

18.650 – Fundamentals of Statistics

1. Introduction and probability

Goals

Goals:

- ▶ To give you a solid introduction to the mathematical theory behind statistical methods;
- ▶ To provide theoretical guarantees for the statistical methods that you may use for certain applications.

At the end of this class, you will be able to

1. From a real-life situation, formulate a statistical problem in mathematical terms
2. Select appropriate statistical methods for your problem
3. Understand the implications and limitations of various methods

Why statistics?

In the press

The New York Times

THE UPSHOT

Nike Says Its \$250 Running Shoes Will Make You Run Much Faster. What if That's Actually True?

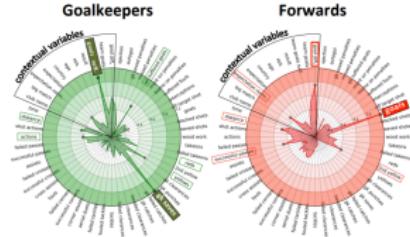
An analysis of nearly 500,000 running times estimates the effect of shoes on race performance.



MIT
Technology
Review

Data Mining Reveals the Way Humans Evaluate Each Other

Vast databases of soccer statistics expose the limited way human observers rate performance and suggest how they can do significantly better.



In businesses

Harvard
Business
Review

How Vineyard Vines Uses Analytics to Win Over Customers

TECHNOLOGY DIGITAL ARTICLE by Dave Sutton

A case study on how personalization is changing retail.

 SAVE  SHARE JUNE 08, 2018



FAST COMPANY

AppNexus is key to AT&T's plans to use HBO for more consumer data

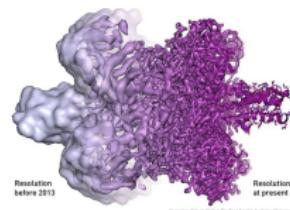
New WarnerMedia CEO John Stankey says HBO is going to "change direction a little bit," and it's all about the advertising.



In science and engineering

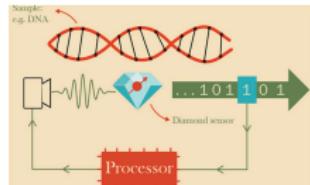
The Guardian

What is cryo-electron microscopy, the Nobel prize-winning technique?



IEEE
SPECTRUM

**Measuring Tiny Magnetic Fields
With an Intelligent Quantum Sensor**



On TV

LAST WEEK TONIGHT WITH JOHN OLIVER



"Last Week Tonight with John Oliver": Scientific Studies



Data Science and the Art of Producing
Entertainment at Netflix



Statistics, Data Science . . . and all that

Statistics, Data Science, Machine Learning, Artificial Intelligence

What's the difference?

Statistics, Data Science . . . and all that

Statistics, Data Science, Machine Learning, Artificial Intelligence

What's the difference?

- ▶ All use data to gather insight and ultimately make decisions
- ▶ Statistics is at the core of the data processing part
- ▶ Nowadays, computational aspects play an important role as data becomes larger

Computational and statistical aspects of data science

- ▶ Computational view: data is a (large) sequence of numbers that needs to be processed by a relatively fast algorithm: approximate nearest neighbors, low dimensional embeddings, spectral methods, distributed optimization, etc.
- ▶ Statistical view: data comes from a **random process**. The goal is to learn how this process works in order to make predictions or to understand what plays a role in it.

To understand randomness, we need PROBABILITY.

Probability

- ▶ Probability studies randomness (hence the prerequisite)
- ▶ Sometimes, the physical process is completely known: dice, cards, roulette, fair coins, . . .

Rolling 1 die:

- ▶ Alice gets \$1 if # of dots ≤ 3
- ▶ Bob gets \$2 if # of dots ≤ 2

Who do you want to be: Alice or Bob?

