

Introduction to probability, statistics and data handling

Tomasz Szumlak, Agnieszka Obłakowska-Mucha
Faculty of Physics and Applied Computer Science

AGH UST Krakow



2

Type I and Type II error

- ❑ A Type I error means rejecting the null hypothesis when it's actually true. It means concluding that results are **statistically significant** when, in reality, they came about purely by chance or because of unrelated factors. This kind of error is called **Type I**.
- ❑ If we accept a hypothesis which is not true, we say we made **Type II** error.

You decide to get tested for pregnancy. There are two errors that could potentially occur:

Type I error (false positive): the test result says you are pregnant, but you actually don't.

Type II error (false negative): the test result says you are not pregnant, but you actually are.

Test Results	Reality	
	Pregnant (=1)	Not Pregnant (=0)
Positive (=1)	Number of True Positives (TP)	Number of False Positives (FP)
Negative (=0)	Number of False Negatives (FN)	Number of True Negatives (TN)

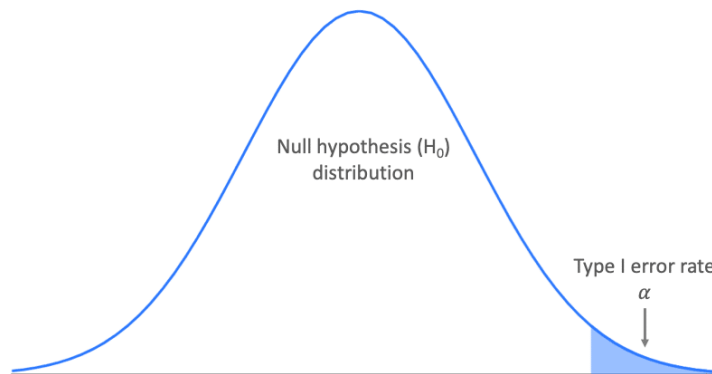


3

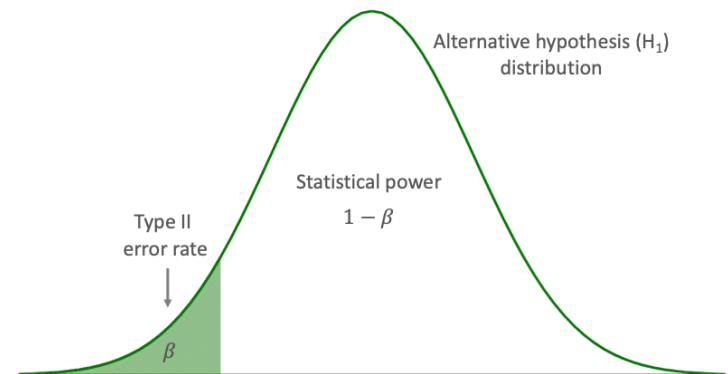
Type I and Type II error

- ❑ Imagine that we rejected a hypothesis and it happens to be true. This kind of error is called **Type I**.
- ❑ If we accepted a hypothesis which is not true, we say we made **Type II** error.

Probability of making a Type I error



Probability of making a Type II error





What can happen?

- ❑ No matter what we do, in real life we are never going to know for certain that we made the right choice or we in fact made one of the errors
- ❑ Remember Type I is to **REJECT** the null when it is true and Type II is **NOT TO REJECT** the null when it is false in reality
- ❑ Let's start from a generic (and easy) example, then we try to make a number of general statements and we get back to some more examples...

- ❑ Say, we are testing (for normal distribution) the following:
$$H_0: \mu = 14 \quad H_1: \mu > 14$$
- ❑ After applying our „procedure“ we decided to reject the hypo at significance level $\alpha = 0.05$
- ❑ Either of the two happened: **the null is wrong** and we did good or the **null is true** and we made a Type I error



What can happen?

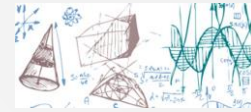
- ❑ After applying our „procedure” we decided to accept the hypo at significance level $\alpha = 0.05$
- ❑ Either of the two happened: **the null is ok** and we make a good decision or the **null is false** and we made a Type II error
- ❑ Is it a bit confusing...? Yes, and we can summarise this using the CONFUSION MATRIX – now you know why we use this name...

