

Introduction to statistics

TODO subtitle

Caio Volpato (caioau)

caioau.keybase.pub → caioauheyuuiavlc.onion

210B C5A4 14FD 9274 6B6A 250E **EFF5 B2E1 80F2 94CE**

All Copylefts are beautiful: licensed under [CC BY-SA 4.0](https://creativecommons.org/licenses/by-sa/4.0/)

Summary:

- Probability concepts
- Discrete distributions.
- Continuous distributions.
- Calculations on the Normal distribution.
- Convergence
- Inference

Motivation



Figure 1: Dados apontam ... (data shows ...)

Basic concepts of probability:

Sample space Ω

It's the set of all the possible outcomes of a experiment, denoted by S or Ω

Event

It's a subset of the sample space.

Properties needed:

Given a sample space Ω the class of events denoted by \mathcal{A} need to satisfy the following properties:

- 1 $\emptyset \in \mathcal{A}$;
- 2 $a_1, a_2, \dots \in \mathcal{A} \implies \bigcup_{i=1}^{\infty} a_i \in \mathcal{A}$
- 3 if $a \in \mathcal{A} \implies a^c \in \mathcal{A}$

Basic concepts of probability:

Probability (Definition):

Given a experiment with a sample space Ω and a class of events \mathcal{A} , the probability denoted by \mathbb{P} is a function which has \mathcal{A} as domain and associate a numerical value between $[0, 1]$ as image.

Probability properties:

- ① $\mathbb{P}(\Omega) = 1$ and $\mathbb{P}(\emptyset) = 0$
- ② $0 \leq \mathbb{P}(A) \leq 1$, for every event A
- ③ For any sequence of mutually exclusive events A_1, A_2, \dots that's events that $A_i \cap A_j$ when $i \neq j$ we have that:

$$\mathbb{P} \left(\bigcup_{i=1}^{\infty} A_i \right) = \sum_{i=1}^{\infty} \mathbb{P}(A_i)$$

Basic concepts of probability:

Event independence:

Two events are independent when the occurrence of the first does not affect the probability of occurrence of the second.

Two events A and B are independent if:

$$\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B)$$

Conditional Events:

The probability of a event A to occur given that the event B occurred is:

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}$$

Basic concepts of probability:

Bayes theorem:

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(B|A)\mathbb{P}(A)}{\mathbb{P}(B)}$$

General case:

$$\mathbb{P}(A_i|B) = \frac{\mathbb{P}(B|A_i)\mathbb{P}(A_i)}{\sum_{j=1}^n \mathbb{P}(B|A_j)\mathbb{P}(A_j)}$$

Bayes example (from Veritasium):

You are feeling sick, so you go to the doctor, there you run a battery of tests. After getting the results you tested positive for a rare disease (affects 0.1% of the population), the test will correctly identify that you have it 99% of the times.

What's the chances that you actually have the disease? 99%?

Bayes example Solution

Let's denote the event of you have the disease H (stands for hypothesis, the prior) and the test been positive denoted by E (stands for evidence), so we have: $\mathbb{P}(H) = 0.001$ and $\mathbb{P}(E|H) = 0.99$

$$\begin{aligned}\mathbb{P}(H|E) &= \frac{\mathbb{P}(E|H)\mathbb{P}(H)}{\mathbb{P}(E)} = \frac{\mathbb{P}(E|H)\mathbb{P}(H)}{\mathbb{P}(H)\mathbb{P}(E|H) + \mathbb{P}(H^C)\mathbb{P}(E|H^C)} = \\ &= \frac{0.99 \cdot 0.001}{0.001 \cdot 0.99 + 0.999 \cdot 0.01} = 0.09 = 9\%\end{aligned}$$

What if you test again and it's also positive? You can just take the posterior probability we just calculated and use as a prior:

$$= \frac{0.99 \cdot 0.09}{0.09 \cdot 0.99 + 0.91 \cdot 0.01} = 0.907 \approx 91\%$$

- Awesome video: [A visual guide to Bayesian thinking](#)

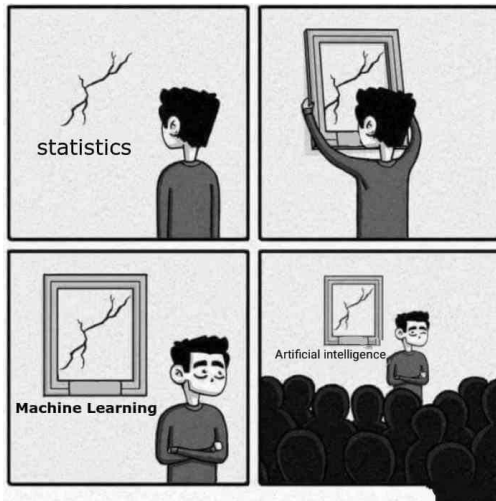


Figure 2: Credits: [sandserifcomics](#)

