

STATISTICAL DATA ANALYSIS IN EXCEL

Part 1

Introduction to Statistics

Dr. Petr Nazarov

petr.nazarov@crp-sante.lu

14-06-2010

The course:

- ◆ Reminds statistical basics
- ◆ Gives the methodological tools for the research
- ◆ Provides practical skill for fast data analysis

Organization

- ◆ 4 sessions (6 hours)
- ◆ Lectures are integrated with practical work
- ◆ **PLEASE: ask questions. Understanding is extremely important for future parts**

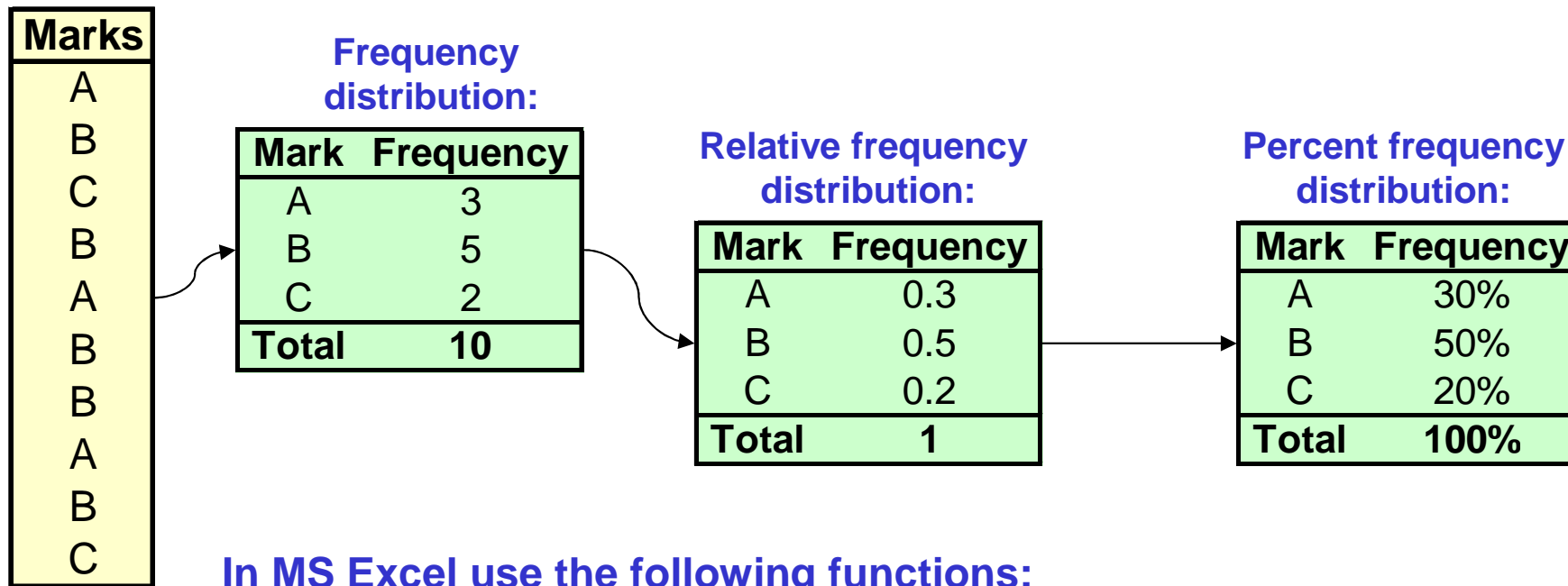
- ◆ **Descriptive statistics**
- ◆ **Exploratory analysis**
- ◆ **Discrete probability distribution**
- ◆ **Continues probability distribution**

Look for the data:

`http://edu.sablab.net/data/xls`

Frequency distribution

A tabular summary of data showing the number (frequency) of items in each of several nonoverlapping classes.



In MS Excel use the following functions:

- ◆ =COUNTIF (data , element) to get number of “elements” found in the “data” area
- ◆ =SUM (data) to get the sum of the values in the “data” area

The role of smoking in the etiology of pancreatitis has been recognized for many years. To provide estimates of the quantitative significance of these factors, a hospital-based study was carried out in eastern Massachusetts and Rhode Island between 1975 and 1979. **53 patients** who had a hospital discharge diagnosis of **pancreatitis** were included in this unmatched case-control study. The **control group** consisted of 217 patients admitted for **diseases other** than those of the pancreas and biliary tract. Risk factor information was obtained from a standardized interview with each subject, conducted by a trained interviewer.

adapted from Chap T. Le, Introductory Biostatistics

pancreatitis.xls

Pancreatitis patients:

Smokers	Ex-smokers	Ex-smokers	Smokers	Smokers	Smokers
Ex-smokers	Smokers	Smokers	Smokers	Smokers	Smokers
Ex-smokers	Smokers	Smokers	Ex-smokers	Smokers	Smokers
Ex-smokers	Ex-smokers	Smokers	Ex-smokers	Smokers	
Smokers	Never	Smokers	Ex-smokers	Ex-smokers	
Smokers	Ex-smokers	Smokers	Smokers	Ex-smokers	
Smokers	Smokers	Smokers	Smokers	Smokers	
Ex-smokers	Smokers	Smokers	Smokers	Smokers	
Smokers	Smokers	Smokers	Smokers	Smokers	
Smokers	Never	Smokers	Smokers	Smokers	

Frequency distribution

A tabular summary of data showing the number (frequency) of items in each of several nonoverlapping classes.

In MS Excel use the following functions:

- ◆ =COUNTIF(data, element) to get number of “elements” found in the “data” area
- ◆ =SUM(data) to get the sum of the values in the “data” area

pancreatitis.xls

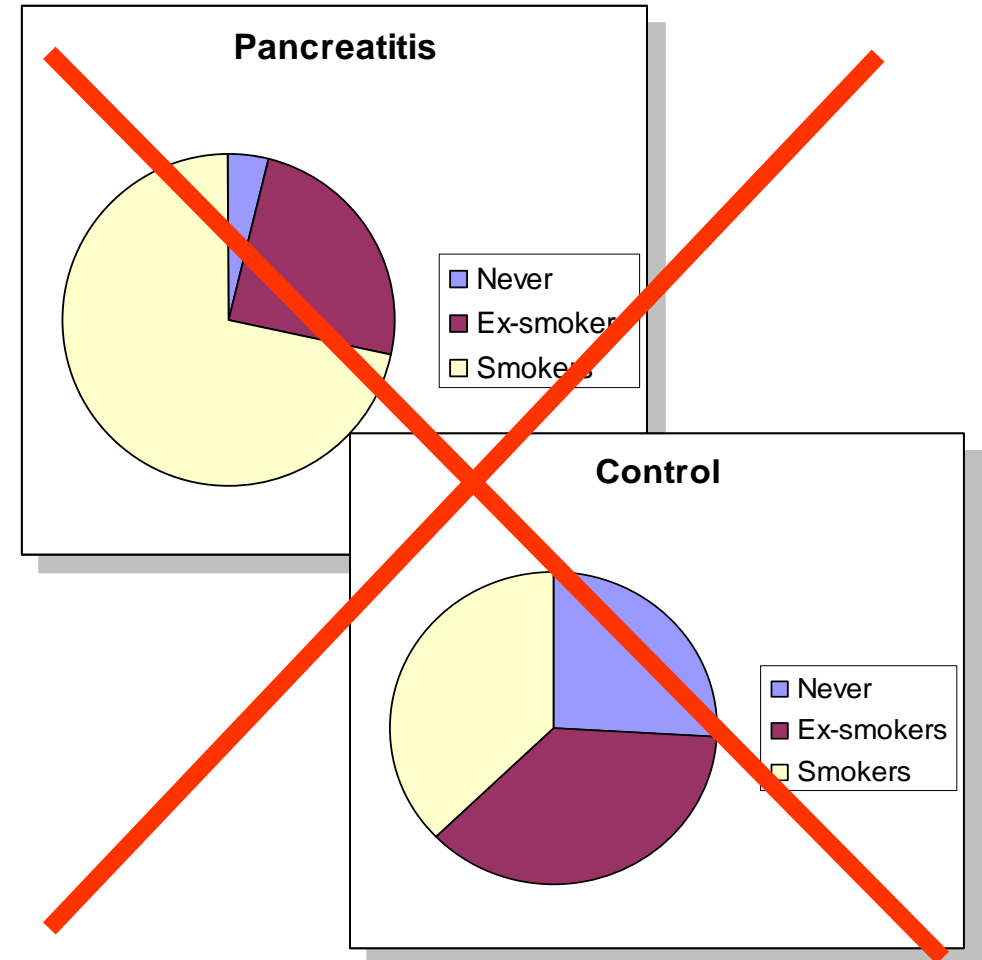
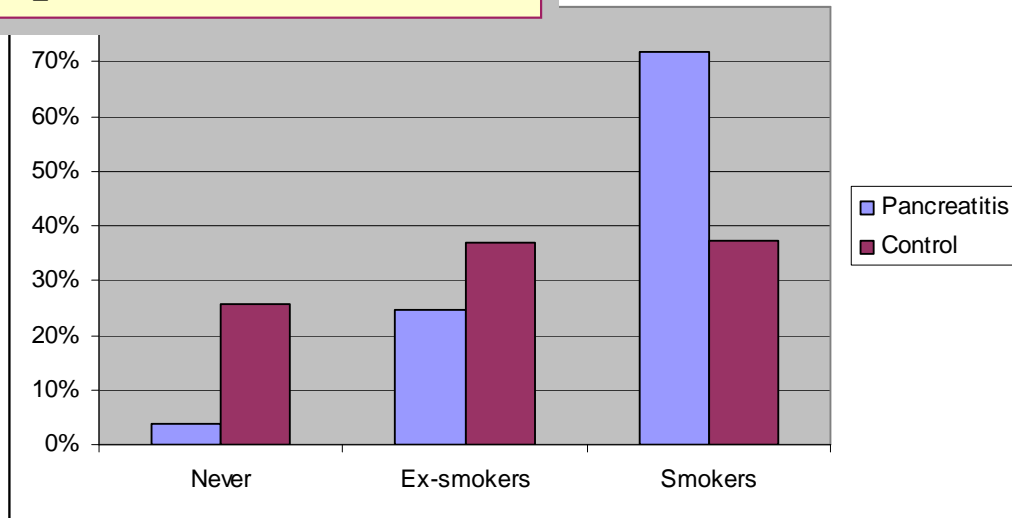
Frequency distribution:

Smoking	Cases	Controls
Never	2	56
Ex-smokers	13	80
Smokers	38	81
Total	53	217

Relative frequency distribution:

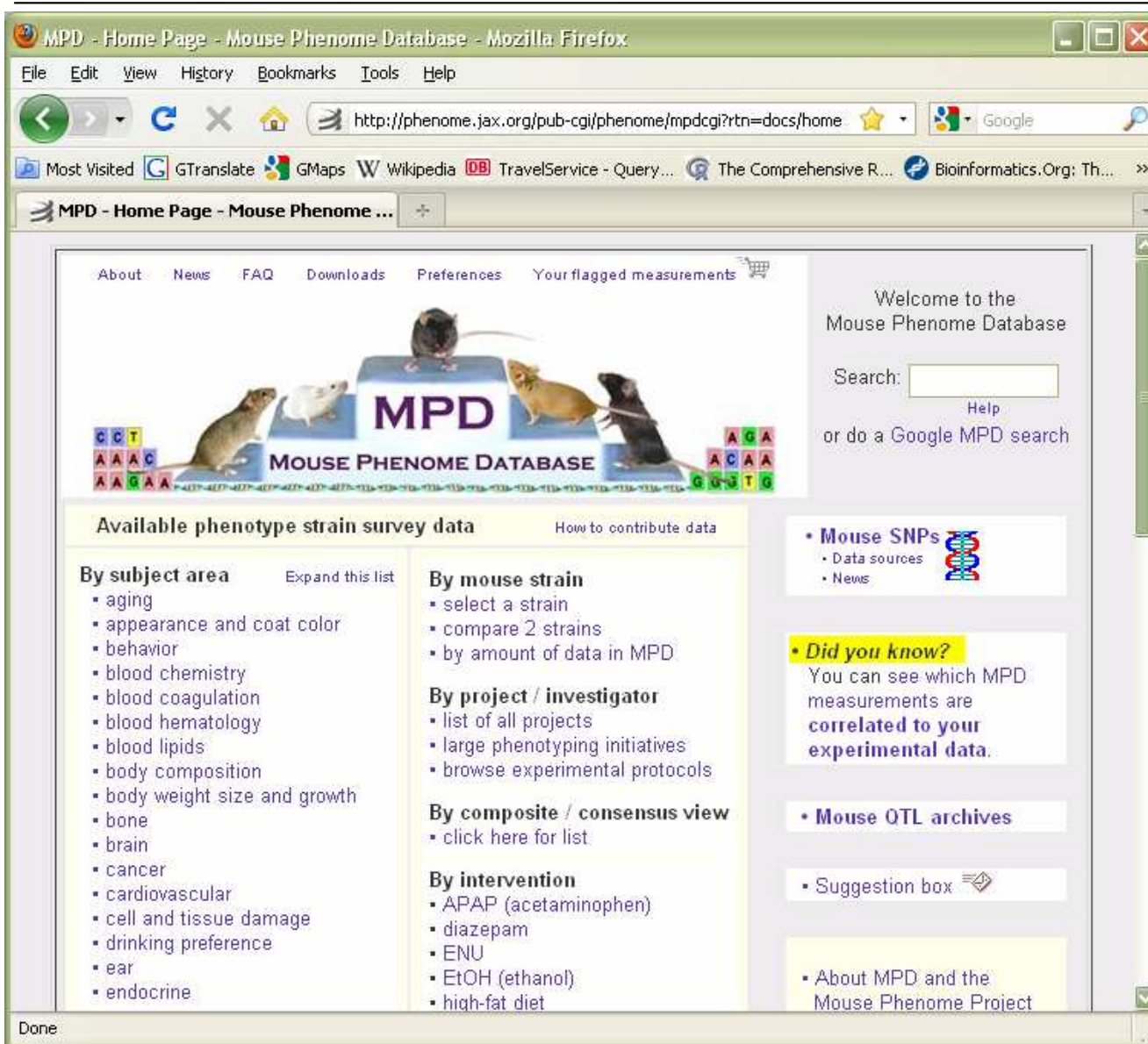
Smoking	Cases	Controls
Never	0.038	0.258
Ex-smokers	0.245	0.369
Smokers	0.717	0.373
Total	1	1

pancreatitis.xls



In MS Excel use the following steps:

- ◆ Chart Wizard → Columns → Set data range (both columns of Percent freq. distribution)
- ◆ Chart Wizard → Pie → Set data range (one column of Percent freq. distribution)



Tordoff MG, Bachmanov AA

Survey of calcium & sodium intake and metabolism with bone and body composition data

Project symbol: **Tordoff3**

Accession number: **MPD:103**

mice.xls

790 mice from different strains

<http://phenome.jax.org>

parameter

Starting age

Ending age

Starting weight

Ending weight

Weight change

Bleeding time

Ionized Ca in blood

Blood pH

Bone mineral density

Lean tissues weight

Fat weight

The following are weights in grams for 970 mice:

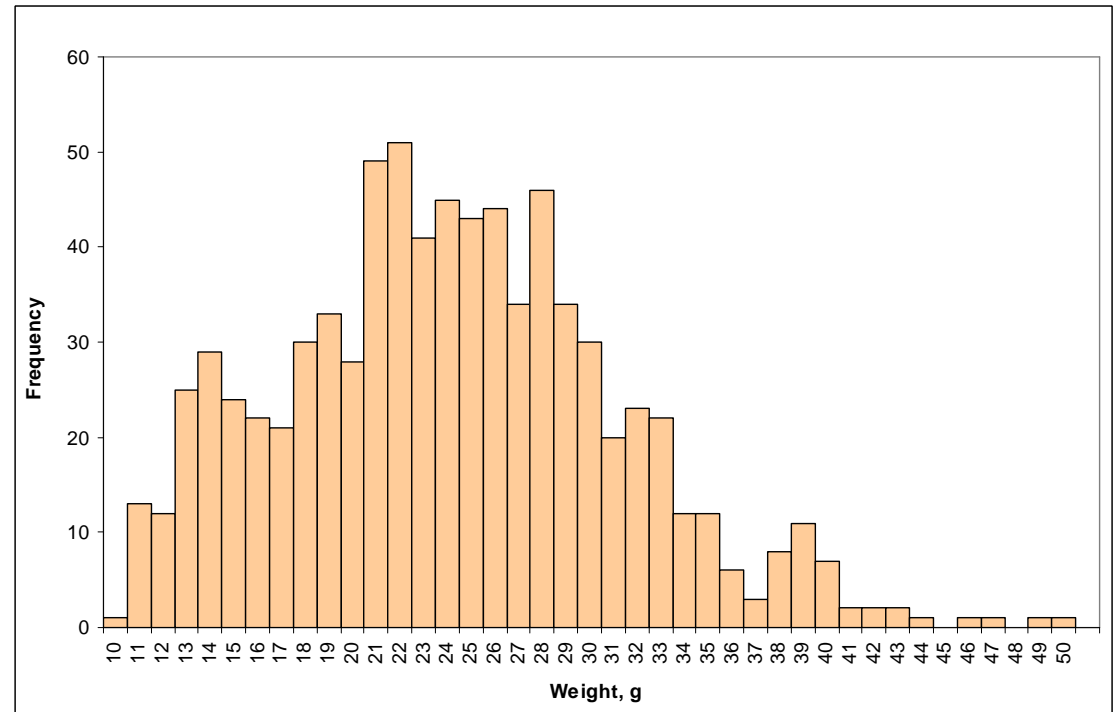
mice.xls	20.5	23.2	24.6	23.5	26	25.9	23.9	22.8	19.9	...
	20.8	22.4	26	23.8	26.5	26	22.8	22.9	20.9	...
	19.8	22.7	31	22.7	26.3	27.1	18.4	21	18.8	...
	21	21.4	25.7	19.7	27	26.2	21.8	22.2	19.2	...
	21.9	22.6	23.7	26.2	26	27.5	25	20.9	20.6	...
	22.1	20	21.1	24.1	28.8	30.2	20.1	24.2	25.8	...
	21.3	21.8	23.7	23.5	28	27.6	21.6	21	21.3	...
	20.1	20.8	24.5	23.8	29.5	21.4	21.5	24	21.1	...
	18.9	19.5	32.3	28	27.1	28.2	22.9	19.9	20.4	...
	21.3	20.6	22.8	25.8	24.1	23.5	24.2	22	20.3	...