Ground truth

	H_0 is false	H_0 is true
Reject	Super	Type I error
Do not reject	Type II error	Super

Conclusions
using data

Note, using data we in fact make calculations, the ground truth remains, however, unknown...



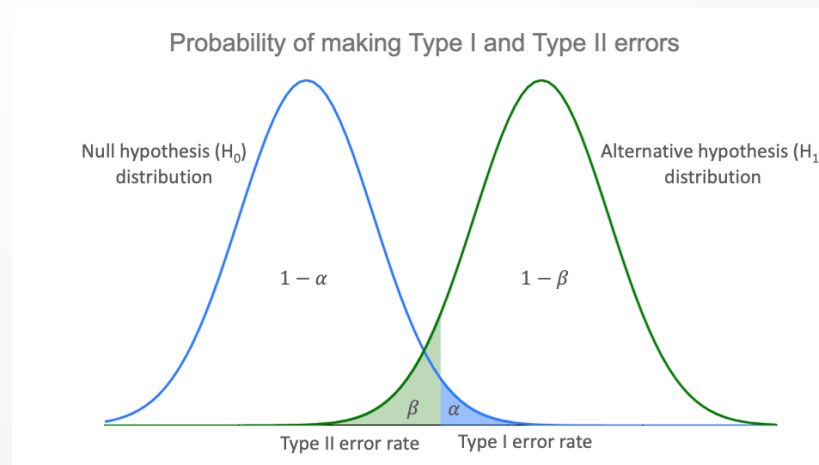
A trial example

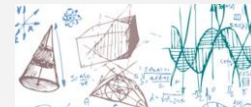
- ❑ This is a classic example: criminal trial

H_0 : not guilty (no crime),

H_1 : guilty (criminal)

- ❑ In a typical jurisdiction system we have a rule „not guilty until proven otherwise”, so we need to have very strong evidences to convict
- ❑ Type I: rule guilty when in reality a person is innocent
- ❑ Type II: not guilty when a person committed a crime
- ❑ Definitely, we want to make sure, that we never convict an innocent person (well, in real life we do...)
- ❑ What we can do – we should decrease the possibility of making Type I error as much as possible





Significance and power

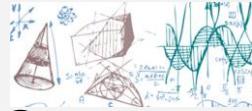
- ❑ The probability of **Type I error is called the significance level** of the statistical test (we know that already), we write this as

$$P(\text{Type I error} | H_0) = \alpha$$

- ❑ Usually we start the test procedure by setting the value of α
- ❑ Now, we call the probability of making a Type II error as β
- ❑ There is a complication – it is much harder to define this error, it depend on the true value of the testing parameter, the size of the data sample and on the significance level itself!
- ❑ The **POWER** of a statistical test is **the probability of rejecting the null hypothesis when it is false**, we write

$$\mathcal{R} = 1 - P(\text{Type II error}) = 1 - \beta$$

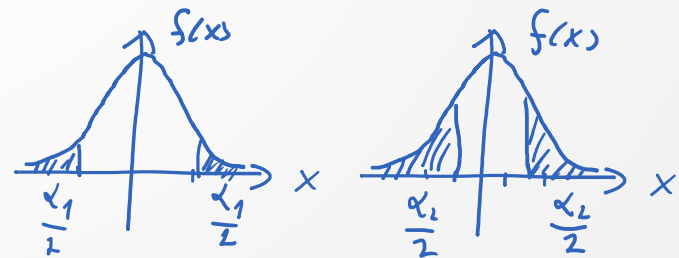
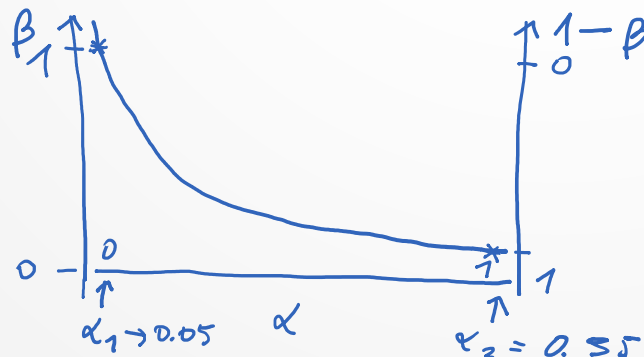
- ❑ The power depends on β – and we still do not know how to calculate it...



