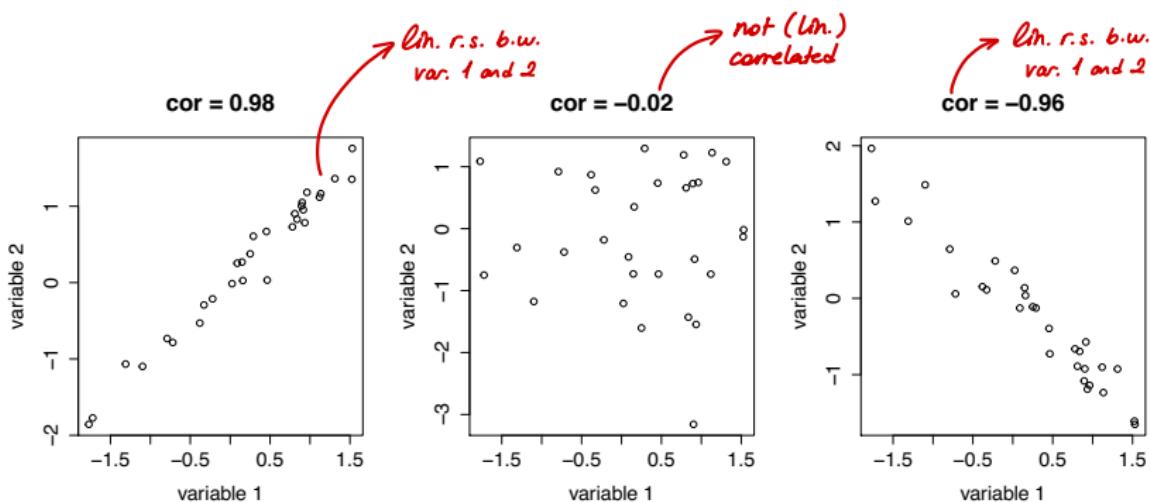
- In practice, the distribution of (X, Y) is almost never known. Instead, one has a sample $(X_1, Y_1), \dots, (X_n, Y_n)$ from the distributions of (X, Y) .
- Then we can compute the sample covariance and sample correlation

$$\hat{c}_{x,y} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}), \quad \hat{\rho} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2 \sum_{i=1}^n (Y_i - \bar{Y})^2}}.$$

estim. of covar. (true covar. can be comp. if dist. is known)

estim. of corr.

Correlation and scatter plot (1)

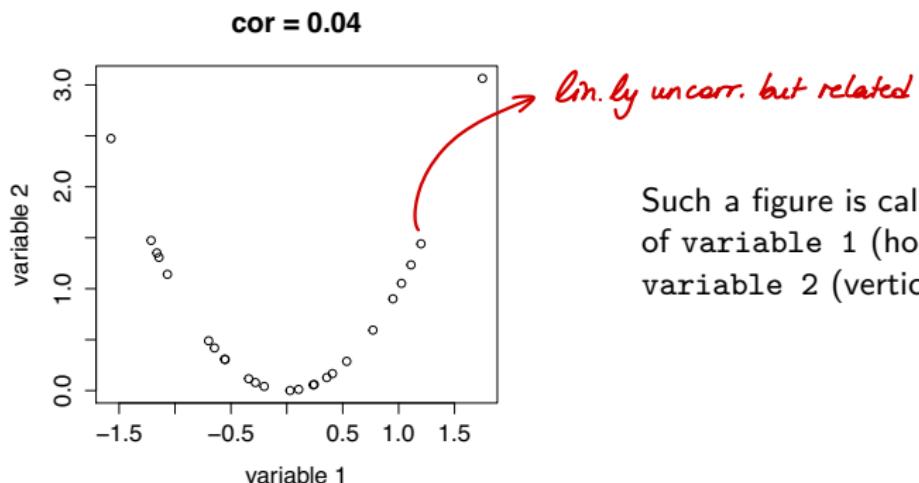


Correlation values:

- ≈ 1 : linear relation (straight line) with positive slope (if $=1$, then perfect linear relation)
- ≈ -1 : linear relation (straight line) with negative slope
- ≈ 0 : no linear relation (but maybe some other relation?)

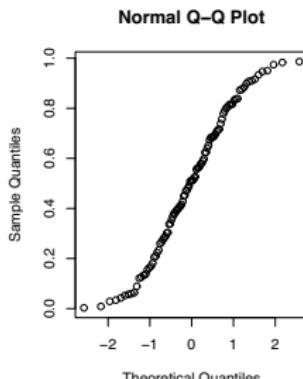
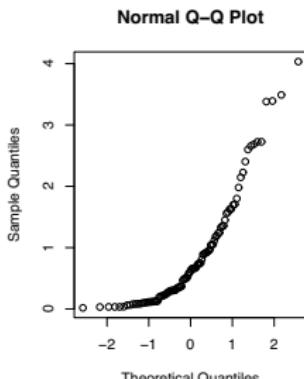
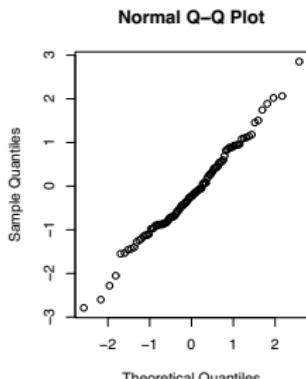
Correlation and scatter plot (2)

Example of two variables that have correlation close to 0, but a clear relation:



QQ-plots

- A QQ-plot can reveal whether data (approximately) follows a certain distribution P (often this is the normal distribution: `qqnorm(x)`).
- It plots the ordered data $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$ versus the quantiles $q_{1/n}, q_{2/n}, \dots, q_{n/n}$ of the distr. P , i.e., $P(X \leq q_\alpha) = \alpha$ for $X \sim P$. (Actually, R uses the quantiles at $\frac{i}{n+1}$ (or another slight adaptation) rather than at $\frac{i}{n}$)
to avoid $q_{n/n} = q_1 = \infty$.
- If $X_i \sim P$, then approx. a fraction $\frac{i}{n}$ of the population should be smaller than the $\frac{i}{n}$ -quantile $q_{i/n}$, i.e., the plot points should follow a straight line.
- If the points are approximately on a straight line, then the data can be assumed to be sampled from P , possibly with different location and scale.



Shapiro-Wilk test for normality

* not reliable when H_0 's not rejecting
 H_0 !

* cannot distinguish v. well if samp. comes from norm. dist.

test whether
this dist. is
normal or not.

Setting: A sample X_1, \dots, X_n from an unknown distribution P .

