## Statistics and learning An introduction: from data to modelling

### Emmanuel Rachelson and Matthieu Vignes

ISAE SupAero

Wednesday 3<sup>rd</sup> September 2013

イロト イポト イヨト イヨト 二日

990

1 / 13

2013

A quick, partial and not very comprehensive overview

Goal of this course: not a recipe cooking handbook, rather a path to mathematical reasoning which leads to dealing with quantitative aspects of decision making from data and still accounting for uncertainty.

A quick, partial and not very comprehensive overview

- Goal of this course: not a recipe cooking handbook, rather a path to mathematical reasoning which leads to dealing with quantitative aspects of decision making from data and still accounting for uncertainty.
- Application of stats (and Machine Learning): pension matters, optimal insurance premium, stock market predictions, clinical trial for a new medicine, marketing research and planning ... Different branches: geostatistics, statistical physics ...

A quick, partial and not very comprehensive overview

- Goal of this course: not a recipe cooking handbook, rather a path to mathematical reasoning which leads to dealing with quantitative aspects of decision making from data and still accounting for uncertainty.
- Application of stats (and Machine Learning): pension matters, optimal insurance premium, stock market predictions, clinical trial for a new medicine, marketing research and planning ... Different branches: geostatistics, statistical physics ...
- Professional AND citizen interest: ad-hoc exploitation of available data and don't be manipulated ?!

A quick, partial and not very comprehensive overview

- Goal of this course: not a recipe cooking handbook, rather a path to mathematical reasoning which leads to dealing with quantitative aspects of decision making from data and still accounting for uncertainty.
- Application of stats (and Machine Learning): pension matters, optimal insurance premium, stock market predictions, clinical trial for a new medicine, marketing research and planning ... Different branches: geostatistics, statistical physics ...
- Professional AND citizen interest: ad-hoc exploitation of available data and don't be manipulated ?!
- ► few prerequisites: basic/intermediate maths and probability calculus.

イロト 不得 トイヨト イヨト 二日

A quick, partial and not very comprehensive overview

- Goal of this course: not a recipe cooking handbook, rather a path to mathematical reasoning which leads to dealing with quantitative aspects of decision making from data and still accounting for uncertainty.
- Application of stats (and Machine Learning): pension matters, optimal insurance premium, stock market predictions, clinical trial for a new medicine, marketing research and planning ... Different branches: geostatistics, statistical physics ...
- Professional AND citizen interest: ad-hoc exploitation of available data and don't be manipulated ?!
- ► few prerequisites: basic/intermediate maths and probability calculus.
- Grail: linking data to mathematical modelling, objectively quantify and interpret conclusions and...awareness of limitations: statistics helps but won't make decision for you !

# Inspiring work / our bibliography



T. Hastie, R. Tibshirani and J. Friedman. *Elements of statistical learning.* Springer, 2nd edition, 2009.



E. Moulines, F. Roueff and J.-L. Pac (and formerly F. Rossi) Statistiques. *Cours TelecomParisTech*, 2008.



A. Garivier

Statistiques avancées. Cours Centrale 2011, 2011.



S. Clémençon.

Apprentissage statistique. Cours TELECOM ParisTech, 2011-2012.



S. Arlot, Francis B., O. Catoni, G. Stolz and G. Obozinski Apprentissage. *Cours ENS*, 2012.



N. Chopin, D. Rosenberg and G. Stolz

Eléments de statistique pour citoyens d'aujourd'hui et managers de demain. *Cours L3 HEC*, 2012–2013.

### And many others we just forgot to mention.

A. Baccini, P. Besse, S. Canu, S. Déjean, B. Laurent, C. Marteau, P. Martin and H. Milhem Wikistat, le cours dont vous êtes le héros. http://wikistat.fr/, 2012.

A.B. Dufour D. Chessel J.R. Lobry, S. Mousset and S. Dray Enseignements de Statistique en Biologie. http://pbil.univ-lyon1.fr/R/, 2012.