Rolling 2 dice:

- ▶ Choose a number between 2 and 12
- ▶ Win \$100 if you chose the sum of the 2 dice

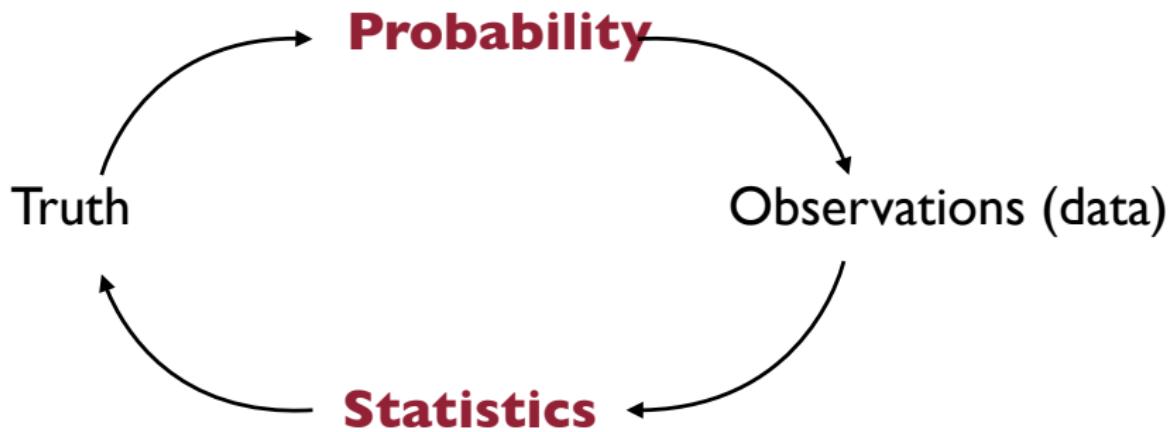
Which number do you choose?

Statistics and modeling

- ▶ Dice are well known random process from physics: 1/6 chance of each side (no need for data!), dice are independent. We can deduce the probability of outcomes, and expected \$ amounts. This is **probability**.
- ▶ How about more complicated processes? Need to estimate parameters from data. This is **statistics**
- ▶ Sometimes real randomness (random student, biased coin, measurement error, . . .)
- ▶ Sometimes deterministic but too complex phenomenon: **statistical modeling**

Complicated process “=” Simple process + random noise

- ▶ (good) Modeling consists in choosing (plausible) simple process **and** noise distribution.



Statistics vs. probability

Probability Previous studies showed that the drug was 80% effective. Then we can anticipate that for a study on 100 patients, in average 80 will be cured and at least 65 will be cured with 99.99% chances.

Statistics Observe that 78/100 patients were cured. We (will be able to) conclude that we are 95% confident that for other studies the drug will be effective on between 69.88% and 86.11% of patients

What this course is about

- ▶ Understand **mathematics** behind statistical methods
- ▶ Justify quantitative statements given modeling assumptions
- ▶ Describe interesting mathematics arising in statistics
- ▶ Provide a math toolbox to extend to other models.

What this course is **not** about

- ▶ Statistical thinking/modeling (e.g., 15.075)
- ▶ Implementation (e.g. IDS.012)
- ▶ Laundry list of methods (e.g. AP stats)

What this course is about

- ▶ Understand **mathematics** behind statistical methods
- ▶ Justify quantitative statements given modeling assumptions
- ▶ Describe interesting mathematics arising in statistics
- ▶ Provide a math toolbox to extend to other models.

What this course is **not** about

- ▶ Statistical thinking/modeling (e.g., 15.075)
- ▶ Implementation (e.g. IDS.012)
- ▶ Laundry list of methods (e.g. AP stats)

Let's do some statistics

The kiss



Le baiser. Auguste Rodin. 1882.

The kiss



Le baiser. Auguste Rodin. 1882.

The kiss

Full text access provided to **Massachusetts Institute of Technology**
by the **MIT Libraries**

 Cart

nature

International weekly journal of science

Search go Advanced search

Journal home > Archive > Brief Communications > Full Text

Journal content

- + Journal home
- + Advance online publication
- + Current issue
- + Nature News
- + Archive**
- + Supplements

Brief Communications

Nature 421, 711 (13 February 2003) | doi:10.1038/421711a

Human behaviour: Adult persistence of head-turning asymmetry

Onur Güntürkün

A neonatal right-side preference makes a surprising romantic reappearance later in life.

