Statistics and learning Tests

Emmanuel Rachelson and Matthieu Vignes

ISAE SupAero

Thursday 24th January 2013

イロト イポト イヨト イヨト

∃ <\0<</p>

1 / 14

2013

WHen could tests be useful ?

► A statistical hypothesis is an assumption on the distribution of a random variable.

500

WHen could tests be useful ?

- ► A statistical hypothesis is an assumption on the distribution of a random variable.
- ► Ex: test whether the average temperature in a holiday ressort is 28°C in the summer.

WHen could tests be useful ?

- ► A statistical hypothesis is an assumption on the distribution of a random variable.
- ► Ex: test whether the average temperature in a holiday ressort is 28°C in the summer.
- A test is a procedure which makes the use of a sample to decide whether we can reject an hypothesis or whether there is nothing wrong with it (it's not really acceptance).

WHen could tests be useful ?

- ► A statistical hypothesis is an assumption on the distribution of a random variable.
- ► Ex: test whether the average temperature in a holiday ressort is 28°C in the summer.
- A test is a procedure which makes the use of a sample to decide whether we can reject an hypothesis or whether there is nothing wrong with it (it's not really acceptance).
- Examples of applications: decide if a new drug can be put on market after adequate clinical trials, decide if items comply with predefined standards, which genes are significantly differentially expressed in pathological cells

・ロト ・ 同ト ・ ヨト ・ ヨト

WHen could tests be useful ?

- ► A statistical hypothesis is an assumption on the distribution of a random variable.
- ► Ex: test whether the average temperature in a holiday ressort is 28°C in the summer.
- ► A test is a procedure which makes the use of a sample to decide whether we can reject an hypothesis or whether there is nothing wrong with it (it's not really acceptance).
- Examples of applications: decide if a new drug can be put on market after adequate clinical trials, decide if items comply with predefined standards, which genes are significantly differentially expressed in pathological cells
- Typically, sources to build hypothesis stem from quality need, values from a previous experiment, a theory that need experimental confirmation or an assumption based on observations.

It's really about **decision making**; don't be fooled, tests shed light on a question, final results heavily depend on a human interpretation !

It's really about **decision making**; don't be fooled, tests shed light on a question, final results heavily depend on a human interpretation !

Today's goals:

▶ introduce basic concepts related to tests through 2 examples

It's really about **decision making**; don't be fooled, tests shed light on a question, final results heavily depend on a human interpretation !

Today's goals:

- ▶ introduce basic concepts related to tests through 2 examples
- ► a general presentation of tests

It's really about **decision making**; don't be fooled, tests shed light on a question, final results heavily depend on a human interpretation !

Today's goals:

- ▶ introduce basic concepts related to tests through 2 examples
- ► a general presentation of tests
- ► some particular cases: one-sample, two-sample, paired tests; Z-tests, t-tests, χ²-tests, F-tests . . .

It's really about **decision making**; don't be fooled, tests shed light on a question, final results heavily depend on a human interpretation !

Today's goals:

- ▶ introduce basic concepts related to tests through 2 examples
- ► a general presentation of tests
- ► some particular cases: one-sample, two-sample, paired tests; Z-tests, t-tests, χ²-tests, F-tests ...

Example 1: cheater detection

To introduce randomness, you are asked to throw a coin 200 times and write down the results. Why would I be suspicious about students that do not exhibit at least one HHHHHH or TTTTTT pattern ? Would I be (totally ?) fair if I was to blame (all of) them ?

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Example 2: rain makers

In a given area of agricultural interest, it usually rains 600mm a year. Suspicious scientists claim that they can locally increase rainfall, when spreading a revolutionary chemical (iodised silver) on clouds. Tests over the 1995-2002 period gave te following results:

Year	1995	1996	1997	1998	1999	2000	2001	2002
Rainfall (mm/year)	606	592	639	598	614	607	616	586

Does this sound correct to you ? Quantify the answer.

Bonus: what would have changed if you wanted to test if the increase was of say 30 mm?

