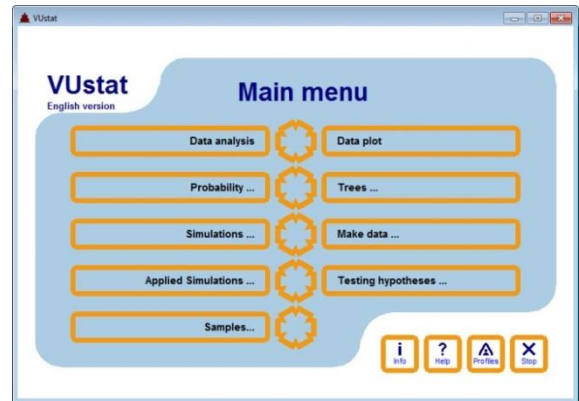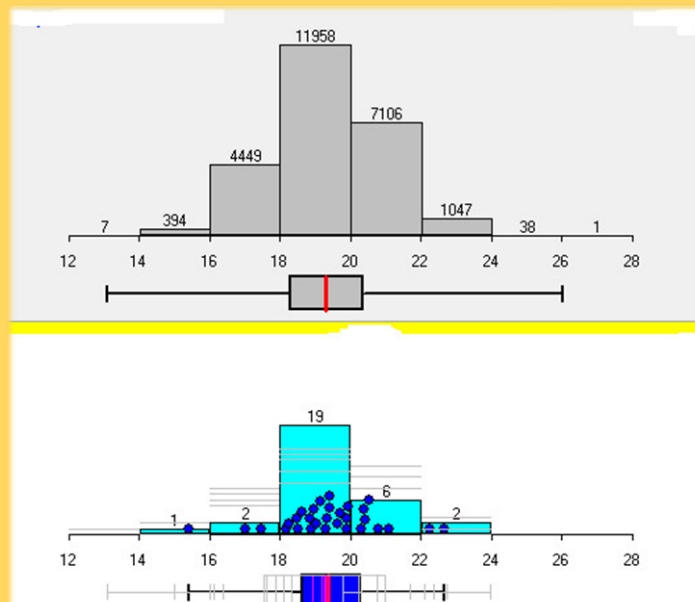# Visual Understanding

## educational software for maths

**VUstat** is a program with a wealth of possibilities to illustrate ideas and concepts in statistics education. VUstat is a statistical package for different levels of education for students and teachers. It contains data, tools and simulations to enhance the ability to teach statistics and probability to students of various levels. VUstat consists of many different modules. Many parts of the program can be turned off by a special tool called *profiles*, so that the student is not confronted with too much complexity at lower levels.
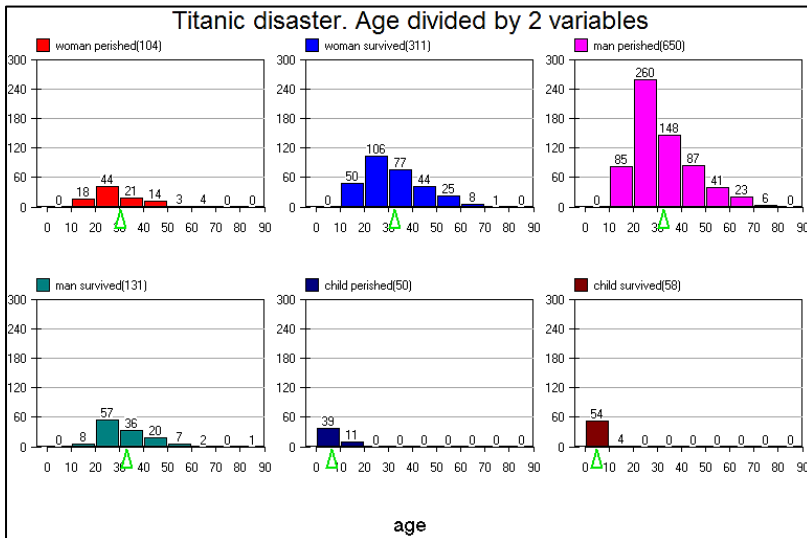
*Visual Understanding software*
*www.vusoft.eu*



*"You don't have to eat the whole ox to know that the meat is tough"*
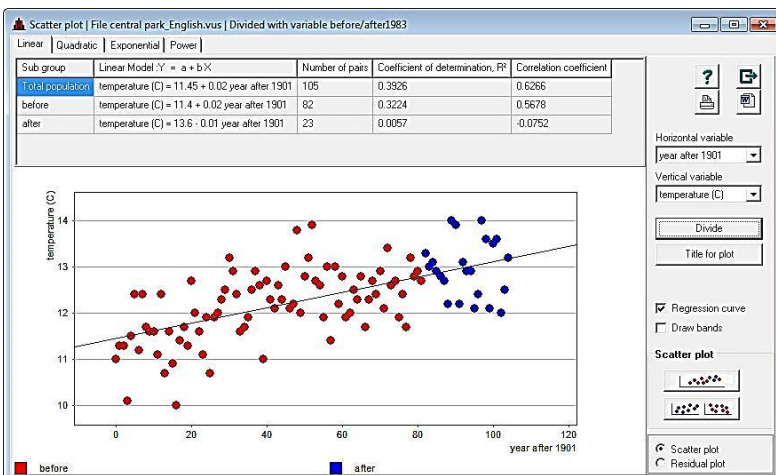*(David Moore about sampling)*

## Data analysis

Data analysis enables you to make the more general tables and graphs of raw data input. It has advanced options, but its main purpose is to make the analysis of data easier to understand for the student who is neither a statistics expert nor a computer whiz. The selection and dividing of data into subgroups is made straightforward.



*Histogram from a dataset of the passengers of the Titanic. The diagrams show passengers age divided by the variables sex and survived/perished. In the diagrams we also show mean values (optional).*

| perished/survived | passenger class | man/woman/child | | | Total |
|---|---|---|---|---|---|
| | | woman | man | child | |
| perished | passenger class 1 | 4 | 118 | 1 | 123 |
| | passenger class 2 | 12 | 146 | . | 158 |
| | passenger class 3 | 88 | 391 | 49 | 528 |
| | Subtotal | 104 | 655 | 50 | 809 |
| survived | passenger class 1 | 139 | 57 | 4 | 200 |
| | passenger class 2 | 81 | 14 | 24 | 119 |
| | passenger class 3 | 91 | 60 | 30 | 181 |
| | Subtotal | 311 | 131 | 58 | 500 |
| Total | | 415 | 786 | 108 | 1309 |

*Contingency table (cross table) divided by variable passenger class. The table can be copied and pasted in text format into other applications, e.g. Microsoft Word®.*
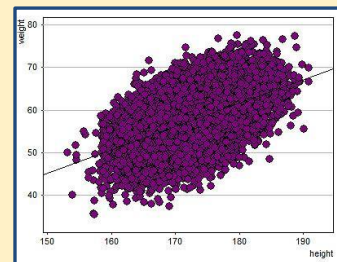


*Climate data for more than 100 years. Scatter plot with data divided into two variables. Regression line and correlation coefficient are calculated.*

---

### Data analysis
✓ Import of data from the Web
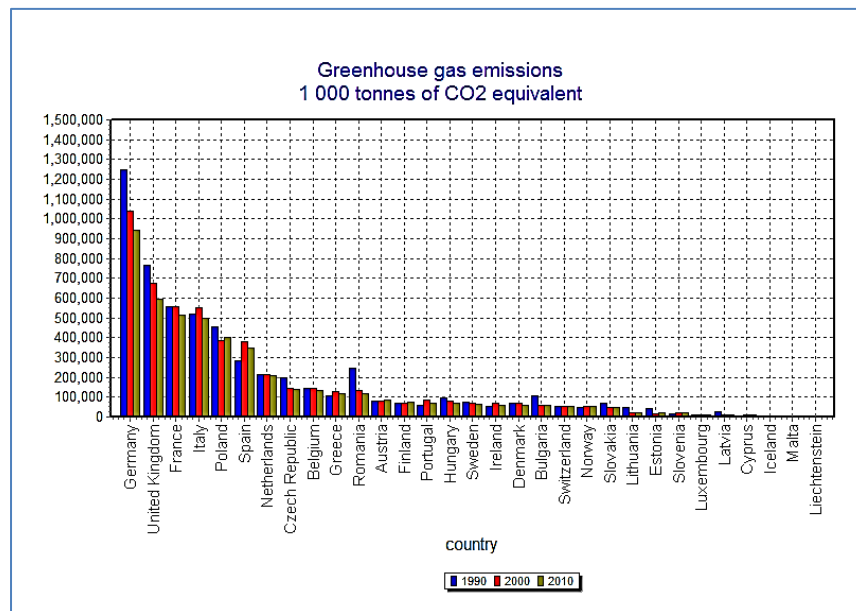✓ Import/export from and to Excel®
✓ Built-in calculator for recalculating data



