

RCDS Statistics 1 - Imperial College London

Introduction to sampling and hypothesis testing

Jesús Martínez Elizari

December 2024

Random sampling and hypothesis testing

Chapter 1. Random events and probability theory

- i) Definition of random events and probability.
- ii) Discrete probability ; "mass" distributions.
- iii) Continuous probability ; "density" distributions.

Chapter 2. Expected values and the central limit theorem

- i) Mean, variance as expected values.
- ii) Parameter estimation , prediction vs inference .
- iii) the law of large numbers (LLN).
- iv) the central limit theorem (CLT).

Chapter 3 . Hypothesis testing

- i) Prediction vs inference . Formulate hypothesis.
- ii) General approach . Significance and p-values.
- iii) Real example
 - { Compare mean to expected value (1 sample T)
 - { Compare two means (2 sample T)

1.1 Random events and probability

* Random events ("stochastic"), from στοχαστικός "to guess")

Something whose output we don't know

* Probability: number $\in [0, 1]$ quantifying certainty / "surprise".

Example: tossing coins (H, T)

$P(H) = 0 \rightarrow$ certain I will never get H

$P(H) = 1 \rightarrow$ always get H

$0 < P(H) < 1 \rightarrow$ level of uncertainty / "surprise"

* Mutuality: The sum of probabilities for all possible outcomes x_i must add up to 1.

$$\sum_{\forall x_i} P(x_i) = 1$$

Example: tossing coins

$$P(H) + P(T) = \frac{1}{2} + \frac{1}{2} = 1 \quad \checkmark$$

Example: rolling dice

$$P(1) + P(2) + P(3) + P(4) + P(5) + P(6) = \frac{1}{6} + \frac{1}{6} + \dots + \frac{1}{6} = 1 \quad \checkmark$$

* A bit of history

i) Girolamo Cardano (1501 - 1576)

Probability in games and dice, Chevalier.

ii) Pierre de Fermat (1601 - 1665)

Blaise Pascal (1623 - 1662)

Foundations of mathematical probability theory.

iii) Pierre S. Laplace (1749 - 1827)

Frequentist definitions of probability.

P as fraction of favorable / possible outcomes.

iv) Andrey Kolmogorov (1903 - 1987)

Axiomatic foundations (1933)

"Foundations of theory of probability"

$P \in [0, 1]$; unitarity. (*)

1.2 Discrete probability distributions

* Discrete: number of possible outcomes is a finite number.

coins
dice
counting

i) Binomial distribution

Jacob Bernoulli; "Ars conjectandi" (1713)

"How many times X I get a specific result in n trials,
if the probability of each success is $p"$

$$\left\{ B(x; n, p) = \binom{n}{x} p^x (1-p)^{n-x} \right\}$$

x : number of successes
 n : number of trials
 p : probability each success.

Example: Probability of 5 times H tossing 10 times a coin

$$B(5; 10, \frac{1}{2}) = \binom{10}{5} \left(\frac{1}{2}\right)^5 \left(1-\frac{1}{2}\right)^{10-5} = 0.246$$

⊗
Exercise

Example: Probability of 3 times a 6 in 10 dice

$$B(3; 10, \frac{1}{6}) = \binom{10}{3} \left(\frac{1}{6}\right)^3 \left(1-\frac{1}{6}\right)^{10-3} = 0.155$$

⊗
Exercise

Example: Probability of passing exam (A, B, C) of 10 quest.

$$B(5; 10, \frac{1}{3}) = \binom{10}{5} \left(\frac{1}{3}\right)^5 \left(1-\frac{1}{3}\right)^{10-5} = 0.137$$

⊗
Exercise

* What if $x=n$?

Example : probability of tossing 10 Heads out of 10

$$\begin{aligned} \mathcal{B}(10; 10, \frac{1}{2}) &= \binom{10}{10} \left(\frac{1}{2}\right)^{10} \left(1-\frac{1}{2}\right)^{0} \\ &= 1 \cdot \left(\frac{1}{2}\right)^{10} \cdot 1 = \left(\frac{1}{2}\right)^{10} \quad \checkmark \end{aligned}$$

* General : recover the individual probability, raised to the number of attempts.

$$\mathcal{B}(n; n, p) = \underbrace{\binom{n}{n}}_1 p^n \underbrace{(1-p)}_1^0 = \underline{\underline{p^n}}$$

iii) Poisson distribution

Simeon D. Poisson; "Research on probability" (1837)

"How many times I observe an event in an interval, know λ"

$$\left\{ \begin{array}{l} P(x; \lambda) = \frac{e^{-\lambda} \lambda^x}{x!} \\ \quad \quad \quad \left\{ \begin{array}{l} x: \text{number of times I observe} \\ \lambda: \text{observed average.} \end{array} \right. \end{array} \right.$$