▲ Top

subscribe to
nature 

FULL TEXT

- + Previous | Next +
- + Table of contents
-  Download PDF

Statistical experiment

"A neonatal right-side preference makes a surprising romantic reappearance later in life."

- ▶ Let p denote the proportion of couples that turn their head to the right when kissing.
- ▶ Let us design a statistical experiment and analyze its outcome.
- ▶ Observe n kissing couples times and collect the value of each outcome (say 1 for RIGHT and 0 for LEFT);
- ▶ Estimate p with the proportion \hat{p} of RIGHT.
- ▶ Study: "Human behaviour: Adult persistence of head-turning asymmetry" (Nature, 2003): $n = 124$ and 80 to the right so

$$\hat{p} = \frac{80}{124} = 64.5\%$$

Random intuition

Back to the data:

- ▶ 64.5% is much larger than 50% so there seems to be a preference for turning right.
- ▶ What if our data was RIGHT, RIGHT, LEFT ($n = 3$). That's 66.7% to the right. Even better?
- ▶ Intuitively, we need a large enough sample size n to make a call. How large?
- ▶ Another way to put the problem: for $n = 124$, what is the minimum number of couple "to the right" would you need to see to be convinced that $p > 50\%$? 63? 72? 75? 80?

We need **mathematical modeling** to understand
the accuracy of this procedure?

A first estimator

Formally, this procedure consists of doing the following:

- ▶ For $i = 1, \dots, n$, define $R_i = 1$ if the i th couple turns to the right **RIGHT**, $R_i = 0$ otherwise.
- ▶ The estimator of p is the sample average

$$\hat{p} = \bar{R}_n = \frac{1}{n} \sum_{i=1}^n R_i.$$

What is the accuracy of this estimator ?

In order to answer this question, we propose a statistical model that describes/approximates well the experiment.

We think of the R_i 's as random variables so that \hat{p} is also a random variable. We need to understand its fluctuation.

Modelling assumptions

Coming up with a model consists of making assumptions on the observations $R_i, i = 1, \dots, n$ in order to draw statistical conclusions. Here are the assumptions we make:

1. Each R_i is a random variable.
2. Each of the r.v. R_i is Bernoulli with parameter p .
3. R_1, \dots, R_n are mutually independent.

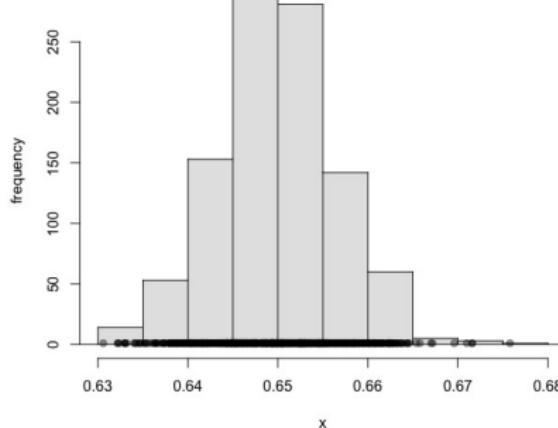
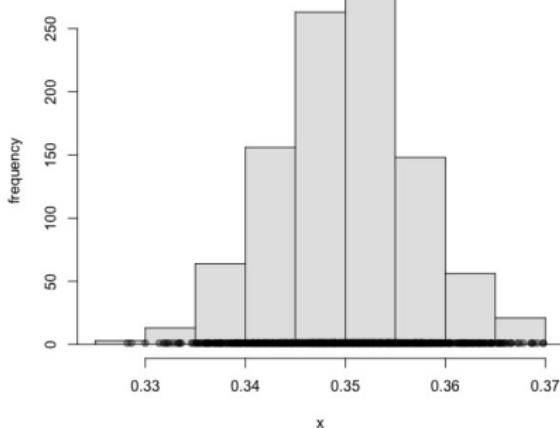
Discussion

Let us discuss these assumptions.