✓ Sorting data
✓ 5 different types of variables: text, integers, real numbers, labels and multiple response
✓ Options for filtering data
✓ Drawing samples from a dataset
✓ 8 different types of diagrams: Bar chart/Histogram, Dotplot, Pie chart, Frequency polygon, Box plot, Scatter plot, Time series, Lorenz curve
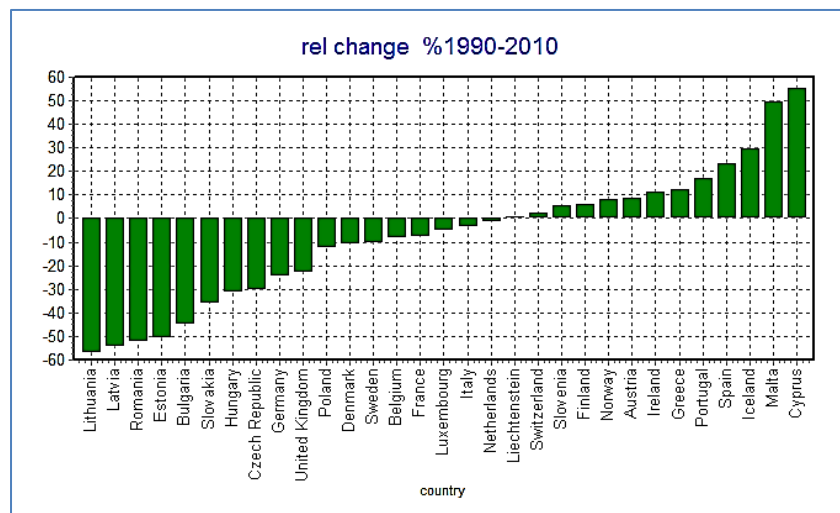


✓ 3 types of tables: Frequency table, Cross table, Stem and leaf plot
✓ Data can be divided for up to two variables.
✓ Aggregation of data

## Data plot

The purpose of this feature is to make graphs of ordered data. There are built-in features for design of diagrams. Some of the tools in the module Data Analysis are here too. Dataplot is especially useful for analyzing *frequency tables*. Traditionally, much teaching in the lower grades is concentrated on frequency tables. The Dataplot module allows students to do more than drawing graphs; it allows them to concentrate on the more useful task of comparing and interpreting data.
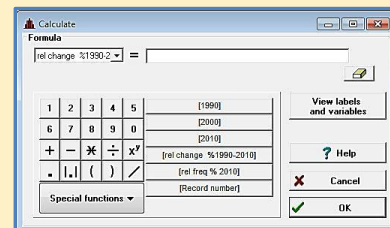


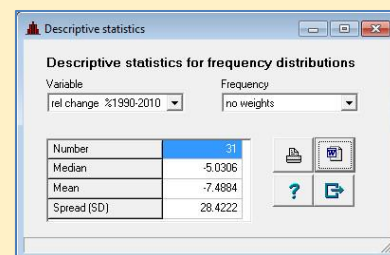*Emissions from all European countries (no figures for Croatia), from highest to lowest.*



Once the files have been opened in modules Data plot and Data analysis you can toggle between them with a shortcut button.
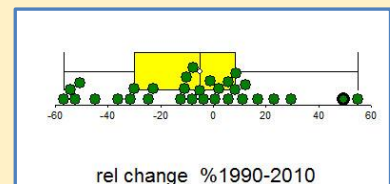


The Data table with the first four variables imported as an Excel file (from Eurostat). The last two variables are calculated with the built-in calculator.
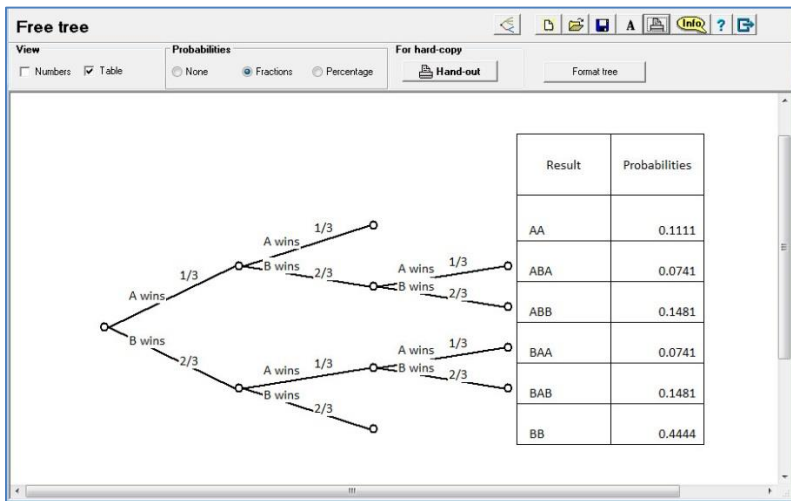


Calculator



Summary statistics for the variable *rel change.*



The file can of course also be opened in the module Data analysis. Here is a Dotplot and a Box plot of the distribution of the variable *rel change*.
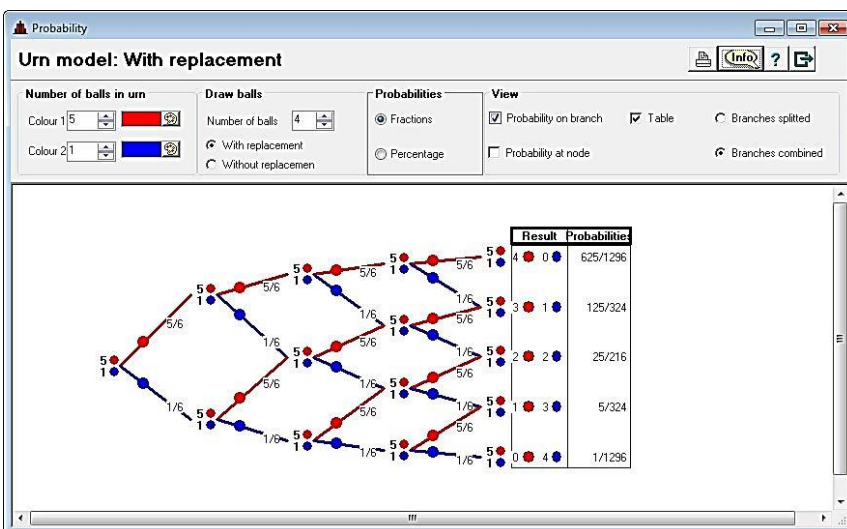
# Trees

A **tree diagram** visualizes a probability tree from different points of view. Both a regular tree and a free tree can be constructed.



*Free tree: Probability tree for a tennis match in 3 sets.*

The module **Urn** shows the probabilities when balls are drawn from an urn with or without replacement. A lot of settings, as you can see on the screen.



*A demonstration of the famous Chevalier De Mere´s problem. The screen shows the possible outcomes and the probability of getting at least one "6" in four throws of a single 6-sided die.*
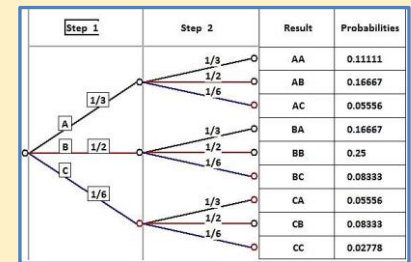
## Table and tree

Table and tree. Data in cross tables (absolute or relative values in rows or columns) can be represented in trees and diagrams in different ways.
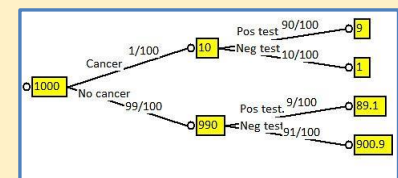




### Trees
✓ Two different types of probability trees: a regular tree and a tree which you can design.



*Regular tree* with three branches and two steps.