Example: Probability of observing 3 new patients, with $\lambda = 5$.

$$P(3; 5) = \frac{e^{-5} 5^3}{3!} = 0.140$$

(*)

Exercise

Example: Probability of observing 5 or less patients in same λ.

$$\begin{aligned} P(x \leq 5; 5) &= P(1) + P(2) + \dots + P(5) \\ &= \sum_{x=1}^5 P(x; \lambda=5) = 0.616 \end{aligned}$$

(*)

Exercise

Cumulative distribution function cdf

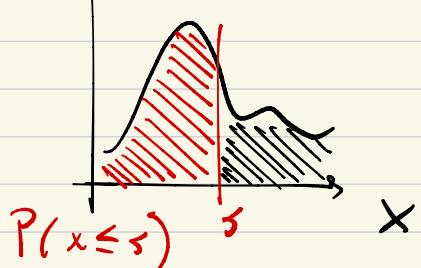
Example: Probability of more than 5

$$P(x > 5; 5) = 1 - \text{cdf}(5) = 0.394$$

\uparrow
Invert

(*)

Exercise



1.3 Continuous distributions

* Continuous: amount of possible outcomes is infinite / uncountable

Case 5: 2 outcomes (H, T) $\rightarrow P = \frac{1}{2}$ Discrete cases "Mass dist."

Dice: 6 outcomes ($1, 2, 3, 4, 5, 6$) $\rightarrow P = \frac{1}{6}$ Frequentist definition /

Continuous variable ($T, h, \text{conc.}$): ∞ outcomes $\rightarrow P = \frac{1}{\infty} = 0$ WTF

↳ Frequentist approach does not work

Need a new mathematical object "Density distribution"

i) Discrete case

Probability $P \in [0, 1]$

Unitarity $\sum_{x_i} P(x_i) = 1$

ii) Continuous case

Density $f(x)$

Unitarity $\int_{-\infty}^{\infty} dx f(x) = 1$

* Gaussian

i) Gaussian distribution

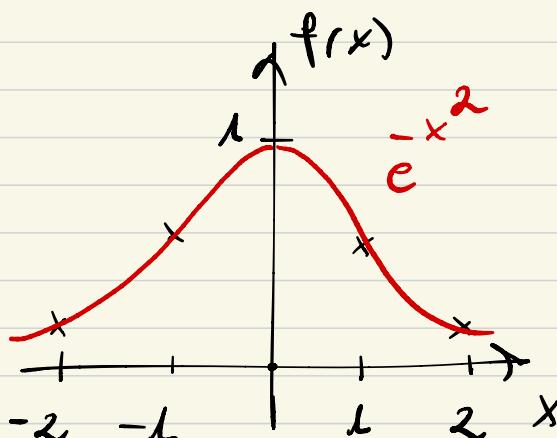
Carl Friedrich Gauss; "theory of movement of celestial bodies" (1809)

"Random variable x centered at some μ and spreaded σ "

$$\left\{ f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} \right. \quad \left. \begin{array}{l} \mu: \text{mean value} \\ \sigma: \text{standard deviation} \end{array} \right.$$

* Consider simplest case $\mu=0$; $\sigma=L$; ignore normalization factor

$$f(x) = e^{-\frac{x^2}{2}} ; \text{ plot same values}$$



$$x = 0; f(0) = e^{-0^2/2} = 1$$

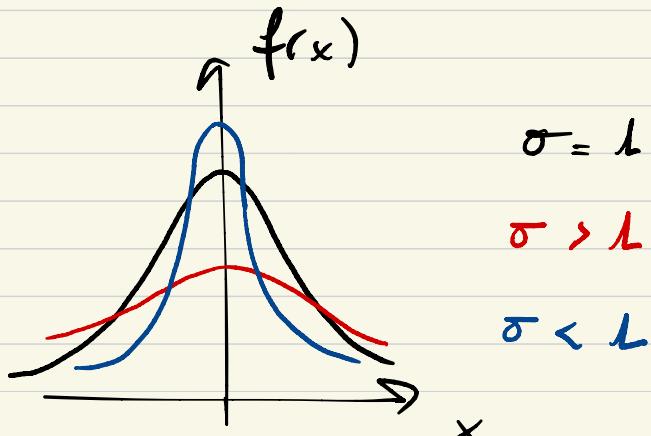
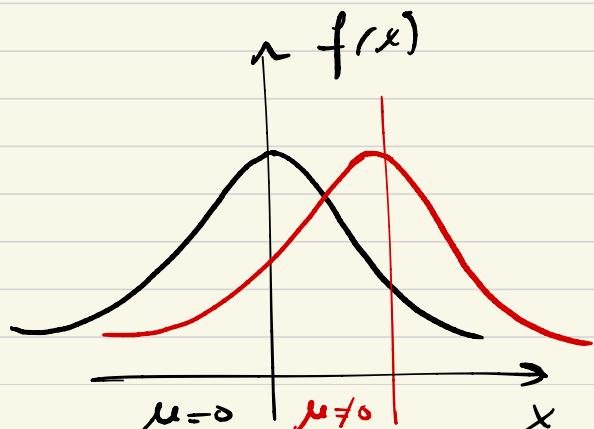
$$x = L; f(L) = e^{-L^2/2} = \frac{1}{e} = 0.37$$