Rain makers et possible errors



(H0) $\theta = \theta_0$ and (H1) $\theta = \theta_1$

Possible situations



990

Possible situations

Real world	(H0)	(H1)	
Decision made			
(H0)	$1 - \alpha$	β	
(H1)	α	$1-\beta$	

Possible situations

Real world	(H0)	(H1)	
Decision made			
(H0)	$1 - \alpha$	β	
(H1)	α	$1-\beta$	

Apply that to "innoncent until proven guilty" and interpret the different situations. How do you want to control α and β ? What about introducing a new drug on the market ??

< □ > < □ > < 豆 > < 豆 > < 豆 > < 豆 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

General methodology

1. Modelling of the problem.

General methodology

- 1. Modelling of the problem.
- 2. Determine alternative hypotheses to test (disjoint but not necessarily exhaustive).

3

590

General methodology

- 1. Modelling of the problem.
- 2. Determine alternative hypotheses to test (disjoint but not necessarily exhaustive).
- 3. Choose of a statistic than (a) can be computed from data and (b) which has a known distribution under (H0).

General methodology

- 1. Modelling of the problem.
- 2. Determine alternative hypotheses to test (disjoint but not necessarily exhaustive).
- 3. Choose of a statistic than (a) can be computed from data and (b) which has a known distribution under (H0).
- 4. Determine the behaviour of statistics under (H1) and buid critical region (where (H0) rejected)

General methodology

- 1. Modelling of the problem.
- 2. Determine alternative hypotheses to test (disjoint but not necessarily exhaustive).
- 3. Choose of a statistic than (a) can be computed from data and (b) which has a known distribution under (H0).
- 4. Determine the behaviour of statistics under (H1) and buid critical region (where (H0) rejected)
- 5. Compute the region at a fixed error I threshold and compare to values obtained from data. Or compute p-value of the test from data.

General methodology

- 1. Modelling of the problem.
- 2. Determine alternative hypotheses to test (disjoint but not necessarily exhaustive).
- 3. Choose of a statistic than (a) can be computed from data and (b) which has a known distribution under (H0).
- 4. Determine the behaviour of statistics under (H1) and buid critical region (where (H0) rejected)
- 5. Compute the region at a fixed error I threshold and compare to values obtained from data. Or compute p-value of the test from data.
- 6. Statistical conclusion: accept or reject (H0). Comment on p-value ?

General methodology

- 1. Modelling of the problem.
- 2. Determine alternative hypotheses to test (disjoint but not necessarily exhaustive).
- 3. Choose of a statistic than (a) can be computed from data and (b) which has a known distribution under (H0).
- 4. Determine the behaviour of statistics under (H1) and buid critical region (where (H0) rejected)
- 5. Compute the region at a fixed error I threshold and compare to values obtained from data. Or compute p-value of the test from data.
- 6. Statistical conclusion: accept or reject (H0). Comment on p-value ?
- opt. Can you say something about the power ?

General methodology

- 1. Modelling of the problem.
- 2. Determine alternative hypotheses to test (disjoint but not necessarily exhaustive).
- 3. Choose of a statistic than (a) can be computed from data and (b) which has a known distribution under (H0).
- 4. Determine the behaviour of statistics under (H1) and buid critical region (where (H0) rejected)
- 5. Compute the region at a fixed error I threshold and compare to values obtained from data. Or compute p-value of the test from data.
- 6. Statistical conclusion: accept or reject (H0). Comment on p-value ?
- opt. Can you say something about the power ?
 - 7. Strategic conclusion: how do YOU decide thanks to the light shed by statistical result ?

Methodology into details

► Hypothesis:= any subset of the family of all considered probability distributions *P*. In practice, hypotheses are often on unknown parameters of distributions → parametric hypotheses, defined by equalities or inequalities: (H0) θ₀ ∈ Θ₀ and (H1) θ₀ ∈ Θ₁. In turn, they can be simple is only one value for the parameters is tested or muliple composite.

Methodology into details

- ► Hypothesis:= any subset of the family of all considered probability distributions *P*. In practice, hypotheses are often on unknown parameters of distributions → parametric hypotheses, defined by equalities or inequalities: (H0) θ₀ ∈ Θ₀ and (H1) θ₀ ∈ Θ₁. In turn, they can be simple is only one value for the parameters is tested or muliple composite.
- ► Choose a test statistic T_n:=a random variable which only depends on (Θ₀, Θ₁) and on obervations of the (X_i)'s. Interesting if the distribution is known given (H0) is true. Note that it is an estimator...depending on (H0) and (H1).