Sorted weights show that the values are in the 10 – 49.6 grams.
 Let us divide the weight into the “bins”

<i>Weight,g</i>	<i>Frequency</i>
>=10	1
10-20	237
20-30	417
30-40	124
40-50	11
More	0

Now, let us use bin-size = 1 gram

<i>Bin</i>	<i>Frequency</i>
10	1
11	13
12	12
13	25
14	29
...	...
46	1
47	1
48	0
49	1
50	1
More	0

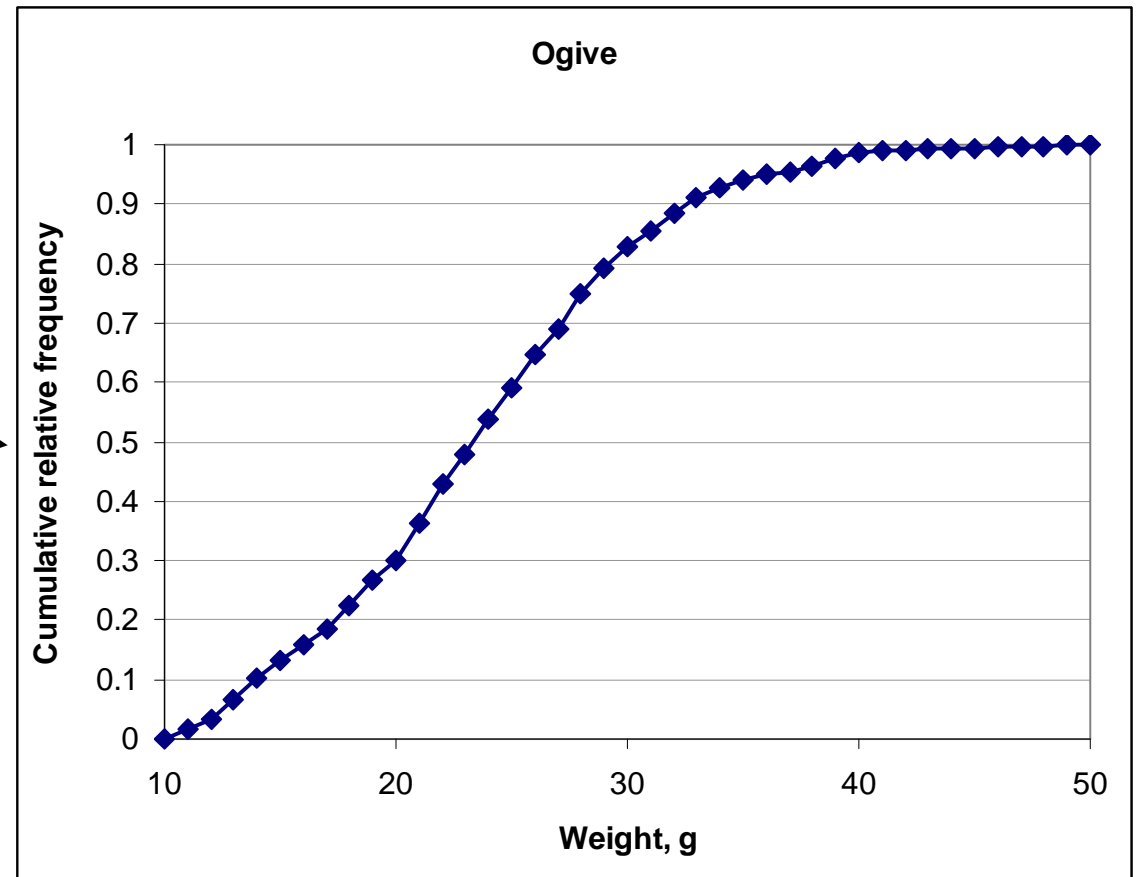
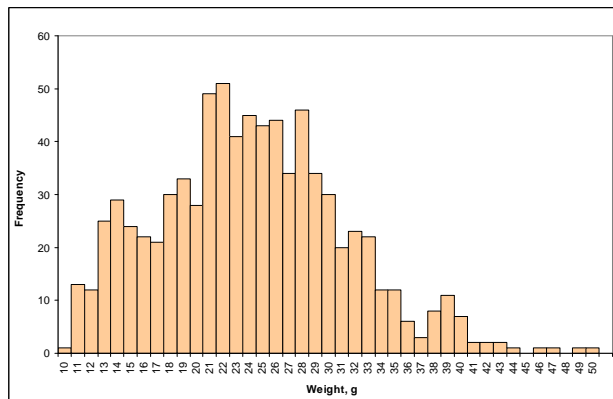


In Excel use the following steps:

- ◆ Specify the column of bins (interval) upper-limits
- ◆ Tools → Data Analysis → Histogram → select the input data, bins, and output (Analysis ToolPak should be installed)
- ◆ use Chart Wizard → Columns to visualize the results

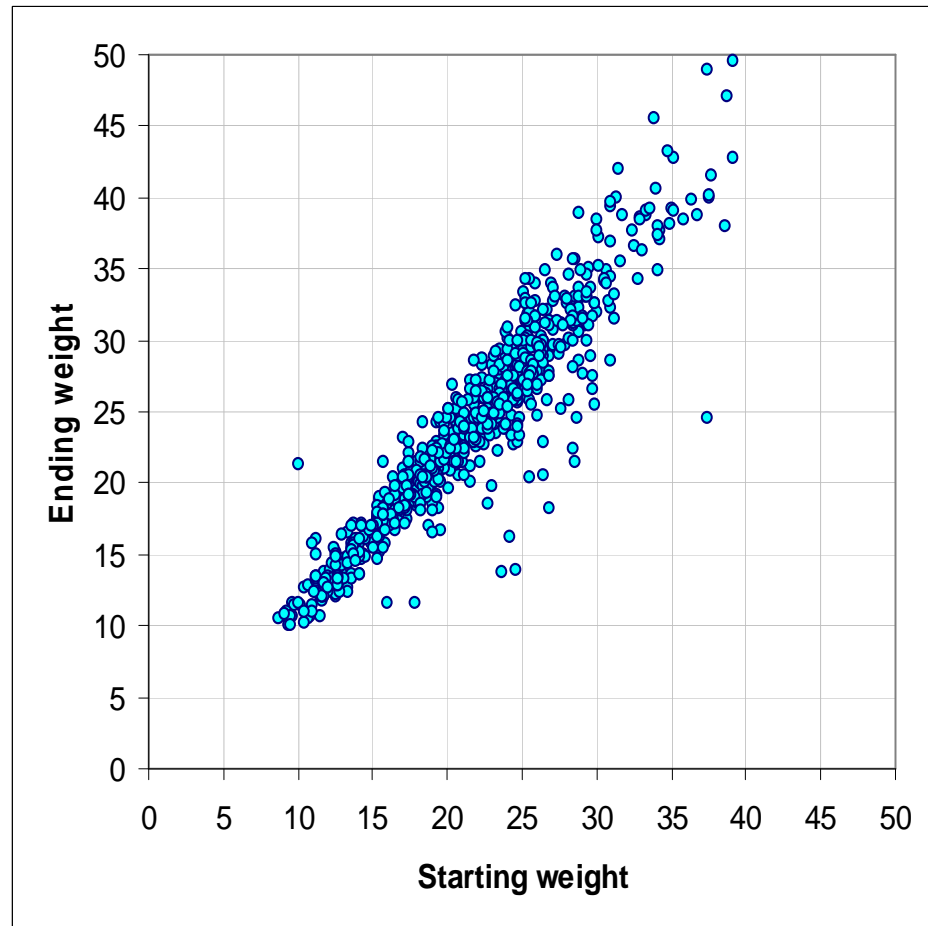
Cumulative frequency distribution

A tabular summary of quantitative data showing the number of items with values less than or equal to the upper class limit of each class.



mice.xls

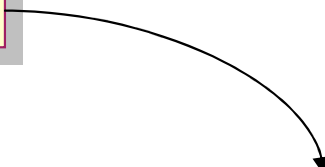
Let us look on mutual dependency of the Starting and Ending weights.



In Excel use the following steps:

- ◆ Select the data region
- ◆ Use Chart Wizard → XY (Scatter)

pancreatitis.xls



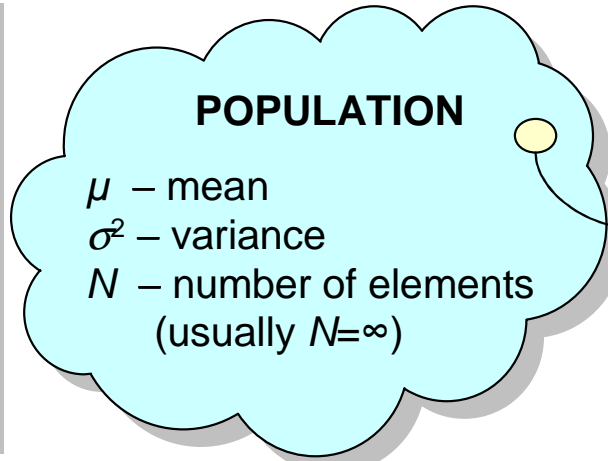
Smoking	Disease		Total
	other	pancreatitis	
Ex-smokers	80	13	93
Never	56	2	58
Smokers	81	38	119
Total	217	53	270

In Excel use the following steps:

- ◆ Data → Pivot Table and PivotChart → MS Office list + Pivot Table
- ◆ Set the range, including the headers of the data
- ◆ Select output and set layout by drag-and-dropping the names into the table

Population parameter

A numerical value used as a summary measure for a population (e.g., the population mean μ , variance σ^2 , standard deviation σ)

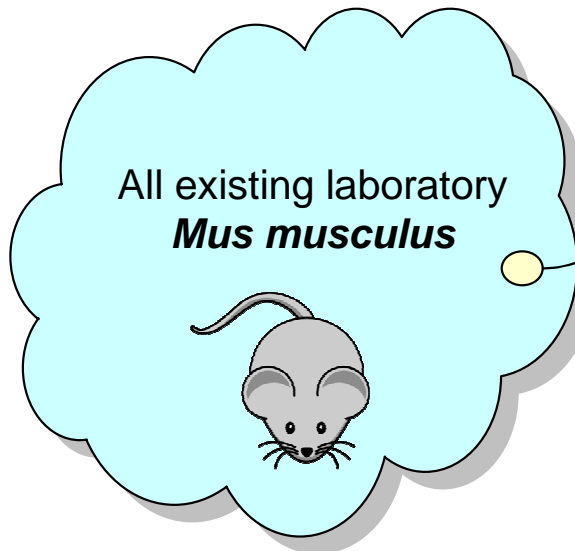


SAMPLE

m, \bar{x} – mean
 s^2 – variance
 n – number of elements

Sample statistic

A numerical value used as a summary measure for a sample (e.g., the sample mean m , the sample variance s^2 , and the sample standard deviation s)



mice.xls

790 mice from different strains

<http://phenome.jax.org>

ID	Strain	Sex	Starting age	Ending age	Starting weight	Ending weight	Weight change	Bleeding time	Ionized Ca in blood	Blood pH	Bone mineral density	Lean tissues weight	Fat weight
1	129S1/SvlmJ	f	66	116	19.3	20.5	1.062	64	1.2	7.24	0.0605	14.5	4.4
2	129S1/SvlmJ	f	66	116	19.1	20.8	1.089	78	1.15	7.27	0.0553	13.9	4.4
3	129S1/SvlmJ	f	66	108	17.9	19.8	1.106	90	1.16	7.26	0.0546	13.8	2.9
368	129S1/SvlmJ	f	72	114	18.3	21	1.148	65	1.26	7.22	0.0599	15.4	4.2
369	129S1/SvlmJ	f	72	115	20.2	21.9	1.084	55	1.23	7.3	0.0623	15.6	4.3
370	129S1/SvlmJ	f	72	116	18.8	22.1	1.176		1.21	7.28	0.0626	16.4	4.3
371	129S1/SvlmJ	f	72	119	19.4	21.3	1.098	49	1.24	7.24	0.0632	16.6	5.4
372	129S1/SvlmJ	f	72	122	18.3	20.1	1.098	73	1.17	7.19	0.0592	16	4.1
4	129S1/SvlmJ	f	66	109	17.2	18.9	1.099	41	1.25	7.29	0.0513	14	3.2
5	129S1/SvlmJ	f	66	112	19.7	21.3	1.081	129	1.14	7.22	0.0501	16.3	5.2
10	129S1/SvlmJ	m	66	112	24.3	24.7	1.016	119	1.13	7.24	0.0533	17.6	6.8
364	129S1/SvlmJ	m	72	114	25.3	27.2	1.075	64	1.25	7.27	0.0596	19.3	5.8
365	129S1/SvlmJ	m	72	115	21.4	23.9	1.117	48	1.25	7.28	0.0563	17.4	5.7
366	129S1/SvlmJ	m	72	118	24.5	26.3	1.073	59	1.25	7.26	0.0609	17.8	7.1
367	129S1/SvlmJ	m	72	122	24	26	1.083	69	1.29	7.26	0.0584	19.2	4.6
6	129S1/SvlmJ	m	66	116	21.6	23.3	1.079	78	1.15	7.27	0.0497	17.2	5.7
7	129S1/SvlmJ	m	66	107	22.7	26.5	1.167	90	1.18	7.28	0.0493	18.7	7
8	129S1/SvlmJ	m	66	108	25.4	27.4	1.079	35	1.24	7.26	0.0538	18.9	7.1
9	129S1/SvlmJ	m	66	109	24.4	27.5	1.127	43	1.29	7.29	0.0539	19.5	7.1

Mean

A measure of central location computed by summing the data values and dividing by the number of observations.

$$\bar{x} = m = \frac{\sum x_i}{n}$$

$$\mu = \frac{\sum x_i}{N}$$

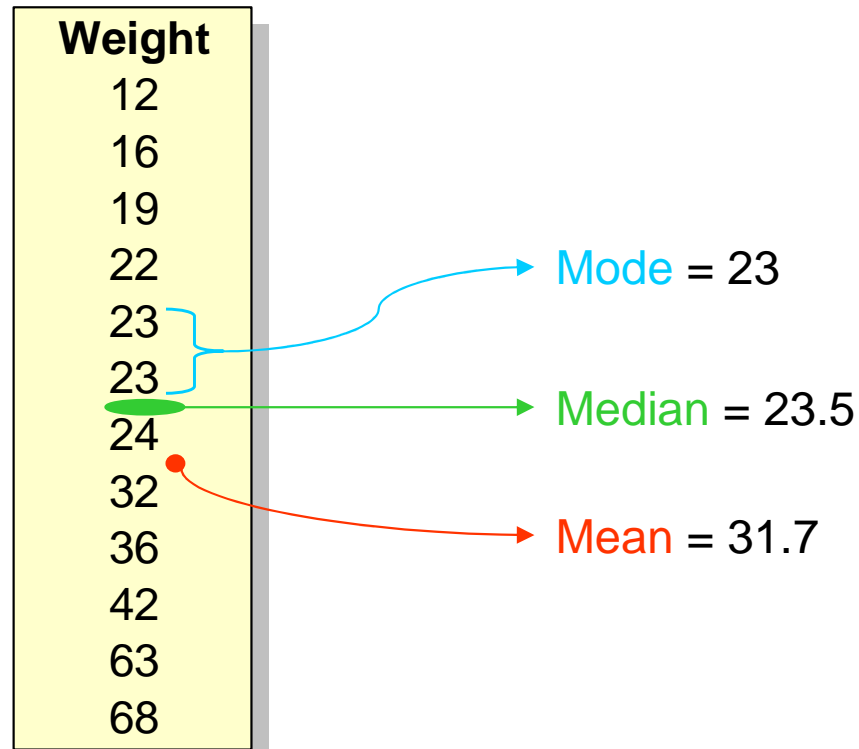
$$p = \frac{\sum (x_i = \text{true})}{n}$$

Median

A measure of central location provided by the value in the middle when the data are arranged in ascending order.

Mode

A measure of location, defined as the value that occurs with greatest frequency.