Hypothesis: $H_0 : P$ is a normal distribution versus $H_1 : P$ is not a normal

Test statistic: for certain constants a_1, \dots, a_n ,

coeff.s are chosen in
such a way that H_0 's poss.
to derive dist of this stat.s under
 H_0 : x_i 's are normal.

$$W = \frac{\left(\sum_{i=1}^n a_i X_{(i)} \right)^2}{\sum_{i=1}^n (X_i - \bar{X})^2}.$$

ordered obs.

$(a_1, \dots, a_n) = m^T V^{-1} / \|V^{-1}m\|$, where m and V are the vector of expected values and covariance matrix of the order statistics of n independent standard normals.

Distribution of W under H_0 : known, but complicated to write down. H_0 is rejected for "small" values of W . It is always the left-sided test.

In R: `shapiro.test(x)`

Note: this test complements the graphical check by a normal QQ-plot.

sample
not
normal
if p 's < 0.05
reject H_0 , else
don't reject H_0
(no strong evid.
against normal.)

Example — expensescrime (1)

The data expensescrime were obtained to determine factors related to state expenditures on fighting criminality (courts, police, etc.). The variables are: state (indicating the state in the USA), expend (state expenditures on fighting criminality in \$1000), bad (number of persons under judicial supervision), crime (crime rate per 100000), lawyers (number of lawyers in the state), employ (number of persons employed in the state) and pop (population of the state in 1000).

```
> expensescrime=read.table("expensescrime.txt",header=TRUE)
> head(expensescrime)

  state  expend   bad  crime  lawyers  employ    pop
1   AK      360   5.1  5877     1749    2796    525
2   AL      498  34.4  3942     6679   13999   4083
3   AR      219  19.2  3585     3741    7227   2388
4   AZ      728  31.3  7116     7535   14755   3386
5   CA     6539 336.2  6518    82001  118149  27663
6   CO      602  25.7  6919    11174   12556   3296
```

Apart from numerical and graphical summaries of the columns separately, we can consider bivariate summaries to see the relation between pairs of columns.

Example — expensescrime (2)

The correlation between all pairs of variables, excluding the first column:

```
> round(cor(expensescrime[,-1]),3)
```

	expend	bad	crime	lawyers	employ	pop
expend	1.000	0.834	0.334	0.968	0.977	0.953
bad	0.834	1.000	0.373	0.832	0.871	0.920
crime	0.334	0.373	1.000	0.375	0.311	0.275
lawyers	0.968	0.832	0.375	1.000	0.966	0.934
employ	0.977	0.871	0.311	0.966	1.000	0.971
pop	0.953	0.920	0.275	0.934	0.971	1.000

lin. corr.s b.w.

- exp and bad → # of prisoners↑ cost↑
- exp and lawyer → # of lawyers↑ cost↑
- exp and employ → # of emp↑ cost↑
- exp and pop → pop↑ cost↑

Ingredients of R-code:

- `expensescrime[,-1]` removes column 1 from `expensescrime`,
- `cor(expensescrime[,-1])` produces pairwise correlations between remaining columns,
- `round(cor(expensescrime[,-1]),3)` rounds the numbers to 3 decimals.

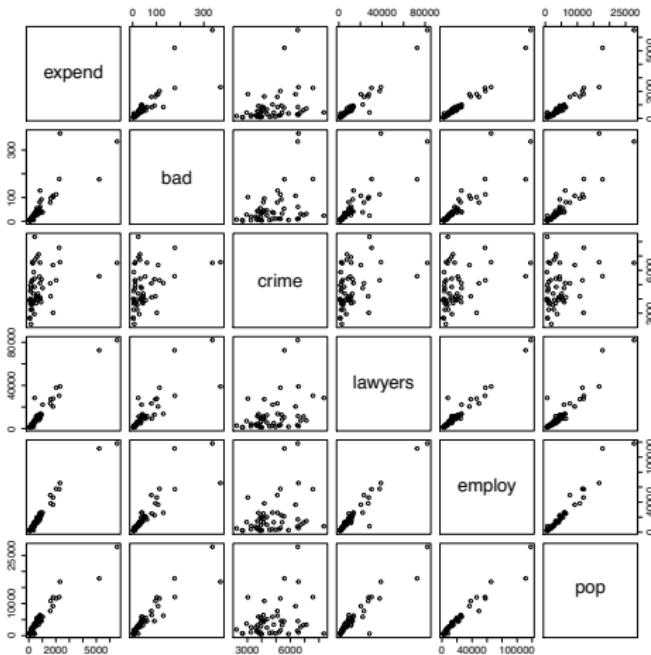
removed
b/c not numeric.

Example — expensescrime (3)

The scatter plots of all pairs of variables, excluding the first column:

```
> pairs(expensescrime[,-1])
```

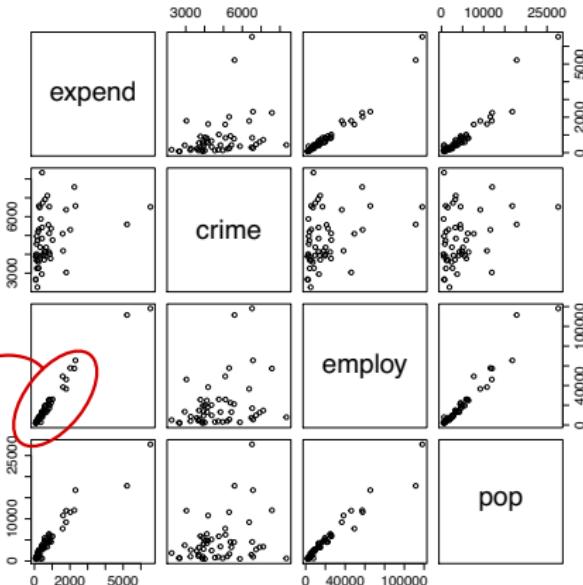
plot 1 col. against the other



Example — expensescrime (4)

The scatter plots of the variables expend, crime, employ, pop:

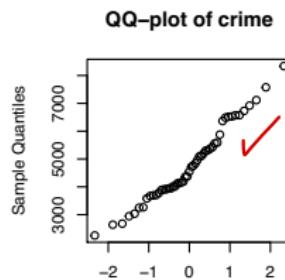
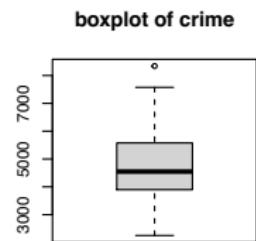
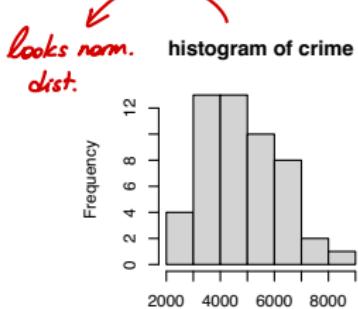
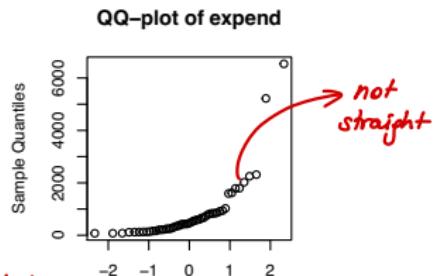
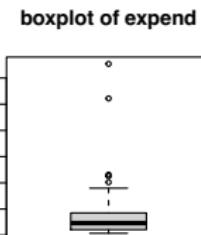
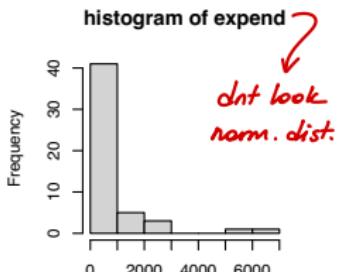
```
> pairs(expensescrime[,c(2,4,6,7)])
```



```
expensescrime[,c(2,4,6,7)] selects columns 2, 4, 6 and 7.
```

Example — expensescrime (5)

Histograms, boxplot and QQ-plots of the two columns (`expend` and `crime`) of the `expensescrime` data. Column `crime` seems to follow a normal distribution.



$$\begin{aligned} & \bullet X \sim N(\mu_1, \sigma_1^2) \\ & Y \sim N(\mu_2, \sigma_2^2) \end{aligned} \quad \left. \begin{array}{l} \text{indep.} \\ \rightarrow aX + bY \sim N(a\mu_1 + b\mu_2, a^2\sigma_1^2 + b^2\sigma_2^2) \end{array} \right\}$$
$$\begin{aligned} E(ax + bY) &= aE(X) + bE(Y) = a\mu_1 + b\mu_2 \\ \text{Var}(a + bx) &= b^2 \text{Var}(x) \end{aligned}$$

Start Lecture 1. Recap basic stat. concepts: estimation, CI, CLT

The sample mean and its distribution, CLT

- The sample mean of a sample X_1, \dots, X_n of sample size n is

estimator $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$; for binomial data $X_1, \dots, X_n \sim \text{Bin}(1, p)$, $\bar{X} = \hat{p}$.

n trials → each try is succ. (1) or fail (0). prob. of succ. is p.

- If $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ -distribution, then $\bar{X} \sim N(\mu, \sigma^2/n)$ exactly.

- When the sample is taken from any distribution with expectation μ and variance σ^2 , \bar{X} still has approximately $N(\mu, \sigma^2/n)$ -distribution (\bar{X} is asymptotically normal). This is the Central Limit Theorem (CLT):

Its dist. converges
to std. norm. dist.

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right), \quad \text{or} \quad \sqrt{n}(\bar{X} - \mu)/\sigma \sim N(0, 1), \quad \text{appr. (for large } n\text{)}.$$

if $n \rightarrow \infty$

- The CLT is a fundamental result of probability theory.
- Example: for binomial data $X_1, \dots, X_n \sim \text{Bin}(1, p)$, $E(X_i) = p$, $\bar{X} = \hat{p}$, $\sigma^2 = \text{Var}(X_i) = p(1-p) \approx \hat{p}(1-\hat{p})$, so that approximately (for large n)

$$\frac{\sqrt{n}(\hat{p}-p)}{\sqrt{\hat{p}(1-\hat{p})}} \sim N(0, 1).$$

*we don't know var.
so we estim. that*

Standardizing the mean

- Any normal random variable $X \sim N(\mu, \sigma^2)$ can be standardized into a standard $N(0, 1)$ -variable by $Z = (X - \mu)/\sigma \sim N(0, 1)$.
- Converse is also true: if $Z \sim N(0, 1)$, then $X = \mu + \sigma Z \sim N(\mu, \sigma^2)$.
- General fact: if $X \sim N(\mu_x, \sigma_x^2)$ and $Y \sim N(\mu_y, \sigma_y^2)$ are independent, then $V = aX + bY + c \sim N(a\mu_x + b\mu_y + c, a^2\sigma_x^2 + b^2\sigma_y^2)$. c disappears
- As $\bar{X} \sim N(\mu, \sigma^2/n)$ (exactly or approximately), standardizing \bar{X} yields

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} = \boxed{\frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \sim N(0, 1)}$$

same as CLT

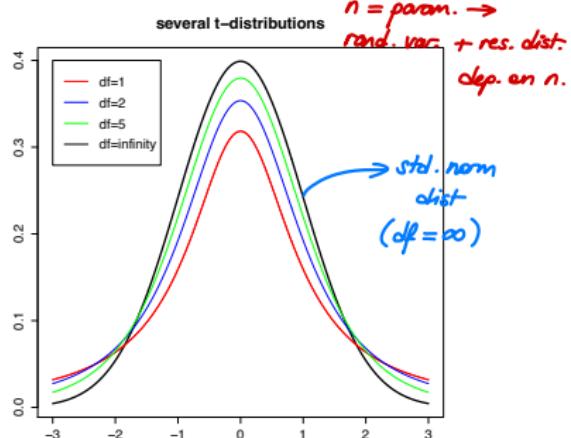
The t -distribution

- In a real data set X_1, \dots, X_n , the population standard deviation σ is **unknown** and needs to be estimated by the **sample standard deviation s** . also σ^2
 - This uncertainty influences the distribution of the resulting statistics $\frac{\bar{X} - \mu}{s/\sqrt{n}} \approx Z \sim N(0, 1)$, which can still be approximated by $N(0, 1)$. CLT holds approx. b/c s converges to σ .
 - If X_1, \dots, X_n is a sample from $N(\mu, \sigma^2)$, then the random variable $T = \frac{\bar{X} - \mu}{s/\sqrt{n}} \sim t_{n-1}$, a t -distribution with $n - 1$ degrees of freedom.
 - $t_{n-1} \neq N(0, 1)$, but $t_{n-1} \approx N(0, 1)$ for **big n** .
- subst.
 s inst. of σ*
- becomes less precise*

For any other generating distribution,

$$T = \frac{\bar{X} - \mu}{s/\sqrt{n}}$$
 is still approximately $N(0, 1)$,

but often t_{n-1} -distribution is used instead. This does no harm to inference on μ as it only leads to more conservative quantiles in testing and confidence intervals.