Assume you conduct breast cancer screening using mammography in a certain region. You know the following information about the women in this region:
  - The probability that a woman has breast cancer is 1 % (prevalence)
  - If a woman has breast cancer, the probability that she tests positive is 90% (sensitivity)
  - If a woman does not have breast cancer, the probability that she nevertheless tests positive is 9% (false-positive rate)



From the tree diagram it is much easier to calculate the probability that a woman who tests positive has breast cancer if we work with *natural frequencies*.

$$P = \frac{9}{9+89} \approx 9\ \%$$

# Probability

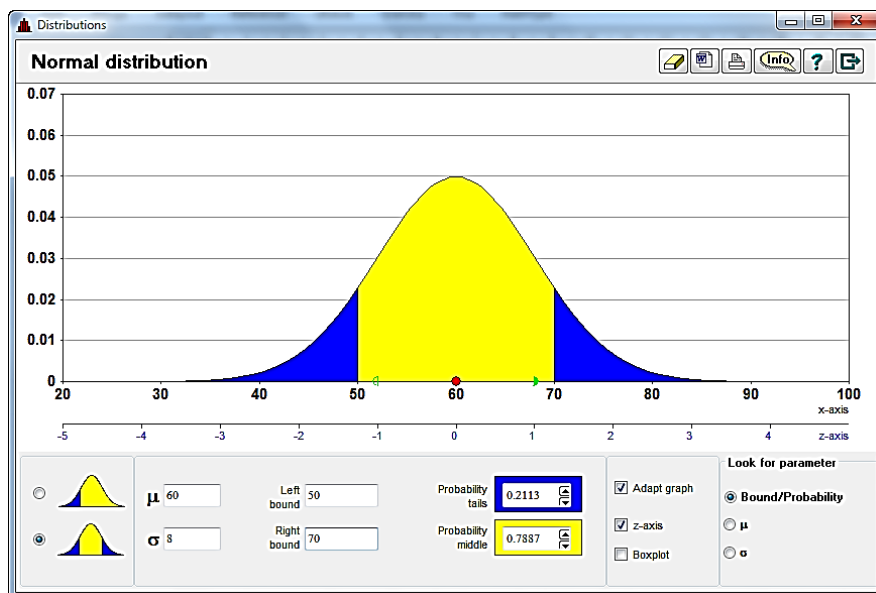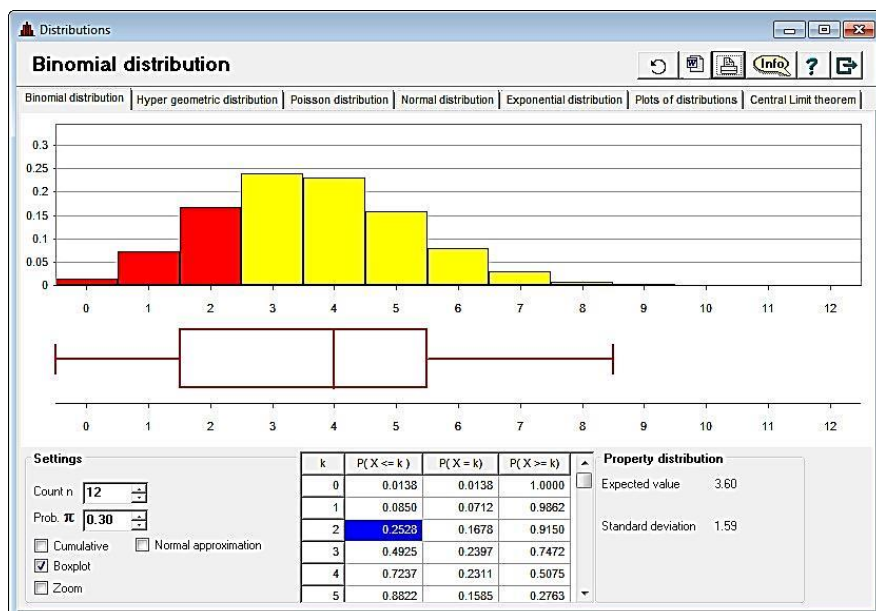This section contains several options that are useful in counting problems and understanding of probability. The modules in this section are Distributions, Galtons board, Grid, Combinatorics and Frequency grid.

## Distributions

The binomial, Poisson, hypergeometric, and normal distributions can easily be manipulated by students using VUstat. The central limit theorem is demonstrated with skewed and normal dice and you can compare plots of different distributions.
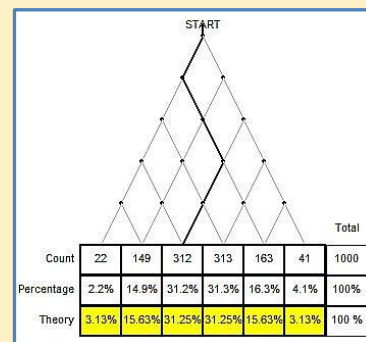


*Normal distribution. You can change parameters by changing values in the frames or by dragging the bounds between the blue and yellow areas, the red points representing μ or the green arrows representing σ.*



*Binomial distribution. The number of trials is 12. The probability of success is 0.3 and we can see the probabilities for success in three different ways in the table. If we highlight a cell in the table we can see the corresponding bins in red color in the diagram. Here we have highlighted P(X≤2).*

In this module there are two small applications for the binomial distribution: *Galton's board* and *Grid*.



Galton's board





In the grid model a "success" is taking a step to the right. The number of trials is the total number of steps taken in the grid to get from start to an end point. You can see the result in different views. Here we show a probability table with formulas.

## Combinatorics

The module Combinatorics deals with things that can be enumerated. For example, how many different ice creams with three flavors there are if you have six different flavors to choose from or how many handshakes take place if ten people meet and everyone shakes hands with everyone else.

We illustrate concepts as variation, permutation and combination in three views and we use colors or characters as objects.

An example: Suppose we have two points in a plane. How many paths are there that pass only through the lines of the grid? Only horizontal (r) or vertical (u) movements are permitted.

All paths between the points have the same length, 3 right + 2 up. We can calculate the number of all the paths between the points as combinations: from five characters, you choose two.

$$\binom{5}{2} = 10$$

The possible ten paths are:
uurrr, ururr, urrur, urrru, rrruu
rruru, rruur, ruurr, rurur, rurru



The path here is rurru.

*Option **Groups** gives you number of combinations. You can see all variations in a tree and all possible groups in a table.*



*If we choose 3 out of 5 we get the same result. Here we have colors as objects.*





From option *Tables*. Here you can show tables for both variations and combinations. The probability that five random cards drawn (out of 52) will give some unique hand where order does not matter is

$$\frac{1}{2\,598\,960}$$

# Simulations

In this module, a lot of simulations can be run. These include classical simulations, such as throwing coins and dice. There are also a lot of simulations with random numbers. In addition, we have a package of *Applied simulations*, for instance the effect of smoking on life expectancy, the effect of overselling of airline tickets on airline profit, the classical game Roulette and waiting times in a fast-food shop. There is also a module where students can create their own experimental data. Theses data can also be further analyzed.
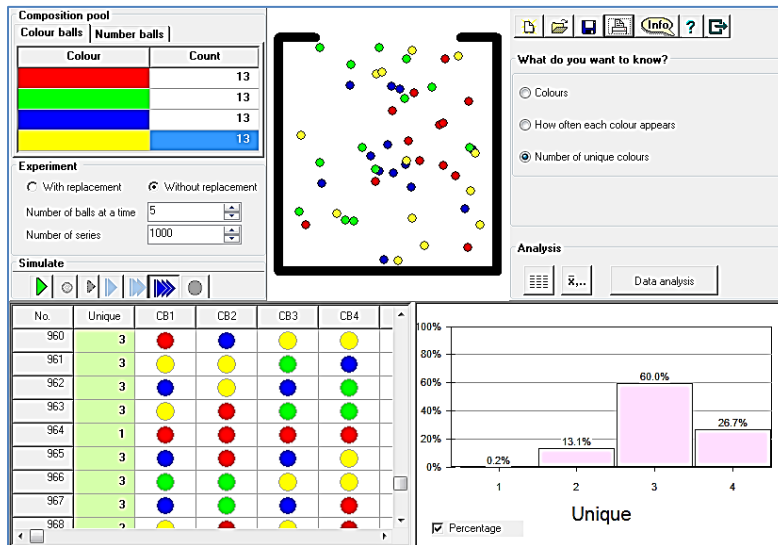
## Ball pool

The ball pool is a random generator based on the functionality of an urn model. You can choose from two options. One model shows a vase with colored balls. In the other model balls are numbered as in a lottery or bingo machine. The simulation process contains several steps.