Methodology into details

- ► Hypothesis:= any subset of the family of all considered probability distributions *P*. In practice, hypotheses are often on unknown parameters of distributions → parametric hypotheses, defined by equalities or inequalities: (H0) θ₀ ∈ Θ₀ and (H1) θ₀ ∈ Θ₁. In turn, they can be simple is only one value for the parameters is tested or muliple composite.
- ► Choose a test statistic T_n:=a random variable which only depends on (Θ₀, Θ₁) and on obervations of the (X_i)'s. Interesting if the distribution is known given (H0) is true. Note that it is an estimator...depending on (H0) and (H1).
- How to choose a good test statistic ? Remember the typology of confidence intervals ? And explore R help ?!

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○○

▶ Determine the rejection region R. Usually of the form (r;∞), (-∞;r) or (-∞;r) ∪ (r';∞). To decide, examine how the test statistic behaves under (H1).

- ▶ Determine the rejection region R. Usually of the form (r;∞), (-∞;r) or (-∞;r) ∪ (r';∞). To decide, examine how the test statistic behaves under (H1).
- type I error:=probability to reject (H0) whilst it is correct. Mathematically:

$$\alpha = \sup_{\theta_0 \in \Theta_0} P\left(T_n \in R | X_1 \dots X_n \text{ iid } \sim P_{\theta_0}\right)$$

- ▶ Determine the rejection region R. Usually of the form (r;∞), (-∞;r) or (-∞;r) ∪ (r';∞). To decide, examine how the test statistic behaves under (H1).
- type I error:=probability to reject (H0) whilst it is correct. Mathematically:

$$\alpha = \sup_{\theta_0 \in \Theta_0} P\left(T_n \in R | X_1 \dots X_n \text{ iid } \sim P_{\theta_0}\right)$$

• Remark: useless to try to get $\alpha = 0$, it is a useless test !

- ▶ Determine the rejection region R. Usually of the form (r;∞), (-∞;r) or (-∞;r) ∪ (r';∞). To decide, examine how the test statistic behaves under (H1).
- type I error:=probability to reject (H0) whilst it is correct. Mathematically:

$$\alpha = \sup_{\theta_0 \in \Theta_0} P\left(T_n \in R | X_1 \dots X_n \text{ iid } \sim P_{\theta_0}\right)$$

- ▶ Remark: useless to try to get $\alpha = 0$, it is a useless test !
- ▶ p-value:=maximal value of α so that the test would accept the observed statistic to be drawn under (H0) ≈ credibility index on (H0). Alternative definition: probability to obtain a test statistic value at least as contradictory to (H0) as the observed value assuming (H0) is true if we repeated the experiment.

dissymetry between (H0) and (H1): (H0) tends to be kept unless good reasons to reject it. (H1) is only used to choose the form of the rejection region, not its bounds ! It is then interesting to look at the

イロト 不得下 イヨト イヨト

- dissymetry between (H0) and (H1): (H0) tends to be kept unless good reasons to reject it. (H1) is only used to choose the form of the rejection region, not its bounds ! It is then interesting to look at the
- ► type II error:=probability to wrongly keep (H0) (while (H1) is true). In mathematical terms:

$$\beta = \sup_{\theta_0 \ in\Theta_1} P\left(T_n \notin R | X_1 \dots X_n \, \mathsf{iid} \sim P_{\theta_0}\right)$$

- dissymetry between (H0) and (H1): (H0) tends to be kept unless good reasons to reject it. (H1) is only used to choose the form of the rejection region, not its bounds ! It is then interesting to look at the
- ► type II error:=probability to wrongly keep (H0) (while (H1) is true). In mathematical terms:

$$\beta = \sup_{\theta_0 \ in\Theta_1} P\left(T_n \notin R | X_1 \dots X_n \, \mathsf{iid} \sim P_{\theta_0}\right)$$

hence (H0) is chosen according to a firmly established theory (you don't want to make a fool of yourself), because caution is needed or...for subjective reasons (consumer choice is not that of manufacturers !)