F. Bertrand Page professionnelle - Enseignements. http://www-irma.u-strasbg.fr/~fbertran/ enseignement/, 2012.

イロト イポト イヨト イヨト

٢

3

# From data to modelling

and back

Two different situations might occur for the same modelling:

 empirical approach to gaining knowledge from an experiment repeated many times,

# From data to modelling

and back

Two different situations might occur for the same modelling:

- empirical approach to gaining knowledge from an experiment repeated many times,
- study of a sample drawn from a population.

イロト 不得下 イヨト イヨト

# From data to modelling

and back

Two different situations might occur for the same modelling:

- empirical approach to gaining knowledge from an experiment repeated many times,
- study of a sample drawn from a population.

Preference between two possible configurations

Consumer ID	1	2	3	4	5	6	
Opinion	А	А	В	А	В	В	

We can denote by  $x_i$  successive opinions taking (binary) values "A" (= 0) or "B" (= 1). Mathematician sees that as realisation of random variables denoted  $X_i$ .

イロト イロト イヨト イヨト 三日

### Randomness...

...arises from the choice of the questioned persons, NOT from in each actual answer.

Incidental reminder: Bernouilli distribution, with parameter 0

▶ laid question: is  $p_0 > 1/2$  or < 1/2 ? This is a **test**.

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

### Randomness...

...arises from the choice of the questioned persons, NOT from in each actual answer.

Incidental reminder: Bernouilli distribution, with parameter 0

- ▶ laid question: is  $p_0 > 1/2$  or < 1/2 ? This is a **test**.
- but hey, what is  $p_0$  actually ??

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 ろの⊙

### Randomness...

...arises from the choice of the questioned persons, NOT from in each actual answer.

Incidental reminder: Bernouilli distribution, with parameter 0

- ▶ laid question: is  $p_0 > 1/2$  or < 1/2 ? This is a **test**.
- but hey, what is  $p_0$  actually ??
- ▶ all consumers cannot be interviewed for obvious reasons !

### Randomness...

...arises from the choice of the questioned persons, NOT from in each actual answer.

Incidental reminder: Bernouilli distribution, with parameter 0

- ▶ laid question: is  $p_0 > 1/2$  or < 1/2 ? This is a **test**.
- but hey, what is  $p_0$  actually ??
- ▶ all consumers cannot be interviewed for obvious reasons !
- Statistics is a sound framework to

### Randomness...

...arises from the choice of the questioned persons, NOT from in each actual answer.

Incidental reminder: Bernouilli distribution, with parameter 0

- ▶ laid question: is  $p_0 > 1/2$  or < 1/2 ? This is a **test**.
- but hey, what is  $p_0$  actually ??
- ▶ all consumers cannot be interviewed for obvious reasons !
- Statistics is a sound framework to
  - 1. describe sample using estimates

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

### Randomness...

...arises from the choice of the questioned persons, NOT from in each actual answer.

Incidental reminder: Bernouilli distribution, with parameter 0

- ▶ laid question: is  $p_0 > 1/2$  or < 1/2 ? This is a **test**.
- but hey, what is  $p_0$  actually ??
- ▶ all consumers cannot be interviewed for obvious reasons !
- Statistics is a sound framework to
  - 1. describe sample using estimates
  - 2. quantitatively answer the question (generalising sample to full population conclusions)

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

a mandatory step

we start by modelling the problem (Sample/population description, variables and their distributions).