*In the simulation above we use a deck of cards where we draw five cards. How many different colors do we get? The result shows that we have "1 unique",which means five cards of the same suit, in 0.2 % of the simulations.*

### 1. Balls in the pool
You choose a model and fill in which balls should fill the pool.
In *Color balls,* you give color and number of each color. In *Number balls,* you fill in the numbers. More times the same number means more of the same balls. You can see the balls displayed in the pool.

### 2. Picking balls
You fill in:
- How the balls are to be picked.
- How many balls at a time: this is actually the size of the sample. You can choose *with* or *without* replacement.
- How many times you want to run the simulation.

### 3. Simulation
First have a small number of simulations (e.g. ten) and establish whether you really obtain what you expect.

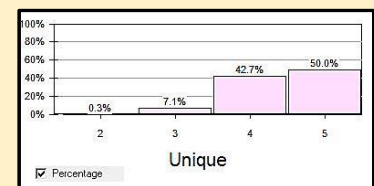### 4. What do you want to know?
Often you want to view the outcome variable regarding a certain problem. Before that, you can choose from several options. With buttons on the right below the chart, you can scroll through the variables.

### 5. Analyze the data using the tools
The tools available are histogram, frequency table, and statistics. For more advanced problems you can send the dataset to the Data analysis module.



A simulation with five balls at a time from this set of numbers (ranks in a deck of cards) gives the result you can see in the bar chart.



In about 43 % of the simulations we have four unique colors, which means that we have a "pair" or two cards of the same rank.

**Another example:**
**What about the surplus of boys?**
In a large hospital 200 children are born each month. In a smaller hospital it´s about 50. Assume that the probability of giving birth to a boy or girl is equally large. *In which hospital is it more likely that more than 60 % of newborns in a given month are boys*?



A simulation for the small hospital gives the result
8.3 %+1.5 % + 0.2 % = 10.5 %
A simulation for the larger hospital gave the result 0.3 %.

## Random numbers

In this module there are a lot of settings for simulations of random numbers. You can see the result in a table and a graph. You can also analyze the result further by pressing the graphs icons or the button "Analyze results" in the upper right corner.

Another option is to use extra formulas to create new variables.



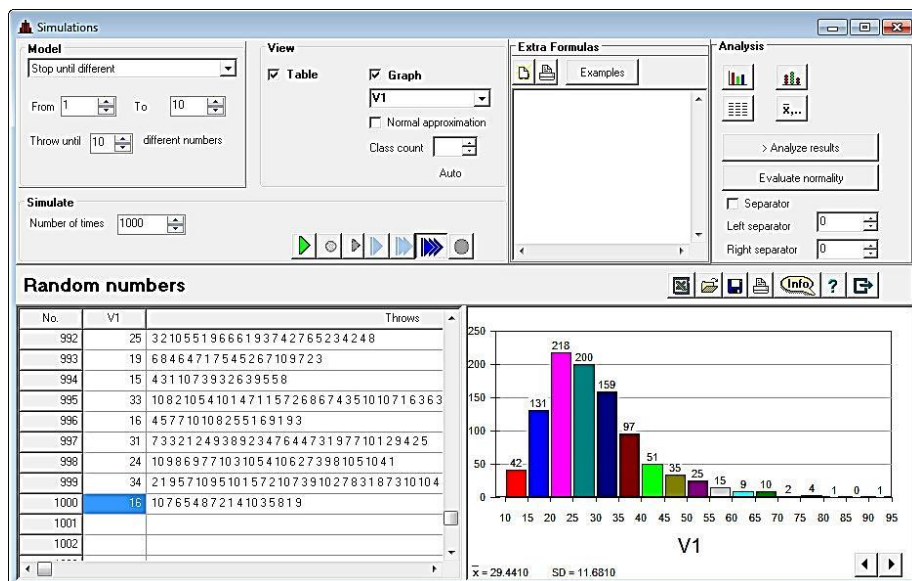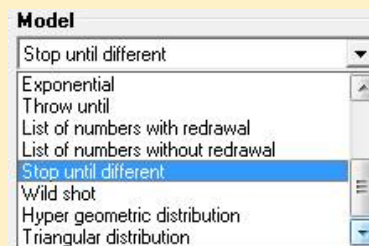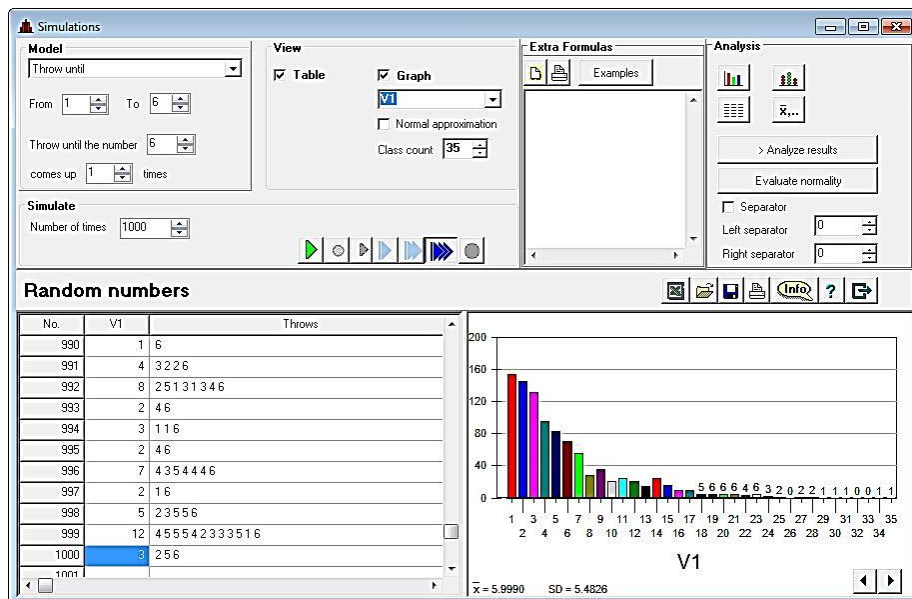*A simulation with Random numbers module and the choice "Stop until .. different". It illustrates the famous "Coupon collectors problem". We simulate a sequence of numbers from 1 to 10, and we want to know how many times we have to draw until we got all 10. The simulation was done 1000 times and we can observe the result in the diagram.*

Another example: You throw a dice. How many times must you on average throw until six comes up? Now, we use the model "Throw until ". We see that the mean becomes 6 times. If the question was "In which throw do you guess that a six comes up?" we would respond differently. Look at the diagram.





There are a lot of options to generate random numbers.



You can also create new variables by creating extra formulas. The diagram below shows the distributions for # of sixes when you throw a dice three times and run the simulation 1000 times.

## Coins

In the module Coins you can throw one or two coins. With one coin you can see how the relative frequency stabilizes after many throws. When you have two coins you will see the frequencies for zero, one or two heads.



## Dice

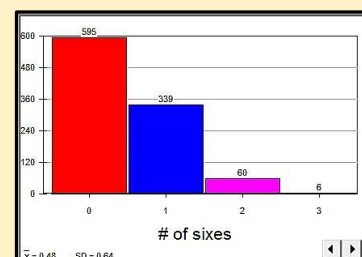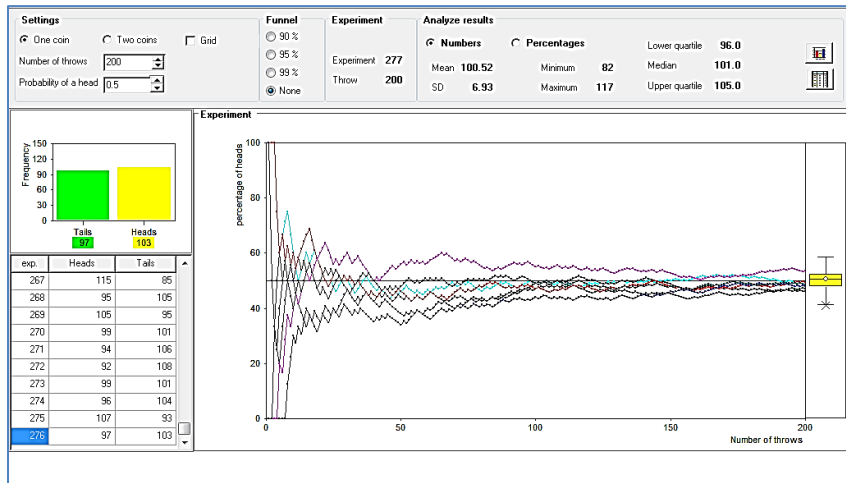This simulation consists of a number of throws with one, two or three dice. If you tick the box "Show sum of all experiments" you also get a graph which shows a summary for all performed experiments. We can see that "7" has the highest frequency.