1. Randomness is a way of modeling lack of information; with perfect information about the conditions of kissing (including what goes on in the kissers' mind), physics or sociology would allow us to predict the outcome.
2. Hence, the R_i 's are necessarily Bernoulli r.v. since $R_i \in \{0, 1\}$. They could still have a different parameter $R_i \sim \text{Ber}(p_i)$ for each couple but we don't have enough information with the data to estimate the p_i 's accurately. So we simply assume that our observations come from the same process: $p_i = p$ for all i
3. Independence is reasonable (people were observed at different locations and different times).

Population vs. Samples

- ▶ Assume that there is a total **population** of 5,000 “airport-kissing” couples
- ▶ Assume for the sake of argument that $p = 35\%$ or that $p = 65\%$.
- ▶ What do **samples** of size 124 look like in each case?



Why probability?

We need to understand probabilistic aspects of the distribution of the random variable:

$$\hat{p} = \bar{R}_n = \frac{1}{n} \sum_{i=1}^n R_i.$$

Specifically, we need to be able to answer questions such as:

- ▶ Is the expected value of \hat{p} close to the unknown p ?
- ▶ Does \hat{p} take values close to p with high probability?
- ▶ Is the variance of \hat{p} large? I.e. does \hat{p} fluctuate a lot?

We need probabilistic tools! Most of them are about **average of independent random variables**.

Probability redux

Averages of random variables: LLN & CLT

Let X, X_1, X_2, \dots, X_n be i.i.d. r.v., $\mu = \mathbb{E}[X]$ and $\sigma^2 = \mathbb{V}[X]$.

- ▶ Laws (weak and strong) of large numbers (LLN):

$$\bar{X}_n := \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow[n \rightarrow \infty]{\mathbb{P}, \text{ a.s.}} \mu.$$

- ▶ Central limit theorem (CLT):

$$\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \xrightarrow[n \rightarrow \infty]{(d)} \mathcal{N}(0, 1).$$

(Equivalently, $\sqrt{n} (\bar{X}_n - \mu) \xrightarrow[n \rightarrow \infty]{(d)} \mathcal{N}(0, \sigma^2)$.)

Another useful tool: Hoeffding's inequality

What if n is not large enough to apply CLT?

Theorem (Hoeffding, 1963)

Let n be a positive integer and X, X_1, \dots, X_n be i.i.d. r.v. such that $\mu = \mathbb{E}[X]$ and

$$X \in [a, b] \quad \text{almost surely} \qquad (a < b \text{ are given numbers})$$

Then,

$$\mathbb{P}[|\bar{X}_n - \mu| \geq \varepsilon] \leq 2e^{-\frac{2n\varepsilon^2}{(b-a)^2}}. \quad \forall \varepsilon > 0$$

This holds even for small sample sizes n .

Consequences

- The LLN's tell us that

$$\bar{R}_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}, \text{ a.s.}} p.$$

(what modeling assumptions did we use?)

- Hence, when the size n of the experiment becomes large, \bar{R}_n is a *good* (say "*consistent*") estimator of p .
- The CLT refines this by quantifying *how good* this estimate is: for n large enough the distribution of \hat{p} is almost:

$$\mathbb{P}(|\bar{R}_n - p| \geq \varepsilon) \simeq \mathbb{P}(|\mathcal{N}(0, \frac{p(1-p)}{n})| > \varepsilon)$$

In the Kiss example, $\mathbb{P}(|\bar{R}_n - p| \geq 0.084) \simeq 5\%$

- Hoeffding's inequality tells us that

$$\mathbb{P}(|\bar{R}_n - p| \geq 0.084) \leq 0.35$$

The Gaussian distribution

Because of the CLT, the Gaussian (a.k.a normal) distribution is ubiquitous in statistics. It is named after German Mathematician Carl Friedrich Gauss (1777–1855) in the context of the method of *least squares* (regression).