3

a mandatory step

- we start by modelling the problem (Sample/population description, variables and their distributions).
- ▶ in the example: "correct model" ∈ (Bern(p))<sub>p</sub>; n realisations (x<sub>i</sub>)<sub>i</sub> of iid (indep. & identical. distr.) random variables (X<sub>i</sub>)<sub>i</sub> ~ Bern(p) are available.

a mandatory step

- we start by modelling the problem (Sample/population description, variables and their distributions).
- ▶ in the example: "correct model" ∈ (Bern(p))<sub>p</sub>; n realisations (x<sub>i</sub>)<sub>i</sub> of iid (indep. & identical. distr.) random variables (X<sub>i</sub>)<sub>i</sub> ~ Bern(p) are available.
- ► remember the Jean Tibéri vs. Lyne Cohen-Solal (+ Ph. Meyer) council election in Paris in 2008 between 20.45 and 21.15 ? At 20.45, (463; 409; 106) but after counting the votes : (11, 044; 11, 269; 2, 730).

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

a mandatory step

- we start by modelling the problem (Sample/population description, variables and their distributions).
- ▶ in the example: "correct model" ∈ (Bern(p))<sub>p</sub>; n realisations (x<sub>i</sub>)<sub>i</sub> of iid (indep. & identical. distr.) random variables (X<sub>i</sub>)<sub>i</sub> ~ Bern(p) are available.
- ► remember the Jean Tibéri vs. Lyne Cohen-Solal (+ Ph. Meyer) council election in Paris in 2008 between 20.45 and 21.15 ? At 20.45, (463; 409; 106) but after counting the votes : (11, 044; 11, 269; 2, 730).
- Construction of **confidence intervals** to answer the question.

## Useful tools

a reminder ?

▶ mean (central tendency):  $E[X] = \overline{X} := \frac{X_1 + \dots + X_n}{n}$ ,

### Useful tools

a reminder ?

- ▶ mean (central tendency):  $E[X] = \overline{X} := \frac{X_1 + \dots + X_n}{n}$ ,
- ► standard deviation (dispersion tendency): =  $\sigma(X) := \sqrt{E[(X - E[X])^2]} = \sqrt{E[X]^2 - E[X^2]}$ ,

▲ロト ▲圖ト ▲ヨト ▲ヨト 三ヨ - のへ⊙

## Useful tools

a reminder ?

- ▶ mean (central tendency):  $E[X] = \overline{X} := \frac{X_1 + \ldots + X_n}{n}$ ,
- ► standard deviation (dispersion tendency): =  $\sigma(X) := \sqrt{E[(X - E[X])^2]} = \sqrt{E[X]^2 - E[X^2]}$ ,
- (almost never use skewness and kurtosis)

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 ろの⊙

Law of large numbers

#### Theorem

Let  $X_1 \dots X_n$  be iid random variables with mean  $\mu$ . Then the empirical mean converges in probability towards  $\mu$ , i.e.:

$$\overline{X_n} := \frac{1}{n} (X_1 + \ldots + X_n) \longrightarrow \mu.$$

<ロト < 回 > < 回 > < 回 > < 回 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Law of large numbers

#### Theorem

Let  $X_1 \dots X_n$  be iid random variables with mean  $\mu$ . Then the empirical mean converges in probability towards  $\mu$ , i.e.:

$$\overline{X_n} := \frac{1}{n} (X_1 + \ldots + X_n) \longrightarrow \mu.$$

In other term, for all  $\epsilon > 0$ ,  $P\left( \mid \overline{X_n} - \mu \mid > \epsilon \right) \to 0$ 

イロト 不得 トイヨト イヨト ヨー のくや

Central limit theorem

#### Theorem

Let  $X_1 \dots X_n$  be iid random variables which admit an order 2 moment. Denote by  $\mu$  and  $\sigma$  the corresponding mean and standard deviation, then:

$$\frac{\sqrt{n}}{\sigma}(\overline{X_n} - \mu) \longrightarrow \mathcal{N}(0, 1).$$

イロト 不得下 イヨト イヨト

Central limit theorem

#### Theorem

Let  $X_1 \dots X_n$  be iid random variables which admit an order 2 moment. Denote by  $\mu$  and  $\sigma$  the corresponding mean and standard deviation, then:

$$\frac{\sqrt{n}}{\sigma}(\overline{X_n} - \mu) \longrightarrow \mathcal{N}(0, 1).$$

In the case of distribution with density functions, this means that

$$P\left(\frac{\sqrt{n}}{\sigma}(\overline{X_n} - \mu) \le x\right) := F_n(x) \longrightarrow P(Z \le x) = \frac{\int_{-\infty}^x e^{-z^2/2} dz}{\sqrt{2\pi}}$$

E. Rachelson & M. Vignes (ISAE)

2013 9 / 13

▶ Back to the "preference" example: A vs. B.