*Summary for a simulation with three dice. 1000 dice were thrown in each experiment and the simulation was performed 80 times. In the table to the right you can see the number of outcomes for different sums. Observe the symmetry in both the diagram and the table.*

## Coins



Simulation of 100 throws with two coins . The bars are frequencies for zero, one and two heads.



Here is a histogram for the distribution of number of tails for 200 throws with one coin.

| # of dots | # of outcomes |
|-----------|---------------|
| 3 | 1 |
| 4 | 3 |
| 5 | 6 |
| 6 | 10 |
| 7 | 15 |
| 8 | 21 |
| 9 | 25 |
| 10 | 27 |
| 11 | 27 |
| 12 | 25 |
| 13 | 21 |
| 14 | 15 |
| 15 | 10 |
| 16 | 6 |
| 17 | 3 |
| 18 | 1 |

**Raindrops**

When it´s raining, the drops hit the ground evenly or unevenly distributed. Rain can be regarded as a random process. You can simulate the process with a surface where rain falls completely randomly. Because the drops fall randomly, they sometimes fall "on each other".

A good application of this process is as folllows: "A baker bakes a large dough which contains 400 raisins. From the dough she then bakes 100 buns. How great is the probability that you pick a bun with no raisins at all? "

This simulation is shown in the screenshot below. The grid has 10 columns and 10 rows, e.g. 100 cells. We repeat the simulation 1000 times.
The mean after 1000 simulations is 1.83. The probability that you will not get any raisins in your bun then becomes 1.84/100 = 0.0184 or about 2 %. The probability that you will get exactly 4 raisins is 0.196. The expected value is of course 4.
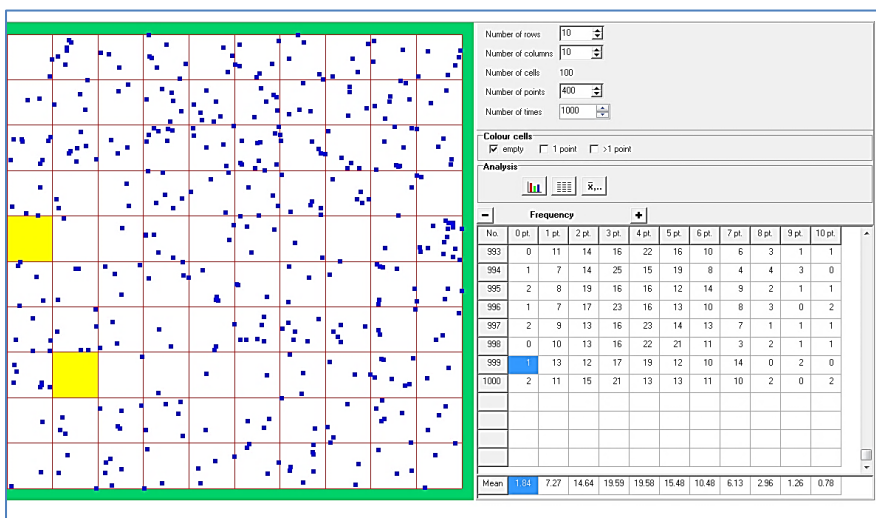
| No. | 0 pt. | 1 pt. | 2 pt. | 3 pt. | 4 pt. | 5 pt. | 6 pt. | 7 pt. | 8 pt. | 9 pt. |
|-----|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 989 | 3 | 7 | 14 | 22 | 15 | 16 | 11 | 5 | 5 | 2 |
| 990 | 1 | 11 | 10 | 26 | 16 | 13 | 11 | 4 | 6 | 1 |
| 991 | 3 | 5 | 14 | 24 | 19 | 11 | 11 | 7 | 5 | 1 |
| 992 | 1 | 10 | 14 | 18 | 23 | 12 | 9 | 7 | 2 | 3 |
| 993 | 0 | 11 | 14 | 16 | 22 | 16 | 10 | 6 | 3 | 1 |
| 994 | 1 | 7 | 14 | 25 | 15 | 19 | 8 | 4 | 4 | 3 |
| 995 | 2 | 8 | 19 | 16 | 16 | 12 | 14 | 9 | 2 | 1 |
| 996 | 1 | 7 | 17 | 23 | 16 | 13 | 10 | 8 | 3 | 0 |
| 997 | 2 | 9 | 13 | 16 | 23 | 14 | 13 | 7 | 1 | 1 |
| 998 | 0 | 10 | 13 | 16 | 22 | 21 | 11 | 3 | 2 | 1 |
| 999 | 1 | 13 | 12 | 17 | 19 | 12 | 10 | 14 | 0 | 2 |
| 1000 | 2 | 11 | 15 | 21 | 13 | 13 | 11 | 10 | 2 | 0 |
| Mean | 1.84 | 7.27 | 14.64 | 19.59 | 19.58 | 15.48 | 10.48 | 6.13 | 2.96 | 1.26 |

| 0 pt. | Freq. |
|-------|-------|
| 0 | 145 |
| 1 | 303 |
| 2 | 270 |
| 3 | 170 |
| 4 | 81 |
| 5 | 24 |
| 6 | 6 |
| 7 | 1 |
| Total | 1000 |

You can analyze all your data in a bar chart, frequency table or as summary statistics. Above is a table for the variable "*0 pt*".
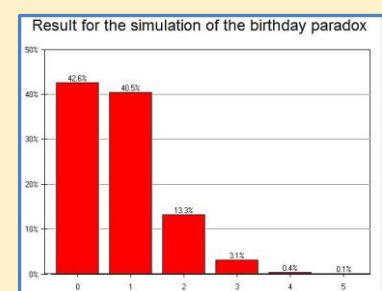


We also show a simulation of the famous *birthday paradox*. The grid is divided into **365** cells (73 × 5). Then let **25** drops fall every time we simulate. We want to do a simulation which calculates the probability that *two or more students* have the same birthday in a group of about 25 students. We repeat this **1000** times.

In the screen to the right we have all the settings. We show the columns 0 pt, 1 pt and > 1pt. It is the *last column* we are interested in here. For the analysis we press the button for the bar chart. We are interested in the variable *> 1pt*. From the diagram, we can see that in 58.6 % of the simulations we get two hits or more. The simulation calculates that the probability that two or more have the same birthday in a class of 25 students is approximately 58.6 %. The theoretical probability is

$$1 - \frac{364}{365} \cdot \frac{363}{365} \cdot \frac{362}{365} \cdot .... \cdot \frac{342}{365} \cdot \frac{341}{365} \approx 0.57$$



*The settings for the birthday paradox*



In the bar chart can we see that in 57.64 % of the simulations we have *two or more* persons who share the same birthday.
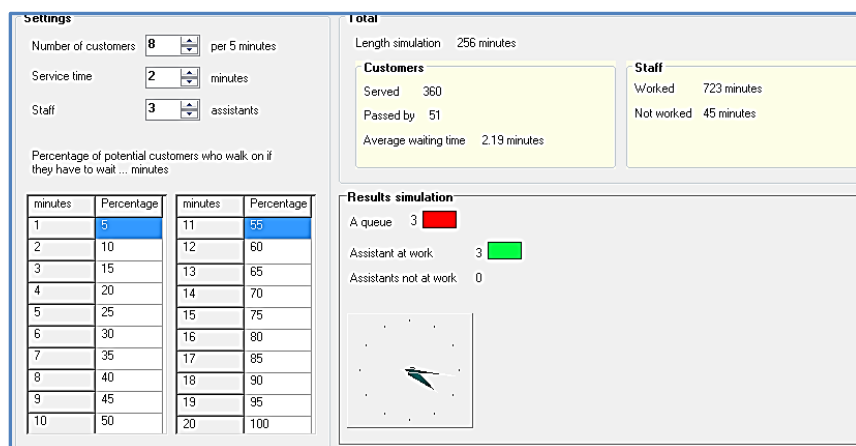
## Applied Simulations

This module contains a number of simulations of specific situations which occur in practice and in games. Examples are for instance the effect of smoking on life expectancy, the effect of overselling of airline tickets on airline profit, the game Roulette and waiting times in a fast-food shop. Here are two examples of these simulations.

