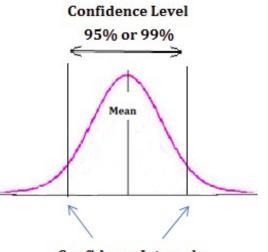
# 7. COMPARISON OF MEANS

# 7.1. Methods for comparison of means

- > There are several *methods*:
- · comparison of a single observed <u>mean</u> with some <u>hypothesized</u> value <sup>1</sup>;
- · comparison of two <u>means</u> arising from <u>paired</u> data <sup>2</sup>;
- comparison of two *means* from *unpaired* data <sup>3</sup>.



- Confidence Intervals
- All can be made by using appropriate *confidence intervals* <sup>4</sup> and *t-tests*.
- The <u>single</u> mean and <u>paired</u> data cases are introduced first.

<sup>&</sup>lt;sup>1</sup> another name "one sample t hypothesis test"

<sup>&</sup>lt;sup>2</sup> e.g. when individual objects are measured **twice** (i.e. once for each type of measurement), or at two different times, etc.: e.g. to study the photosynthetic performance of ten plants in two environments in a greenhouse (<u>shady</u> and <u>sunny</u>), we could <u>measure</u> each individual plant <u>twice</u>, once in the <u>shade</u> and once in the <u>sun</u> - the measures are <u>paired</u> by belonging to the <u>same</u> individual plant.

<sup>&</sup>lt;sup>3</sup> another name "two sample t hypothesis test" and is probably the most common;

<sup>&</sup>lt;sup>4</sup> a range of values we are fairly sure our true value lies in.

# 7.2. Comparison of a single mean with a hypothesized value

Not very *common* in practice but it may be used (table) <sup>1</sup>.

E.g. these are the haemoglobin concentrations of 15 UK adult males admitted into an intensive care unit (ICU) <sup>2</sup>.

The population <u>mean</u> haemoglobin concentration in UK males is **15.0** g/dl <sup>3</sup>.

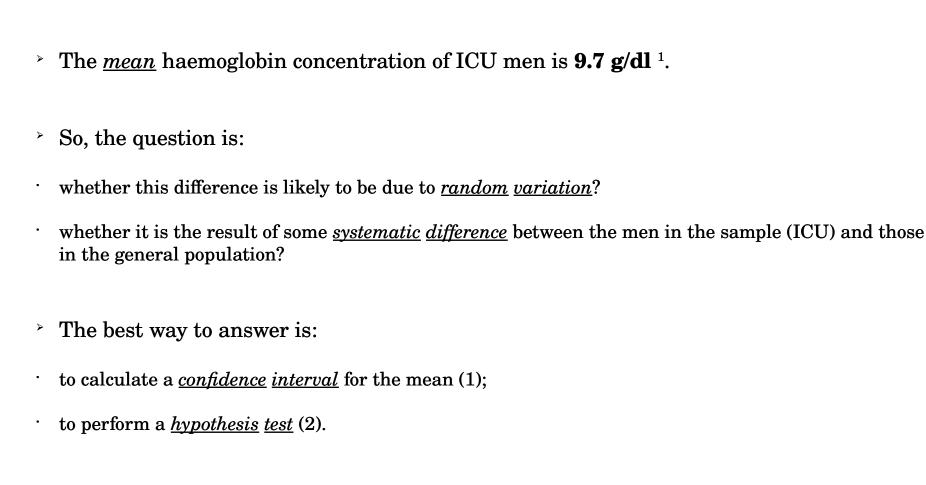
8,1	10,1	12,3
9,7	11,7	11,3
11,9	9,3	13
10,5	8,3	8,8
9,4	6,4	5,4

> Is there any *evidence* that critical illness is associated with an acute anemia?

<sup>&</sup>lt;sup>1</sup> to compare a *mean* value from a *sample* with some *hypothesized* value (e.g. with some standards).

<sup>&</sup>lt;sup>2</sup> Adapted from Whitley F, Ball J. Statistics review 5: Comparison of means. Crit. Care. 2002;6(5): 424–428.

<sup>&</sup>lt;sup>3</sup> i.e. *hypothesized* value.



<sup>1</sup> In practice any sample of 15 ICU men would be unlikely to have a mean haemoglobin of exactly 15.0 g/dl.

The **SD** of these data = 2.2 g/dl.