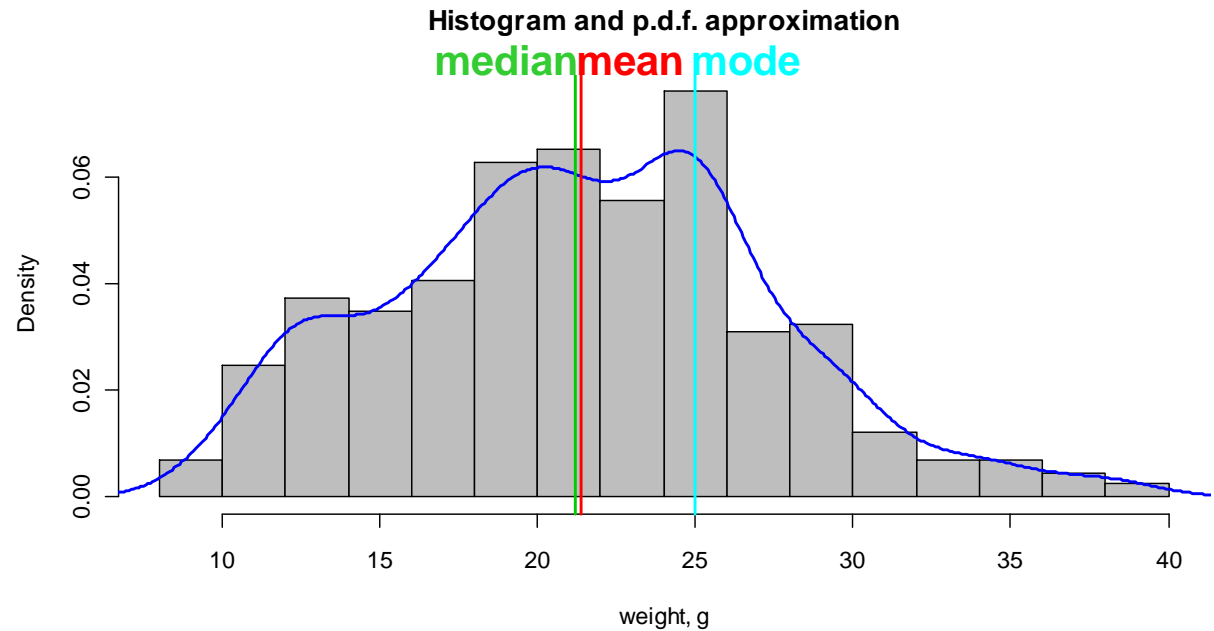


mice.xls

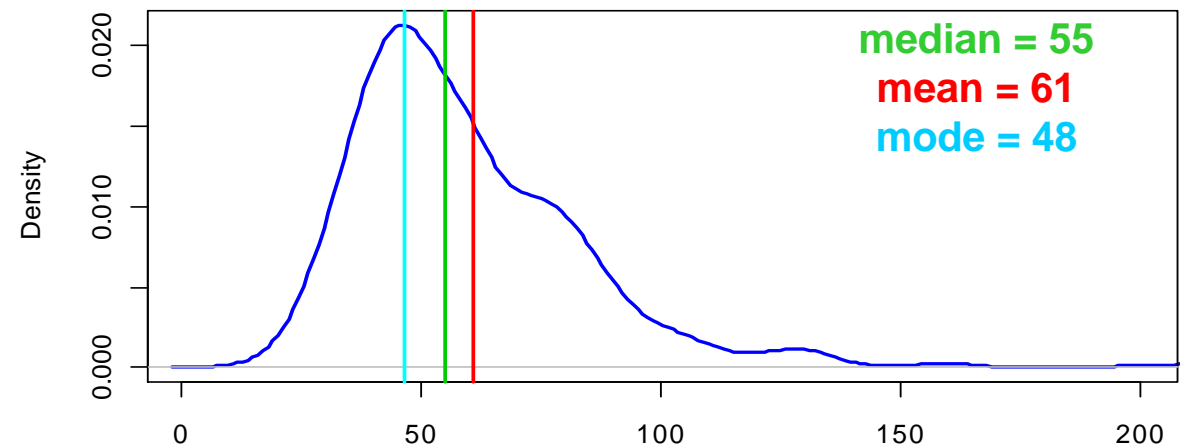
Female proportion
 $p_f = 0.501$

In Excel use the following functions:

- ◆ = AVERAGE(data)
- ◆ = MEDIAN(data)
- ◆ = MODE(data)



Bleeding time



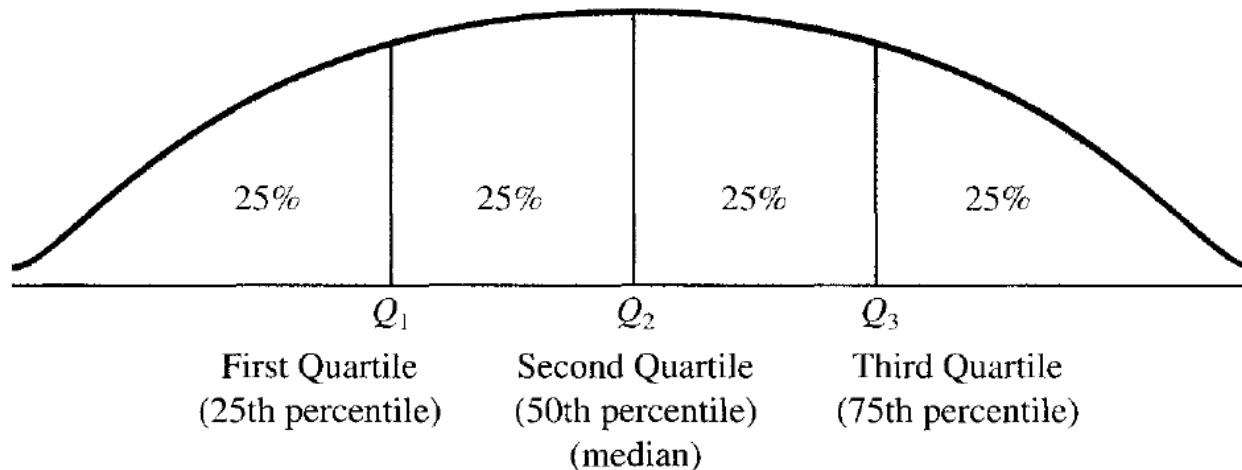
N = 760 Bandwidth = 5.347

Percentile

A value such that at least $p\%$ of the observations are less than or equal to this value, and at least $(100-p)\%$ of the observations are greater than or equal to this value. The 50-th percentile is the **median**.

Quartiles

The 25th, 50th, and 75th percentiles, referred to as the first quartile, the second quartile (median), and third quartile, respectively.



In Excel use the following functions:

◆ =PERCENTILE(data, p)

Weight	12	16	19	22	23	23	24	32	36	42	63	68
---------------	----	----	----	----	----	----	----	----	----	----	----	----

$Q_1 = 21$

$Q_2 = 23.5$

$Q_3 = 39$

Interquartile range (IQR)

A measure of variability, defined to be the difference between the third and first quartiles.

$$IQR = Q_3 - Q_1$$

Variance

A measure of variability based on the squared deviations of the data values about the mean.

population

$$\sigma^2 = \frac{\sum (x_i - \mu)^2}{N}$$

sample

$$s^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$$

Standard deviation

A measure of variability computed by taking the positive square root of the variance.

$$\text{Sample standard deviation} = s = \sqrt{s^2}$$

$$\text{Population standard deviation} = \sigma = \sqrt{\sigma^2}$$

Weight	12	16	19	22	23	23	24	32	36	42	63	68
---------------	----	----	----	----	----	----	----	----	----	----	----	----

IQR = 18

Variance = 320.2

St. dev. = 17.9

In Excel use the following functions:

◆ =VAR(data) , =STDEV(data)

Coefficient of variation

A measure of relative variability computed by dividing the standard deviation by the mean.

Weight	12	16	19	22	23	23	24	32	36	42	63	68
--------	----	----	----	----	----	----	----	----	----	----	----	----

$$\left(\frac{\text{Standard deviation}}{\text{Mean}} \times 100 \right) \%$$

CV = 57%

Median absolute deviation (MAD)

MAD is a robust measure of the variability of a univariate sample of quantitative data.

$$MAD = \text{median} \left(|x_i - \text{median}(x)| \right)$$

Set 1	Set 2
23	23
12	12
22	22
12	12
21	21
18	81
22	22
20	20
12	12
19	19
14	14
13	13
17	17

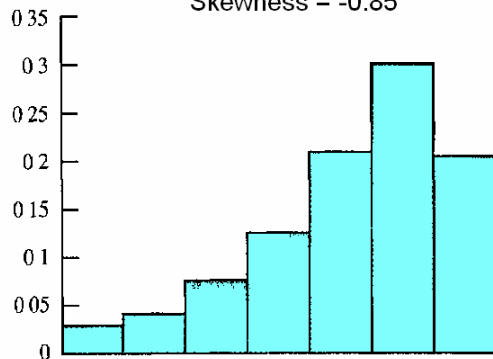
	Set 1	Set 2
Mean	17.3	22.2
Median	18	19
St.dev.	4.23	18.18
MAD	5.93	5.93

Skewness

A measure of the shape of a data distribution. Data skewed to the left result in negative skewness; a symmetric data distribution results in zero skewness; and data skewed to the right result in positive skewness.

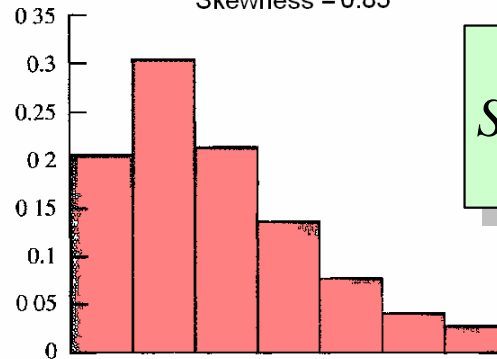
Panel A: Moderately Skewed Left

Skewness = -0.85



Panel B: Moderately Skewed Right

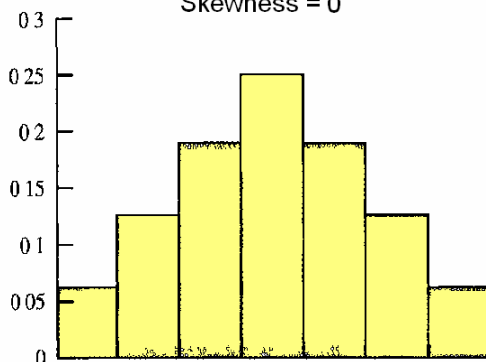
Skewness = 0.85



$$Skewness = \frac{n}{(n-1)(n-2)} \sum_i \left(\frac{x_i - \bar{x}}{s} \right)^3$$

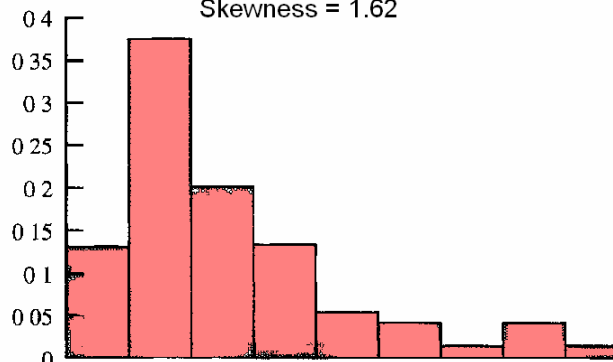
Panel C: Symmetric

Skewness = 0



Panel D: Highly Skewed Right

Skewness = 1.62



adapted from Anderson et al Statistics for Business and Economics

Z-score

A value computed by dividing the deviation about the mean ($x_i - \bar{x}$) by the standard deviation s . A **z-score** is referred to as a standardized value and denotes the number of standard deviations x_i is from the mean.

$$z_i = \frac{x_i - \bar{x}}{s}$$

Chebyshev's theorem

For **any data set**, at least $(1 - 1/z^2)$ of the data values must be within **z** standard deviations from the mean, where z – any value > 1 .

Weight	z-score
12	-1.10
16	-0.88
19	-0.71
22	-0.54
23	-0.48
23	-0.48
24	-0.43
32	0.02
36	0.24
42	0.58
63	1.75
68	2.03

For ANY distribution:

- ◆ At least **75 %** of the values are within **z = 2** standard deviations from the mean
- ◆ At least **89 %** of the values are within **z = 3** standard deviations from the mean
- ◆ At least **94 %** of the values are within **z = 4** standard deviations from the mean
- ◆ At least **96%** of the values are within **z = 5** standard deviations from the mean

For bell-shaped distributions:

- ◆ Approximately 68 % of the values are within 1 st.dev. from mean
- ◆ Approximately 95 % of the values are within 2 st.dev. from mean
- ◆ Almost all data points are inside 3 st.dev. from mean

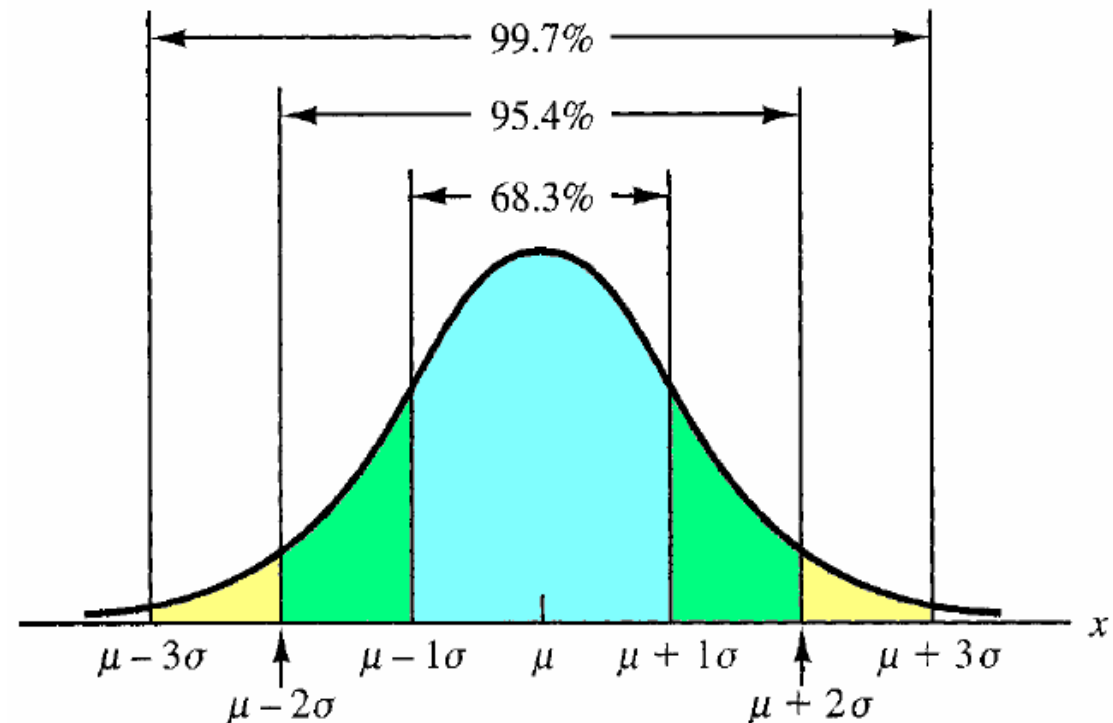
Outlier

An unusually small or unusually large data value.

For bell-shaped distributions data points with $|z| > 3$ can be considered as outliers.

Weight	z-score
23	0.04
12	-0.53
22	-0.01
12	-0.53
21	-0.06
81	3.10
22	-0.01
20	-0.11
12	-0.53
19	-0.17
14	-0.43
13	-0.48
17	-0.27

Example: Gaussian distribution



Five-number summary

An exploratory data analysis technique that uses five numbers to summarize the data: smallest value, first quartile, median, third quartile, and largest value

children.xls

Min. :	12
Q ₁ :	25
Median:	32
Q ₃ :	46
Max. :	79

In Excel use:

◆ Tool → Data Analysis → Descriptive Statistics

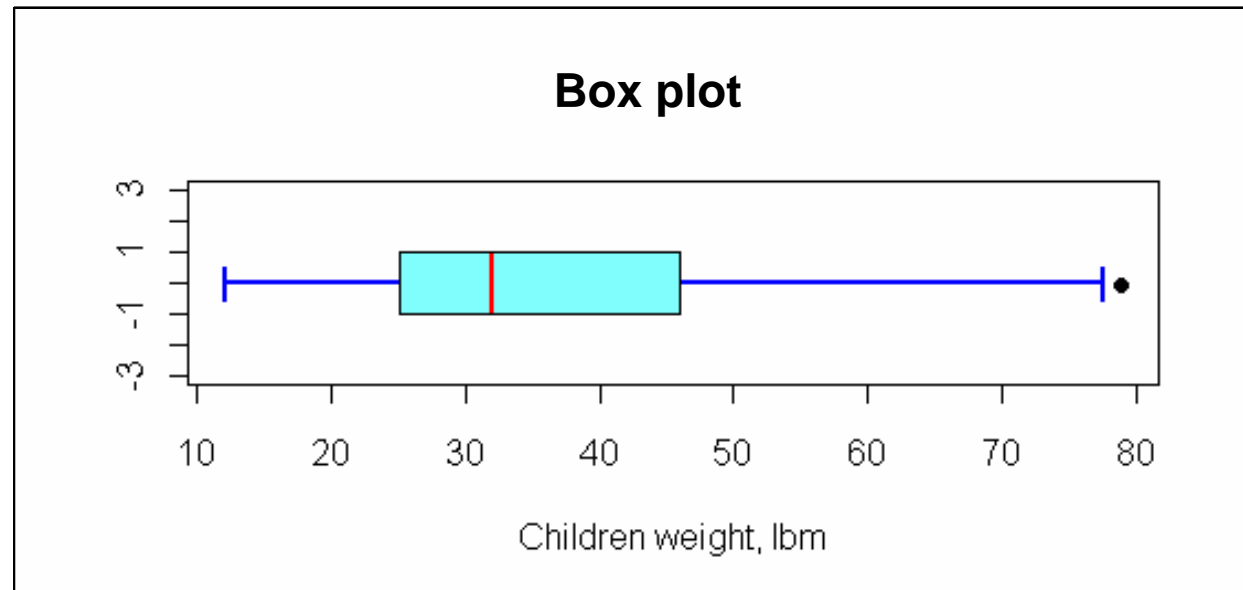
Box plot

A graphical summary of data based on a five-number summary

In Excel use (indirect):

◆ Chart Wizard → Stock → Open-high-low-close

open	Q3
high	Q3+1.5*IQR
low	Q1-1.5*IQR
close	Q1



Example

Build a box plot for weights of male and female mice

`mice.xls`

1. Build 5 number summaries for males and females

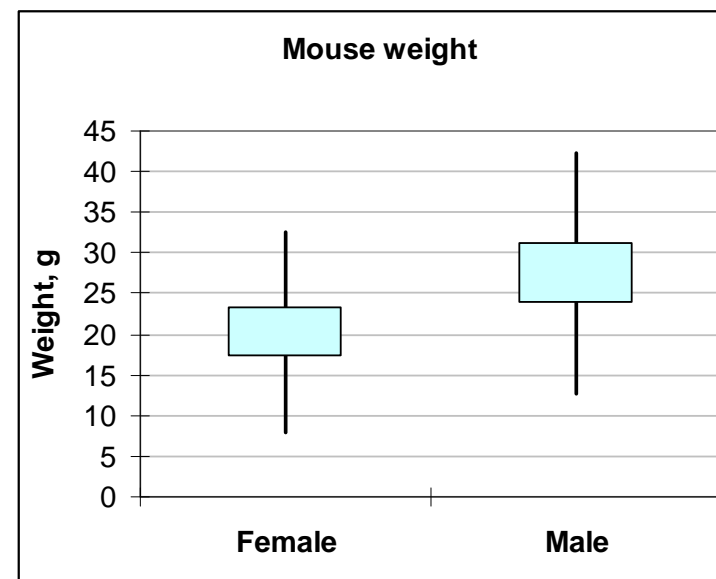
	Female	Male
Min	10.0	12.0
Q1	17.2	23.8
Q2	20.7	27.1
Q3	23.3	31.2
Max	41.5	49.6