8

Intricate interplay of α and β

- There is some intuitive understanding about β .
- So, why not pick-up some very, very tiny value for the significance? The critical region would be very small, thus we would make almost impossible to reject the null hypo. We also say, **that the efficiency of the selection is very high**, but the probability of making the Type II error will also be very high (or we say, **the purity of the selection is low**)
- Conversely, for the large values of α the probability of Type II errors will go down
- Also, from this discussion you should start to understand, that the dependency is not trivial (i.e., it is not just some linear function or such)





9

Examples!

- Let's start calculating something!
- Ex. 1** Say, we consider a sampling from a normal distribution with known variance $\sigma = 18$, the size of the data sample $n = 36$. Our hypothesis

$$H_0: \mu = 40, H_1: \mu < 40$$

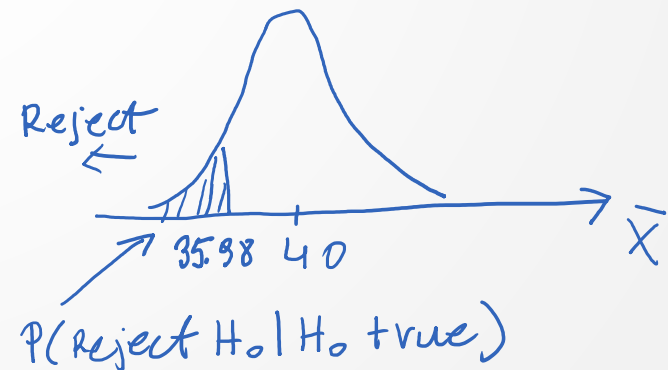
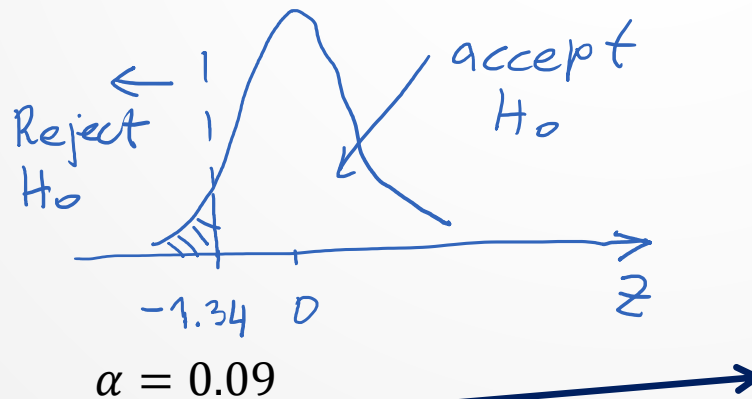
- We set the significance of the test: $\alpha = 0.09$ and our test statistics is

$$Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}}$$

$$\bar{X} = \mu_0 + Z \frac{\sigma}{\sqrt{n}} = 40 + \frac{18}{\sqrt{36}} (-1.34)$$

- The critical region

$$\bar{X} = 35.98$$





10

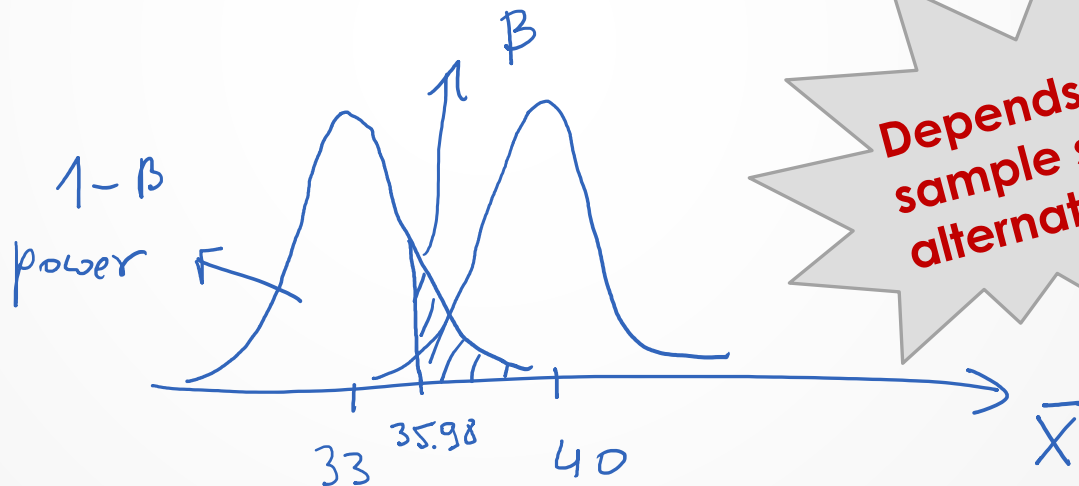
Examples!