イロト イポト イヨト イヨト

990

► Back to the "preference" example: A vs. B.

Preference between two possible configurations

Consumer ID	1	2	3	4	5	6	
Opinion	А	А	В	А	В	В	

A = A = A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A

► Back to the "preference" example: A vs. B.

Preference between two possible configurations

Consumer ID	1	2	3	4	5	6	
Opinion	А	А	В	А	В	В	

• We compute  $\overline{x_{100}} = \frac{x_1 + \dots + x_{100}}{100} = 0.42.$ 

► Back to the "preference" example: A vs. B.

Preference between two possible configurations

Consumer ID	1	2	3	4	5	6	
Opinion	А	А	В	А	В	В	

- We compute  $\overline{x_{100}} = \frac{x_1 + \dots + x_{100}}{100} = 0.42.$
- Our intuition and the LLN tell us that  $p_0$  is "close" to 0.42.

► Back to the "preference" example: A vs. B.

Preference between two possible configurations

Consumer ID	1	2	3	4	5	6	
Opinion	А	А	В	А	В	В	

- We compute  $\overline{x_{100}} = \frac{x_1 + \dots + x_{100}}{100} = 0.42.$
- Our intuition and the LLN tell us that  $p_0$  is "close" to 0.42.
- ► Can we conclude ? Is this **estimate** enough ?

Let's play around the Central limit theorem...

at the price of a slight risk

▶ we directly derive 
$$\frac{\sqrt{n}}{\sqrt{p_0(1-p_0)}}(\overline{X_n}-p_0) \rightarrow \mathcal{N}(0,1)$$
 so

E. Rachelson & M. Vignes (ISAE)

at the price of a slight risk

▶ we directly derive 
$$\frac{\sqrt{n}}{\sqrt{p_0(1-p_0)}}(\overline{X_n}-p_0) \to \mathcal{N}(0,1)$$
 so

▶ with a confidence level (what is that ??) of 95%:  $\frac{\sqrt{n}}{\sqrt{p_0(1-p_0)}} |\overline{X_n} - p_0| \le u = 1.96.$  This is equivalent to:

《日》《聞》《臣》《臣》 三臣

at the price of a slight risk

▶ we directly derive 
$$\frac{\sqrt{n}}{\sqrt{p_0(1-p_0)}}(\overline{X_n}-p_0) \to \mathcal{N}(0,1)$$
 so

E. Rachelson & M. Vignes (ISAE)

at the price of a slight risk

▶ we directly derive 
$$\frac{\sqrt{n}}{\sqrt{p_0(1-p_0)}}(\overline{X_n}-p_0) \to \mathcal{N}(0,1)$$
 so

▶ with a confidence level (what is that ??) of 95%:  $\frac{\sqrt{n}}{\sqrt{p_0(1-p_0)}} |\overline{X_n} - p_0| \le u = 1.96.$  This is equivalent to: ▶  $p_0 \in I_n := \left[\overline{X_n} - 1.96\frac{\sqrt{p_0(1-p_0)}}{\overline{x_n}} : \overline{X_n} + 1.96\frac{\sqrt{p_0(1-p_0)}}{\overline{x_n}}\right]$ 

• 
$$p_0 \in I_n := \left[ X_n - 1.96 \frac{\sqrt{p_0(1-p_0)}}{\sqrt{n}}; X_n + 1.96 \frac{\sqrt{p_0(1-p_0)}}{\sqrt{n}} \right],$$

▶ again, does this help ?