2. Combine the numbers into the following order

open	Q3
high	$Q3 + \min(1.5 * (Q3 - Q1), \text{Max})$
low	$Q1 - \max(1.5 * (Q3 - Q1), \text{Min})$
close	Q1

In Excel use:

- ◆ Chart Wizard → Stock → Open-high-low-close
- ◆ Put “series-in-rows”
- ◆ Adjust colors, etc



Measure of Association between 2 Variables

Covariance

A measure of linear association between two variables. Positive values indicate a positive relationship; negative values indicate a negative relationship.

population

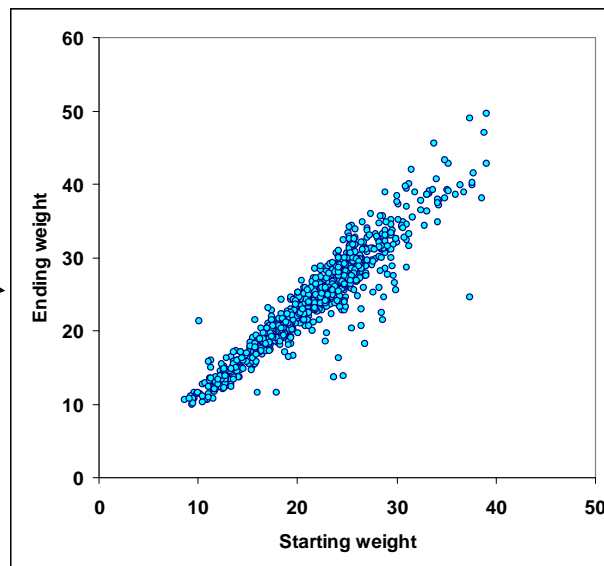
$$\sigma_{xy} = \frac{\sum (x_i - \mu_x)(y_i - \mu_y)}{N}$$

sample

$$s_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

mice.xls

Ending weight vs.
Starting weight



In Excel use function:

◆ =COVAR(data)

$$s_{xy} = 39.8$$

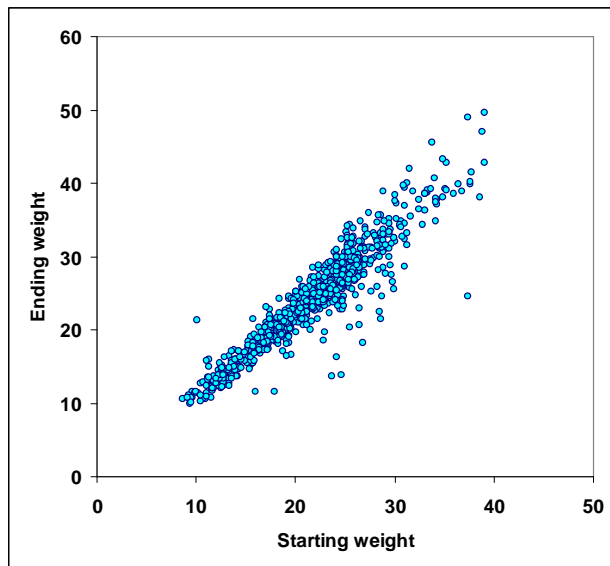
hard to interpret

Correlation (Pearson product moment correlation coefficient)

A measure of linear association between two variables that takes on values between -1 and +1. Values near +1 indicate a strong positive linear relationship, values near -1 indicate a strong negative linear relationship; and values near zero indicate the lack of a linear relationship.

population

$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sigma_x \sigma_y N}$$



sample

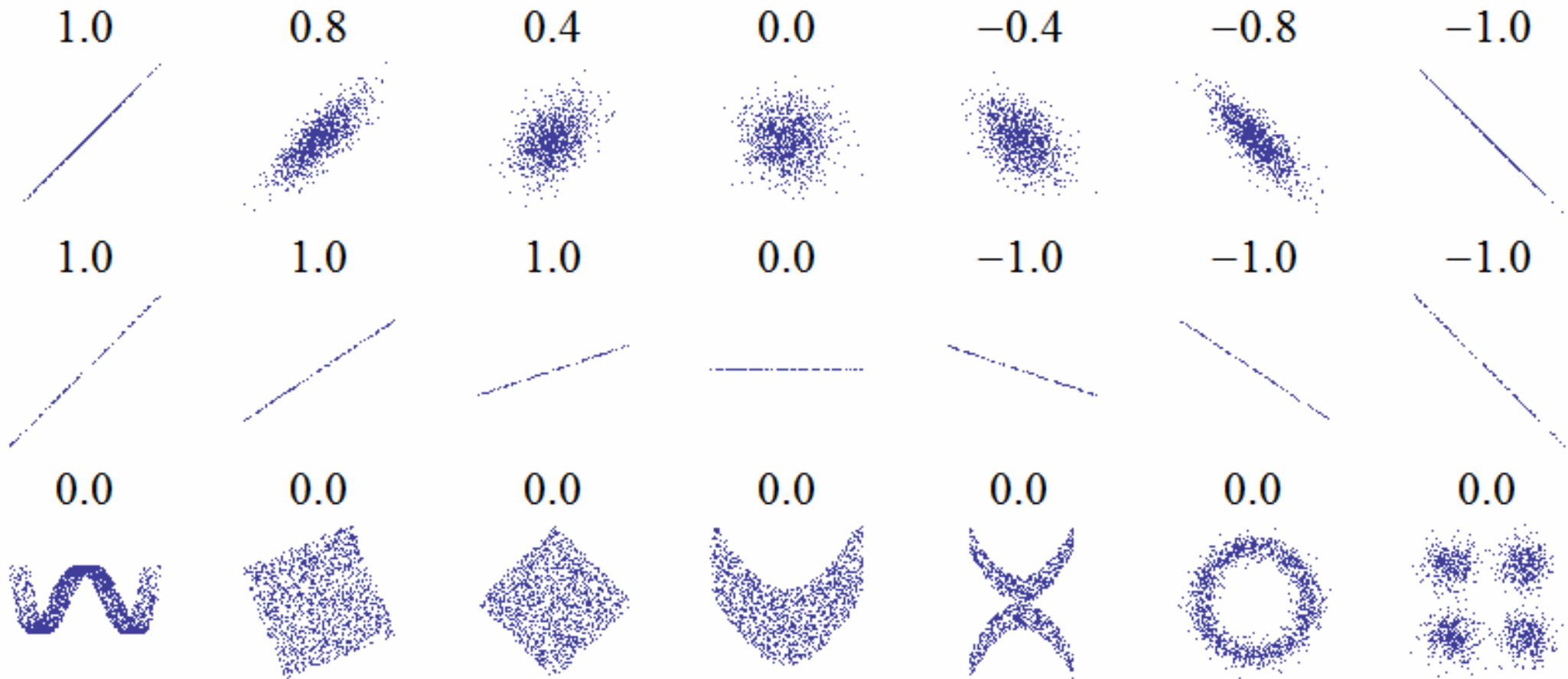
$$r_{xy} = \frac{s_{xy}}{s_x s_y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{s_x s_y (n - 1)}$$

In Excel use function:

◆ =CORREL(data)

$$r_{xy} = 0.94$$

mice.xls



Wikipedia

?

If we have only 2 data points in x and y datasets, what values would you expect for correlation b/w x and y ?

◆ Discrete and continuous probability distributions

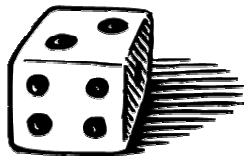
- ◆ discrete probability distribution
- ◆ continuous probability distribution
- ◆ normal probability distribution

Random variable

A numerical description of the outcome of an experiment.

A random variable is always a numerical measure.

Roll a die



Discrete random variable

A random variable that may assume either a finite number of values or an infinite sequence of values.

Continuous random variable

A random variable that may assume any numerical value in an interval or collection of intervals.

Number of calls to a reception per hour



Time between calls to a reception



Volume of a sample in a tube



Weight, height, blood pressure, etc



Probability distribution

A description of how the probabilities are distributed over the values of the random variable.

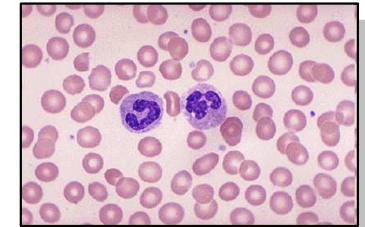
Probability function

A function, denoted by $f(x)$, that provides the probability that x assumes a particular value for a discrete random variable.

Number of cells under microscope

Random variable X :

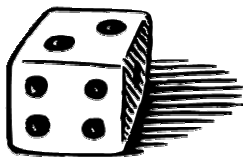
- $x = 0$
- $x = 1$
- $x = 2$
- $x = 3$
- ...



Roll a die

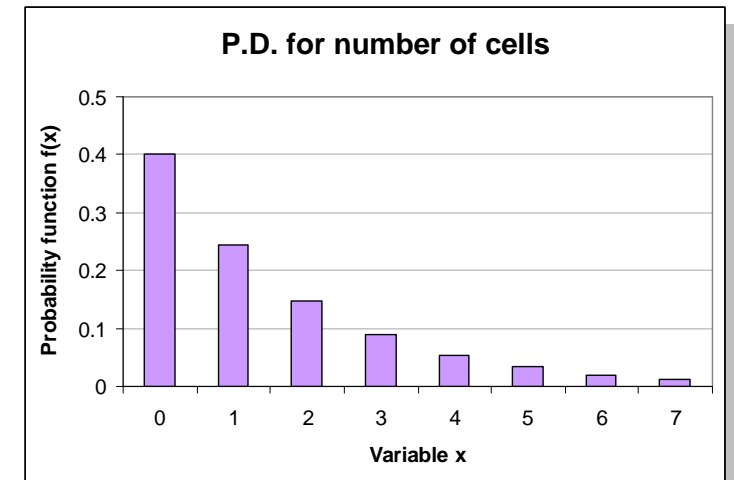
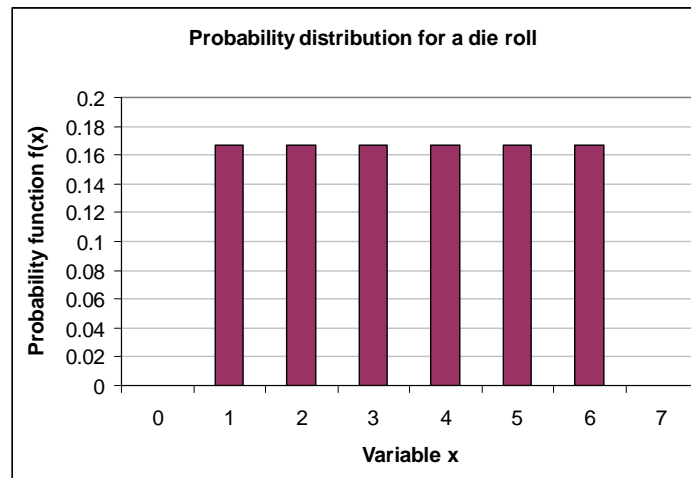
Random variable X :

- $x = 1$
- $x = 2$
- $x = 3$
- $x = 4$
- $x = 5$
- $x = 6$



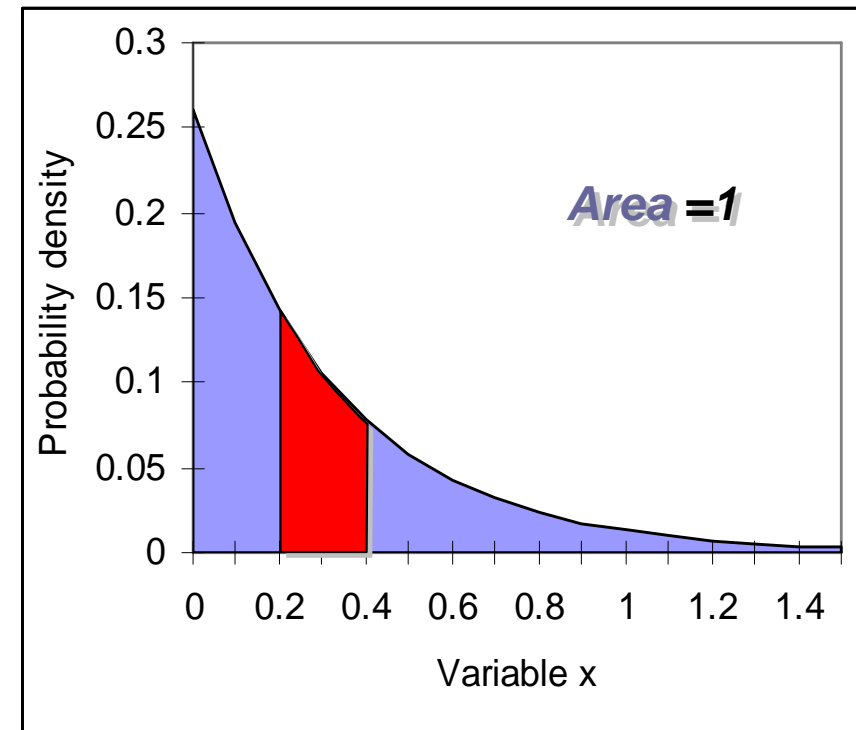
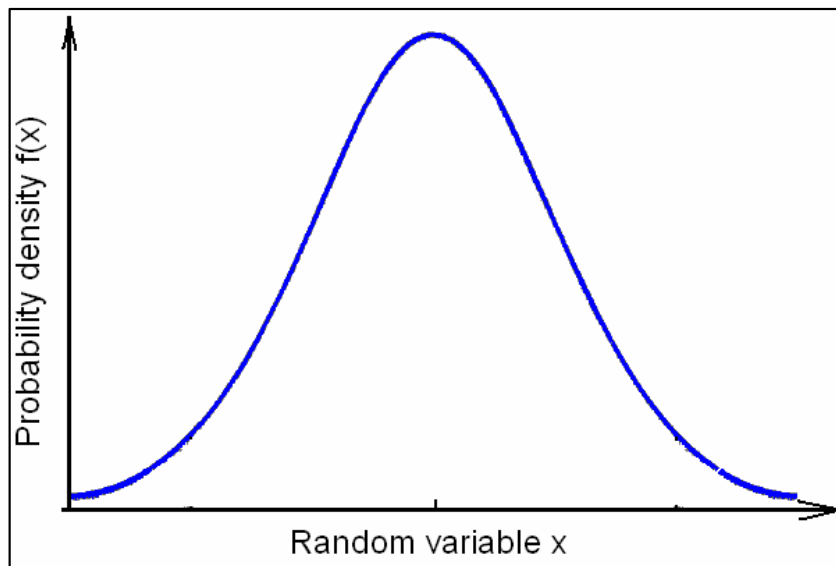
$$f(x) \geq 0$$

$$\sum f(x) = 1$$



Probability density function

A function used to compute probabilities for a continuous random variable. The area under the graph of a probability density function over an interval represents probability.

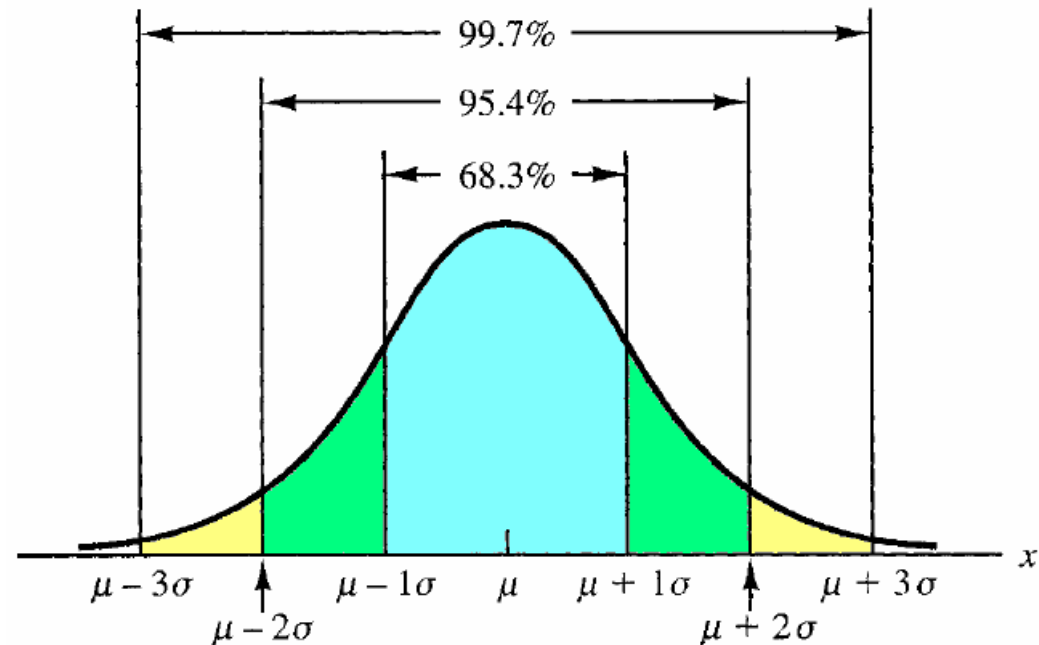


$$\int_x f(x) = 1$$

Normal probability distribution

A continuous probability distribution. Its probability density function is bell shaped and determined by its mean μ and standard deviation σ .