- Now, what if the true mean is $\mu = 33$? What is the Type II error in this case?

$$P(\text{Accept } H_0 | \mu = 33)$$

- We reject the H_0 hypo if the mean is less than 35.98, so the Type II error will occur with the probability

$$P(\bar{X} > 35.98 | H_1) = P(\bar{X} > 35.98 | \mu = 33)$$



Depends on the sample size and alternative!

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} = \frac{35.98 - 33}{18/6} = \frac{2.98}{3} = 0.99 \rightarrow P(Z > 0.99 | H_1) = 0.16$$

β



Sample variance

- If $\{X_1, X_2, \dots, X_n\}$ denote R.Vs for a random sample of size n , the R.V. giving the variance of the sample (the sample variance) is defined as:

$$S^2 = \frac{1}{n} [(X_1 - \bar{X})^2 + (X_2 - \bar{X})^2 + \dots + (X_n - \bar{X})^2]$$

- We already know, that $E[\bar{X}] = \mu$, is this the same for $E[S^2] = \sigma^2$?
 - A little digression – whenever the expected value of a statistics is **equal** to the corresponding **population parameter**, we call this statistics **an unbiased estimator**. Its value is then an unbiased estimate of the respective parameter
- Unfortunately, it can be proved that for the sample variance, we have some bias:

$$E[S^2] = \mu_{S^2} = \frac{n-1}{n} \sigma^2$$

- However, an unbiased variance estimator $\hat{\sigma}^2$, based on data, is easy to find:

$$\hat{S}^2 \equiv \hat{\sigma}^2 = \frac{n}{n-1} S^2 = \frac{1}{n-1} [(X_1 - \bar{X})^2 + (X_2 - \bar{X})^2 + \dots + (X_n - \bar{X})^2]$$

for large samples the difference between S^2 and \hat{S}^2 is small.

S^2 is unbiased for normal variance,



Sampling dist. of variances (lab)

- With such unbiased estimator, we have:

$$E[\hat{S}^2] = \sigma^2$$

- In order to create the sampling distribution of variances we do exactly the same as for the SDoM, we take all the possible samples of size n , that can be drawn from a population and calculate their variances
- One change is, that instead of looking directly at the distribution of the sample variance, we look at the R.V.:

$$\chi^2 \equiv \frac{nS^2}{\sigma^2} = \frac{(n-1)\hat{S}^2}{\sigma^2} = \frac{(X_1 - \bar{X})^2 + (X_2 - \bar{X})^2 + \dots + (X_n - \bar{X})^2}{\sigma^2}$$

new RV: $\chi^2 = (\text{coefficient}) \times \text{sum of random variables}$

what is the μ and σ of this new RV?

- **Theorem 6.** If a random samples of size n are taken from a population having a normal distribution, than the sampling variable $\frac{nS^2}{\sigma^2}$ has a χ^2 **distribution** with $n - 1$ degrees of freedom

χ^2 for fitting



- ❑ χ^2 distribution should be always associated with a RV which describes the dispersion of the square of the deviations of a an RV around a fixed point.
- ❑ Now: what if the point is the „TRUE” expected value of X , i.e. μ_X ?