$$x = -L; f(-L) = e^{-(-L)^2/2} = \frac{1}{e} = 0.37$$

$$x = \pm 2; f(\pm 2) = e^{-4/2} = \frac{1}{e^2} = 0.018$$

* Switch back on μ and σ $\left\{ \begin{array}{l} \mu \text{ will displace the central value} \\ \sigma \text{ will drive sharper/wider shape} \end{array} \right.$

$$f(x) = e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$



* Add normalization factor

i) Remember discrete probability : Normality $\sum_{x_i} P(x_i) = 1$

ii) Continuous case. Density distributions $\int_{-\infty}^{+\infty} dx f(x) = 1$

$$\int_{-\infty}^{\infty} dx e^{-\frac{1}{2\sigma^2}(x-\mu)^2} = \sigma \sqrt{2\pi}$$

*
Exercise

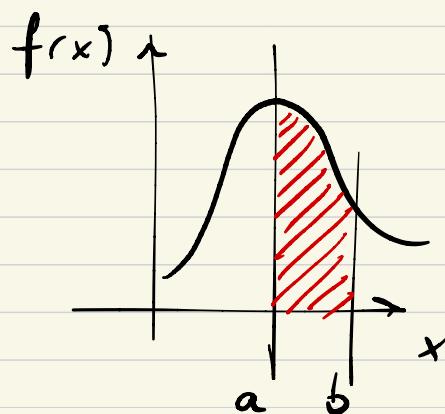
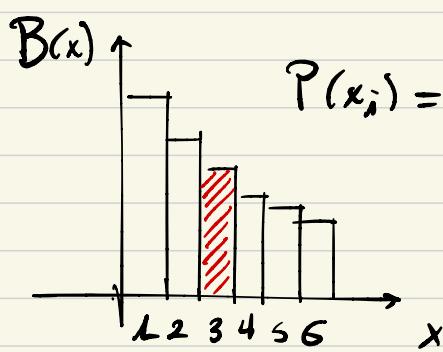
* for $f(x)$ to behave as a probability, we add a normalization factor

$$\left\{ f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} \right\} \left\{ \begin{array}{l} \text{Gaussian distribution ; } \int_{-\infty}^{\infty} dx f(x) = 1 \\ (\text{Normal distribution}) \end{array} \right.$$

* Particular case $\mu=0; \sigma=1$

$$\left\{ f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} \right\} \left\{ \begin{array}{l} \text{Standard Normal distribution } N(0,1) \\ (\text{Z distribution ; } Z(0,1)) \end{array} \right.$$

* In continuous distributions, we can only compute probability in a given range. Never a single value.



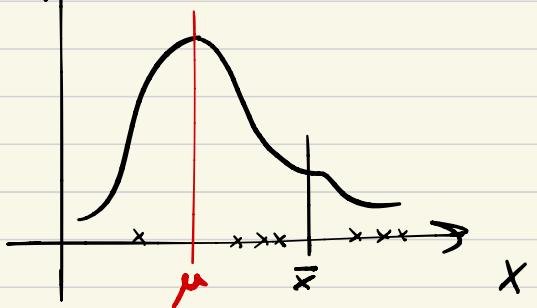
$$P(a < x < b) = \text{cdf}(b) - \text{cdf}(a)$$

2.1 Parameter estimation

We have a set of observations $X = \{x_1, x_2, \dots, x_n\}$

We assume they came from a certain distribution with same true μ and same true σ^2

$f(x)$: $B(x)$, $P(x)$, Gaussian(x) ...



I don't know the true μ, σ^2

I only have a discrete set of observations. I can build some "informative quantities" from data(x)

* Observed coverage ("sample mean")

$$\bar{x} = \frac{1}{n} (x_1 + x_2 + \dots + x_n) = \frac{1}{n} \overline{\sum_{i=1}^n x_i}$$

* Observed variance ("sample variance")

$$s^2 = \frac{1}{n-1} \left\{ (x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2 \right\} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

\hookrightarrow Bessel's factor (*)

$\left\{ \begin{array}{l} \text{Sample mean } \bar{x} \text{ is an estimator of } \mu \\ \text{Sample variance } s^2 \text{ is an estimator of } \sigma^2 \end{array} \right\}$