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



In Excel use the function:

◆ = **NORMDIST(x,m,s,false)** for probability density function

◆ = **NORMDIST(x,m,s,true)** for cumulative probability function of normal distribution (area from left to x)

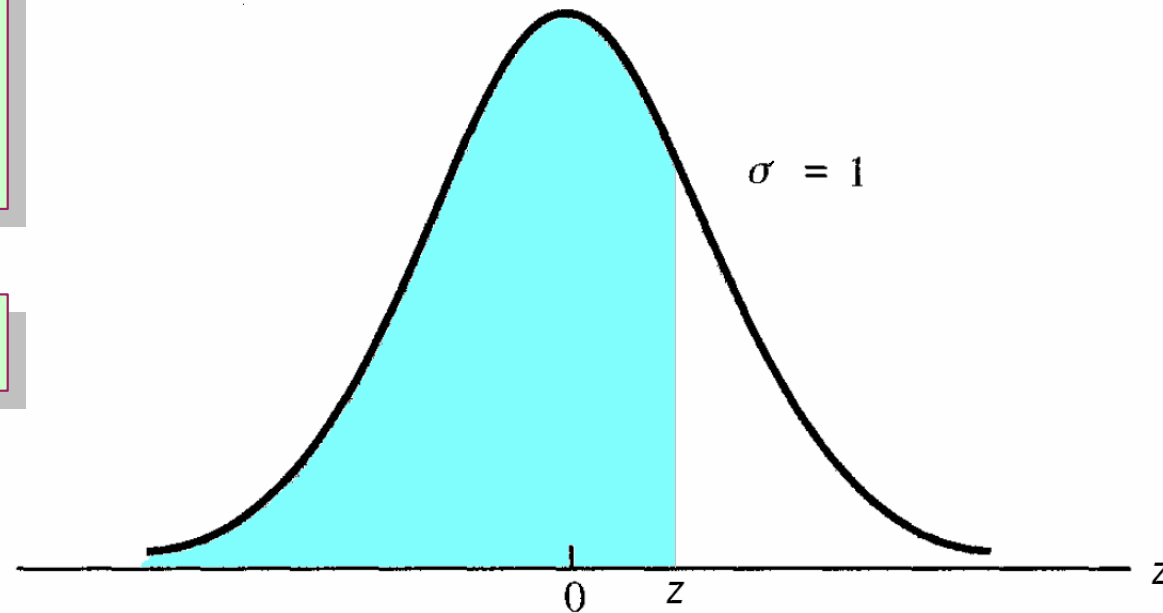
Standard normal probability distribution

A normal distribution with a mean of zero and a standard deviation of one.

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

$$z = \frac{x - \mu}{\sigma}$$

$$x = z\sigma + \mu$$



In Excel use the function:

◆ = NORMSDIST(z)

Example

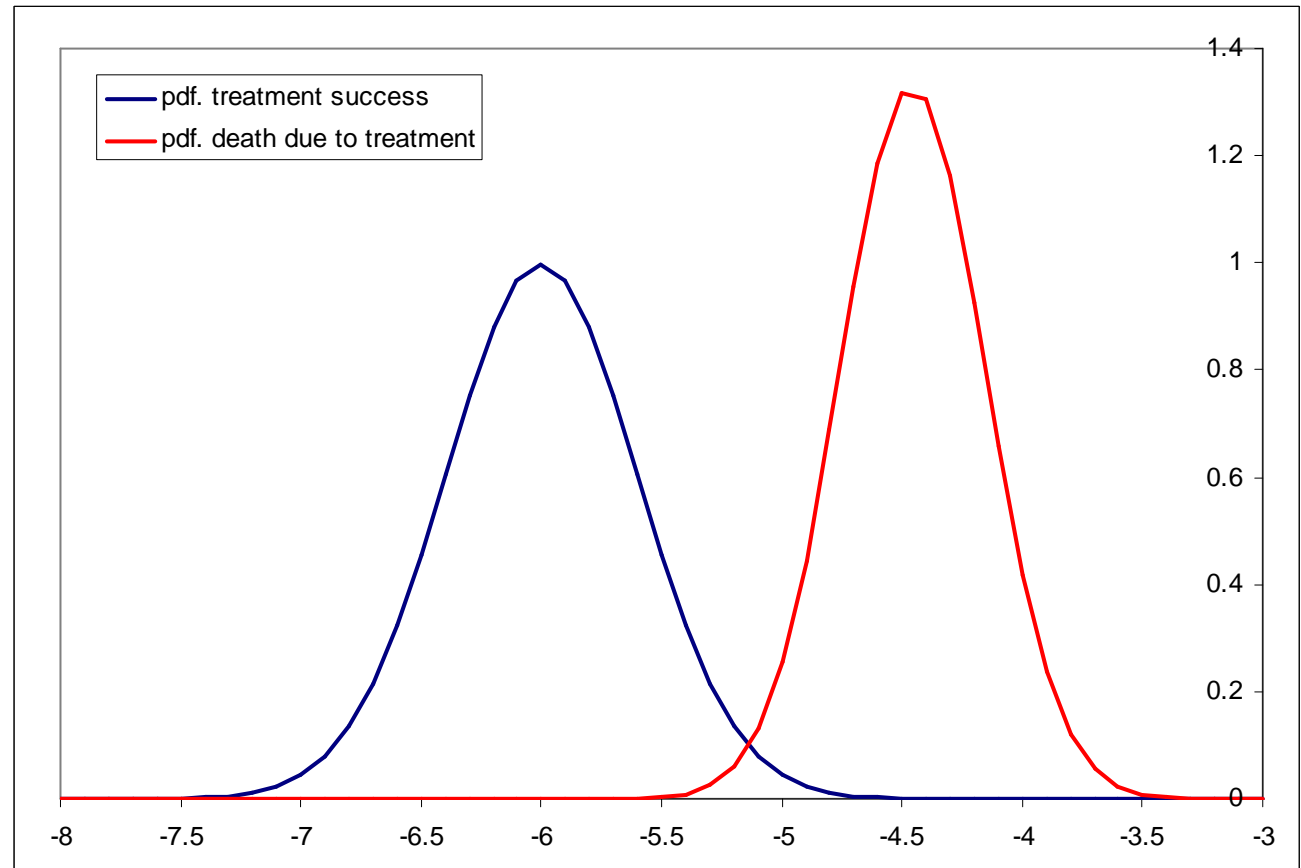
Assume that you have developed an extremely efficient chemical treatment for glioblastoma. During tests on animal models it was found that the substance X, which you use, is able to kill all tumor cells (theoretically), but being given at high concentration it leads to the death of a patient due to intoxication. As the survived cancer cells fast evolve into resistant form, the efficiency of the treatment is significantly reduced if the second course is given. Therefore the treatment should be performed in one injection.

The experimental data suggest that the average concentration needed for the positive treatment is 1 $\mu\text{g}/\text{kg}$. The concentration needed for effective treatment is, of course, a random variable. Being presented in \log_{10} scale and in g/kg , it can be approximated by a normal random variable with mean of -6 and standard deviation of 0.4 .

The 50% lethal dose for human is 35 $\mu\text{g}/\text{kg}$. And the tests on animals suggest that in \log_{10} scale it has a normal distribution as well with the standard deviation of 0.3 .

parameter	ug/kg	log scale
mean positive treatment	1	-6
std positive treatment	x	0.4
mean lethal dose	35	-4.456
std lethal dose	x	0.3

parameter	ug/kg	log scale
mean positive treatment	1	-6
std positive treatment	x	0.4
mean lethal dose	35	-4.456
std lethal dose	x	0.3

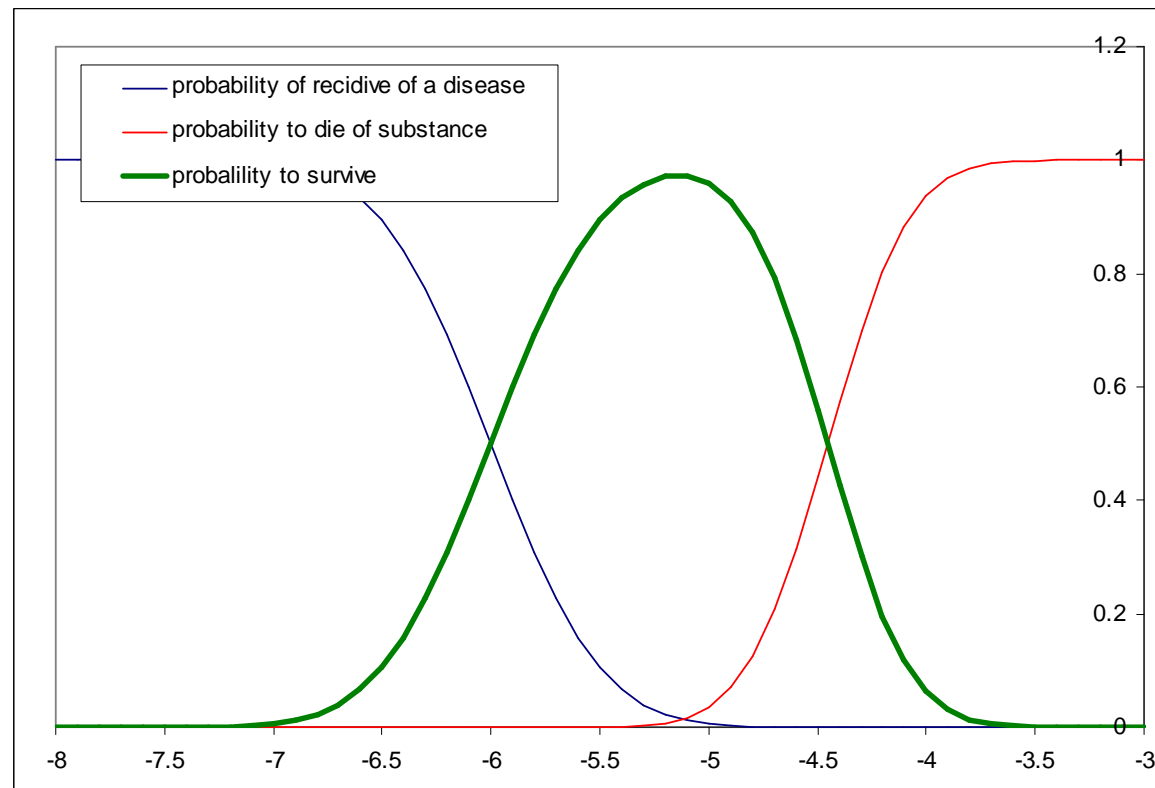


In Excel use the function:

◆ = `NORMDIST(x,mean,std,FALSE)`

- ◆ Probability to die from disease = inverse probability to treat
- ◆ Over-dose and disease behaviors are independent =>

$$P(\text{ survive }) = P(\text{ heal disease }) * P(\text{ survive treatment })$$



In Excel use the function:

◆ = `NORMDIST(x,mean,std,TRUE)`

◆ Sampling distribution

- ◆ sample and population and their parameters
- ◆ central limit theorem
- ◆ types of sampling

Population parameter

A numerical value used as a summary measure for a population (e.g., the mean μ , variance σ^2 , standard deviation σ , proportion π)

POPULATION

μ – mean
 σ^2 – variance
 N – number of elements
 (usually $N=\infty$)

SAMPLE

m, \bar{x} – mean
 s^2 – variance
 n – number of elements

Sample statistic

A numerical value used as a summary measure for a sample (e.g., the sample mean m , the sample variance s^2 , and the sample standard deviation s)

All existing laboratory
Mus musculus



mice.xls

790 mice from different strains

<http://phenome.jax.org>

ID	Strain	Sex	Starting age	Ending age	Starting weight	Ending weight	Weight change	Bleeding time	Ionized Ca in blood	Blood pH	Bone mineral density	Lean tissues weight	Fat weight
1	129S1/SvlmJ	f	66	116	19.3	20.5	1.062	64	1.2	7.24	0.0605	14.5	4.4
2	129S1/SvlmJ	f	66	116	19.1	20.8	1.089	78	1.15	7.27	0.0553	13.9	4.4
3	129S1/SvlmJ	f	66	108	17.9	19.8	1.106	90	1.16	7.26	0.0546	13.8	2.9
368	129S1/SvlmJ	f	72	114	18.3	21	1.148	65	1.26	7.22	0.0599	15.4	4.2
369	129S1/SvlmJ	f	72	115	20.2	21.9	1.084	55	1.23	7.3	0.0623	15.6	4.3
370	129S1/SvlmJ	f	72	116	18.8	22.1	1.176		1.21	7.28	0.0626	16.4	4.3
371	129S1/SvlmJ	f	72	119	19.4	21.3	1.098	49	1.24	7.24	0.0632	16.6	5.4
372	129S1/SvlmJ	f	72	122	18.3	20.1	1.098	73	1.17	7.19	0.0592	16	4.1
4	129S1/SvlmJ	f	66	109	17.2	18.9	1.099	41	1.25	7.29	0.0513	14	3.2
5	129S1/SvlmJ	f	66	112	19.7	21.3	1.081	129	1.14	7.22	0.0501	16.3	5.2
10	129S1/SvlmJ	m	66	112	24.3	24.7	1.016	119	1.13	7.24	0.0533	17.6	6.8
364	129S1/SvlmJ	m	72	114	25.3	27.2	1.075	64	1.25	7.27	0.0596	19.3	5.8
365	129S1/SvlmJ	m	72	115	21.4	23.9	1.117	48	1.25	7.28	0.0563	17.4	5.7
366	129S1/SvlmJ	m	72	118	24.5	26.3	1.073	59	1.25	7.26	0.0609	17.8	7.1
367	129S1/SvlmJ	m	72	122	24	26	1.083	69	1.29	7.26	0.0584	19.2	4.6
6	129S1/SvlmJ	m	66	116	21.6	23.3	1.079	78	1.15	7.27	0.0497	17.2	5.7
7	129S1/SvlmJ	m	66	107	22.7	26.5	1.167	90	1.18	7.28	0.0493	18.7	7
8	129S1/SvlmJ	m	66	108	25.4	27.4	1.079	35	1.24	7.26	0.0538	18.9	7.1
9	129S1/SvlmJ	m	66	109	24.4	27.5	1.127	43	1.29	7.29	0.0539	19.5	7.1

mice.xls

790 mice from different strains

<http://phenome.jax.org>

ID	Strain	Sex	Starting age	Ending age	Starting weight	Ending weight	Weight change	Bleeding time	Ionized Ca in blood	Blood pH	Bone mineral density	Lean tissues weight	Fat weight
1	129S1/SvlmJ	f	66	116	19.3	20.5	1.062	64	1.2	7.24	0.0605	14.5	4.4
2	129S1/SvlmJ	f	66	116	19.1	20.8	1.089	78	1.15	7.27	0.0553	13.9	4.4
3	129S1/SvlmJ	f	66	108	17.9	19.8	1.106	90	1.16	7.26	0.0546	13.8	2.9
368	129S1/SvlmJ	f	72	114	18.3	21	1.148	65	1.26	7.22	0.0599	15.4	4.2
369	129S1/SvlmJ	f	72	115	20.2	21.9	1.084	55	1.23	7.3	0.0623	15.6	4.3
370	129S1/SvlmJ	f	72	116	18.8	22.1	1.176		1.21	7.28	0.0626	16.4	4.3
371	129S1/SvlmJ	f	72	119	19.4	21.3	1.098	49	1.24	7.24	0.0632	16.6	5.4
372	129S1/SvlmJ	f	72	122	18.3	20.1	1.098	73	1.17	7.19	0.0592	16	4.1
4	129S1/SvlmJ	f	66	109	17.2	18.9	1.099	41	1.25	7.29	0.0513	14	3.2
5	129S1/SvlmJ	f	66	112	19.7	21.3	1.081	129	1.14	7.22	0.0501	16.3	5.2
10	129S1/SvlmJ	m	66	112	24.3	24.7	1.016	119	1.13	7.24	0.0533	17.6	6.8
364	129S1/SvlmJ	m	72	114	25.3	27.2	1.075	64	1.25	7.27	0.0596	19.3	5.8
365	129S1/SvlmJ	m	72	115	21.4	23.9	1.117	48	1.25	7.28	0.0563	17.4	5.7
366	129S1/SvlmJ	m	72	118	24.5	26.3	1.073	59	1.25	7.26	0.0609	17.8	7.1
367	129S1/SvlmJ	m	72	122	24	26	1.083	69	1.29	7.26	0.0584	19.2	4.6
6	129S1/SvlmJ	m	66	116	21.6	23.3	1.079	78	1.15	7.27	0.0497	17.2	5.7
7	129S1/SvlmJ	m	66	107	22.7	26.5	1.167	90	1.18	7.28	0.0493	18.7	7
8	129S1/SvlmJ	m	66	108	25.4	27.4	1.079	35	1.24	7.26	0.0538	18.9	7.1
9	129S1/SvlmJ	m	66	109	24.4	27.5	1.127	43	1.29	7.29	0.0539	19.5	7.1

1. Add a column to the table
2. Fill it with `=RAND()`
3. Sort all the table by this column

4. Assume that these mice is a population with size $N=790$. Build 3 samples with $n=20$
5. Calculate m , s for ending weight and p – proportion of males for each sample

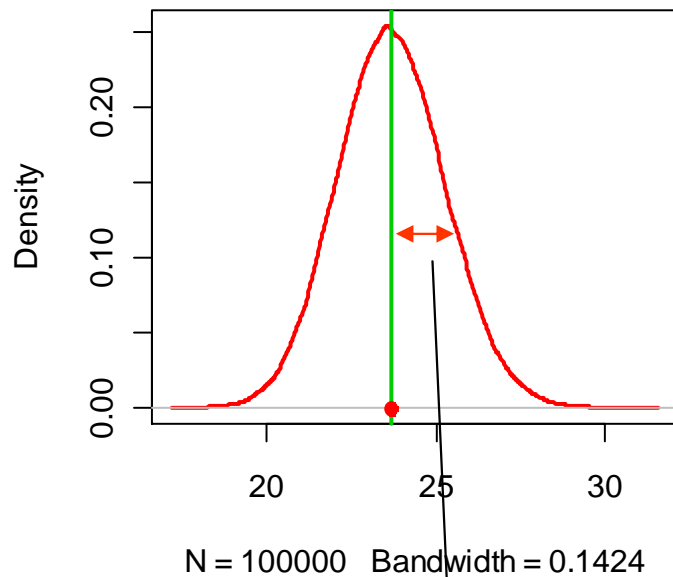
Point estimator

The sample statistic, such as m , s , or p , that provides the point estimation the population parameters μ , σ , π .