It means that the variable:

$$\chi^2 = \sum_{i=1}^n \frac{(X_i - E\{X\})^2}{\sigma^2} = \sum_{i=1}^n \frac{(X_i - \mu_X)^2}{\sigma^2}$$

has a χ^2 distribution with $\nu = n$ (!) degrees of freedom

- ❑ One can use this to determine wheather the data fit a particular distribution (goodness-of-fit test).

the test distribution is χ^2 distribution with a test statistic for goodness-of-fit test:

$$\chi^2 = \sum_k \frac{(O_k - E_k)^2}{E_k}$$

E_k - TRUE value (or expected), for instance from theory or expected based on other very precise measurements



Differences of means

- **Ex. 4** Two classes attended the same course and took an exam. The mean mark for the first class (40 students) was 74 points with a standard deviation of 8, the second (50 students) scored the mean of 78 points with the standard deviation of 7. Can we claim that one of the class is significantly better than other? We set the $\alpha = 0.05$
- $H_0: \mu_1 = \mu_2$ - no significant difference, just fluctuation
- $H_1: \mu_1 \neq \mu_2$ - the second class is better



Differences of means

- Ex. 4** Two classes attended the same course and took an exam. The mean mark for the first class (40 students) was 74 points with a standard deviation of 8, the second (50 students) scored the mean of 78 points with the standard deviation of 7. Can we claim that one of the class is significantly better than other? We set the $\alpha = 0.05$

- $H_0: \mu_1 = \mu_2$ - no significant difference, just fluctuation
- $H_1: \mu_1 \neq \mu_2$ - the second class is better
- The test statistics is:



Differences of means

- **Ex. 4** Two classes attended the same course and took an exam. The mean mark for the first class (40 students) was 74 points with a standard deviation of 8, the second (50 students) scored the mean of 78 points with the standard deviation of 7. Can we claim that one of the class is significantly better than other? We set the $\alpha = 0.05$
- $H_0: \mu_1 = \mu_2$ - no significant difference, just fluctuation
- $H_1: \mu_1 \neq \mu_2$ - the second class is better
- The test statistics is: the difference of means

$$\bar{X}_1 \rightarrow N\left(\mu, \frac{\sigma_1}{\sqrt{n_1}}\right) \quad \bar{X}_2 \rightarrow N\left(\mu, \frac{\sigma_2}{\sqrt{n_2}}\right)$$

$$\bar{X}_1 - \bar{X}_2 \rightarrow N(0, \sigma(\bar{X}_1 - \bar{X}_2)) \quad \sigma(\bar{X}_1 - \bar{X}_2) = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_{\bar{X}_1} - \mu_{\bar{X}_2})}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

$$Z_{st} = \frac{RV - \mu}{\sigma}$$



Differences of means

- **Ex. 4** Two classes attended the same course and took an exam. The mean mark for the first class (40 students) was 74 points with a standard deviation of 8, the second (50 students) scored the mean of 78 points with the standard deviation of 7. Can we claim that one of the class is significantly better than other? We set the $\alpha = 0.05$
- $H_0: \mu_1 = \mu_2$ - no significant difference, just fluctuation
- $H_1: \mu_1 \neq \mu_2$ - the second class is better
- The test statistics is the difference of means

$$\mu_{\bar{X}_1 - \bar{X}_2} = 0, \sigma_{\bar{X}_1 - \bar{X}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} = \sqrt{\frac{64}{40} + \frac{49}{50}} = 1.606$$

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_{\bar{X}_1} - \mu_{\bar{X}_2})}{\sigma_{\bar{X}_1 - \bar{X}_2}} = \frac{74 - 78}{1.606} = -2.49$$

$$z_{1,2} = \pm 1.96 \rightarrow \text{rejected } H_0 \text{ @ } 5\%$$

- p-value: $P(Z \leq -2.49) + P(Z \geq 2.49) = 0.0128$

we can use z-scores
for normal distribution
(big sample)

Examples – C.I.



- **Ex. 6** A sample of 150 brand A light bulbs showed a mean life-time of 1400 hours and a standard dev. of 120 h. A sample of 200 brand B light bulbs showed the corresponding values of 1200 h and 80 h. Find (a) 95% and (b) 99% *C.I.* for the difference of the mean life-times of the populations of brands A and B.