2.2 The law of large numbers.

Bernoulli, Chebyshev, Borel (1713)

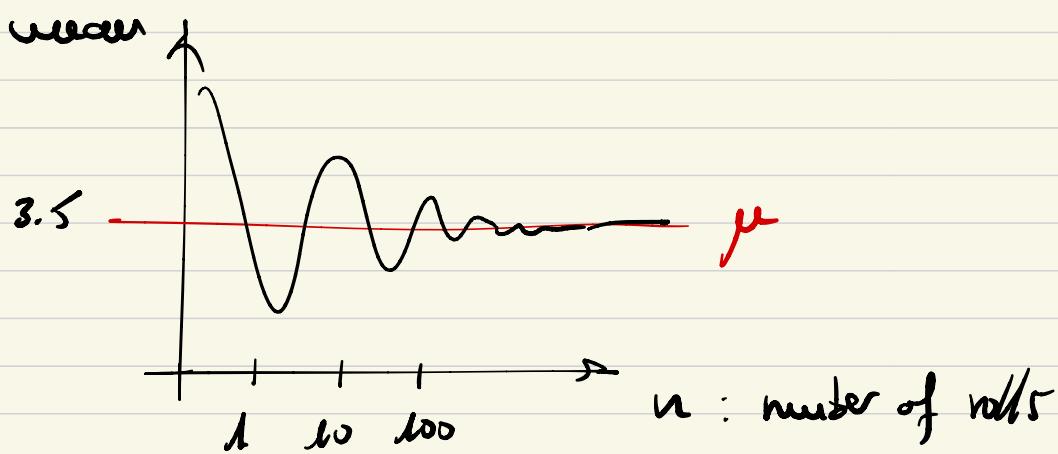
"Let $\{x_1, x_2, \dots, x_n\}$ independent and identically distributed (iid) random variables, with same true mean μ (expected), the sample mean \bar{x} converges to μ as n increases"

$$\left\{ \lim_{n \rightarrow \infty} \bar{x} = \mu ; \text{ being } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \right\}$$

Example : rolling dice.

i) True / expected mean $\mu = \frac{1}{6}(1+2+3+4+5+6) = 3.5$

ii) As we increase n rolls, sample mean \bar{x}
will converge to the true / expected value μ .



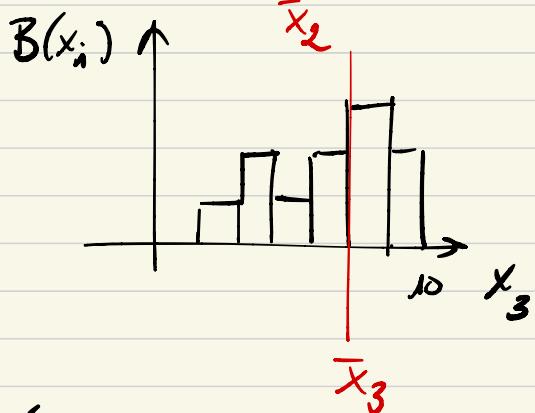
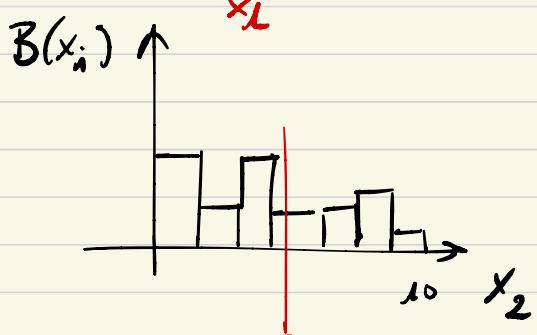
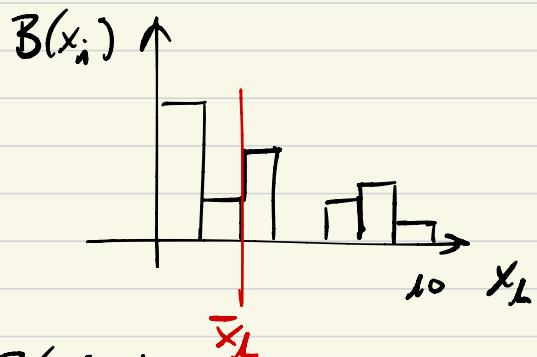
2.3) The Central Limit Theorem