Estimation – the concepts

- Suppose we assume that our population of interest has a certain distribution with an unknown parameter, e.g., its mean μ or a fraction p .
- A point estimate for the unknown parameter is a function of only the observed data (X_1, \dots, X_n) , seen as a random variable.
- We denote estimators by a hat: $\hat{\mu}$, \hat{p} , etc.
- Examples of point estimates: $\hat{\mu} = \bar{X}$, the sample proportion \hat{p} .
- A confidence interval (CI) of level $1 - \alpha$ for the unknown parameter is a random interval based only on the observed data (X_1, \dots, X_n) that contains the true value of the parameter with probability at least $1 - \alpha$.

prob. that
our rand. interval
contains our true val.

any func.
of this data
will deliv. a
pt. estim. of
smth.

Estimating the mean, CI

var. became smaller = our dist. for \bar{X} is concentr. around μ .

- Recall that $\bar{X} \sim N(\mu, \sigma^2/n)$ for X_1, \dots, X_n from $N(\mu, \sigma^2)$ distribution.
- The upper quantile z_α of the $N(0, 1)$ -distribution is such z_α that $P(Z \geq z_\alpha) = \alpha$ for $Z \sim N(0, 1)$, (in R: $z_\alpha = qnorm(1-\text{alpha})$). Then

$$\begin{aligned} 1 - \alpha &= P(|Z| \leq z_{\alpha/2}) = P\left(\frac{|\bar{X} - \mu|}{\sigma/\sqrt{n}} \leq z_{\alpha/2}\right) \xrightarrow{1-\alpha} \\ &= P\left(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right). \end{aligned}$$

fixed, unknown

- In other words, $\bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}} = [\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}]$ is the confidence interval of μ of level $1 - \alpha$.
- If the standard deviation σ is unknown, we estimate it by s and the confidence interval is based on a t -distribution and the upper t -quantile $t_\alpha = qt(1-\text{alpha}, df=n-1)$ (i.e., $P(T \geq t_\alpha) = \alpha$ for $T \sim t_{n-1}$).
- The t -confidence interval of level $1 - \alpha$ for μ then becomes

$$\bar{X} \pm t_{\alpha/2} \frac{s}{\sqrt{n}} = \left[\bar{X} - t_{\alpha/2} \frac{s}{\sqrt{n}}, \bar{X} + t_{\alpha/2} \frac{s}{\sqrt{n}} \right].$$

*our int. becomes larger
b/c t quant. is bigger than z quant.*

The t -CI's are (nearly) always used, since σ is almost never known in practice. In view of CLT, this can be used also for non-normal data.

Margin of error for the mean

so we want to know how many obs in our sample to ensure that CI is at cert. level.

- The $(1 - \alpha)$ -confidence interval for μ

$$\bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

or

$$\bar{X} \pm t_{\alpha/2} \frac{s}{\sqrt{n}}.$$

bigger n , smaller
CI = more certainty
in CI

- The margin of error is thus $E = z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$ or $E = t_{\alpha/2} \frac{s}{\sqrt{n}}$.
- Remark 1. If we take larger n , the confidence interval will be smaller (shorter), i.e., gaining more accuracy at the same confidence level.
- Remark 2. If σ (or s) is smaller, the confidence interval will be shorter, again yielding more accuracy at the same confidence level.
- Remark 3. If we take bigger α , the confidence interval will be shorter.
Warning: more accuracy at the cost of a lower confidence level.

less
uncert.
in CI

b/c $1 - \alpha$ becomes
smaller

Minimal sample size

- Question: how big should the sample size be in order to obtain a margin of error at most E ? (This is the same as having the CI length at most $2E$.)
- Answer: n must satisfy $z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq E$ or $t_{\alpha/2} \frac{s}{\sqrt{n}} \leq E$, or equivalently

$$\sqrt{n} \geq \frac{z_{\alpha/2}\sigma}{E} \quad \text{or} \quad \sqrt{n} \geq \frac{t_{\alpha/2}s}{E}, \quad \text{so that}$$

we can compute this
if we know the
quantities

$$n \geq \frac{(z_{\alpha/2})^2 \sigma^2}{E^2} \quad \text{or} \quad n \geq \frac{(t_{\alpha/2})^2 s^2}{E^2} \approx \frac{(z_{\alpha/2})^2 s^2}{E^2}.$$

can't compute this
b/c t quant. dep. on
 n . → we subst. z
quant. b/c if n is big,
they are close.

- Remark. For large n we have $t_{\alpha/2} \approx z_{\alpha/2}$ and $s \approx \sigma$. Actually, it makes sense to use $z_{\alpha/2}$ in the second formula instead of $t_{\alpha/2}$, because $t_{\alpha/2}$ depends on (unknown) n as well.

Estimating a proportion, CI, minimal sample size

- We want to estimate a population proportion p , based on a sample $X_1, \dots, X_n \sim \text{Bin}(1, p)$. The point estimate for p is $\hat{p} = \bar{X}$. approx.
- The $(1 - \alpha)$ -confidence interval for p is $\hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$ (based on CLT).
- To ensure a margin of error at most E , the minimal sample size must satisfy $z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq E$ or $n \geq z_{\alpha/2}^2 \hat{p}(1 - \hat{p}) / E^2$.
- Example: trains in time. We want take a sample trains to estimate the fraction p of trains that arrive in time. This fraction was estimated as 0.95. We want a 98% confidence interval for p with length at most 3% (0.03). Question: how many trains should we have in the sample? Answer. A CI length of 3% means $2E = 0.03$ so that $E = 0.015$. Next, $\hat{p} = 0.95$ and $1 - \hat{p} = 0.05$. For a 98% interval we have $z_{\alpha/2} = \text{qnorm}(0.99) = 2.326$. Hence, the minimal sample size must satisfy

$$\text{upper quant. } n \geq \frac{z_{\alpha/2}^2 \hat{p}(1 - \hat{p})}{E^2} = \frac{(2.326)^2 \times 0.95 \times 0.05}{(0.015)^2} = 1142.5.$$

$$1 - 0.98 = 0.02$$

$$0.02 \div 2 = 0.01$$

$$1 - 0.01 = 0.99$$

which is found in R by `qnorm(0.99)^2 * 0.95 * 0.05 / (0.015)^2`. In words: we need at least 1143 trains to ensure a 98%-CI of length at most 0.03.