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 ろの⊙

at the price of a slight risk

▶ we directly derive 
$$\frac{\sqrt{n}}{\sqrt{p_0(1-p_0)}}(\overline{X_n}-p_0) \to \mathcal{N}(0,1)$$
 so

▶ with a confidence level (what is that ??) of 95%:  $\frac{\sqrt{n}}{\sqrt{p_0(1-p_0)}} |\overline{X_n} - p_0| \le u = 1.96.$  This is equivalent to:  $\left[ -\frac{\sqrt{p_0(1-p_0)}}{\sqrt{p_0(1-p_0)}} - \frac{\sqrt{p_0(1-p_0)}}{\sqrt{p_0(1-p_0)}} \right]$ 

• 
$$p_0 \in I_n := \left[ \overline{X_n} - 1.96 \frac{\sqrt{p_0(1-p_0)}}{\sqrt{n}}; \overline{X_n} + 1.96 \frac{\sqrt{p_0(1-p_0)}}{\sqrt{n}} \right],$$

- again, does this help ?
- ▶ yes, using the simplifying trick  $\sqrt{x(1-x)} \le 1/2$ :  $I_n \subseteq [\bar{X}_n - 1/\sqrt{n}; \bar{X}_n + 1/\sqrt{n}] = [32\%; 52\%]$  in our scenario. Your final conclusion ? Again what is 95% here ?

at the price of a slight risk

▶ we directly derive 
$$\frac{\sqrt{n}}{\sqrt{p_0(1-p_0)}}(\overline{X_n}-p_0) \to \mathcal{N}(0,1)$$
 so

▶ with a confidence level (what is that ??) of 95%:  $\frac{\sqrt{n}}{\sqrt{p_0(1-p_0)}} |\overline{X_n} - p_0| \le u = 1.96.$  This is equivalent to:

• 
$$p_0 \in I_n := \left[ \overline{X_n} - 1.96 \frac{\sqrt{p_0(1-p_0)}}{\sqrt{n}}; \overline{X_n} + 1.96 \frac{\sqrt{p_0(1-p_0)}}{\sqrt{n}} \right],$$

- ► again, does this help ?
- ▶ yes, using the simplifying trick  $\sqrt{x(1-x)} \le 1/2$ :  $I_n \subseteq [\bar{X}_n - 1/\sqrt{n}; \bar{X}_n + 1/\sqrt{n}] = [32\%; 52\%]$  in our scenario. Your final conclusion ? Again what is 95% here ?
- is the conclusion similar if n = 1,000 ?

Note: 95% could have been replaced by 99%. How could this have affected the conclusion ? What about 100% ?

2013 11 / 13

decision depends on what is tested ?!

► (H<sub>0</sub>) is the basic hypothesis. It will be rejected iif data strongly supports it (e.g. dangerous drug or alleged innocent). Then

decision depends on what is tested ?!

- ► (H<sub>0</sub>) is the basic hypothesis. It will be rejected iif data strongly supports it (e.g. dangerous drug or alleged innocent). Then
- ► (H<sub>1</sub>) is the alternative hypothesis which would be accepted iif (H<sub>0</sub>) was regognised to be unacceptable.

decision depends on what is tested ?!

- ► (H<sub>0</sub>) is the basic hypothesis. It will be rejected iif data strongly supports it (e.g. dangerous drug or alleged innocent). Then
- ► (H<sub>1</sub>) is the alternative hypothesis which would be accepted iif (H<sub>0</sub>) was regognised to be unacceptable.
- ► Back on the example: do you want to test (H<sub>0</sub>) p<sub>0</sub> ≥ 1/2 ? Or (H<sub>0</sub>') p<sub>0</sub> ≤ 1/2 ? Which one is the sensible/aggressive choice ?

decision depends on what is tested ?!