A 95% <u>confidence interval</u> for the mean can be calculated using the **SE** = **SD**/ $\sqrt{n}$  <sup>1</sup>:

$$SE = 2.2/\sqrt{15} = 0.56$$

The corresponding 95% *confidence interval* is as follows: <sup>2</sup>

$$9.7 \pm 2.14 \times 0.56 = 9.7 \pm 1.19 = (8.5, 10.9)^{3}$$

So, assuming that this sample is <u>representative</u>, it is likely that the true mean haemoglobin in the population of ICU adult male patients is between 8.5 and 10.9 g/dl.

 $^{1}$  SD and SE both measure <u>variability</u>: high values of both indicate more dispersion, however, SD and SE are not the same:

 ${\tt SD}$  quantifies the <u>variability</u> of data points around the <u>mean</u> in a given dataset (it tells us, on average, how far each data point is away from the mean).

 ${\tt SE}$  quantifies the  ${\tt variability}$  between samples drawn from the same population (indicates how different the sample mean is likely to be from the population mean).

 $^2$  a <u>confidence interval</u> for  $\mu$  (when the population  $\sigma$  is unknown) can be calculated as:  $\bar{\chi} \pm t \frac{s}{\sqrt{n}}$ 

 $^3$  the multiplier, in this case **2.14**, comes from the **t-distribution** because the sample size is <u>small</u>.

$$t = \frac{sample mean - hypothesised mean}{SE of sample mean^{1}}$$

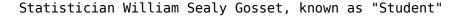
- > The associated **P** value is obtained by <u>comparison</u> with the <u>t-distribution</u> <sup>2</sup> (regardless of sign) corresponding to smaller **P** values.
- > The *t-statistic* for the haemoglobin example is as follows:

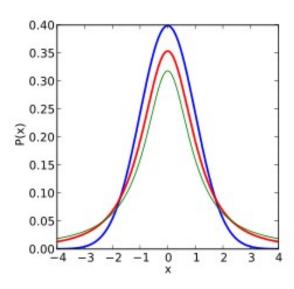
$$t = \frac{9.7 - 15}{0.56} = \frac{-5.3}{0.56} = -9.54$$

 $^{1}$  it is the number of **SE**s that separate the sample <u>mean</u> from the hypothesized <u>value</u>.

 $^2$  the **t-distribution** (Student's t-distribution) is a **probability** distribution that is used to estimate population parameters when the sample size is **small** and/or when the population **variance** is unknown.

The shape of the t distribution is determined by the  $\underline{df}$ , which, in the case of the one sample t-test, is equal to the sample size minus 1, i.e.  $\underline{14}$  (15-1).





Density of the t-distribution: 2 df – red; the standard normal distribution – blue; 1 df - green.

- > The <u>observed</u> (t) mean haemoglobin concentration is **9.54**, which is <u>below</u> the <u>hypothesized</u> mean.
- $rac{Tabulated}{Tabulated}$  values indicate how likely this is to <u>occur</u>: for a sample size of **14** (15-1) <u>df</u> the <u>P value</u> is < than 0.0001 <sup>1</sup> (see next slide).

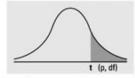
#### Conclusion:

- · it is *extremely unlikely* that the mean haemoglobin in this sample would differ from that in the general population to this extent by *chance* alone;
- this may *indicate* that there is a *truly* difference in haemoglobin concentrations in ICU men <sup>2</sup>.

the <u>P value</u> gives no indication of the <u>size</u> of any <u>difference</u>; it merely indicates the probability that the difference arose by chance. In order to assess the <u>magnitude</u> of any difference, it is essential also to have the <u>confidence interval</u> calculated above.

<sup>&</sup>lt;sup>2</sup> It is important to know how this sample of men was selected and whether they are representative of all UK men admitted to ICUs.