Random Variable (RV)

Consider an experiment with a sample space Ω associated with it. A function that maps each element $\omega \in \Omega$ to a Real number such that $[w/leq X]$ it's called random variable (RV) ($X : \Omega \rightarrow \mathbb{R}$)

- Example: Imagine an experiment that consists of 3 consecutive fair coin tosses, so the sample space of this experiment is:
 $S = \{(H,H,H), (H,H,T), \dots, (T,T,T)\}$. Now we want to create a random variable X that counts the number of heads in each outcome, so $X((H,H,H)) = 3$ and $X((H,H,T)) = 2$.

Random Variable:

Probability Mass Function (PMF):

$$f_X(x) = \mathbb{P}[X = x] = \mathbb{P}[\{\omega \in \Omega : X(\omega) = x\}]$$

Probability Density Function (PDF)

$$\mathbb{P}[a \leq X \leq b] = \int_a^b f(x) dx$$

Cumulative Distribution Function (CDF)

$$F_X(x) = \mathbb{P}[X \leq x]$$

Expectation:

- Discrete : $\mathbb{E}[X] = \sum x\mathbb{P}(X = x)$
- Continuous: $\mathbb{E}[X] = \int_{-\infty}^{\infty} xf(x)dx$

Variance:

$$\mathbb{V}[X] = \sigma_X^2 = \mathbb{E}[X^2] - \mathbb{E}^2[X]$$

Sample mean:

$$\overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

Sample variance and standard deviation:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \overline{X})^2$$

Standard deviation = s

Discrete distributions

Bernoulli:

Consider an experiment with two possible outcomes: success ($X=1$, with probability p) or failure ($X=0$), this random variable is called Bernoulli, the PMF is:

$$\mathbb{P}(X = k) = p^k(1 - p)^{1-k}$$

Binomial:

Now consider a Bernoulli experiment conducted n times, let X be the random variable that represents the number of successes, X is called Binomial, the PMF is:

$$\mathbb{P}(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

Discrete distributions

Geometric:

Again consider a Bernoulli experiment conducted n times, but the first $n-1$ are failures and the last n th is a success. Let X be number of tries, which is called Geometric, the PMF is:

$$\mathbb{P}(X = k) = (1 - p)^k p$$

- A important property is that Geometric distribution is **memoryless** (TODO definir)

Poisson:

A random variable which value can assume $0, 1, 2, \dots$ is called Poisson with $\lambda > 0$ parameter if your PMF is:

$$\mathbb{P}(X = k) = \frac{e^{-\lambda} \lambda^k}{k!}$$

Continuous distributions

- Normal
- Exponential memoryless, sunk costs
(<https://youarenotsosmart.com/2011/03/25/the-sunk-cost-fallacy/>)
- Pareto, principio de pareto (80-20)

meme know distributions

Calculations on the Normal distribution

tabela e calcular python, excel, normalizar

Assumptions on distribution choice

antes se eh discreta ou continua

descrever a natureza da normal, exp, pareto

Discretas:

Bernouli: resultado dicotomico , exemplo: moeda, homem ou mulher, sim ou nao. Voto em 2o turno

Binom: quantidade de sucesso dado um numero fixo de experimento independente; Dados 20 dispositivos independentes, depois de muitas horas, qual a prob de 15 apresentarem defeito.

Geometrica: Numero de falhas ate primeiro sucesso. Exemplo

loteria: Dado p ser 1 em 1 milhao qual a prob de ganhar depois de 3 tentativas?

Poisson: Contagem de pessoas inscritas em algum programa que desistem.

Continuas

Normal : Sem restricao de valores (pode ser positivo ou negativo).

Exemplo: Altura de criancas do mesmo sexo e idade (funciona bem pra qui quadrado).

order statistics

defs, min, max, median, q_1, q_3 , IQR, p_q ? estat robusta, boxplot

Convergence

defs, lei dos grandes numeros, teorema do limite central

Inference

metodo da maxima verosimilhanca, e grafico qxq
(`scipy.stats.probplot`) (colocar exemplo com distribuicao errada: t
student e fitar com a normal)

max veross: find the most likely parameter value, given data. That
is, given a prob description of data, find the optimum value for that
data (derivatives).

Further reading:

- Portal action (pt)
- stat cookbook
- havard youtube
(<https://www.youtube.com/playlist?list=PL2SOU6wwxB0uwwH80KT>)
- ross, barry james, meyer
- khan academy
- <http://www.randomservices.org/random/>
- divulgacao: pizza de dados, senhora toma cha, andar do bebado