- ► (H<sub>0</sub>) is the basic hypothesis. It will be rejected iif data strongly supports it (e.g. dangerous drug or alleged innocent). Then
- ► (H<sub>1</sub>) is the alternative hypothesis which would be accepted iif (H<sub>0</sub>) was regognised to be unacceptable.
- ▶ Back on the example: do you want to test (H<sub>0</sub>)  $p_0 \ge 1/2$  ? Or (H<sub>0</sub>')  $p_0 \le 1/2$  ? Which one is the sensible/aggressive choice ?
- lessons from this: tests are not reductible to confidence intervals and...don't be fooled by an obscure choice of hypotheses !

From 18<sup>th</sup> September 2013 to 7<sup>th</sup> February 2014, you will hear about:

 descriptive statistics and modelling: estimation, dispersion measure, confidence intervals, PCA and perhaps more (CA, clustering, risk function).

▶ Projects with some stats content from Nov. 2013 'til May 2014...tbc

From 18<sup>th</sup> September 2013 to 7<sup>th</sup> February 2014, you will hear about:

- descriptive statistics and modelling: estimation, dispersion measure, confidence intervals, PCA and perhaps more (CA, clustering, risk function).
- ► Tests, ANOVA (one, several factors, analysis of covariance)

▶ Projects with some stats content from Nov. 2013 'til May 2014...tbc

From 18<sup>th</sup> September 2013 to 7<sup>th</sup> February 2014, you will hear about:

- descriptive statistics and modelling: estimation, dispersion measure, confidence intervals, PCA and perhaps more (CA, clustering, risk function).
- ► Tests, ANOVA (one, several factors, analysis of covariance)
- ► (linear) Regressions

▶ Projects with some stats content from Nov. 2013 'til May 2014...tbc

From 18<sup>th</sup> September 2013 to 7<sup>th</sup> February 2014, you will hear about:

- descriptive statistics and modelling: estimation, dispersion measure, confidence intervals, PCA and perhaps more (CA, clustering, risk function).
- ► Tests, ANOVA (one, several factors, analysis of covariance)
- ► (linear) Regressions
- Trees and ensemble methods

### ▶ Projects with some stats content from Nov. 2013 'til May 2014...tbc

From 18<sup>th</sup> September 2013 to 7<sup>th</sup> February 2014, you will hear about:

- descriptive statistics and modelling: estimation, dispersion measure, confidence intervals, PCA and perhaps more (CA, clustering, risk function).
- ► Tests, ANOVA (one, several factors, analysis of covariance)
- ► (linear) Regressions
- Trees and ensemble methods
- Kernels and parcimony

▶ Projects with some stats content from Nov. 2013 'til May 2014...tbc

From 18<sup>th</sup> September 2013 to 7<sup>th</sup> February 2014, you will hear about:

- descriptive statistics and modelling: estimation, dispersion measure, confidence intervals, PCA and perhaps more (CA, clustering, risk function).
- ► Tests, ANOVA (one, several factors, analysis of covariance)
- ► (linear) Regressions
- Trees and ensemble methods
- Kernels and parcimony
- Markov Chain Monte Carlo methods
- ▶ Projects with some stats content from Nov. 2013 'til May 2014...tbc

- 4 伺 ト 4 ヨ ト 4 ヨ ト

From 18<sup>th</sup> September 2013 to 7<sup>th</sup> February 2014, you will hear about:

- descriptive statistics and modelling: estimation, dispersion measure, confidence intervals, PCA and perhaps more (CA, clustering, risk function).
- ► Tests, ANOVA (one, several factors, analysis of covariance)
- ► (linear) Regressions
- Trees and ensemble methods
- Kernels and parcimony
- Markov Chain Monte Carlo methods
- ▶ ...and lots of R ;) !
- ▶ Projects with some stats content from Nov. 2013 'til May 2014...tbc