Numbers in each row of the table are values on a t-distribution with (df) degrees of freedom for selected right-tail (greater-than) probabilities (p).



df/p	0.40	0.25	0.10	0.05	0.025	0.01	0.005	0.0005
1	0.324920	1.000000	3.077684	6.313752	12.70620	31.82052	63.65674	636.6192
2	0.288675	0.816497	1.885618	2.919986	4.30265	6.96456	9.92484	31.5991
3	0.276671	0.764892	1.637744	2.353363	3.18245	4.54070	5.84091	12.9240
4	0.270722	0.740697	1.533206	2.131847	2.77645	3.74695	4.60409	8.6102
5	0.267181	0.726687	1.475884	2.015048	2.57058	3.36493	4.03214	0.8688
6	0.264835	0.717558	1.439756	1.943180	2.44691	3.14267	3,70,743	5.9588
7	0.263167	0.711142	1.414924	1.894579	2.36462	2.99795	3.49948	5.4079
8	0.261921	0.706387	1.396815	1.859548	2.30600	2.89646	3.35539	5.0413
9	0.260955	0.702722	1.383029	1.833113	2.26216	2.82144	3.24984	4.7809
10	0.260185	0.699812	1.372184	1.812461	2.22814	2.76377	3.16927	4 5009
11	0.259556	0.697445	1.363430	1.795885	2.20099	2.71808	3.10591	4.4370
12	0.259033	0.695483	1.356217	1.782288	2.17881	2.68100	3.05454	43178
13	0.258591	0.693829	1.350171	1.770933	2.16037	2.65031	3.01228	4.2208
14	0.258213	0.692417	1.345030	1.761310	2.14479	2.62449	2.97684	4.1405
15	0.257885	0.691197	1.340606	1.753050	2.13145	2.60248	2.94671	4.0728
16	0.257599	0.690132	1.336757	1.745884	2.11991	2.58349	2.92078	4.0150
17	0.257347	0.689195	1.333379	1.739607	2.10982	2.56693	2.89823	3.9651
18	0.257123	0.688364	1.330391	1.734064	2.10092	2.55238	2.87844	3.9216
19	0.256923	0.687621	1.327728	1.729133	2.09302	2.53948	2.86093	3.8834
20	0.256743	0.686954	1.325341	1.724718	2.08596	2.52798	2.84534	3.8495
21	0.256580	0.686352	1.323188	1.720743	2.07961	2.51765	2.83136	3.8193
22	0.256432	0.685805	1.321237	1.717144	2.07387	2.50832	2.81876	3.7921
23	0.256297	0.685306	1.319460	1.713872	2.06866	2.49987	2.80734	3.7676
24	0.256173	0.684850	1.317836	1.710882	2.06390	2.49216	2.79694	3.7454
25	0.256060	0.684430	1.316345	1.708141	2.05954	2.48511	2.78744	3.7251
26	0.255955	0.684043	1.314972	1.705618	2.05553	2.47863	2.77871	3700
27	0.255858	0.683685	1.313703	1.703288	2.05183	2.47266	2.77068	3.6896
28	0.255768	0.683353	1.312527	1.701131	2.04841	2.46714	2 70326	3.6739
29	0.255684	0.683044	1.311434	1.699127	2.04523	2.46202	2.75639	3.6594
30	0.255605	0.682756	1.310415	1.697261	2.04227	2.45726	2.75000	3.6460
z	0.253347	0.674490	1.281552	1.644854	1.95996	2.32635	2.57583	3.2905
CI	s <del></del> 8	S	80%	90%	95%	98%	99%	99.9%

2.26

2.14

1.96

#### > Conclusion:

- the haemoglobin concentration in the *general* population of adult men in the UK is well *outside* this range;
- · ICU men may *truly* have haemoglobin concentrations *lower* than the national average.
- Question: how likely it is that this <u>difference</u> is due to chance?
- Answer: we need a *hypothesis* test (one sample t-test) <sup>1</sup>.
- $\rightarrow$  H<sub>0</sub> is that <u>sample</u> mean (9.7 g/dl) is the <u>same</u> as <u>hypothesized</u> mean (15.0 g/dl).
- The *t statistic*, from which a *P value* is derived, is as follows.

<sup>&</sup>lt;sup>1</sup> t-test formally examines how *far* the estimated mean haemoglobin of ICU men (9.7 g/dl), lies from the hypothesized value (15.0 g/dl).