- a) The confidence interval for the difference in means can be written as:

$$\bar{L}_A - \bar{L}_B \pm z_c \sqrt{\frac{\sigma_A^2}{n_A} + \frac{\sigma_B^2}{n_B}} = 1400 - 1200 \pm 1.96 \sqrt{\frac{120^2}{150} + \frac{80^2}{100}} = 200 \pm 24.8 h$$

We can be 95% confident that the difference in population means lies between 175 and 225 h

- b) And for the 99% *C.I.*:



Differences of means

- What would be a change if σ is unknown?
- The test statistics is the same: the difference of means

$$\bar{X}_1 - \bar{X}_2 \rightarrow N(0, \sigma(\bar{X}_1 - \bar{X}_2))$$

- but now $\sigma(\bar{X}_1 - \bar{X}_2)$ is estimated:

$$\sigma(\bar{X}_1 - \bar{X}_2) \rightarrow \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

$$s_1^2 = \frac{1}{n_1} \sum_{i=1}^{n_1} (X_i - \bar{X})^2 \frac{1}{n_1}$$

- and test statistic is t-score:

$$t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_{\bar{X}_1} - \mu_{\bar{X}_2})}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

$$v = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\frac{1}{n_1 - 1} \frac{s_1^2}{n_1} + \frac{1}{n_2 - 1} \frac{s_2^2}{n_2}}$$

- and we need to use t -Student distribution with ν degrees of freedom



Differences of means

- ❑ **How we can compare two independent Population Proportion?**
- ❑ There are a few assumptions that must be fulfilled:
 - samples should be random and independent
 - number of successes > 5

This enables using the normal dist, approximation.

The pooled proportion RV:

$$P = \frac{x_A + x_B}{n_A + n_B}$$

the difference:

$$\Delta_P = P'_A - P'_B \quad \Delta_P \rightarrow N(0, \sigma_\Delta)$$

$$\sigma_\Delta = \sqrt{P(1 - P) \left(\frac{1}{n_A} + \frac{1}{n_B} \right)}$$

and test statistics (z-score):

$$Z = \frac{(P'_A - P'_B) - (P_A - P_B)}{\sqrt{P(1 - P) \left(\frac{1}{n_A} + \frac{1}{n_B} \right)}}$$

$$H_0: P_A = P_B, \text{ so } \Delta_P = 0$$

$$H_a: P_A \neq P_B, \text{ so } \Delta_P \neq 0$$



Differences of means: Population Proportion

- Example: Two types of valves are being tested to determine if there is a difference in pressure tolerances. Fifteen out of a random sample of 100 of Valve A cracked under 4,500 psi. Six out of a random sample of 100 of Valve B cracked under 4,500 psi. Test at a 5% level of significance.

$$P = \frac{x_A + x_B}{n_A + n_B} = \frac{15 + 6}{100 + 100} = 0.105 \quad P'_A = \frac{15}{100} \quad P'_B = \frac{6}{100}$$

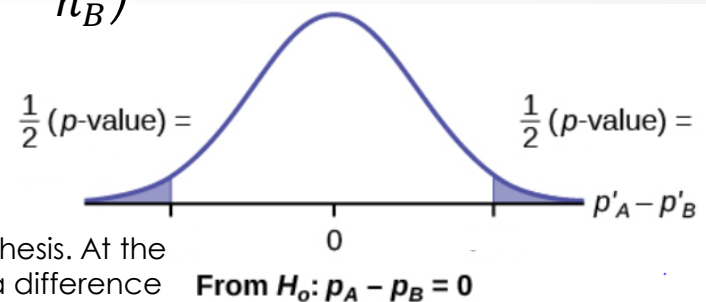
the difference: $\Delta_P = P'_A - P'_B = 0.21 \quad \Delta_P \rightarrow N(0, \sigma_\Delta)$

and test statistics (z-score):

$$Z = \frac{(P'_A - P'_B) - (P_A - P_B)}{\sqrt{P(1 - P) \left(\frac{1}{n_A} + \frac{1}{n_B} \right)}} \quad \sigma_\Delta = \sqrt{P(1 - P) \left(\frac{1}{n_A} + \frac{1}{n_B} \right)}$$

$H_0: P_A = P_B, \text{ so } \Delta_P = 0$

$H_a: P_A \neq P_B, \text{ so } \Delta_P \neq 0$

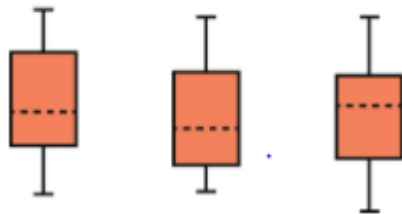


The p-value is 0.0379, so we can reject the null hypothesis. At the 5% significance level, the data support that there is a difference in the pressure tolerances between the two valves.