Exp. design
oooo

Recap probab. theory
oooooooooooooooooooo

Summarizing data
oooooooooooooooooooo

Recap basic stat. concepts
oooooooo●oooooooo

Recap: examples in R
oooooooooooooooooooo

Recap basic stat. concepts: hypothesis testing

Hypothesis testing: the concepts

- Null hypothesis H_0 and alternative hypothesis H_1 about the world.
- A statistical test based on the observed data $X = (X_1, \dots, X_n)$ chooses between H_0 and H_1 . The claim of interest is usually represented by H_1 .
- Precisely, for some test statistic $T = T(X)$ and critical region K , we reject H_0 (and accept H_1) if $\{T(X) \in K\}$ (the strong outcome), otherwise do not reject H_0 (the weak outcome).
- A test statistic $T = T(X)$ summarizes the data $X = (X_1, \dots, X_n)$ in a relevant way. Critical region K is chosen in such a way that $T(X)$ is hardly ever expected to take values in K if H_0 were true.
- In general, to construct a good K we need to know the distribution of $T(X)$ under H_0 . This is usually the main difficulty in constructing tests.
- The test (and test statistic) is not unique. Different tests are possible for the same pair of hypothesis H_0, H_1 , with different performances.

you only need to say H_0 or H_1 is true.

you apply cert func.
to data →
by looking at val.s of this func. you decide which one to confirm/rej.
stat. test

you decide to confirm/rej. H_0 or H_1 by obs. where T stats end up.

if H_0 is true, it should not be in critical region.

Hypothesis testing: p -values

but you compute it w/o // 2nd one: you don't have to know data. // fix α before. but you have to have data to compute T

diff: you fix α beforehand, then you compute crit. region before perf. test (1st one)

- 3 ways to test, say $H_0 : \mu \geq \mu_0$ vs. $H_1 : \mu < \mu_0$, test stat. $T(X)$, level α :

1-1 corresp.

- by checking whether $T(X) \in K_\alpha = \{T(X) < t_\alpha\}$ or not; check if test stat. is in crit. reg.
- by comparing the p -value to α : $p = P(T(X) \leq t) \leq \alpha$ or not;
- by checking whether μ_0 is in the (relevant) $(1 - \alpha)$ -CI for μ or not.

best way

By using p -values is the most common way: e.g., for the realized value $T(x) = t$ and $T \sim t_{n-1}$, check whether $p = P(T \leq t) \leq \alpha$ or not.

p val. =
what you
obs. →

you look at
it from persp.
of H_0 +

if the prob.
of what you
obs. is small
you reject H_0
b/c that means
 H_0 is not true.

- Given observed value t of the test statistic, the p -value is the probability under H_0 of observing a value for T that is at least as extreme as t . A small p -value indicates that the observed data is unlikely if H_0 were true.

- When the p -value is below the chosen significance level α (e.g., 0.05), reject H_0 (strong outcome), otherwise do not reject H_0 (weak outcome).
- If H_0 is rejected, the data are said to be statistically significant at level α .

- By construction, under H_0 , the p -value is like a uniform draw from $[0, 1]$.

Let us show this for our example. Let $p(t) = P(T \leq t) = F_T(t)$ for $T \sim t_{n-1}$, then the (random) p -value is $\tilde{p} = p(T(X)) = F_T(T(X))$, and for any $\alpha \in (0, 1)$, $P(\tilde{p} \leq \alpha) = P(F(T(X)) \leq \alpha) = P(T(X) \leq F_T^{-1}(\alpha)) = F(F^{-1}(\alpha)) = \alpha$.

Example: the one sample t -test(s)

- Data $X_1, \dots, X_n \sim N(\mu, \sigma^2)$. The **t -test** is for testing about μ .

- 1-side.* {
1. $H_0 : \mu \leq \mu_0$ versus $H_1 : \mu > \mu_0$ (`t.test(data, mu=μ₀, alt="g")`)
 2. $H_0 : \mu \geq \mu_0$ vs. $H_1 : \mu < \mu_0$ (`t.test(data, mu=μ₀, alt="l")`)
- 2-side.* ← 3. $H_0 : \mu = \mu_0$ vs. $H_1 : \mu \neq \mu_0$ (`t.test(data, mu=μ₀)`)

- In all 3 cases, at the border of H_0 and H_1 (i.e. for $\mu = \mu_0$), the

test statistic $T = \frac{\bar{X} - \mu_0}{s/\sqrt{n}}$ has t -distribution with $n - 1$ degrees of freedom.

- The **p -value** for observed value $T(x) = t$ of the test statistic is

1. $p = P(T \geq t)$ under H_0 (i.e., assuming that $T \sim t_{n-1}$);
2. $p = P(T \leq t)$ under H_0 ;
3. $p = P(|T| \geq |t|) = 2 \min\{P(T \geq t), P(T \leq t)\}$ under H_0 .

- For testing, say, situation 3, $H_0 : \mu = \mu_0$ vs. $H_1 : \mu \neq \mu_0$, we reject H_0 if

either $|T(x)| > |t_{\alpha/2}|$,

or $p = P(|T| \geq |t|) < \alpha$ under H_0 ,

or μ_0 does not belong to the CI $\bar{X} \pm t_{\alpha/2} \frac{s}{\sqrt{n}}$.

Hypothesis testing: types of errors, power of the test

- Statistical tests $\psi = 1\{T(X) \in K\} \in \{0, 1\}$ make two types of errors:
 - Error of the first kind (type I error): rejecting H_0 while it is true.
 - Error of the second kind (type II error): not rejecting H_0 while it is false.
- It is desirable to construct tests with small probability of type I error P_{H_0} (type I error) ($\leq 5\%$). P_{H_0} (type I error) is called the level of this test.
- P_{H_1} (type II error) depends (among others) on the amount of data.
- The probability of correct (i.e., when H_0 is not true) rejecting H_0 is called the power of the test. Under H_1 , power = $1 - P_{H_1}$ (type II error).
- Different test statistics can yield different statistical power of the test.
- Higher sample sizes typically yield higher power.
- The Neyman-Pearson paradigm: tests with high statistical power are preferred, while controlling the level of the test by a fixed margin (5%).
we are concern. w/ type I error = we want it bound, we are concern. w/ tests where prob. of type I err is control. and of those tests, we choose one that has smallest prob. of type 2 err.

The power of a test is specified for each possibility under H_1 . E.g., if $H_0 : \mu \leq 0$ then the power can be calculated in each $\mu > 0$. A good test (that is, a test based on a good test statistic) has high power in all positive μ -values, relative to other tests.