How does the pharmaceutical industry proceed to approve the commercialisation of a new drug ? Basically two possibilities

test again a placebo; (H0) the new drug is better than the placebo. Do you like it ?

2013

11 / 14

How does the pharmaceutical industry proceed to approve the commercialisation of a new drug ? Basically two possibilities

- ► test again a placebo; (H0) the new drug is better than the placebo. Do you like it ?
- I don't: it's not difficult to find a chemical compound which makes better than empty pills (or sugar) ?!

How does the pharmaceutical industry proceed to approve the commercialisation of a new drug ? Basically two possibilities

- ► test again a placebo; (H0) the new drug is better than the placebo. Do you like it ?
- ► I don't: it's not difficult to find a chemical compound which makes better than empty pills (or sugar) ?!
- you can also test again an existing drug. But then (H0) can be "the new drug is at least as efficient as the old one" (good for the compagny).

How does the pharmaceutical industry proceed to approve the commercialisation of a new drug ? Basically two possibilities

- test again a placebo; (H0) the new drug is better than the placebo. Do you like it ?
- I don't: it's not difficult to find a chemical compound which makes better than empty pills (or sugar) ?!
- you can also test again an existing drug. But then (H0) can be "the new drug is at least as efficient as the old one" (good for the compagny).
- ► if the social healthcare hired me, I would test (H0) "the new drug does not improve over existing ones".

How does the pharmaceutical industry proceed to approve the commercialisation of a new drug ? Basically two possibilities

- ► test again a placebo; (H0) the new drug is better than the placebo. Do you like it ?
- ► I don't: it's not difficult to find a chemical compound which makes better than empty pills (or sugar) ?!
- you can also test again an existing drug. But then (H0) can be "the new drug is at least as efficient as the old one" (good for the compagny).
- ► if the social healthcare hired me, I would test (H0) "the new drug does not improve over existing ones".
- ► Sadly enough, it's the forst option that is used ??!! For fairness between new and existing molecules...

How does the pharmaceutical industry proceed to approve the commercialisation of a new drug ? Basically two possibilities

- ► test again a placebo; (H0) the new drug is better than the placebo. Do you like it ?
- I don't: it's not difficult to find a chemical compound which makes better than empty pills (or sugar) ?!
- you can also test again an existing drug. But then (H0) can be "the new drug is at least as efficient as the old one" (good for the compagny).
- ► if the social healthcare hired me, I would test (H0) "the new drug does not improve over existing ones".
- Sadly enough, it's the forst option that is used ??!! For fairness between new and existing molecules...

Historical notes: statistics were of great help in modern medicine.

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○○

Tests you need to know

and we shall see during the next session

- ▶ parametric tests (observations drawn from N or large samples so that C.L.Th. applies)
 - \blacktriangleright one sample: comparing the empirical mean to a theoretical value \rightarrow Z-test or t-test
 - ► two independent samples: t-test, F-test
 - paired samples: paired t-test
 - ► several samples: ANOVA, not today !
- \blacktriangleright adequation tests: $\chi^2\text{-test.}$ Normality check: Kolmogorov or Shapiro-Wilks.
- non-parametric tests (when small samples or non Gaussian distributions)
 - ► comparing 2 medians from independent samples: Mann-Whitney test.
 - ► two paired samples: Wilcoxon test on differences.
 - ► several samples: Kruskal-Wallis

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 ろの⊙

Exercises

Poisson arrival at motorway tolls

For two hours, at a motorway toll, we write down the number of cars arriving during each 2 minute intervals. We obtain:

Evolution of purchasing power

In 2004, the total amount spent on products which are not essentials (*e.g.* travels, shows ... as opposed to food, hoosing ...) was 632 euros per month per household accoring to the INSEE during a partial survey over millions of households. In 2008, from a sample of 2000 interviewed by telephone, 1837 answers were obtained and the declared mean value was 598 euros (with sd 254 euros). If you assume a 2% inflation per year, would you say that the amount spent on non-essentials has significantly decreased ?

E. Rachelson & M. Vignes (ISAE)



Next time: more tests and analysis of variance (ANOVA)