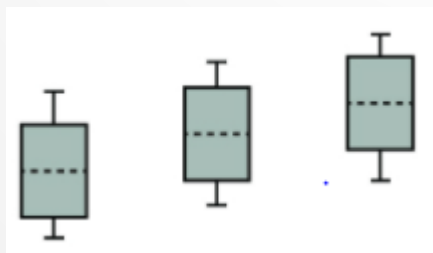


C.I. for **variance**

- ❑ In case of comparing the mean values of two populations we introduced how the sampling distribution of differences (and sums...) of means can be obtained. And what about the distribution of **differences of variances** that can be denote as: $S_1^2 - S_2^2$?
- ❑ For hypothesis tests comparing averages between more than two groups: "Analysis of Variance" (ANOVA).
- ❑ The null hypothesis is simply that all the group population means are the same. The alternative hypothesis is that at least one pair of means is different.



$H_0: \mu_1 = \mu_2 = \mu_3$ and the three populations have the same distribution if the null hypothesis is true. Differences are due to random variation.



If the null hypothesis is false, then the variance of the combined data is larger which is caused by the different means. Differences are too large to be due to random variation.



C.I. for **variance ratios**

- It turns out, that such distribution is **surprisingly** troublesome! Instead we can use **another statistic**: $F \propto S_1^2/S_2^2$. We conclude that if the ratio is small/large the difference between variance is large, conversely if the ratio is close to unity the difference should be small, and we have a theorem...
- **Theorem 1.** Consider, we have **two random and independent** samples, of size n and m respectively, drawn from two normal populations with variances σ_1^2 and σ_2^2 . It can be shown that the statistic

$$F = \frac{\frac{mS_1^2}{(m-1)\sigma_1^2}}{\frac{nS_2^2}{(n-1)\sigma_2^2}} = \frac{\frac{\hat{S}_1^2}{\sigma_1^2}}{\frac{\hat{S}_2^2}{\sigma_2^2}}$$

has the F distribution with $m - 1, n - 1$ degrees of freedom



C.I. for variance ratios

- Let's repeat our scheme again. We denote by $F_{0.01}$ and $F_{0.99}$ the values of F R.V. for which 1% of the area lies in each tail of the F distribution, then with 98% confidence we have the following

$$F_{0.01} \leq \frac{\frac{\hat{S}_1^2}{\sigma_1^2}}{\frac{\hat{S}_2^2}{\sigma_2^2}} \leq F_{0.99}$$

- So, a 98% C.I. for the variance ratio σ_1^2/σ_2^2 of two populations is given by:

$$\frac{1}{F_{0.99}} \frac{\hat{S}_1^2}{\hat{S}_2^2} \leq \frac{\sigma_1^2}{\sigma_2^2} \leq \frac{1}{F_{0.01}} \frac{\hat{S}_1^2}{\hat{S}_2^2}$$

- The respective critical values are tabulated. The way we read the tables are a bit difficult but we do that during the tutorial, so don't worry!



Examples – C.I.

- **Ex. 8** Two samples of sizes 16 and 10 are drawn from two normal populations. If their variances are found to be 24 and 18, respectively, find (a) 98% and (b) 90% confidence limits for the ratio of the variances
 - First we can find the modified sample variances:

$$\hat{s}_1^2 = \frac{m}{m-1} s_1^2 = \frac{16}{15} \cdot 24 = 25.2$$

$$\hat{s}_2^2 = \frac{n}{n-1} s_2^2 = \frac{10}{9} \cdot 18 = 20$$

- We need to use F distribution. For the 98% confidence interval we have: $F_{0.99} = 4.96$ for $v_1 = 16 - 1 = 15$ and $v_2 = 10 - 1 = 9$. And $F_{0.01} = 1/3.89$ for the same values of degrees of freedom. So, the c.i.:

$$C.I._{0.98}^{s^2} = \frac{1}{4.96} \frac{25.2}{20} < \frac{\sigma_1^2}{\sigma_2^2} < 3.89 \frac{25.2}{20} \quad 0.28 < \frac{\sigma_1^2}{\sigma_2^2} < 4.90$$

Examples – C.I.