7.3. Comparison of two means arising	Subject	On admission	6 h after admission	Difference (%)
from paired data	1	39.7	52.9	13.2
> <u>Paired</u> data:	2	59.1	56.7	-2,4
<ul> <li>each data set has the <u>same number</u> of data points;</li> <li>each data point in one data set is related to</li> </ul>	3	56.1	61.9	5.8
<ul> <li>each <u>data</u> point in one data set is <u>related</u> to <u>one</u>, and only <u>one</u>, data point in the other data set <sup>1</sup>.</li> </ul>	4	57.7	71.4	13.7
E.g. table shows central venous $O_2$ saturation in 10 patients on admission	5	60.6	67.7	7.1
and 6 hours after admission to an ICU <sup>2</sup>	6	37.8	50	12.2
<sup>1</sup> e.g. <i>before-after</i> drug test: you record the blood pressure of each subject in the study, <i>before</i> and	7	58.2	60.7	2.5
<pre>after a drug is administered - each "before" measure is related only to the "after" measure from the same subject.</pre>	8	33.6	51.3	17.7
<sup>2</sup> Whitley F, Ball J. Statistics review 5: Comparison of means. Crit. Care. 2002;6(5): 424–428.	9	56	59.5	3.5
	10	65.3	59.8	-5.5

Mean

52.4

59.2

6.8

M. Vinichuk

•
a

Subject		
1		
2		
3		
4		
5		
6		
7		
8		
9		
10		

Mean

6 h after

admission

52.9

56.7

61.9

71.4

67.7

50

60.7

51.3

59.5

59.8

59.2

Difference

(%)

13.2

-2,4

5.8

13.7

7.1

12.2

2.5

17.7

3.5

-5.5

6.8

On

admission

39.7

59.1

56.1

57.7

60.6

37.8

58.2

33.6

56

65.3

**52.4** 

The data are <u>paired</u> 1, and it is <u>important</u>	Subject	On admission	6 h after admission	Difference (%)
to account for this pairing in the analysis.	1	39.7	52.9	13.2
<ul> <li>How to do this?</li> <li>To <u>concentrate</u> on the <u>differences between</u></li> </ul>	2	59.1	56.7	-2,4
the <i>pairs</i> of measurements rather than on the measurements themselves.	3	56.1	61.9	5.8
$H_0$ : the <u>mean</u> of the <u>differences</u> in central	4	57.7	71.4	13.7
venous oxygen saturation = <b>0</b> .  > <u>T-test</u> therefore <u>compares</u> the observed	5	60.6	67.7	7.1
$\underline{mean}$ of the differences with a hypothesized value of 0 $^2$ .	6	37.8	50	12.2
	7	58.2	60.7	2.5
<pre>¹ the two sets of observations are not independent of each other;</pre>	8	33.6	51.3	17.7
<sup>2</sup> i.e. the <i>paired t-test</i> is a special case of the single sample <i>t-test</i> described above.	9	56	59.5	3.5
	10	65.3	59.8	-5.5
M. Vinichuk	Mean	52.4	59.2	6.8

> The <u>t-statistic</u> :	Subject	On admission	6 h after admission	Difference (%)
$t = \frac{\text{sample mean of differences} - 0}{\text{SE of sample mean of differences}}$	1	39.7	52.9	13.2
$t = \frac{\text{sample mean of differences}}{\text{SE of sample mean of differences}}$	2	59.1	56.7	-2,4
, , , ,,	3	56.1	61.9	5.8
The <b>SD</b> of the differences is <b>7.5</b> , and this corresponds to a <b>SE</b> of $7.5/\sqrt{10} =$ <b>2.4</b> .	4	57.7	71.4	13.7
The t-statistic: $t = 6.8/2.4 = 2.87$ , and	5	60.6	67.7	7.1
this corresponds to a $\underline{P \ value}$ of $0.01^{-1}$ .	6	37.8	50	12.2
	7	58.2	60.7	2.5
based on a t-distribution with 10-1=9 df.	8	33.6	51.3	17.7
i. e. there is some evidence to suggest that admission to ICU and subsequent treatment may <u>increase</u> central venous oxygen saturation <u>beyond</u> the level expected by chance.	9	56	59.5	3.5
	10	65.3	59.8	-5.5

Mean

52.4

6.8

59.2

M. Vinichuk

> Remember: <i>P value</i> gives no info. about	Subject	On admission	6 h after admission	Difference (%)
the likely <u>size</u> of any <u>effect</u> .	1	39.7	52.9	13.2
This may be done by calculating a 95% <u>confidence interval</u> from the <u>mean</u> and <b>SE</b> of the differences:	2	59.1	56.7	-2,4
	3	56.1	61.9	5.8
$6.8 \pm 2.26 \times 2.4 = 6.8 \pm 5.34 = (1.4, 12.2)^{1}$	4	57.7	71.4	13.7
So, the <u>true increase</u> in central venous O <sub>2</sub> saturation is probably <u>between</u> 1.4% and 12.2%.	5	60.6	67.7	7.1
> Although the <i>increase</i> may be small	6	37.8	50	12.2
(1.4%), it is <b>unlikely</b> that the effect is to decrease saturation.	7	58.2	60.7	2.5
	8	33.6	51.3	17.7
<sup>1</sup> the multiplier, in this case <b>2.26</b> , comes from the t-distribution because the sample size is <i>small</i> . If n>30, z-tables are better.	9	56	59.5	3.5
	10	65.3	59.8	-5.5
M. Vinichuk	Mean	52.4	59.2	6.8