Sampling distribution

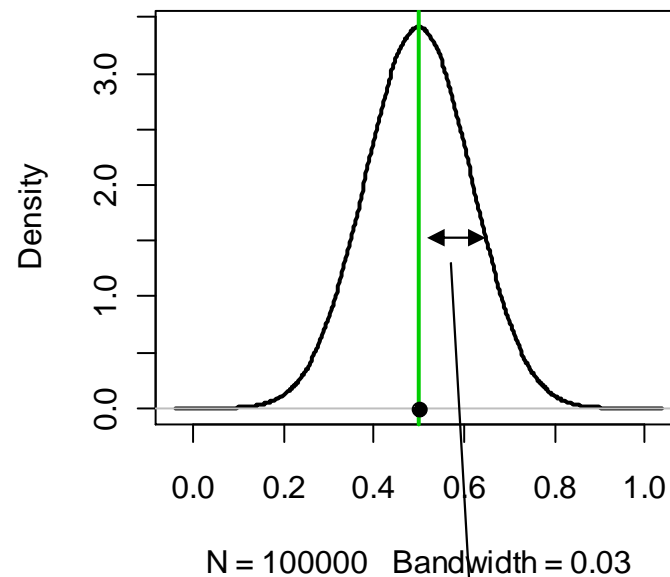
A probability distribution consisting of all possible values of a sample statistic.

Distribution of m



$$\sigma_m = \frac{\sigma}{\sqrt{n}}$$

Distribution of p



$$\sigma_p = \sqrt{\frac{\pi(1-\pi)}{n}}$$

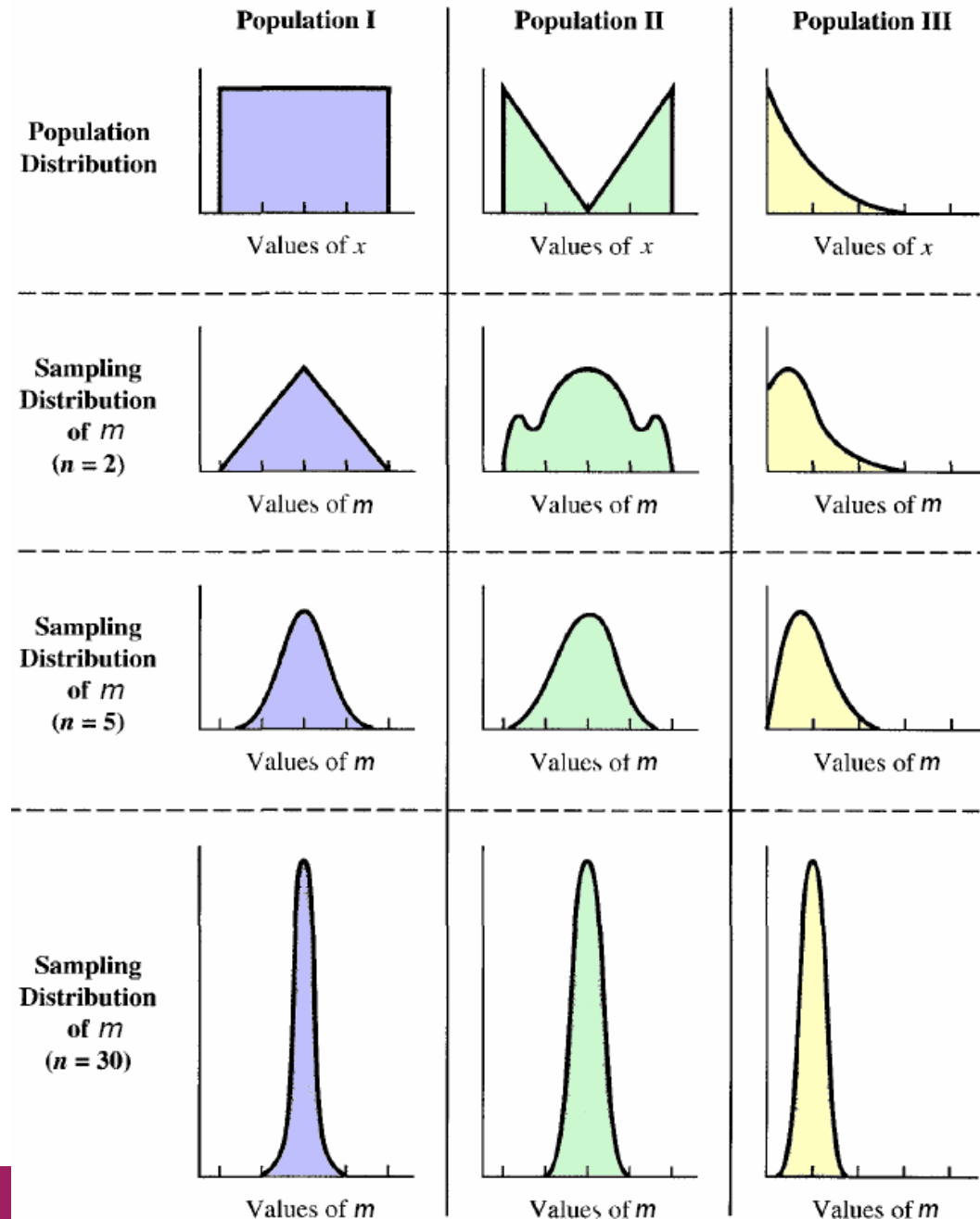
$$E(m) = \mu$$

$$E(p) = \pi$$

Central limit theorem

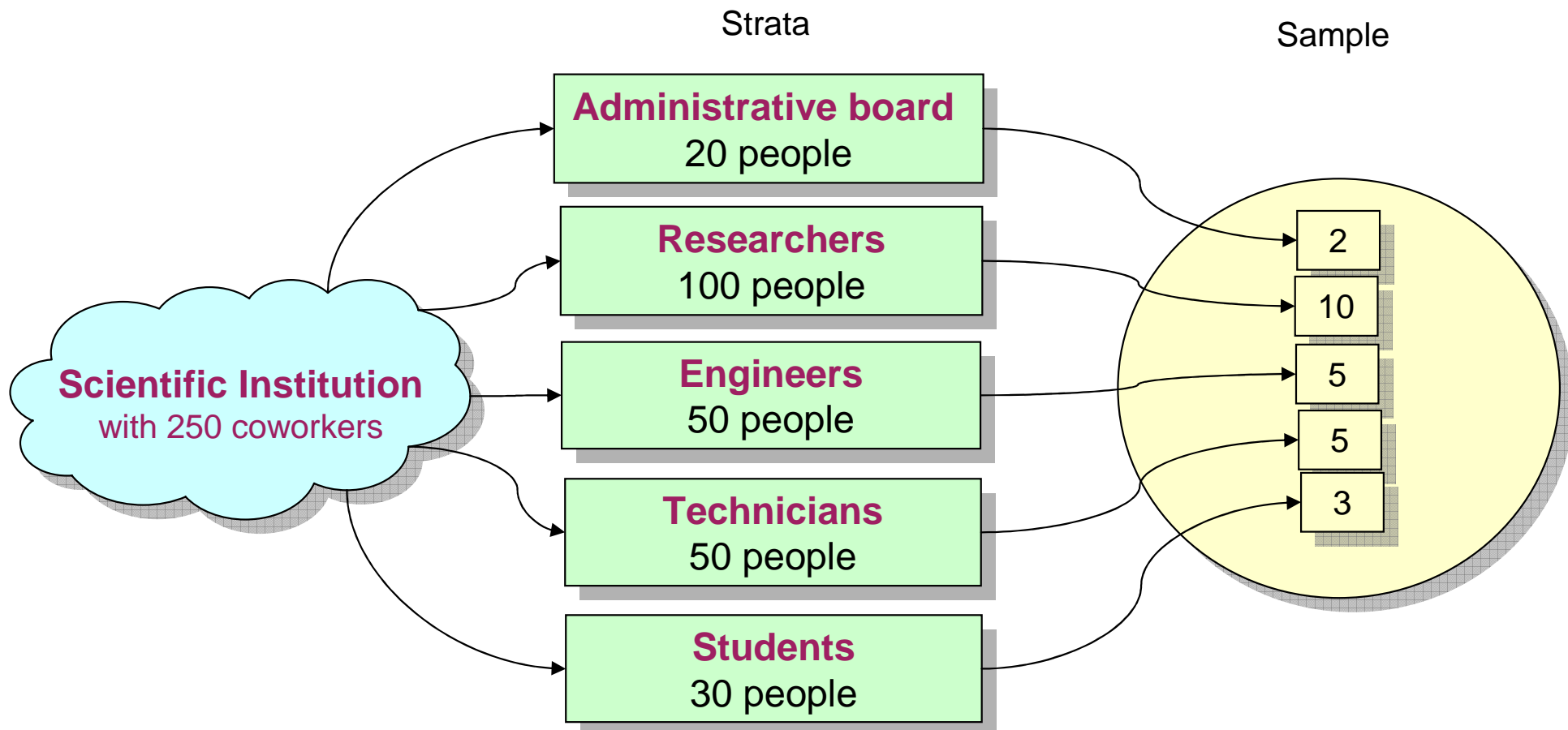
In selecting simple random sample of size n from a population, the **sampling distribution of the sample mean m can be approximated by a normal distribution** as the sample size becomes large

In practice if the sample size is $n > 30$, the normal distribution is a good approximation for the sample mean for any initial distribution.



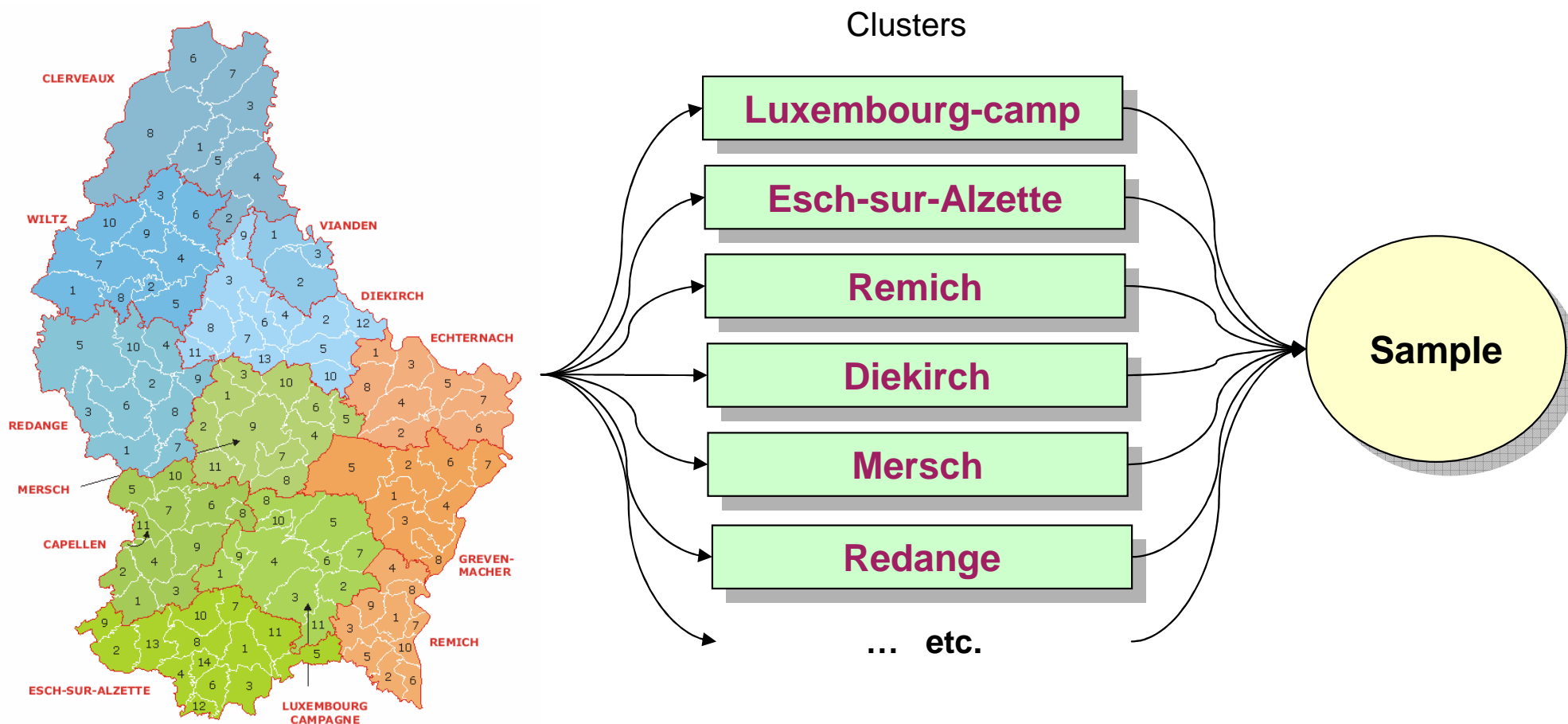
Stratified random sampling

A probability sampling method in which the population is first divided into strata and a simple random sample is then taken from each stratum.



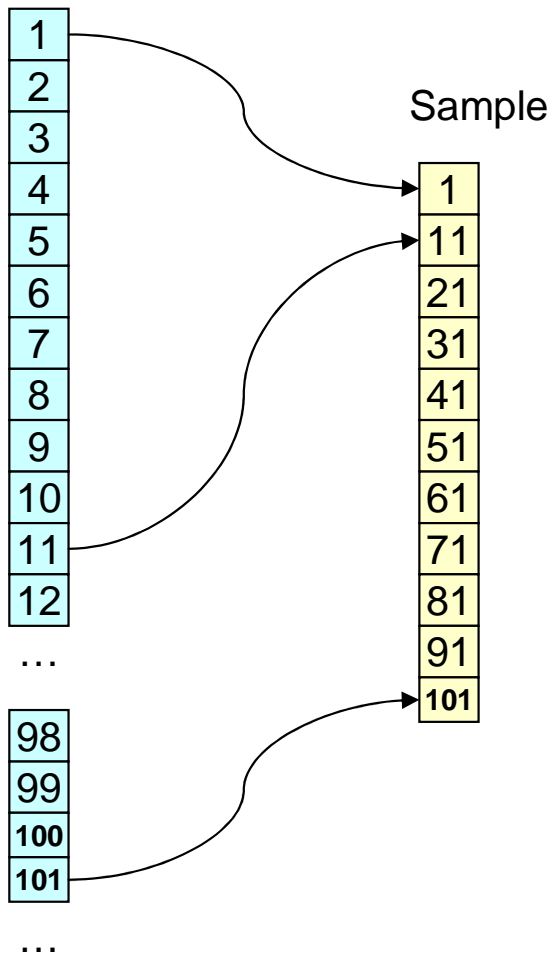
Cluster sampling

A probability sampling method in which the population is first divided into clusters and then a simple random sample of the clusters is taken.



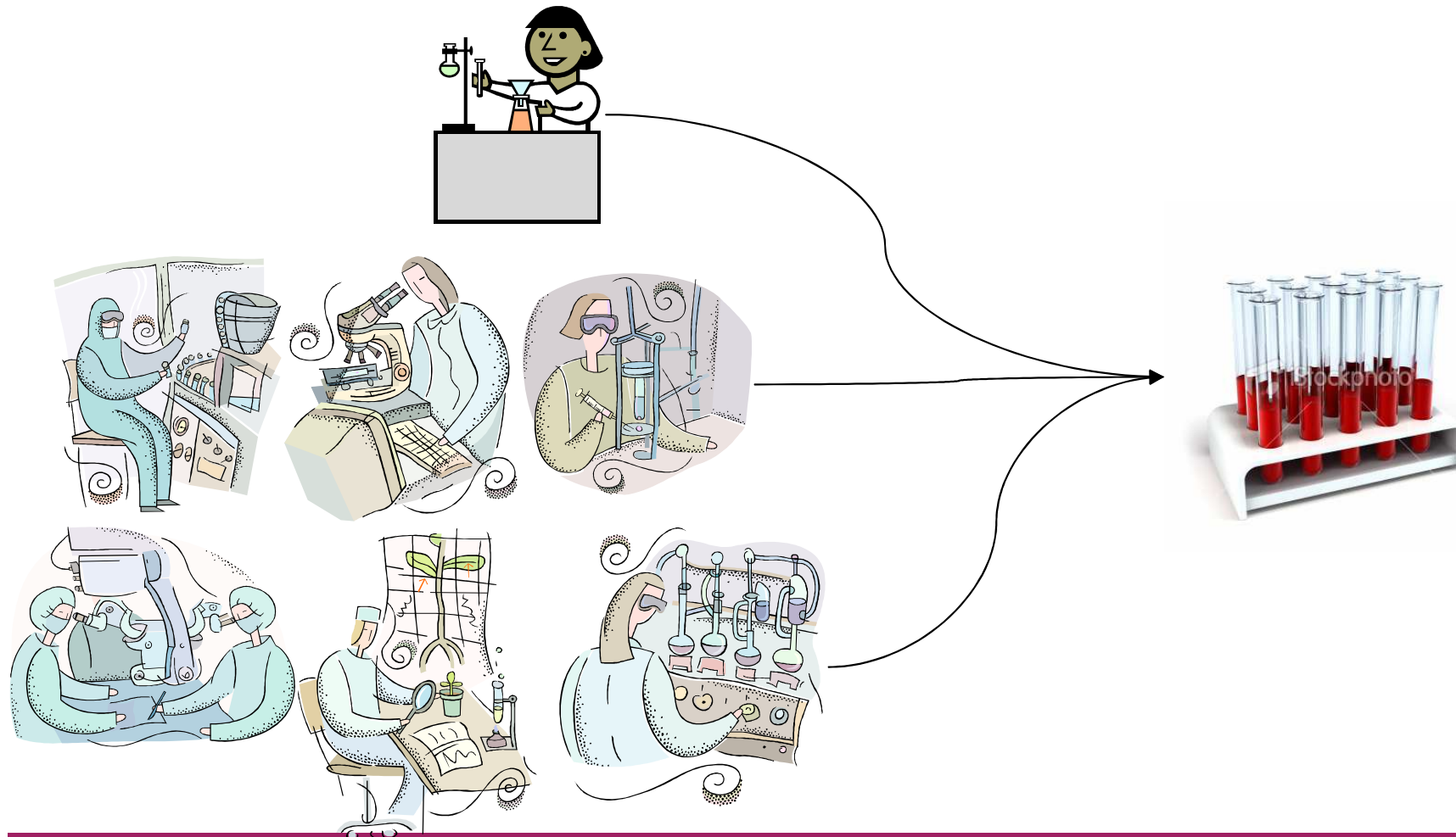
Systematic sampling

A probability sampling method in which we randomly select one of the first k elements and then select every k -th element thereafter.