### Example Fast-food shop

Everyone has experienced having to wait in a queue at a checkout or at a booking office. Queuing is a problem that has to do with capacity (e.g. the number of assistants) and costs. If the capacity is large enough, nobody has to wait. But this is an expensive solution. If the capacity is too low, potential customers walk away, and this costs money too.

This simulation gives you the opportunity to study this problem. In this simulation customers arrive and are served inside by the staff, or they have to wait. If customers have to wait too long, some will walk away.
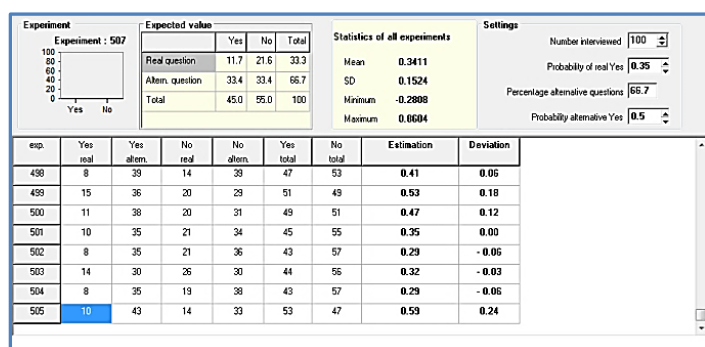
There are several variables that can be chosen. In the simulation you can see how many customers are served and for how much of the time the staff is at work.



### Example Randomized response

Randomized response is a technique which tries to obtain answers to sensitive questions like "*Have you ever taken illegal drugs*". The simulation shows how this technique works for different values of the parameters. The parameters are:

- The probability of "yes" with the real question; in practice this is the value you are looking for.
- The number of persons interviewed
- The percentage of people answering the alternative question
- The probability of "yes" for the alternative question





*Simulations with the same settings give different results.*

**How Randomized response works:**
You get a bag with three cards. Then you randomly draw a card, which only you can see, and respond in accordance with what you can read on the card:

First card: "say YES"
Second card: "say NO"
Third card: "Do you use illegal drugs?"

On average,If you are honest, one-third will answer YES, one-third will answer NO and one-third will respond if they use illegal drugs.

| Expected value | Yes | No | Total |
|---|---|---|---|
| Real question | 11.7 | 21.6 | 33.3 |
| Altern. question | 33.4 | 33.4 | 66.7 |
| Total | 45.0 | 55.0 | 100 |

In the table for Expected values, you can see that with 100 persons in the simulation 45 percent said YES, and 55 percent said NO. The figures tell you what percentage of persons have taken illegal drugs!

## Samples

Sampling is to gain information about the whole by examining a part. That is why sampling is a basic issue for statistics. We need to understand the process of drawing samples to know what you can do with the results. Important concepts include variability, uncertainty, sampling distribution, randomization, confidence, square root *n* law, central limit theorem .. Statistics is not mathematics.
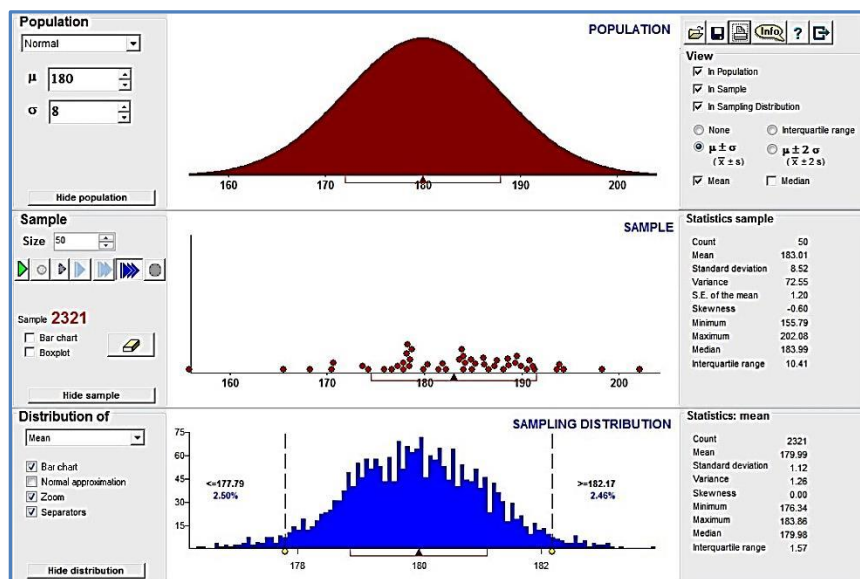
Simulations are excellent to visualize the concepts and so lead to understanding. This will enhance statistical reasoning with insight.

### Sampling distribution

The idea with this module is that students should get some insight in the the sample process and then be able to interpret the result you get from the sample distribution. You have a lot of settings for the population, the sample and for the sample distribution and can view the result numerically and graphically in different ways.

If you want to get information about a large unknown population, you need to draw one or more samples. The data you collect can be used to make estimates of mean and median and the distribution of the population.

Of course, these estimates are uncertain due to random fluctuations in the data collected. This uncertainty will be less if you draw a large sample, but in practice such sampling is often too costly. *The idea behind this module is to gain insight into the sampling process and be able to interpret the results you get from the sampling distribution.*

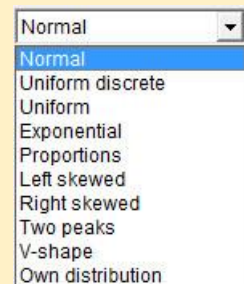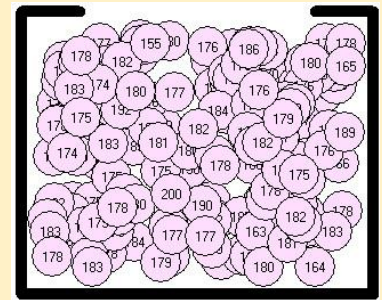*Sampling from a normal distribution with zooming. Sample size is 50.*

One can show that in about 19 cases out of 20, e.g. 95 % of all samples, the sample mean will be in the range

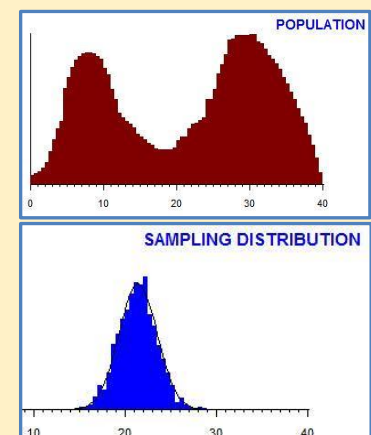$$\overline{x} \pm 2 \cdot \frac{s}{\sqrt{n}}$$

where $\overline{x}$ is the mean of the sampling distribution, *s* is the standard deviation and *n* is the sample size. In this example we have the interval

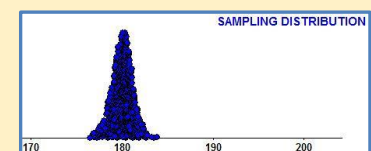$$180.0 \pm 2 \cdot \frac{8}{\sqrt{50}} \approx 180.0 \pm 2.3$$

We can check this in the histogram above. About 95 % of the total area of the blue bars is in the range 177.7 - 182.3. We call this a confidence interval at 95 % confidence level.

For the population you can choose different models. You can also make your own distribution with the help of the mouse.

It is surprising to see is that the sampling distribution of the mean is more or less a normal distribution, whatever the population distribution may be.