## 7.4. Comparison of two means arising from unpaired data

> Comparison of two means arising from <u>unpaired</u> <u>data</u> <sup>1</sup> is most <u>common</u>.

E.g. lets compare <u>early goal</u>-<u>directed</u> <u>therapy</u> (**EGDT**) with <u>standard</u> <u>therapy</u> (**ST**) in the

treatment of septic shock <sup>2</sup>.

A total of 263 patients were randomized and 236 *completed* 6 hours of treatment.

The <u>mean</u> arterial pressures after 6 hours of treatment in the **ST** and **EGDT** groups are shown in table <sup>3</sup>.

	Mean arterial pressure (mmHg)			
	ST EGDT			
Number of patients	119	117		
Mean	81	95		
St. Dev.	18	19		

<sup>&</sup>lt;sup>1</sup> i.e. comparison of data from two **independent** groups;

<sup>&</sup>lt;sup>2</sup> is a potentially fatal *medical condition* when organ injury or damage in response to infection leads to dangerously low blood pressure and abnormalities in cellular metabolism;

<sup>&</sup>lt;sup>3</sup> Adapted from E. Whitley, and J. Ball. **Statistics review 5: Comparison of means.** Rivers E,et al. Early goal-directed therapy in the treatment of severe sepsis and septic shock. N Engl J Med 2001;345:1368–77.

The <u>mean</u> arterial pressure was 14 mmHg <u>higher</u> in the EGDT.

> The 95% *confidence intervals* for the mean arterial pressure in the two groups:

$$ST: 81 \pm 1.96 \times \frac{18}{\sqrt{119}} = 81 \pm 3.23 = [77.8;84.2]$$

EGDT: 95 ± 1.96 × 
$$\frac{19}{\sqrt{117}}$$
 = 95 ± 3.4 = [91.6;98.4]

There is no <u>overlap</u> between the two <u>confidence intervals</u> and there **may be** a <u>difference</u> between the two groups.

	Mean arterial pressure (mmHg)			
	ST EGDT			
Number of patients	119	117		
Mean	81	95		
St. Dev.	18	19		

- Let us *estimate* the *size* of the difference.
- The only <u>difference</u> is in the <u>calculation</u> of the **SE**.
- In the *paired* case attention is focused on the *mean* of the *differences*.
- In the <u>unpaired</u> case interest is in the <u>difference</u> of the <u>means</u> 1.

	Mean arterial pressure (mmHg)		
	ST	EGDT	
Number of patients	119	117	
Mean	81	95	
St. Dev.	18	19	

<sup>&</sup>lt;sup>1</sup> Because the <u>sample sizes</u> in the <u>unpaired</u> case may be (and indeed usually are) <u>different</u>, the combined **SE** takes this into account and gives more weight to the <u>larger sample size</u> because this is likely to be more <u>reliable</u>.

The *pooled* **SD** for the *difference* in means:

$$SD_{difference}$$
: =  $\frac{\sqrt{(n_1 - 1) \times SD_1^2 + (n_2 - 1) \times SD_2^2}}{(n_1 + n_2 - 2)}$ 

where  $SD_1$  and  $SD_2$  are the SDs in the two groups and  $n_1$  and  $n_2$  are the two sample sizes.

The *pooled* **SE** for the *difference* in means:

$$SE_{difference} = SD_{difference} \times \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

This **SE** can now be used to calculate a <u>confidence interval</u> for the difference in means and to perform an <u>unpaired</u> t-test, as above.