Convenience sampling

A nonprobability method of sampling whereby elements are selected for the sample on the basis of convenience.



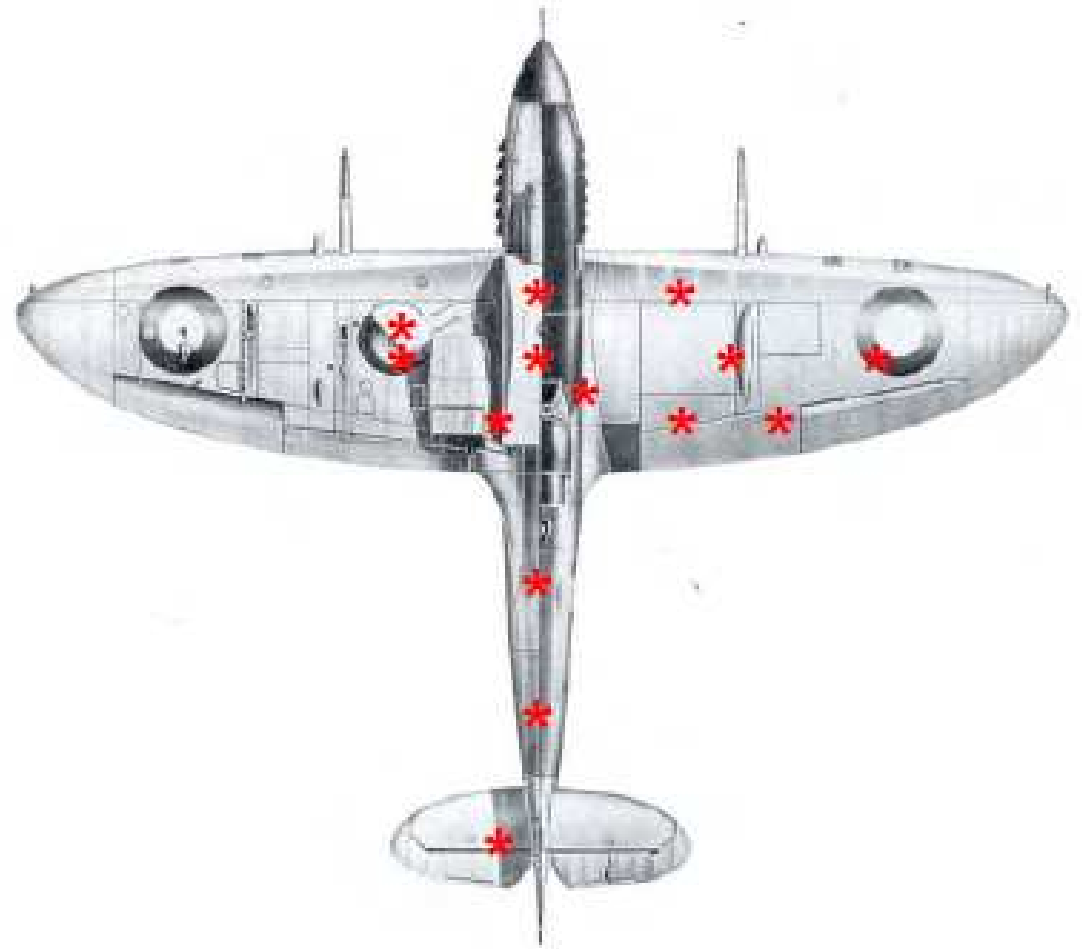
Judgment sampling

A nonprobability method of sampling whereby elements are selected for the sample based on the judgment of the person doing the study.



Perform of a selection of most confident or most experienced experts.

Spitfire: analysis of the damage



Were to put additional protection?

◆ Interval estimation

- ◆ interval estimation
- ◆ population mean: σ known
- ◆ population proportion
- ◆ population mean: σ unknown
- ◆ Student's distribution
- ◆ estimation the size of a sample

Population parameter

A numerical value used as a summary measure for a population (e.g., the mean μ , variance σ^2 , standard deviation σ , proportion π)

POPULATION

μ – mean
 σ^2 – variance
 N – number of elements (usually $N=\infty$)

SAMPLE

m, \bar{x} – mean
 s^2 – variance
 n – number of elements

Sample statistic

A numerical value used as a summary measure for a sample (e.g., the sample mean m , the sample variance s^2 , and the sample standard deviation s)

All existing laboratory
Mus musculus



mice.txt

790 mice from different strains

<http://phenome.jax.org>

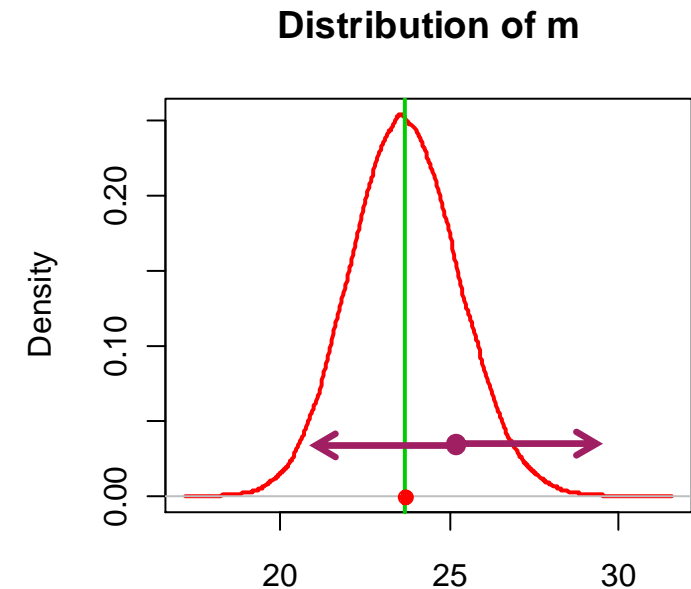
ID	Strain	Sex	Starting age	Ending age	Starting weight	Ending weight	Weight change	Bleeding time	Ionized Ca in blood	Blood pH	Bone mineral density	Lean tissues weight	Fat weight
1	129S1/SvlmJ	f	66	116	19.3	20.5	1.062	64	1.2	7.24	0.0605	14.5	4.4
2	129S1/SvlmJ	f	66	116	19.1	20.8	1.089	78	1.15	7.27	0.0553	13.9	4.4
3	129S1/SvlmJ	f	66	108	17.9	19.8	1.106	90	1.16	7.26	0.0546	13.8	2.9
368	129S1/SvlmJ	f	72	114	18.3	21	1.148	65	1.26	7.22	0.0599	15.4	4.2
369	129S1/SvlmJ	f	72	115	20.2	21.9	1.084	55	1.23	7.3	0.0623	15.6	4.3
370	129S1/SvlmJ	f	72	116	18.8	22.1	1.176		1.21	7.28	0.0626	16.4	4.3
371	129S1/SvlmJ	f	72	119	19.4	21.3	1.098	49	1.24	7.24	0.0632	16.6	5.4
372	129S1/SvlmJ	f	72	122	18.3	20.1	1.098	73	1.17	7.19	0.0592	16	4.1
4	129S1/SvlmJ	f	66	109	17.2	18.9	1.099	41	1.25	7.29	0.0513	14	3.2
5	129S1/SvlmJ	f	66	112	19.7	21.3	1.081	129	1.14	7.22	0.0501	16.3	5.2
10	129S1/SvlmJ	m	66	112	24.3	24.7	1.016	119	1.13	7.24	0.0533	17.6	6.8
364	129S1/SvlmJ	m	72	114	25.3	27.2	1.075	64	1.25	7.27	0.0596	19.3	5.8
365	129S1/SvlmJ	m	72	115	21.4	23.9	1.117	48	1.25	7.28	0.0563	17.4	5.7
366	129S1/SvlmJ	m	72	118	24.5	26.3	1.073	59	1.25	7.26	0.0609	17.8	7.1
367	129S1/SvlmJ	m	72	122	24	26	1.083	69	1.29	7.26	0.0584	19.2	4.6
6	129S1/SvlmJ	m	66	116	21.6	23.3	1.079	78	1.15	7.27	0.0497	17.2	5.7
7	129S1/SvlmJ	m	66	107	22.7	26.5	1.167	90	1.18	7.28	0.0493	18.7	7
8	129S1/SvlmJ	m	66	108	25.4	27.4	1.079	35	1.24	7.26	0.0538	18.9	7.1
9	129S1/SvlmJ	m	66	109	24.4	27.5	1.127	43	1.29	7.29	0.0539	19.5	7.1

Interval estimate

An estimate of a population parameter that provides an interval believed to contain the value of the parameter. For the interval estimates in this chapter, it has the form: point estimate \pm margin of error.

Margin of error

The \pm value added to and subtracted from a point estimate in order to develop an interval estimate of a population parameter.



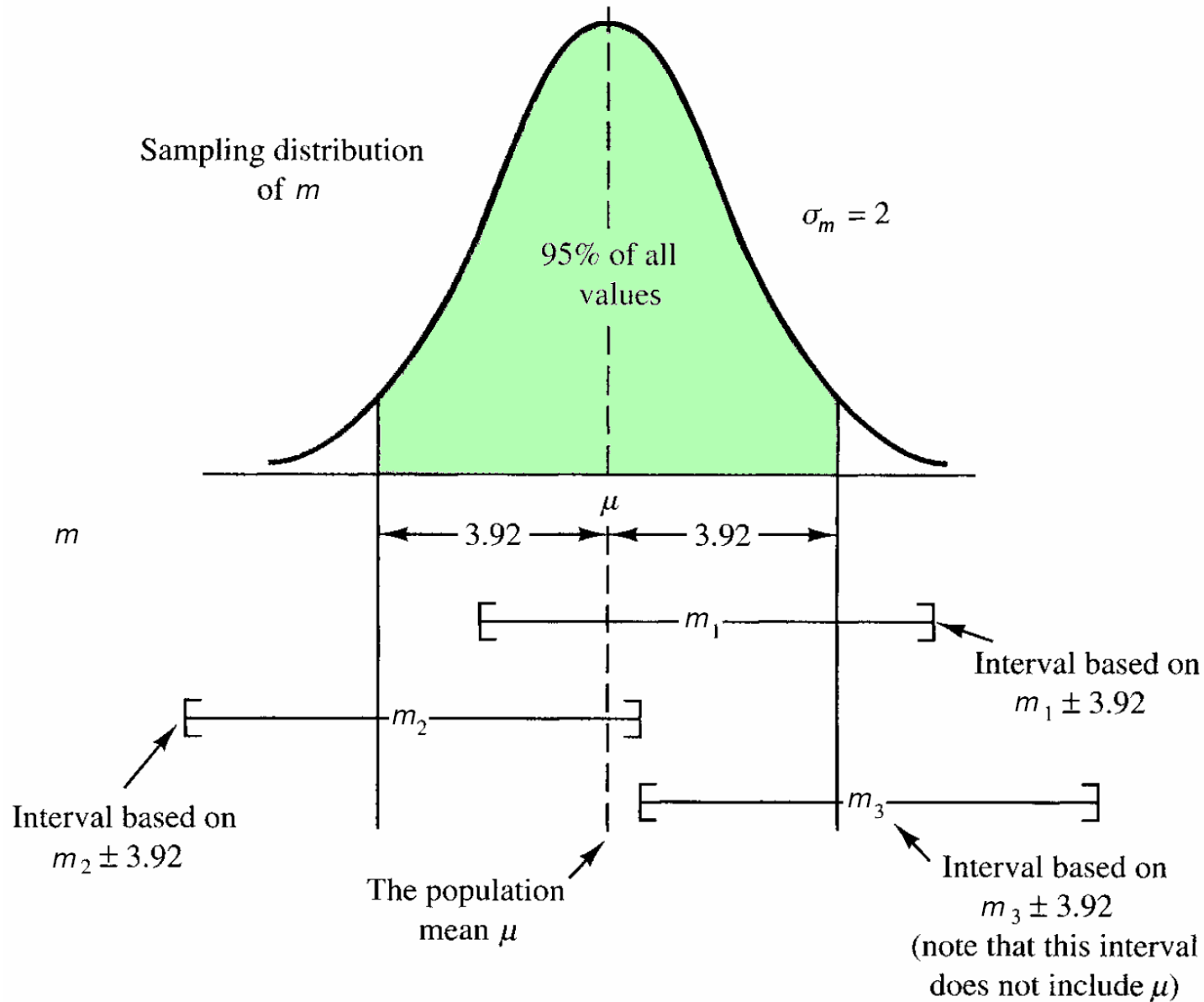
$$\mu = m \pm \text{margin of error}$$

σ known

The condition existing when historical data or other information provides a good value for the population standard deviation prior to taking a sample. The interval estimation procedure uses this known value of σ in computing the margin of error.

σ unknown

The condition existing when no good basis exists for estimating the population standard deviation prior to taking the sample. The interval estimation procedure uses the sample standard deviation s in computing the margin of error.



Confidence level

The confidence associated with an interval estimate. For example, if an interval estimation procedure provides intervals such that 95% of the intervals formed using the procedure will include the population parameter, the interval estimate is said to be constructed at the 95% confidence level.

Confidence interval

Another name for an interval estimate.

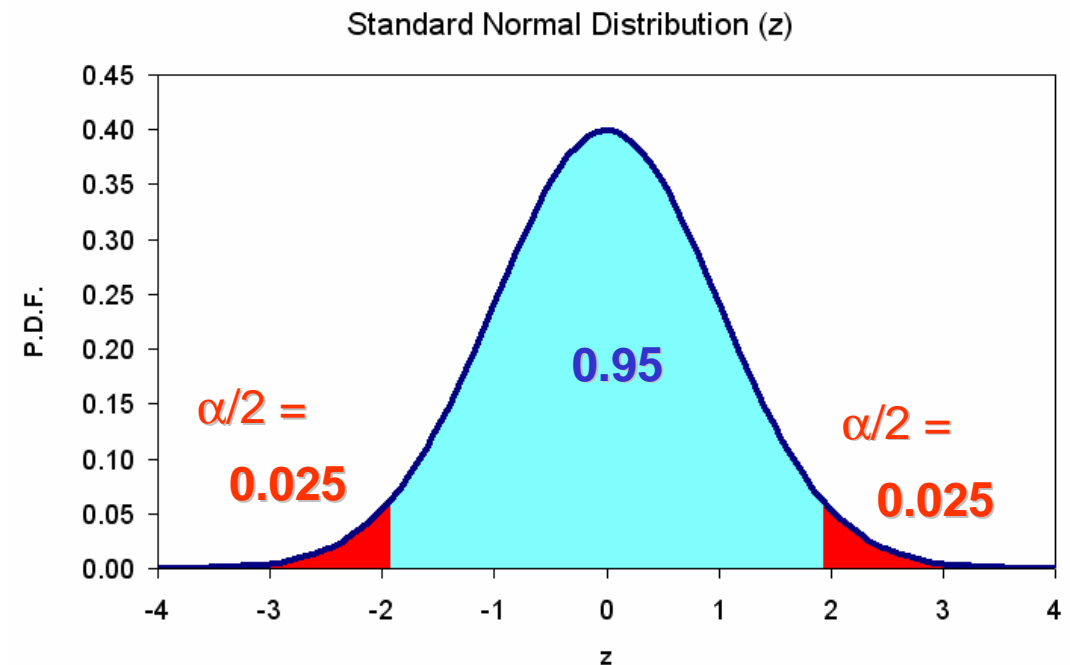
$$\mu = m \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

For 95 % confidence $\alpha = 0.05$, which means that in each tail we have 0.025. Corresponding $z_{\alpha/2} = 1.96$

In Excel use one of the following functions:

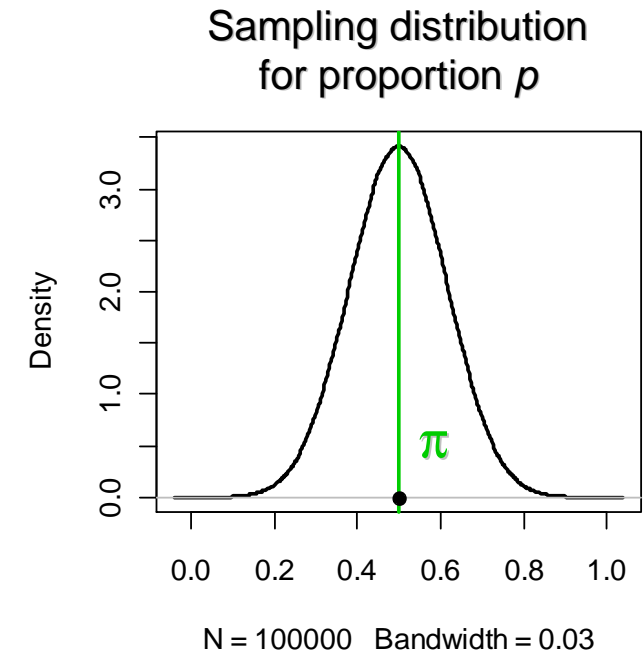
◆ = CONFIDENCE(alpha, σ , n)