- ▶ $X \sim \mathcal{N}(\mu, \sigma^2)$
- ▶ $\mathbb{E}[X] = \mu$
- ▶ $\text{var}(X) = \sigma^2 > 0$



Gaussian density (pdf)

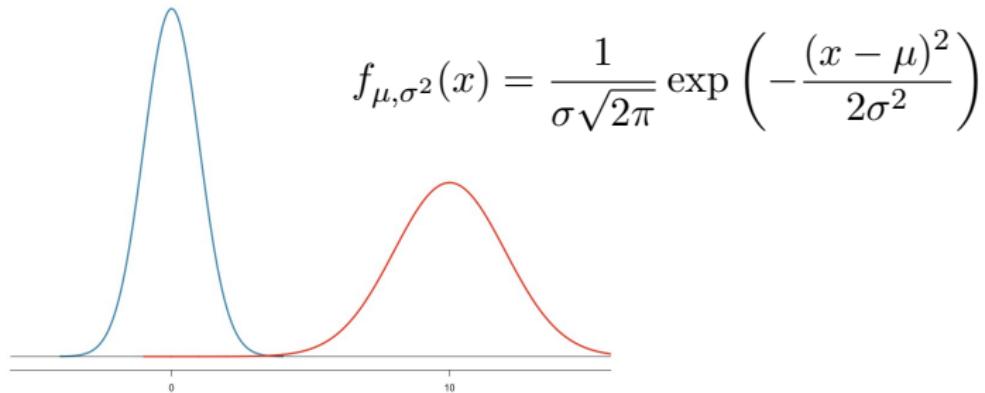


Figure 1: Two pdfs: $\mathcal{N}(0, 1)$ and $\mathcal{N}(10, 4)$

- ▶ Tails decay very fast (like $e^{-\frac{x^2}{2\sigma^2}}$): almost in finite interval.
- ▶ There is no closed form for their cumulative distribution function (CDF). We use tables (or computers):

$$F_{\mu,\sigma^2}(x) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x \exp\left(-\frac{(t-\mu)^2}{2\sigma^2}\right) dt$$

Some useful properties of Gaussians

Perhaps the most useful property of the Gaussian family is that it's *invariant under affine transformation*:

- $X \sim \mathcal{N}(\mu, \sigma^2)$, then for any $a, b \in \mathbb{R}$,

$$a \cdot X + b \sim \mathcal{N}(a \cdot \mu + b, a^2 \sigma^2)$$

- **Standardization** (a.k.a Normalization/Z-score): If $X \sim \mathcal{N}(\mu, \sigma^2)$, then

$$Z = \frac{X - \mu}{\sigma} \sim \mathcal{N}(0, 1)$$

Useful to compute probabilities from CDF of $Z \sim \mathcal{N}(0, 1)$:

$$\mathbb{P}(u \leq X \leq v) = \mathbb{P}\left(\frac{u - \mu}{\sigma} \leq Z \leq \frac{v - \mu}{\sigma}\right)$$

- **symmetry**: If $X \sim \mathcal{N}(\mathbf{0}, \sigma^2)$ then $-X \sim \mathcal{N}(\mathbf{0}, \sigma^2)$: If $x > 0$

$$\mathbb{P}(|X| > x) = 2\mathbb{P}(X > x)$$

Gaussian probability tables



Z	Second decimal place of Z									
	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.5000	0.5040	0.5080	0.5120	0.5160	0.5199	0.5239	0.5279	0.5319	0.5359
0.1	0.5398	0.5438	0.5478	0.5517	0.5557	0.5596	0.5636	0.5675	0.5714	0.5753
0.2	0.5793	0.5832	0.5871	0.5910	0.5948	0.5987	0.6026	0.6064	0.6103	0.6141
0.3	0.6179	0.6217	0.6255	0.6293	0.6331	0.6368	0.6406	0.6443	0.6480	0.6517
0.4	0.6554	0.6591	0.6628	0.6664	0.6700	0.6736	0.6772	0.6808	0.6844	0.6879
0.5	0.6915	0.6950	0.6985	0.7019	0.7054	0.7088	0.7123	0.7157	0.7190	0.7224
0.6	0.7257	0.7291	0.7324	0.7357	0.7389	0.7422	0.7454	0.7486	0.7517	0.7549
0.7	0.7580	0.7611	0.7642	0.7673	0.7704	0.7734	0.7764	0.7794	0.7823	0.7852
0.8	0.7881	0.7910	0.7939	0.7967	0.7995	0.8023	0.8051	0.8078	0.8106	0.8133
0.9	0.8159	0.8186	0.8212	0.8238	0.8264	0.8289	0.8315	0.8340	0.8365	0.8389
1.0	0.8413	0.8438	0.8461	0.8485	0.8508	0.8531	0.8554	0.8577	0.8599	0.8621
1.1	0.8643	0.8665	0.8686	0.8708	0.8729	0.8749	0.8770	0.8790	0.8810	0.8830
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