	Mean arterial pressure (mmHg)				
	ST EGDT				
Number of patients	119	117			
Mean	81	95			
St. Dev.	18	19			

# The pooled **SD** is:

$$SD_{difference}$$
: =  $\frac{\sqrt{(119-1) \times 18^2 + (117-1) \times 19^2}}{(119 + 117 - 2)}$ 

$$SD_{difference}$$
: =  $\sqrt{\frac{38.232 + 41.876}{234}} = \sqrt{342.34} = 18.50$ 

#### and the corresponding pooled **SE** is:

$$SD_{difference} = 18.50 \times \sqrt{\frac{1}{119} + \frac{1}{117}}$$
  
=  $18.50 \times \sqrt{0.008 + 0.009}$   
=  $18.50 \times 0.13 = 2.41$ 

	Mean arterial pressure (mmHg)		
	ST	EGDT	
Number of patients	119	117	
Mean	81	95	
St. Dev.	18	19	

The <u>difference</u> in <u>mean</u> arterial pressure between **EGDT** and **ST** groups is **14 mmHg**, with a corresponding 95% <u>confidence interval</u> of  $14 \pm 1.96 \times 2.41 = (9.3, 18.7)$  mmHg.

The <u>confidence interval</u> suggests that the <u>true</u> <u>difference</u> is likely to be between 9.3 and 18.7 mmHg.

To explore the likely role of chance in explaining this difference, an unpaired <u>t-test</u> can be performed.

	Mean arterial pressure (mmHg)		
	ST	EGDT	
Number of patients	119	117	
Mean	81	95	
St. Dev.	18	19	

- $\rightarrow$  H<sub>0</sub> is that the difference in means = 0.
- As for the previous two cases, a *t-statistic* is calculated.

$$t = \frac{\text{difference of sample means}}{\text{SE of difference of sample means}}$$

- > A <u>P value</u> may be obtained by comparison with the t-distribution on  $n_1 + n_2 2$  df <sup>1</sup>.
- t = 14/2.41 = 5.81, with a corresponding *P value* < 0.0001 <sup>2</sup>.
- Conclusion: there may be a <u>truly</u><u>difference</u> between the two groups.

 $^{\ 1}$  the larger the t-statistic, the smaller the  $\underline{\textit{P}\ \textit{value}}$  will be;

	Mean arterial pressure (mmHg)		
	ST	EGDT	
Number of patients	119	117	
Mean	81	95	
St. Dev.	18	19	

<sup>&</sup>lt;sup>2</sup> **14** mmHg is the *difference* in *mean* arterial pressure; it is extremely *unlikely* that a difference in mean arterial pressure of this magnitude would be observed just by chance.

## 7.5. Comparison of two means: assumptions and limitations

- The t-tests require certain <u>assumptions</u>.
- The <u>one sample</u> t-test: the data have an <u>approximately</u> Normal distribution.
- The *paired* t-test: the distribution of the **differences** are *approximately* Normal.
- > The *unpaired* t-test: the data from the two samples are both *Normally* distributed <sup>1</sup>.
- There are tests to <u>examine</u> whether a set of data are Normal or whether two **SD**s (or, two variances) are <u>equal</u> <sup>2</sup>.

<sup>&</sup>lt;sup>1</sup> and the **SDs** from the two samples are approximately **equal**;

<sup>&</sup>lt;sup>2</sup> the Kolmogorov-Smirnov test (K-S) and Shapiro-Wilk (S-W) test are designed to test normality by comparing your data to a normal distribution.

## What to do if Normality is violated?

- The appropriate <u>transformation</u> <sup>1</sup> of the data may be used before performing any calculations ("Let the data decide").
- For "big" sample size (30, or even better 50), an impact to validity from non-normal data usually "small".
- Fig. 1 If not possible, alternative tests can be used, e.g. nonparametric tests 2.
- To explore differences in <u>means</u> across <u>three</u> or <u>more</u> groups, an analysis of <u>variance</u> (**ANOVA**) can be used.

Can be performed well with <u>skewed</u> and <u>nonnormal</u> distributions. E.g. <u>Mann-Whitney</u> test instead of <u>2-sample t test</u>; or <u>Kruskal-Wallis</u>, <u>Mood's median test</u> instead of <u>One-Way ANOVA</u>, etc.

<sup>&</sup>lt;sup>1</sup> e.g. to apply a square root to each value or data may be log-transformed: your data may now be normal, but interpreting that data may be much more difficult;

 $<sup>^2</sup>$  it is like a <u>parallel</u> universe to parametric tests, also called <u>distribution</u>-<u>free</u> tests: don't assume that your data follow a specific <u>distribution</u>.