When H_0 is true, power = P_{H_0} (do not reject H_0) = P_{H_0} (type I error) $\leq \alpha$.

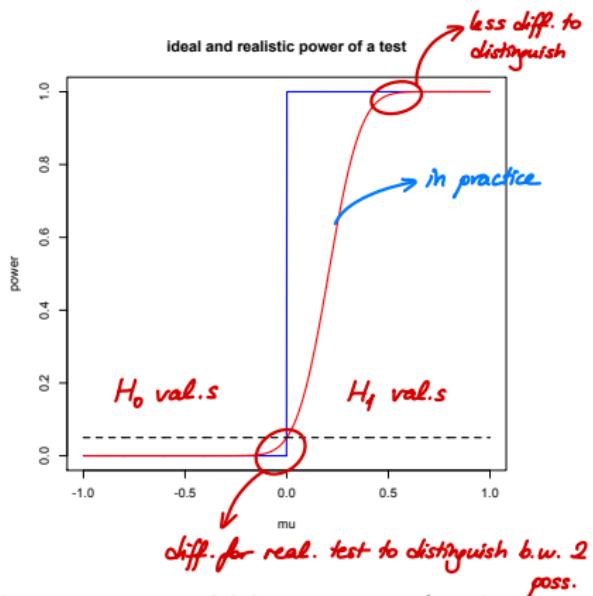
Ideal test and realistic test

The ideal test ψ_{ideal} makes no errors:

- never falsely reject (no error of type I): $\psi_{ideal} = 0$ on H_0 ;
- always reject when H_1 is true (no error of type II): $\psi_{ideal} = 1$ on H_1 .

The power of the ideal test and a realistic test for $H_0 : \mu \leq 0$ vs. $H_1 : \mu > \mu_0$. The dashed line is the level of the test, here 0.05.

One can think of probability of type I error as the proportion of false positives (or the **false positive rate**), and probability of type II error as the proportion of false positives (or the **false negative rate**) in binary classification.



Practical significance

- Statistical significance is about generalization: an observed effect is not due to chance, it should be observed again if a new experiment were performed.
- In practice, this boils down to practical significance which is about the relation between the size of the effect and the available information.

small
num.s

EXAMPLE Suppose that a coin has probabilities $1/2 - 10^{-10}$ and $1/2 + 10^{-10}$ to land HEAD or TAIL.

If we use the coin to decide who will kick-off in a soccer game, then TAIL has a slight advantage, but the difference is negligible. A statistical test based on observing 100 tosses of the coin will not reject H_0 , but a test based on observing 10^{21} coin tosses almost certainly will.

→ not enough data to disting. b.w.
2 very close prob.s

Exp. design
oooo

Recap probab. theory
oooooooooooooooooooo

Summarizing data
oooooooooooooooooooo

Recap basic stat. concepts
oooooooooooooooooooo

Recap: examples in R
●oooooooooooooooooooo

Recap: examples in R

Example of one sample right-sided t -test – crime rate

We want to test whether the mean crime rate (recall column `crime` from the dataset `expencescrime`) is bigger than 4500. Use `t.test` to do the t -test in R:

```
> x=expensescrime$crime; n=length(x); t.test(x,mu=4500,alt="g")
   One Sample t-test
   data: x
   t = 1.5583, df = 50, p-value = 0.06273 > 0.05
   alternative hypothesis: true mean is greater than 4500
   95 percent confidence interval:
   4477.224      Inf
   sample estimates:
   mean of x
   4801.843
```

check ggplot = looks normal

of dof of dist. of stat.

H_0

$H_1: \mu > \mu_0$

we don't reject H_0

** CI = acceptance region → we accept H_0*

if μ is in IT.

vs.

crit. region (rej. reg.)

μ_0 has to be in IT!

test stat.

The R-output gives $\bar{X} = 4801.843$, the value of the test statistics $t = 1.5583$ (or $t=(\text{mean}(x)-4500)/(\text{sd}(x)/\sqrt{n})$), the p -value $p \approx 0.063$ (or $1-\text{pt}(t,n-1)$). Conclude that the mean crime rate is not greater than 4500.

Interestingly, also confidence interval $[4477.224, +\infty)$ is given in the R-output.

But why is Inf in it?

one, right-sided CI

b/c we are testing 1-side hypot.

Point and interval estimation, one sample two-sided t -test

Given a random sample X_1, \dots, X_N from a population with mean μ and unknown variance σ^2 , we wish to estimate μ , construct a CI for it, and to test $H_0 : \mu = \mu_0$ against $H_1 : \mu \neq \mu_0$ for some given number μ_0 , e.g., $\mu_0 = 0$.

2-sided
hypot.

> mu=0.2; x=rnorm(50,mu,1) # creating artificial data
> t.test(x,mu=0)
One Sample t-test

data: x
t = 2.4211, df = 49, p-value = 0.01922 < 0.05 → reject H_0
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
0.05219746 0.56202370
sample estimates:
mean of x
0.3071106

① if you don't put alt. param., by def. you get 2-side. hyp.

if μ_0 belongs to this reg., H_0 is not rejected.

$\mu_0 = 0$ is not in interval.

$$\begin{aligned}1 - 0.95 &= 0.05 \\0.05 &\div 2 = 0.025 \\1 - 0.025 &= 0.975\end{aligned}$$

For this (synthetic) data $X_1, \dots, X_n \sim N(0.2, 1)$, we read off from the above R-output the estimate $\bar{X} \approx 0.31$, the 95% CI [0.052, 0.562], p -value ≈ 0.019 .

$H_0 : \mu = 0$ is rejected because 1) $|t| = 2.42 > |t_{\alpha/2}| = qt(0.975, 49) \approx 2.01$, or because 2) p -value = 0.01922 < 0.05, or because 3) $0 \notin [0.052, 0.562]$.

Standard error and confidence interval

The standard error $\frac{\sigma}{\sqrt{n}} \approx \frac{s}{\sqrt{n}}$ of the estimator \bar{X} is a measure of its precision.
 By CLT, this estimator is approximately normally distributed, hence

Estimate $\pm 1.96 \times \text{Std.Error}$ gives an approx. 95% CI.

s/\sqrt{n}
 The bigger the sample size n , the smaller the standard error and the confidence intervals. The estimates get more precise, as more information is available.

Generate estimates \bar{X} from standard normal samples (i.e., the true $\mu = 0$):

sample size	Estimate	Std.Error
10	0.3564	0.3604
50	0.2198	0.1510
100	0.1098	0.1067
1000	-0.007433	0.031466

In all cases the true value 0 is in the 95% confidence interval

more data you have = better perform.