A. de Moivre, Laplace, Gauss (1733)

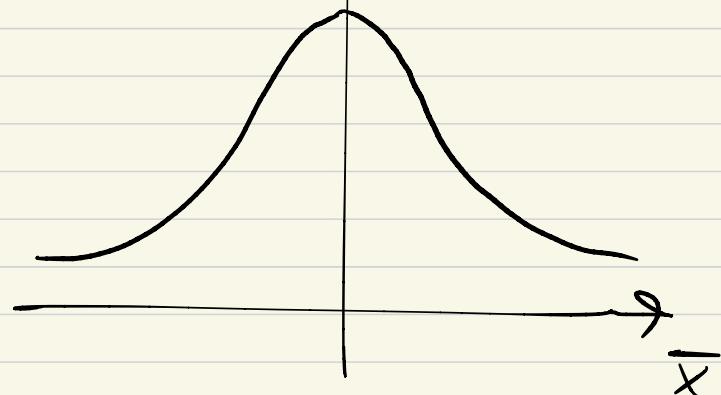
"For a set $\{X_1, \dots, X_n\}$ of n independent and identically distributed (iid) random variables, the distribution of sample means tends to gaussian as n increases, regardless of how are X_i distributed themselves"

i) Consider 10 rolls of a dice. How many times I observe a 6.

Binomial distributed $B(x; 10, 1/6)$. Repeat N times.



* As the sample size n increases,
the distribution of means \bar{X}_i
tends to a gaussian distribution
 $f(\bar{x}) \sim \text{Gaussian}(\bar{x})$



{ * Sample mean
 $\bar{X} \rightarrow N(\mu, \sigma^2/n)$
* Standard error (SE)
 $SE = \frac{\sigma}{\sqrt{n}}$; as increase n , \sqrt{SE}

ii) Build new random variable

$$\bar{X} = \{\bar{X}_1, \bar{X}_2, \dots, \bar{X}_N\}$$

2.2 Confidence Intervals.

* See how many standard deviations an observation is away from the mean. Z-score (*)

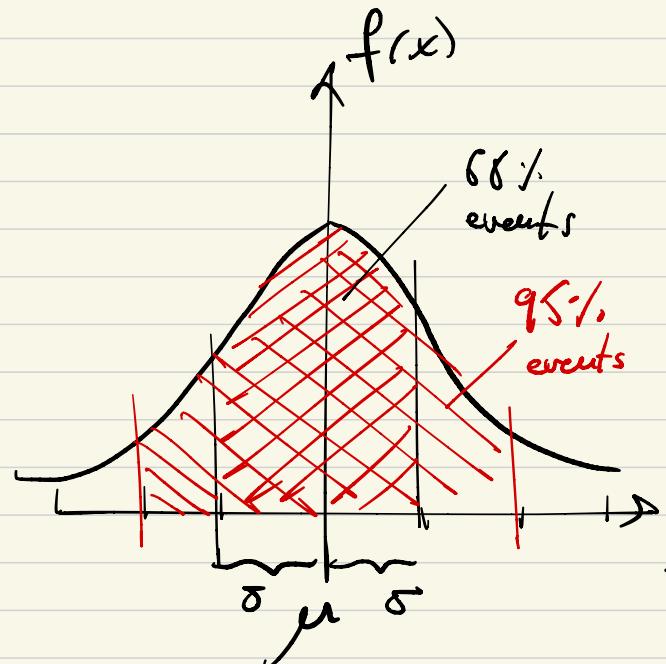
$$\left\{ z = \frac{x-\mu}{\sigma} ; \text{ if dealing with samples } z = \frac{x-\bar{x}}{s} \right\}$$

* When comparing sample means themselves, use δE instead of σ of individual observations.

$$\left\{ z = \frac{x-\mu}{\delta/\sqrt{n}} ; \text{ dealing with samples } z = \frac{x-\bar{x}}{\delta/\sqrt{n}} \right\}$$

* Confidence interval $P_{\text{critical}} = (x+z, x-z)$

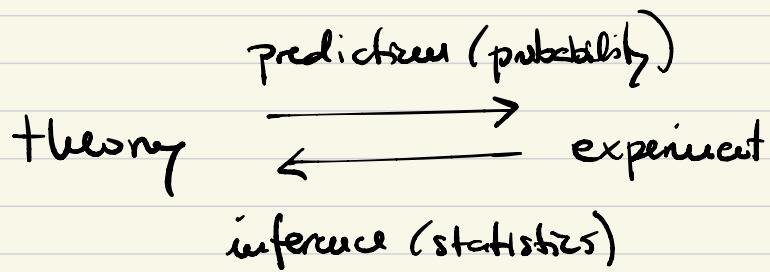
i) Consider a Gaussian distribution $f(x; \mu, \sigma)$



- * 68% of events are contained in the $(\mu-\sigma, \mu+\sigma)$ interval. 1 σ conf. int.
- * 95% of events are contained in the $(\mu-2\sigma, \mu+2\sigma)$ interval. 2 σ conf. int.
- * 99% of events are contained in the $(\mu-3\sigma, \mu+3\sigma)$ interval. 3 σ conf. int.