We can see without zooming that the spread is significantly less.
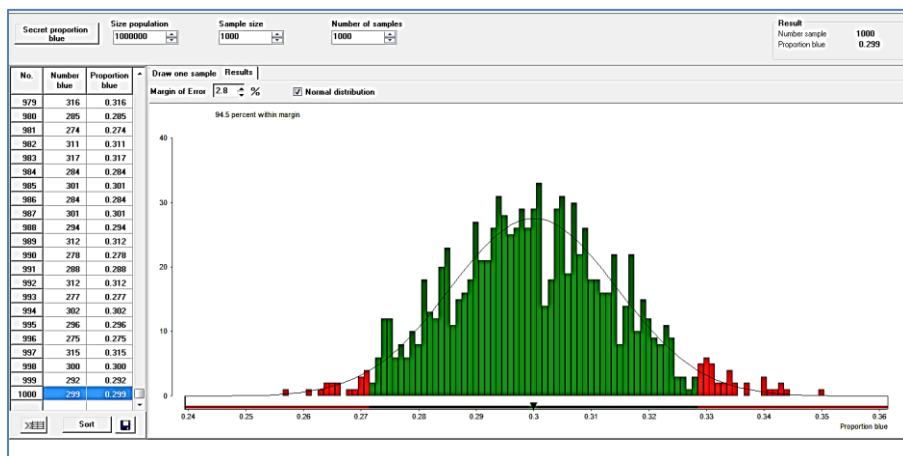
## Sampling for proportion

It surprises many that rather small samples can give a good picture of a population. Think of polls. Facing parliamentary elections polling institutes usually ask about 1000 people which party they would vote for if the election were today. People say that they don't believe in polls because they have never been asked about their opinion! In this simulation, you can study the relationship between *sample size*, *population size* and *population characteristics*.



The population consists of red and blue squares. They represent the proponents and the opponents. Population characteristics of the colours are displayed. Before the sample is drawn, the proportion in the "Secret proportion blue" is entered. In reality, this is an *unknown parameter* of the population and it is the one to be estimated.

You can now enter the size of the population and the sample. Because you can draw many samples, you can then make assessments of the reliability of the estimate of the proportion.

The result of all trials are displayed on the left. The result is also displayed with a bar chart for each sample. You can also see a preview of the distribution of the proportion in each sample.



If you press the button **Results** you see the distribution of the proportions. This is called the *Sampling distribution*. If you change *Margin of error* the diagram is updated and you can see how many are within the margin by the green bars in the diagram. This value is also shown as text in the upper left corner.
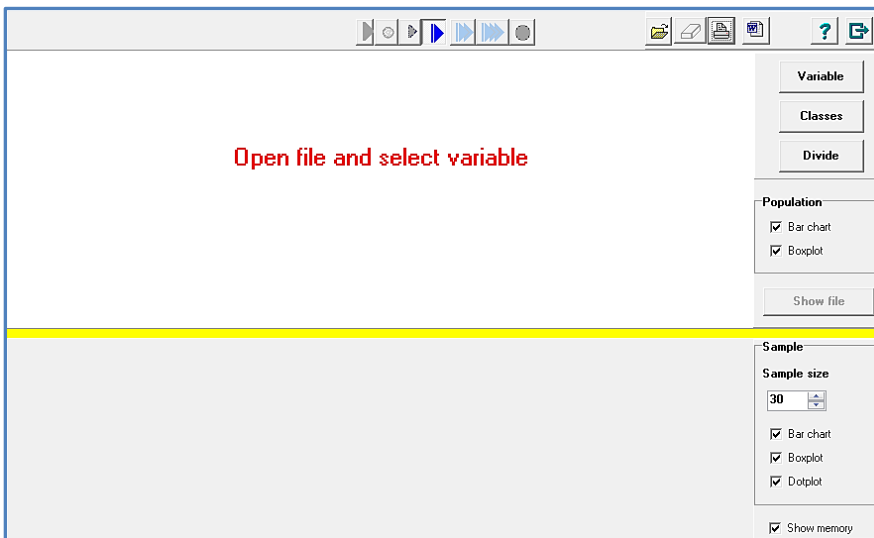
# Dancing samples

The purpose of this module is to demonstrate the influence of the size of the sample on the sampling results. To see the effect of this module, it makes sense to use a large file. The file in the example was created for a nationwide survey of students, with 50 000 participants. In this module it is possible to draw many samples.
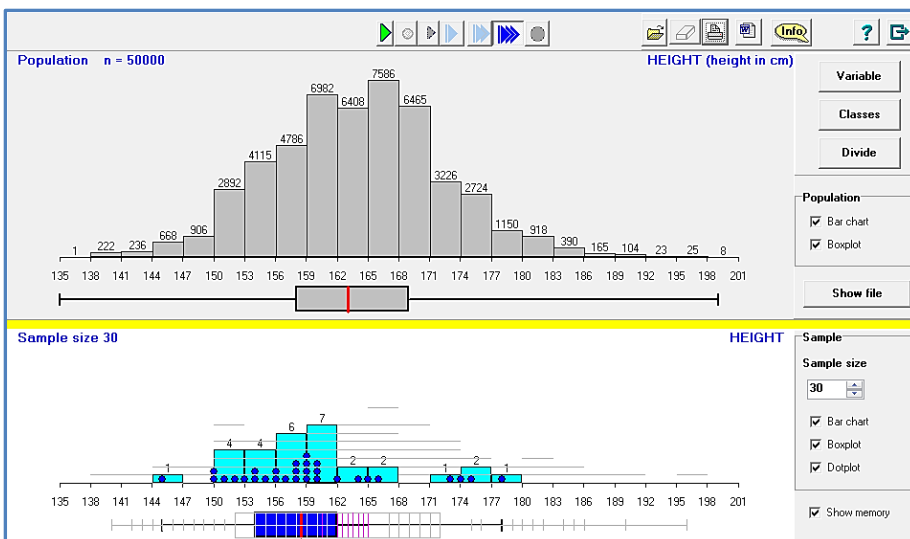
In the top half of the screen is the population. The bottom half shows the samples drawn. The results are presented using a box plot, histogram or dotplot. If memory is turned on, boxplot and histogram show the results of the last 30 sample drawings. It is clear that with a small sample of 30 individuals the variation is very large. With a sample size of 300 there is much less variation.

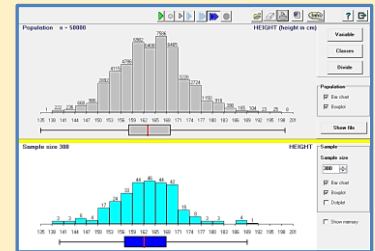The idea of this module comes from Journal of the Royal Statistical Society:Towards more accessible conceptions of statistical inference. Authors: C. J. Wild, M. Pfannkuch, M. Regan, N. J. Horton.
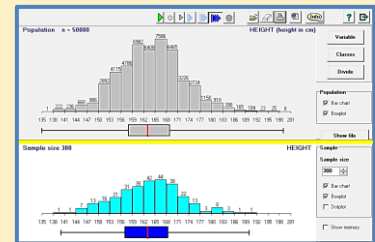


*You start by opening a file. Then you select a variable which can be divided. You can also set the number of classes in the histogram.*



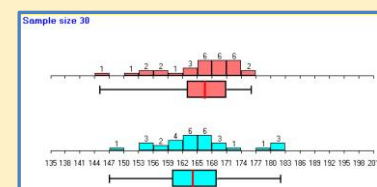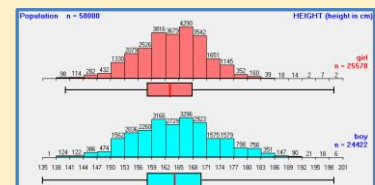*Sample size here is 30. Large variation between samples.*



*sample 1*



*sample 2*
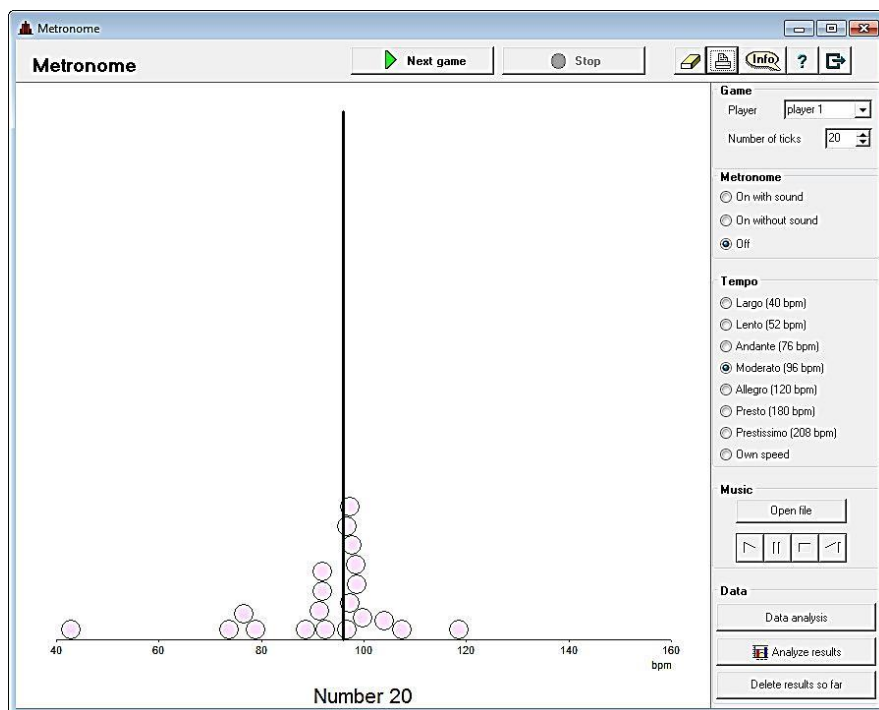With larger samples the variation is much less.