Estimate $\pm 1.96 \times \text{Std.Error}$.

The margin $m = 1.96 \times \text{Std.Error}$ is based on the asymptotic normality of \bar{X} and the fact that s is a good estimator of σ . If in the CI we use the upper t -quantile $t_{0.025, n-1}$ instead of $z_{0.025} \approx 1.96$, the CI will be bigger (i.e., more "conservative") because always $t_{\alpha, n-1} > z_{\alpha}$, but $t_{\alpha, n-1} \rightarrow z_{\alpha}$ as $n \rightarrow \infty$.

Recap binomial and (appr.) normal tests for a proportion

Setting: $X \sim \text{Bin}(n, p)$, e.g., the number of successes in n trials, p is the success proportion (or the probability of success). We want to test about p .

Hypotheses: $H_0 : p \left\{ \begin{array}{l} \leq \\ \geq \end{array} \right\} p_0$ versus $H_1 : p \left\{ \begin{array}{l} \neq \\ < \end{array} \right\} p_0$ → complement of H_0

Test statistic: X or $T = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$, where $\hat{p} = \frac{X}{n}$.

Distribution under H_0 : $X \sim \text{Bin}(n, p_0)$ (exactly) or $T \sim N(0, 1)$ (approx.)

In R: `binom.test(x, n, p=p0, alt=...)` `prop.test(x, n, p=p0, alt=...)`

exact test = more pref.

Testing for binomial data: example trains on time

- In a (fictive) sample of 100 trains arriving at Amsterdam Central station, we observe a sample proportion $\hat{p} = 0.89$ (89/100) trains arriving in time.
- We want to test whether this is significantly lower than the reported 95% for the Netherlands. Hence, we test $H_0 : p \geq 0.95$ versus $H_1 : p < 0.95$.
- This is a binomial sample with $n = 100$ and p unknown. One can use the exact binomial test binom.test or the proportion prop.test.

The exact binomial test:

```
> binom.test(89,100,p=0.95,alt="l")  
[ some output is deleted ]  
p-value = 0.01147
```

of succ.

less

The approximate proportion test:

```
> prop.test(89,100,p=0.95,alt="l")  
[ some output is deleted ]  
p-value = 0.005808
```

we are
interested
in this

The p -values in both tests < 0.05 (although different). Conclusion: reject H_0 .

Example continued: trains on time

* large d.set w/small dev. → small violation of assump. can lead to significance

Now perform the two-sided test $H_0 : p = 0.95$ versus $H_1 : p \neq 0.95$.

The exact binomial test:

```
> binom.test(89,100,p=0.95)
[ some output is deleted ]
p-value = 0.01739
```

The approximate proportion test:

```
> prop.test(89,100,p=0.95)
[ some output is deleted ]
p-value = 0.01162
```

The p -values in both tests < 0.05 (although different). Conclusion?

we still reject b/c of small p

The influence of the sample size: suppose we had found 890 trains arriving in time amongst 1000 trains:

larger n but same prop.

The exact binomial test:

```
> binom.test(890,1000,p=0.95)
[ some output is deleted ]
p-value = 3.786e-14
```

The approximate proportion test:

```
> prop.test(890,1000,p=0.95)
[ some output is deleted ]
p-value < 2.2e-16
```

$e^{-14} = 10^{-14} = 0.00000000000001$, $3.786e-14 = 0.0000000000003786$. The same deviation from H_0 in more data yields a lower p -value.

H_0 gets rej. w/strong. evid.

It becomes more signif. w/ more data

Tests for a difference in proportions

Setting: X_1 successes in a sample of size n_1 taken from population 1 and X_2 successes in a sample of size n_2 from population 2. We want to test about the difference in population success proportion p_1 and p_2 .

Hypotheses: $H_0 : p_1 - p_2 \left\{ \begin{array}{l} \stackrel{=} {\leq} \\ \stackrel{>} {\geq} \end{array} \right\} 0$ versus $H_1 : p_1 - p_2 \left\{ \begin{array}{l} \stackrel{\neq} {>} \\ < \end{array} \right\} 0$. } take diff. and comp. to 0

Test statistic: $T = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\bar{p}\bar{q}\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$, where $\hat{p}_1 = \frac{x_1}{n_1}$, $\hat{p}_2 = \frac{x_2}{n_2}$, $\bar{q} = 1 - \bar{p}$,

b1nom. test is not avail. $\bar{p} = \frac{x_1+x_2}{n_1+n_2}$ is the pooled sample fraction (the best estimate of p under $H_0 : p_1 = p_2 = p$).

Distribution of T under H_0 : $N(0, 1)$ (approximately).

b/c $X_1 - X_2$ is not binom dist. In R: `prop.test(c(x1,x2),c(n1,n2),alt=...)`

default
//
2 sided, l or g

pair of obs.

of succ.
as pair

Example: compare two proportions of defective items

We test whether the proportions of defective items in two manufacturing processes are (significantly) different. In a sample of 1000 items in process A we find 20 defective items, and in a sample of 1500 items in process B we find 19 defective ones. Question: is there a significant difference in (population) proportions p_A and p_B of defective items for processes A and B?

Thus the sample proportions are $\hat{p}_A = \frac{20}{1000} = 0.02$ and $\hat{p}_B = \frac{19}{1500} = 0.013$, but are they significantly different? We apply the approximate proportion test:

```
> prop.test(c(20,19),c(1000,1500))  
[ some output is deleted ]  
p-value = 0.1989 → relatively big = don't  
reject  $H_0$ 
```

Conclusion? Do not reject
 $H_0 : p_A = p_B$.

Suppose we found the same sample proportions but in larger samples:

```
> prop.test(c(200,190),c(10000,15000))  
[ some output is deleted ]  
p-value = 5.85e-06 → p becomes signif.
```

Now we do **reject** $H_0 : p_A = p_B$.
Why? because samp. size became
larger, p val. became small = significant.

More information (estimates, CI's) can be extracted from the complete R-output.

Two sample t -test

- Given two populations with means μ and ν , we wish to test

$H_0 : \mu \left\{ \begin{array}{l} = \\ \leq \\ \geq \end{array} \right\} \nu$ versus $H_1 : \mu \left\{ \begin{array}{l} \neq \\ > \\ < \end{array} \right\} \nu$. Take a sample X_1, \dots, X_M from the first population and, independently, Y_1, \dots, Y_N from the second.