③ Hypothesis testing

* Prediction vs inference



i) Check if an observation is compatible with a given hypothesis H_0

* Statistic test : a function / number computed out of our data
(t-test, χ^2 test, Fisher test, Wald test, ...)

* p-value : probability of obtaining a result at least as extreme as the one actually observed, under the assumption that our null hypothesis H_0 was correct.

(How likely / unlikely was to observe this result)

- { i) \uparrow p-value : it was likely to obtain this result. **Accept H_0 .**
ii) \downarrow p-value : // unlikely // . **Reject H_0 .**

ii) General approach

- { 1. Formulate null hypothesis H_0 and significance level α .
2. Make measurement / observation.
3. Compute statistic test
4. Compute p-value
5. If $p\text{-value} < \alpha$, reject H_0 . Otherwise, accept H_0 .

* A bit of history

i) Thomas Bayes (1701-1761)

Bayes rule; probability as "updated degree of belief"

ii) Karl Pearson (1857-1936)

Correlation, regression, P-values

iii) Ronald Fisher (1890-1962)

Null hypothesis, analysis of variance

iv) Jerzy Neyman (1890-1981)

Egon Pearson (1895-1980)

False discovery rate, type I/II errors.

3.1 One-sided t-test.

* Compare mean of observed die rolls to the expected mean of a fair die.

* For a series of observations $X = \{x_1, x_2, \dots, x_n\}$

$$\left\{ t = \frac{\bar{x} - \mu}{s/\sqrt{n}} \right\}$$

\bar{x} : observed average ("sample mean")

μ : expected average under H_0 ("population mean")

s : observed standard deviation ("sample standard dev.")

n : length of collected data ("sample size")

Example: Check if a die is biased. (1)

* You suspect your opponent is using a loaded die

i) Formulate null hypothesis $H_0: P = \frac{1}{6}$

ii) alternative // $H_1: P > \frac{1}{6}$

Choose significance level $\alpha = 0.01$

iii) Collect data. Roll die 15 times.



10 times out of 100
we get a 6. Suspicious!

iv) Perform 1-sample t-test.

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{5 - 3.5}{1.6 / \sqrt{15}} \approx 3.6$$

μ : expected average is a fair die
(all faces with same probability)

$$\mu = \frac{1}{6} (1 + 2 + 3 + 4 + 5 + 6) = 3.5$$

$$\bar{x} = 5 ; s = 1.6 ; n = 15$$

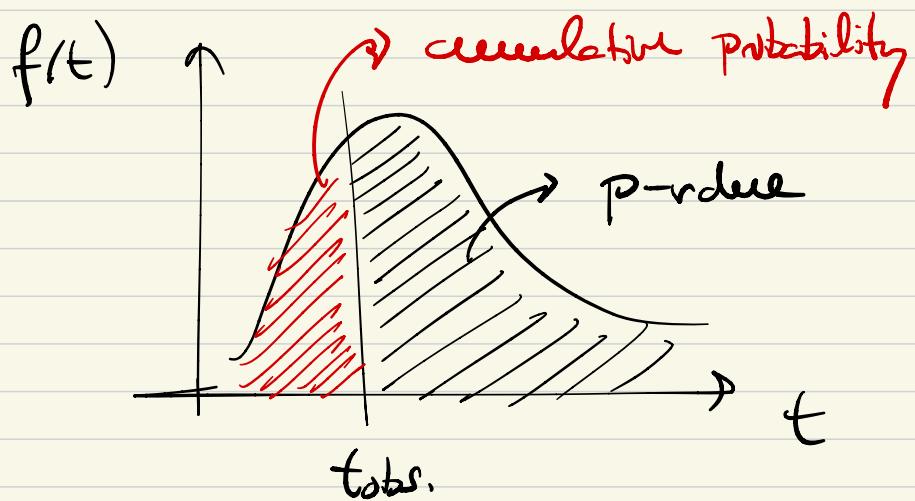
(*)

(*)

iv) P-value. Probability of obtaining a value at least
as extreme as the one we obtained for our t-variable,
given H_0 was true.

If H_0 is true, then t variable follows

a "Student's t distribution"



$$p\text{-value} = P(t \geq t_{\text{obs}}) = 1 - \text{cdf}(t_{\text{obs}}) \quad (*)$$

$\left. \begin{array}{c} \text{Student's t tables} \\ \text{Computer simulations} \end{array} \right\}$

Example: Check if a die is biased. (11)

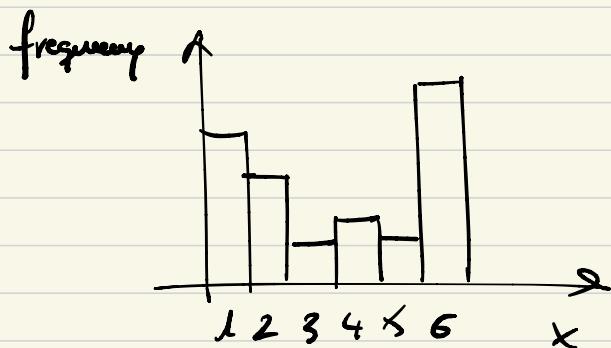
* You suspect your opponent is using a loaded die

i) Formulate null hypothesis $H_0: p = \frac{1}{6}$

ii) alternative $H_1: p > \frac{1}{6}$

Choose significance level $\alpha = 0.01$

iii) Collect data. Roll die 100 times



43 times out of 100
we get a 6.

iv) Compute probability of 43 or less 6's, given H_0 .