Examples

Assume that $Z \sim \mathcal{N}(0, 1)$ and compute

- ▶ $\mathbb{P}(Z \leq 1)$
- ▶ $\mathbb{P}(Z \geq -1)$
- ▶ $\mathbb{P}(|Z| > 1)$

Assume that the score distribution for a final exam is approximately $X \sim \mathcal{N}(85, 4)$, compute

- ▶ $\mathbb{P}(X > 90)$
- ▶ $\mathbb{P}(80 < X < 90)$

More complicated: what is x such that $\mathbb{P}(X < x) = 90\%$ (85-th percentile?). For that we need to read the table backwards.

Quantiles

Definition

Let α in $(0, 1)$. The quantile of order $1 - \alpha$ of a random variable X is the number q_α such that

$$\mathbb{P}(X \leq q_\alpha) = 1 - \alpha$$

Let F denote the CDF of X :

- ▶ $F(q_\alpha) = 1 - \alpha$
- ▶ If F is invertible, then $q_\alpha = F^{-1}(1 - \alpha)$
- ▶ $\mathbb{P}(X > q_\alpha) = \alpha$
- ▶ If $X = Z \sim \mathcal{N}(0, 1)$: $\mathbb{P}(|Z| > q_{\alpha/2}) = \alpha$

Some important quantiles of the $Z \sim \mathcal{N}(0, 1)$ are:

α	2.5%	5%	10%
q_α	1.96	1.65	1.28

We get that $\mathbb{P}(|Z| > 1.96) = 5\%$

Three types of convergence

- $(T_n)_{n \geq 1}$ is a sequence of random variables
- T is a random variable (T may be deterministic).

- Almost surely (a.s.) convergence:

$$T_n \xrightarrow[n \rightarrow \infty]{\text{a.s.}} T \quad \text{iff} \quad \mathbb{P} \left[\left\{ \omega : T_n(\omega) \xrightarrow{n \rightarrow \infty} T(\omega) \right\} \right] = 1.$$

- Convergence in probability:

$$T_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}} T \quad \text{iff} \quad \mathbb{P} [|T_n - T| \geq \varepsilon] \xrightarrow{n \rightarrow \infty} 0, \quad \forall \varepsilon > 0.$$

- Convergence in distribution:

$$T_n \xrightarrow[n \rightarrow \infty]{(d)} T \quad \text{iff} \quad \mathbb{E}[f(T_n)] \xrightarrow{n \rightarrow \infty} \mathbb{E}[f(T)]$$

for all continuous and bounded function f .

Properties

- ▶ If $(T_n)_{n \geq 1}$ converges a.s., then it also converges in probability, and the two limits are equal a.s.
- ▶ If $(T_n)_{n \geq 1}$ converges in probability, then it also converges in distribution
- ▶ **Convergence in distribution** implies convergence of probabilities if the limit has a density (e.g. Gaussian):

$$T_n \xrightarrow[n \rightarrow \infty]{(d)} T \quad \Rightarrow \quad \mathbb{P}(a \leq T_n \leq b) \xrightarrow{n \rightarrow \infty} \mathbb{P}(a \leq T \leq b)$$

Exercises

a) Is the following statement correct? "If $(T_n)_{n \geq 1}$ converges in probability, then it also converges a.s"

1. Yes
2. No

Let $\{X_1, X_2, \dots, X_n\}$ be a sequence of r.v. such that

$X_n \sim \text{Ber}(\frac{1}{n})$. Exercises b), c) and d) are about this sequence.

b) Let $0 < \epsilon < 1$, $n \geq 1$. What is the value of $P(|X_n| > \epsilon)$?