Data for height could also be divided to a variable. The division variable here is gender.
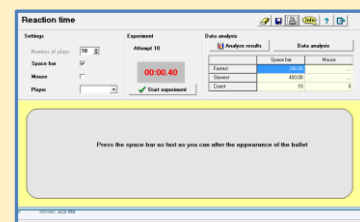
## Make data

VUstat has three modules with which one can create one´s own data from *experiments*.

- In **Reaction time** we measure your response speed at the keyboard. Data are number of milliseconds which can be analysed directly or in Data analysis. Several players can compete against each other.
- Every human being shows unconsciously regularities in his/her behavior. If these regularities can be found, then you can predict his or her behavior based on this knowledge. This is the case in module **Artificial intelligence** when typing zero and one on a computer. In this simulation the computer tries to predict your behavior from the zeroes and ones you type.
- **Metronome:** a player needs to maintain a certain metronome tempo as far as possible and does this by holding the tempo by tapping the space bar. After each tap, BPM (beats per minute) is calculated and displayed as a dot. The screen displays a dotplot of the scores. Several players can compete against each other and data can be analyzed in different ways.
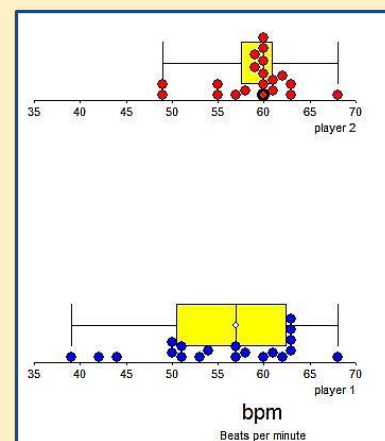
**Metronome**



*For each game round, you should turn on the metronome with sound, so that new players can listen in and learn the tempo. There are 8 different tempos.*



Reaction time. Here we test the response speed at the keyboard. There are three settings for the experiment. You use just the spacebar (easy), just the mouse (a little more difficult) or a mixture (very difficult).





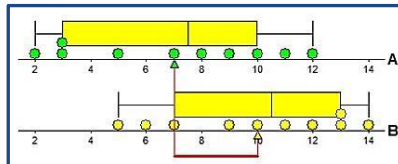| Player | player 2 | player 1 | |
|--------|----------|----------|-------|
| bpm | Freq. | Freq. | Total |
| 39 - 41 | 0 | 1 | 1 |
| 42 - 44 | 0 | 2 | 2 |
| 45 - 47 | 0 | 0 | 0 |
| 48 - 50 | 2 | 2 | 4 |
| 51 - 53 | 0 | 3 | 3 |
| 54 - 56 | 2 | 1 | 3 |
| 57 - 59 | 4 | 3 | 7 |
| 60 - 62 | 9 | 3 | 12 |
| 63 - 65 | 2 | 4 | 6 |
| 66 - 68 | 1 | 1 | 2 |
| **Total** | **20** | **20** | **40** |

You can compare the results for two or more players in different ways.

## Testing hypotheses

The classic way of introducing statistical inference runs mostly through techniques of hypothesis testing. The level of mathematics and abstract language involved makes this approach difficult for students. With simulations, an intuitive inference approach is possible, which focuses on basic concepts such as sampling and chance. The *randomization test* is very suitable for such an intuitive introduction into statistical inference. The classic ways of hypotheses testing can be also be demonstrated through visualization in modules such as binomial test, *z*-test and *t*-test.

### Randomization test - Which battery brand is best?

For two brands of batteries, A and B, 10 lifetimes are listed. Though the mean of sample B is 3 hours ahead of the mean of brand A, the dotplots of the two samples show too much overlap and contain too few data to provide a basis for any conclusion that brand B is better.
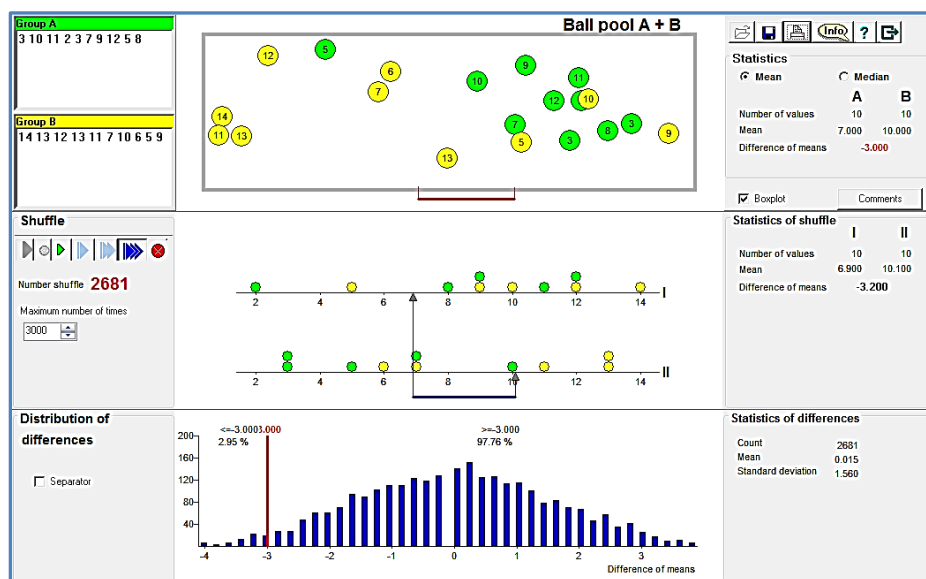


*Here is the distribution of the samples for Brand A and Brand B*

We put both samples in a ball pool of A+B.
A sample of 10 from this pool results in a remainder of 10. These two samples I and II have also a difference between the means. Since the data values were randomly divided into two groups this difference between the means can be attributed to chance. This procedure can be repeated a lot of times with test statistic: the difference between the means. The observed difference is 3 hours in favor of brand B. *The question is whether this difference is significant.*

The distribution of differences indicates that about 3 % of the simulations show a greater difference than the observed difference. Then the conclusion is that brand B is significantly better than brand A.
Although the mean is most commonly used it can be unstable because of the small dataset. The median is more stable. However, when the median is the key, the difference is not significant. This can be seen easily by switching from mean to median.
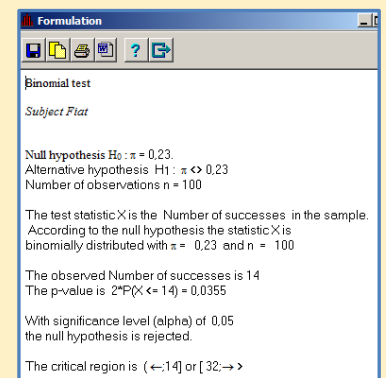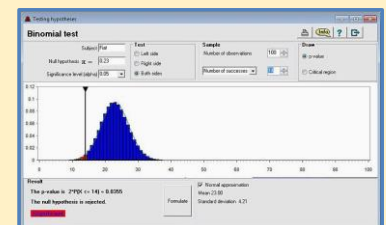
**Example of Binomial test**
A Fiat car dealer claims that Fiat's share of the car market is about 23 %. A consumer organization is to consider whether this claim is correct. A random sample of 100 car owners are asked to mark the car they drive. How many Fiat owners will be needed to prove that the dealer is wrong at a significance level of 5 %?





You can display a full list of test data with comments, copy it and paste into other applications.