Given by Binomial distribution with $p = 1/6$

$$B(43; 100, \frac{1}{6}) \sim L_{\text{calf}}$$



Exercise

v) Compute p-value. Probability of obtaining at least a value as extreme as 43.

$$\text{p-value} = 1 - B_{\text{calf}}(43; 100, \frac{1}{6}) = 1^{-10}$$

↪ p-value $< \alpha$. → evidence to reject H_0 .

Example: check if coin is fair

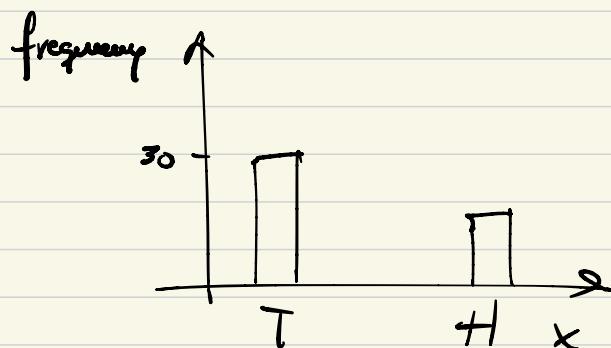
* You suspect your opponent is using a biased coin

i) Formulate null hypothesis $H_0: p = \frac{1}{2}$

ii) alternative " $H_1: p \neq \frac{1}{2}$

Choose significance level $\alpha = 0.05$

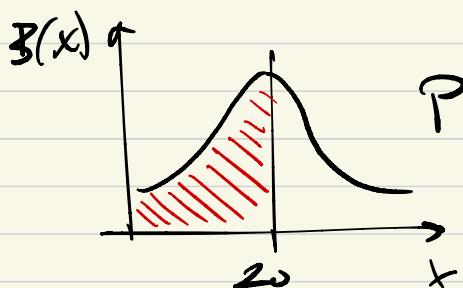
iii) Collect data. Roll die 100 times



20 heads in 50
flips of a coin

iv) Compute probability of 20 or less H, given H_0 .

Given by Binomial distribution with $p = 1/2$



$$P(x=20) = B_{\text{cdf}}(20; 50, \frac{1}{2}) = 0.1$$



v) Compute p-value.

Now H_1 is just $p \neq \frac{1}{2}$, either larger or smaller.

$p\text{-value} = 2 \cdot \text{cdf}(20) = 0.2 < \alpha$ \nexists evidence
to reject H_0 .

3.3 Comparing means. t-test

* Compare if two samples have significantly different means

Given two sets of observations x_1 and x_2 , of sample size

n_1 and n_2 , and sample means \bar{x}_1 and \bar{x}_2 .

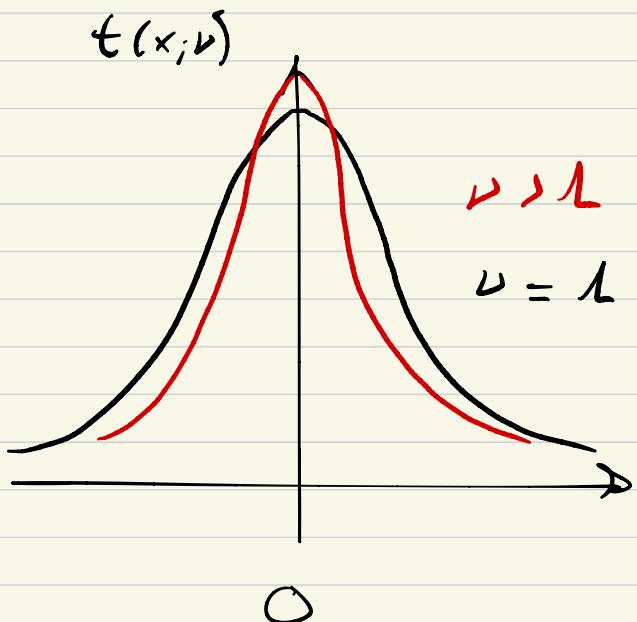
* Compute t variable (t "statistic")

$$t = \frac{\bar{x}_1 - \bar{x}_2}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \quad \left\{ \begin{array}{l} \\ \end{array} \right.$$

being s_p the pooled / mixed standard deviation

$$s_p = \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}}$$

* Under null hypothesis of equal means, the t statistic follows a Student's t distribution with (n_1+n_2-2) dof.