(answer: $\frac{1}{n}$)

c) Does $\{X_n\}$ converges in probability?

1. Yes
2. No

Exercises

d) Denote by X the limit of $\{X_n\}$ (if it exists) (that is, $X_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}} X$). What is the value of X ?

1. X does not exist
2. 0
3. 1
4. None of the above

e) Does $\{X_n\}$ converge in distribution?

1. Yes
2. No

f) What is the limit of the sequence $\mathbb{E}[\cos(X_n)]$ as n tends to infinity?

Addition, multiplication, division

... only for a.s. and \mathbb{P} ...

Assume

$$T_n \xrightarrow[n \rightarrow \infty]{\text{a.s./}\mathbb{P}} T \quad \text{and} \quad U_n \xrightarrow[n \rightarrow \infty]{\text{a.s./}\mathbb{P}} U$$

Then,

- ▶ $T_n + U_n \xrightarrow[n \rightarrow \infty]{\text{a.s./}\mathbb{P}} T + U,$
- ▶ $T_n U_n \xrightarrow[n \rightarrow \infty]{\text{a.s./}\mathbb{P}} TU,$
- ▶ If in addition, $U \neq 0$ a.s., then $\frac{T_n}{U_n} \xrightarrow[n \rightarrow \infty]{\text{a.s./}\mathbb{P}} \frac{T}{U}.$



In general, these rules **do not** apply to convergence (d).

Slutsky's theorem

Some partial results exist for convergence in distribution on the form of *Slutsky's theorem*.

Let $(X_n), (Y_n)$ be two sequences of r.v., such that:

$$\text{(i)} \quad T_n \xrightarrow[n \rightarrow \infty]{(d)} T \quad \text{and} \quad \text{(ii)} \quad U_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}} u$$

where T is a r.v. and u is a given real number (deterministic limit: $\mathbb{P}(U = u) = 1$). Then,

- $T_n + U_n \xrightarrow[n \rightarrow \infty]{(d)} T + u,$
- $T_n U_n \xrightarrow[n \rightarrow \infty]{(d)} Tu,$
- If in addition, $u \neq 0$, then $\frac{T_n}{U_n} \xrightarrow[n \rightarrow \infty]{(d)} \frac{T}{u}.$

...

Taking functions

Continuous functions (for all three types) . If f is a continuous function:

$$T_n \xrightarrow[n \rightarrow \infty]{\text{a.s./P/(d)}} T \Rightarrow f(T_n) \xrightarrow[n \rightarrow \infty]{\text{a.s./P/(d)}} f(T).$$

Example: Recall that by LLN, $\bar{R}_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}, \text{ a.s.}} p$. Therefore

$$f(\bar{R}_n) \xrightarrow[n \rightarrow \infty]{\mathbb{P}, \text{ a.s.}} f(p) \quad \text{for any continuous } f$$

(Only need f to be continuous around p : $f(x)=1/x$ works if $p > 0$)

We also have by CLT: $\sqrt{n} \frac{\bar{R}_n - p}{\sqrt{p(1-p)}} \xrightarrow[n \rightarrow \infty]{(d)} Z$, $Z \sim \mathcal{N}(0, 1)$. So

$$f \left(\sqrt{n} \frac{\bar{R}_n - p}{p(1-p)} \right) \xrightarrow[n \rightarrow \infty]{(d)} f(Y) \quad Y \sim \mathcal{N}(0, p(1-p))$$

⚠ not the limit of $\sqrt{n}[f(\bar{R}_n) - f(p)]$!! (see Delta-method)

Recap

- ▶ Averages of random variables occur naturally in statistics
- ▶ We make modeling assumptions to apply probability results
- ▶ For large sample size they are *consistent* (LLN) and we know their distribution (CLT)
- ▶ CLT gives the (weakest) convergence in distribution but is enough to compute probabilities
- ▶ We use standardization and Gaussian tables to compute probabilities and quantiles
- ▶ We can make operations (addition, multiplication, continuous functions) on sequences of random variables