◆ = -NORMINV(alpha/2, 0, 1) * σ / SQRT(n)



$$\sigma_p = \sqrt{\frac{\pi(1-\pi)}{n}} \rightarrow \sigma_p = \sqrt{\frac{p(1-p)}{n}}$$

$$\pi = p \pm z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} \quad \text{if } np \geq 5 \text{ and } n(1-p) \geq 5$$



Practical Work

pancreatitis.txt

n= 270
p(never)= 0.214815
sp= 0.024994
E= 0.048988

Define a 95% confidence interval for **never-smoking** proportion of people coming to a hospital

for 95% confidence $z_{0.025} = 1.96$

$$\pi = 21.5 \pm 4.9 \%$$

$$\pi = p \pm z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}}$$

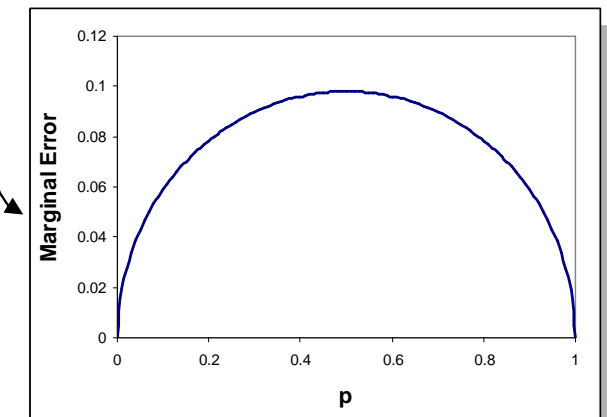
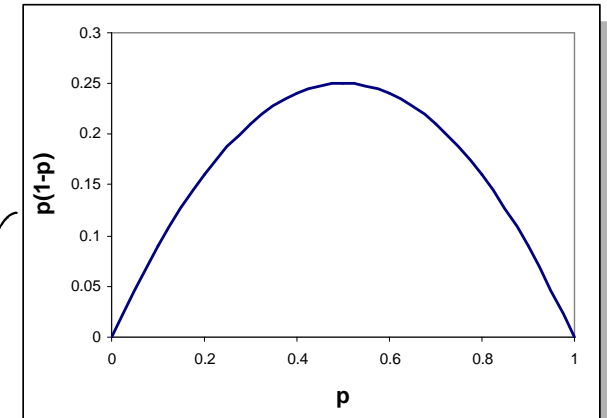
1. The normal distribution is applicable only when enough data points are observed. The rule of thumb is: $np \geq 5$ and $n(1-p) \geq 5$

2. The maximal marginal error is observed when $p=0.5$

3. The estimation of the sample size can be obtained:

$$n = \frac{z_{\alpha/2}^2 p(1-p)}{E^2}$$

$np \geq 5$ and $n(1-p) \geq 5$



where p is a best guess for π or the result of a preliminary study

Assume that we have a sample of 20 mice and would like to estimate an average size of a mice in population.

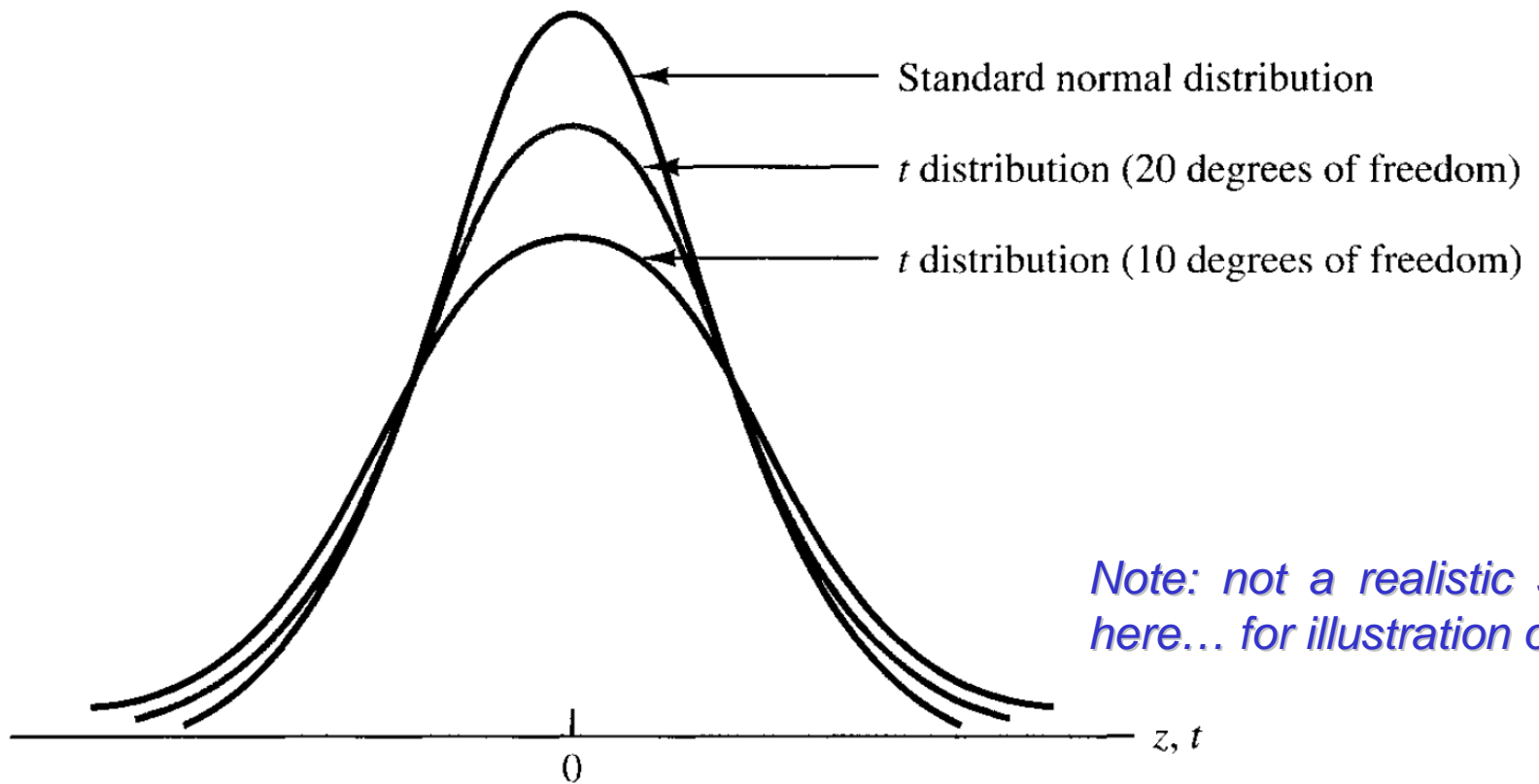
Weight
39.9
19.8
32.4
21
27.5
20.8
21.3
40
10.7
22.6
27
10.8
20.9
14.7
31.4
17.2
11.4
19.1
31.3
14.8

$$m = 22.73$$

$$s = 8.84$$

$$\sigma_m = \frac{\sigma}{\sqrt{n}} \approx \frac{s}{\sqrt{n}}$$

As we replace $\sigma \rightarrow s$, we introduce an additional error and this change the distribution from z to t (Student)



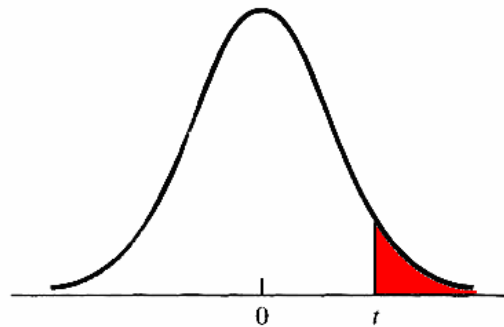
Note: not a realistic scale here... for illustration only

***t*-distribution**

A family of probability distributions that can be used to develop an interval estimate of a population mean whenever the population standard deviation σ is unknown and is estimated by the sample standard deviation s .

Degrees of freedom

A parameter of the t -distribution. When the t distribution is used in the computation of an interval estimate of a population mean, the appropriate t distribution has $n - 1$ degrees of freedom, where n is the size of the simple random sample.



Degrees of Freedom	Area in Upper Tail					
	.20	.10	.05	.025	.01	.005
1	1.376	3.078	6.314	12.706	31.821	63.656
2	1.061	1.886	2.920	4.303	6.965	9.925
3	.978	1.638	2.353	3.182	4.541	5.841
4	.941	1.533	2.132	2.776	3.747	4.604

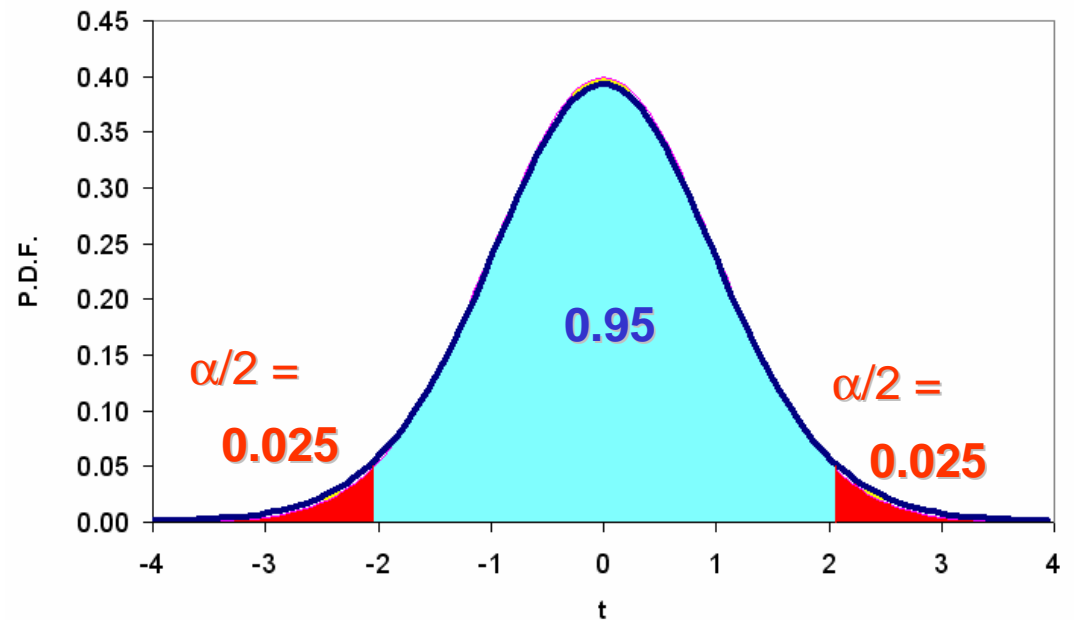
Weight
39.9
19.8
32.4
21
27.5
20.8
21.3
40
10.7
22.6
27
10.8
20.9
14.7
31.4
17.2
11.4
19.1
31.3
14.8

$m = 22.73$
 $s = 8.84$

$s(m) = 1.98$
 $t = 2.09$
 $m.e. = 4.14$

$$\mu = m \pm t_{\alpha/2}^{(n-1)} \frac{s}{\sqrt{n}}$$

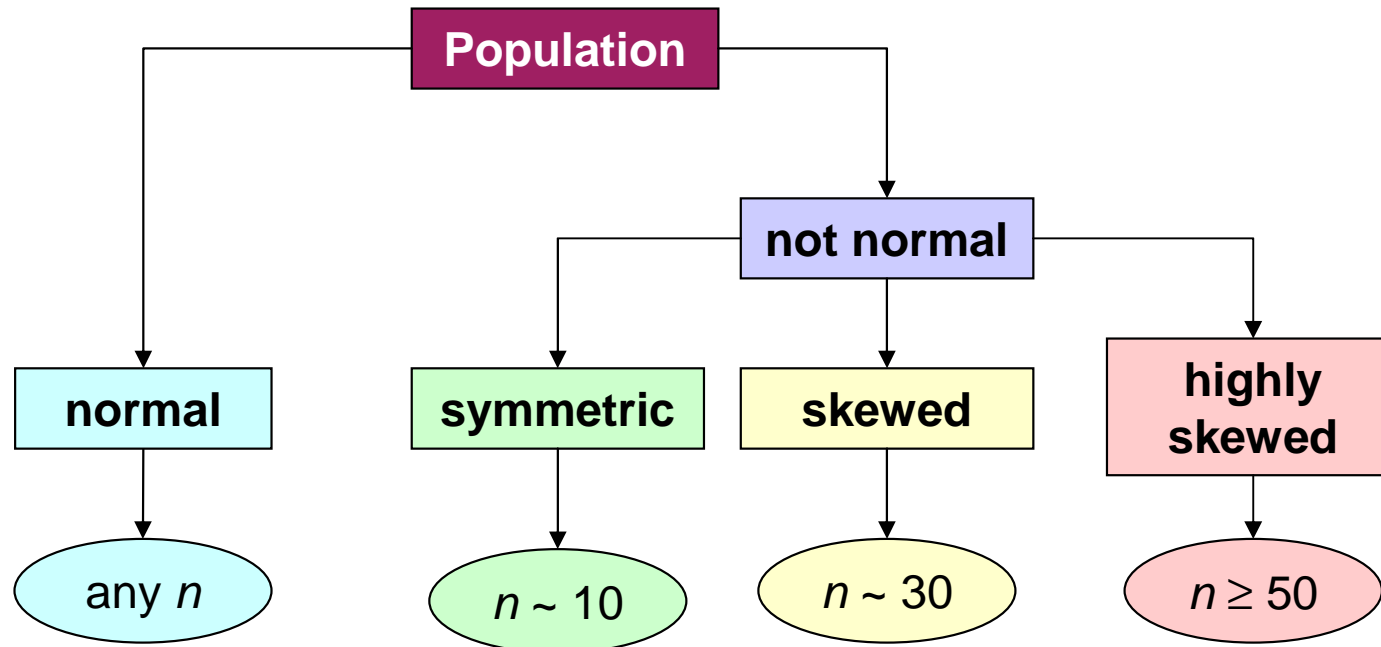
Student Distribution (t), df=19



In Excel use:

◆ = `TINV(alpha, degree-of-freedom)` !!!

Advice 1



$$\mu = m \pm t_{\alpha/2}^{(n-1)} \frac{s}{\sqrt{n}}$$

Advice 2

if $n > 100$ you can use z-statistics instead of t-statistics (error will be $< 1.5\%$)

Let's focus on another aspect: how to select a proper number of experiments.

$$\mu = m \pm E(n, \sigma)$$

$$E(n, \sigma) = E$$

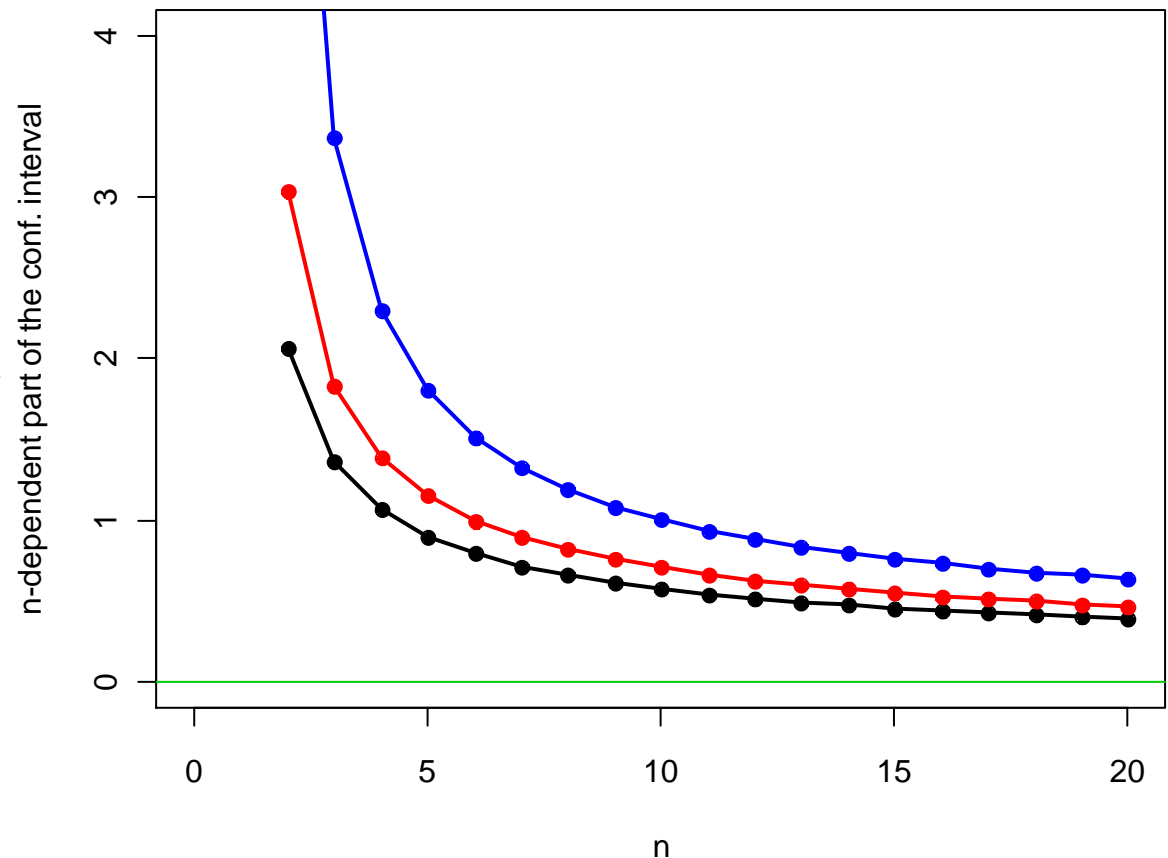
$$n - ?$$

$$E = z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$$n = \frac{z_{\alpha/2}^2 \sigma^2}{E^2}$$

$$n = \frac{z_{\alpha/2}^2 \sigma^2}{E^2}$$

Effect of the Sample Size



Thank you for your attention



to be continued...