i) For $v=1$; $t(x; v)$ is referred

to as the Cauchy distribution

$$\left\{ f(x, x_0, \gamma) = \frac{1}{\pi} \frac{\gamma}{(x-x_0)^2 + \gamma^2} \right\}$$

ii) For v large; $t(x; v)$

\times tends to a normal distribution

* A bit of history

1876. t-distribution derived first by Helmert and Bünnell

1895. General form from Pearson, as Pearson type IV distribution

1908. William Sealy Gosset publishes under name "Student"

1925. Popularized by Fisher as "Student's distribution"

Example : Comparing means of 2 populations

* Check if the mean of birth weights are significantly different between smokers and no smokers.

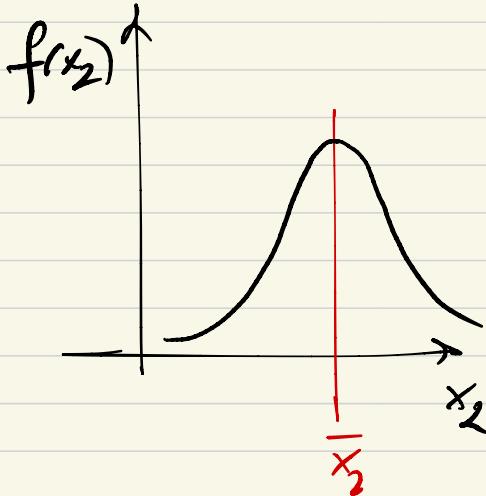
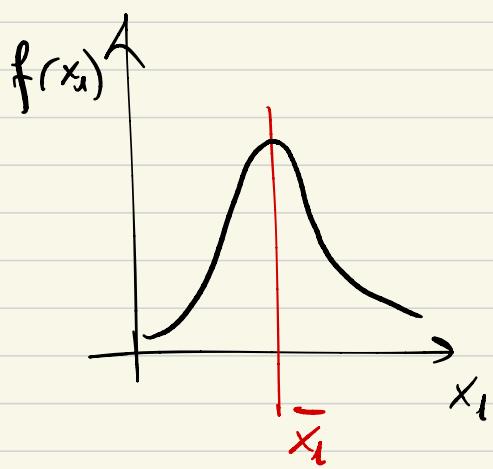
i) Formulate null hypothesis

$$H_0: \text{No difference b/w the means. } d = \bar{x}_1 - \bar{x}_2 = 0$$

$$H_1: \exists \text{ } d \neq 0$$

Choose significance level α

ii) Collect data . Smokers and non smokers .



n_1 non smokers

Sample mean \bar{x}_1

n_2 heavy smokers

Sample mean \bar{x}_2

standard deviation s_1

standard deviation s_2

iii) Compute t statistic iv) Compute d.o.f

$$t = \frac{\bar{x}_1 - \bar{x}_2}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

$$D = n_1 + n_2 - 2$$

v) Check if t falls inside the
95% confidence interval.

Example: Comparing means of two populations.

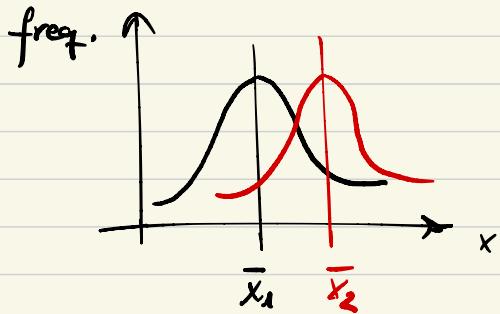
i) Formulate H_0, H_1

H_0 : both samples have the same mean $\bar{x}_1 = \bar{x}_2 = \mu$

H_1 : the samples have different means $\bar{x}_1 \neq \bar{x}_2$

ii) Collect data

Example: luminosity of 2 different types of stars.



$\{x_1\}$: Type IV stars ("Supergiant")

$\{x_2\}$: type V stars ("hypergiant")

iii) Compute t "statistic" (t variable)

$$\left\{ t = \frac{\bar{x}_1 - \bar{x}_2}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \right\} \quad \text{if } \bar{x}_1 = \bar{x}_2 \Rightarrow t=0 \quad \checkmark$$

$$S_p = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2} \quad \text{"pooled" standard deviation.}$$

iv) Compute p-value.

"What was the probability of obtaining a value at least as extreme as the one we got for t."

* If $\bar{x}_1 = \bar{x}_2$, the t variable follows a Student's t distribution with n_1+n_2-2 degrees of freedom. (*)