- The test is based on $\bar{X}_M - \bar{Y}_N$ which is a reasonable estimate for $\mu - \nu$. If it deviates from 0 too much (in the relevant direction), we reject H_0 .
- How different? $\bar{X}_M - \bar{Y}_N$ will not exactly be $\mu - \nu$. The estimation error depends on M and N and the standard deviations of the populations.
- T-statistic: $\bar{X}_M - \bar{Y}_N$ is divided by an estimate $S_{M,N}$ of its standard error.

we have
to normalize
diff.

$$\text{under } H_0, \quad T = \frac{\bar{X}_M - \bar{Y}_N}{S_{M,N}} \sim t_{M+N-2}, \quad S_{M,N}^2 = S_{X,Y}^2 \left(\frac{1}{M} + \frac{1}{N} \right).$$

where $S_{X,Y}^2 = \frac{1}{M+N-2} \left(\sum_{i=1}^M (X_i - \bar{X}_M)^2 + \sum_{j=1}^N (Y_j - \bar{Y}_N)^2 \right)$.

- Then T is compared to the critical value (quantile from t_{M+N-2} -distrib.), or the p-value (computed by using t_{M+N-2} -distribution) is compared to α .

The standard t-test assumes that the two populations are (approx.) normal. If the sample sizes M and N are large, then the test performs well even without this assumption, but the test is unreliable for M, N less than 20.

The quantity $S_{M,N}^2$ is called pooled sample variance.

Two sample two-sided t -test: implementing in R

For example, we test $H_0 : \mu = \nu$ against $H_1 : \mu \neq \nu$ by the two sample t-test:

```
> mu=0;nu=0.5      ↗n
> x=rnorm(50,mu,1);y=rnorm(50,nu,1) #creating artificial data
> t.test(x,y)
    Welch Two Sample t-test
data: x and y
t = -2.4339, df = 96.574, p-value = 0.01677
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-0.85202520 -0.08659066
sample estimates:
mean of x mean of y
0.06552453 0.53483246 }
```

(Handwritten annotations: red arrows point from n to the sample sizes, from the p-value to $< 0.05 \rightarrow \text{reject } H_0$, from the confidence interval to "doesn't contain 0 $\rightarrow \text{reject } H_0$!", and from the sample estimates to "subtr. mean x from mean y = estim. of diff.")

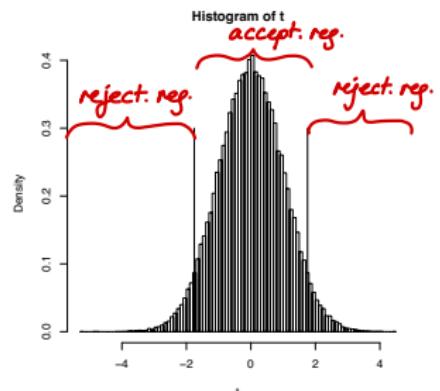
The observed $t = -2.4339$, so that the corresponding p -value is

$$P(|T| > |t|) = 2P(T > |t|) = 2(1 - P(T \leq |t|)) \approx 0.0167.$$

This can be found in the above output, and we could also compute this directly in R as $2*(1-pt(2.4339,98))=0.01674788$. We thus reject H_0 as p -value $\approx 0.017 < 0.05$.

p -value for two sample t -test

We can also evaluate this p -value by simulating from the null hypothesis.



```
> mu=nu=0; t=numeric(100000)  
> for(i in 1:100000){x=rnorm(50,mu,1);y=rnorm(50,nu,1);t[i]=t.test(x,y)[[1]]}  
> sum(abs(t)>=abs(-2.4339))/length(t) ##cf. 2*(1-pt(2.4339,98))=0.01674788  
[1] 0.01744
```

simulations are used when you don't have analytical formulas

We generate a population of T -values under H_0 by repeating the sampling. The p -value of the observed value t is approximately the fraction of this population that is bigger than $|t|$ or smaller than $-|t|$.

we compute
prop. of times we
end up in
ref. reg. = p -val.
very close to p -val.

perform t -test many times, each w/a new sample

Different test statistics

EXAMPLE The **t-test** is for testing the population mean μ of a **normal** population, $H_0 : \mu = \mu_0$. Given a sample X_1, \dots, X_n , the test statistic is

$$T = \frac{\bar{X} - \mu_0}{S_X / \sqrt{n}}, \quad \text{where } \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i, \quad S_X^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

When T is **very different** from 0, reject H_0 . The **critical value** for T that acts as border between rejecting and not rejecting H_0 is based on the distribution of T under H_0 . For t-test, this distribution is the t_{n-1} -distribution.

EXAMPLE For testing $H_0 : \mu = 0$ we can as well use the **sign test**. Given a sample X_1, \dots, X_n from the population, the test statistic for the sign test is

$$T = \#(X_i < 0).$$

① if $\mu = 0 \rightarrow$ approx. half of our sample val.s should be < 0 .

If T is very different from $\frac{n}{2}$, we reject H_0 . The critical value comes from the $\text{Bin}(n, \frac{1}{2})$ -distribution, the distribution of number of heads in throwing n times a fair coin.

Comparing powers of different tests

Assume we have a normal sample and test $H_0 : \mu = 0$ using the t-test and the sign test. We can compare the power in $\mu = 0.5$ of the two tests by simulation.

Recall that always power = $P_{H_1}(\text{reject } H_0)$.

we compute the prob. of rej. H_0
under assump. that actual situation is H_1 .

```

> B=1000; n=50 sample size under assump. that actual
# of times
> psign=numeric(B) ## will contain p-values of sign test
> pttest=numeric(B) ## will contain p-values of t-test
to
> for(i in 1:B) {
+   x=rnorm(n,mean=0.5,sd=1) ## generate data under H1 with mu=0.5
+   pttest[i]=t.test(x)[[3]] ## extract p-value
+   psign[i]=binom.test(sum(x>0),n,p=0.5)[[3]] } ## extract p-value
> sum(psign<0.05)/B power of tests → test is good if power is high/close to 1.
[1] 0.746
> sum(pttest<0.05)/B
[1] 0.937

```

The power in $\mu = 0.5$ for the t-test (0.937) is higher than for the sign test (0.746) when we reject for p -values smaller than the level 0.05. Why? Because for normal data, the t-test has better performance than the sign test.

To finish

We discussed

- ① course organization
- ② experimental design
- ③ recap probability theory and basic statistics
- ④ recap: examples in R

Study the exam to **test your prerequisite knowledge** and Assignment 0 (not to be submitted) to learn how to make assignment reports to submit.

Next time bootstrap methods, one sample tests.