DATA ANALYSIS

CERN School of Computing 2022, Krakow, Poland

Toni Šćulac Faculty of Science, University of Split, Croatia

LECTURES OUTLINE

- Introduction to Data Analysis 1)
- Probability density functions and Monte Carlo methods 2)
- 3) Parameter estimation and Confidence intervals
- 4) Hypothesis testing and p-value



INTRODUCTION TO DATA ANALYSIS



EVENT DISTRIBUTION





EVENT DISTRIBUTION



- We can not tell until we can compare to the expected distribution Is there any place on where
 - data does not agree with the expectation? Where? How significant?





EVENT DISTRIBUTION







WHAT IS DATA ANALYSIS?

"Data analysis is a process for obtaining raw data and converting it into information useful for decision-making by users. Data are collected and analyzed to answer questions, test hypotheses or disprove theories."



- of data
- A mathematical foundation for statistics is the probability theory



USABLE INFORMATION

• Data analysis uses statistics for presentation and interpretation (explanation)



DATA ANALYSIS IN THE INDUSTRY

DATA ANALYSIS



(search string₁,location₁)^{user 1} (search string₂,location₂)^{user 1}

(search string_n, location_n)^{user 1} (search string₁,location₁)^{user 2}

(search string_m,location_m)^{user 2} (search string₁,location₁)^{user 3}

(search string₁,location₁)^{user k}

 $\bullet \bullet \bullet$

Maximum Likelihood fit Significance **Hypothesis testing P-value Neural Networks**

USABLE INFORMATION











DATA ANALYSIS IN HEP





CMS *Preliminary*

H→ZZ→4I **ggH,b**bH 0.97^{+0.09}_{-0.09}(stat.) ^{+0.09}_{-0.07}(syst.) m_H profiled $\mu_{inclusive}$ =0.94 $^{+0.11}_{-0.10}$ **Maximum Likelihood fit VBF** 0.64^{+0.45}_{-0.36}(stat.) ^{+0.16}_{-0.09}(syst.) Significance **Hypothesis testing** VH **P-value** 1.15^{+0.89}_{-0.72}(stat.) ^{+0.26}_{-0.16}(syst.) **Neural Networks** tīH,tH

0.13^{+0.92}_{-0.13}(stat.) ^{+0.11}_{-0.00}(syst.)



9



DATA ANALYSIS IN HEP

Main goals are:

- estimate (measure) the parameters
- quantify the uncertainty of the parameter estimates • test the extent to which the predictions of a theory are in agreement with the data
- Use of statistics for presentation and interpretation (explanation) of data
- A mathematical foundation for statistics is the probability theory

Why is statistics even needed?

- theory predictions in quantum mechanics are not deterministic
- finite size of data sample
- imperfection of the measurement



DATA ANALYSIS GENERAL PICTURE



Sampling reality

EXPERIMENT

Data sample

 $x = (x_1, x_2, ..., x_N)$

x is a multivariate random variable



Described by PDFs, depending on unknown parameters with true values $\theta^{true} = (m_H^{true}, \Gamma_H^{true}, \dots, \sigma^{true})$



PROBABILITY DEFINITION

What is probability anyway?

"Unfortunately, statisticians do not agree on basic principles." - Fred James

- Mathematical (axiomatic) definition
 - **Classical definition**
 - **Frequentist definition**
 - **Bayesian** (subjective) definition





MATHEMATICAL DEFINITION

- Probability"
- \bullet Define an exclusive set of all possible elementary events x_i • Exclusive means the occurrence of one of them implies that none of the others occurs
- For every event x_i , there is a probability $P(x_i)$ which is a real number satisfying the Kolmogorov Axioms of Probability: $P(x_i) \ge 0$ II) $P(x_i \text{ or } x_j) = P(x_i) + P(x_j)$ III) $\sum P(x_i) = 1$
- From these properties more complex probability expressions can be deduced • For non-elementary events, i.e. set of elementary events • For non-exclusive events, i.e. overlapping sets of elementary events
- In Entirely free of meaning, does not tell what probability is about

• Developed in 1933 by Kolmogovor in his "Foundations of the Theory of







CLASSICAL DEFINITION

"The theory of chance consists in reducing all the events of the same kind to a certain number of cases equally possible, that is to say, to such as we may be equally undecided about in regard to their existence, and in determining the number of cases favourable to the event whose probability is sought. The ratio of this number to that of all the cases possible is the measure of this probability, which is thus simply a fraction whose numerator is the number of favourable cases and whose denominator is the number of all the cases possible."

- Pierre-Simon Laplace, A Philosophical Essay on Probabilities







FREQUENTIST DEFINITION

- Experiment performed N times, outcome x occurs N(x) times
- Define probability:

$$P(x) = \lim_{N \to \infty} \frac{N(x)}{N}$$

- Such a probability has big restrictions: • depends on the sample, not just a property of the event experiment must be repeatable under identical conditions • For example one can't define a probability that it'll snow tomorrow

- Probably the one you're implicitly using in everyday life
- Frequentist statistics is often associated with the names of Jerzy Neyman and Egon Pearson





BAYESIAN DEFINITION

• Define probability: P(x) = degree of belief that x is true

- It can be quantified with betting odds:
 - the event

• What's amount of money one's willing to bet based on their belief on the future occurrence of

In particle physics frequency interpretation often most useful, but Bayesian probability can provide more natural treatment of non-repeatable phenomena





BAYES' THEOREM

- Define conditional probability: $P(A|B) = P(A \cap B)/P(B)$
 - o probability of A happening given B happened
 - for independent events P(A|B) = P(A), hence $P(A \cap B) = P(A)P(B)$
- From the definition of conditional probability Bayes' theorem states:

- T is a theory and D is the data
- seen
- P(DIT) is called the likelihood.
- P(D) is the marginal probability of D.
 - \odot P(D) is the prior probability of witnessing the data D under all possible theories
- and the previous state of belief about the theory

В

• P(T) is the prior probability of T: the probability that T is correct before the data D was

• P(DIT) is the conditional probability of seeing the data D given that the theory T is true.

• P(TID) is the posterior probability: the probability that the theory is true, given the data









EXAMPLE

cards

• This situation happens many times in the following days. What is the probability that your friend cheats if you end up paying wins consecutive times?*

• You assume:

• P(cheat) = 5% and P(honest) = 95% (surely an old friend is an unlikely cheater...)

- P(wins | cheat) = 1 and P(wins | honest) = 2^{-wins}
- Bayesian solution:

$$P(cheat | wins) = \frac{P(v)}{P(wins | cheat)P(v)}$$

$$P(cheat \mid 0) = \frac{1 \cdot P(cheat)}{1 \cdot P(cheat) + 2^{-0}P(cheat)}$$

$$P(cheat \mid 5) = \frac{1 \cdot P(cheat)}{1 \cdot P(cheat) + 2^{-5}P(h)}$$

*taken from G.D'Agostini, Bayesian Reasoning in HEP, Principles and Applications, CERN-99-03, 1999

• You meet an old friend in a pub. He proposes that the next round should be payed by whoever of the two extracts the card of lower value from a pack of

P(wins | cheat)P(cheat) cheat) + P(wins | honest)P(honest) $\frac{0.05}{(honest)} = \frac{0.05}{0.05 + 0.95} = 5\%$ $\frac{0.05}{honest} = \frac{0.05}{0.05 + 0.03} = 63\%$



LEARNING BY EXPERIENCE

- The process of updating the probability when new experimental data becomes available can be follow easily if we insert
 - P(cheat) = P(cheat | wins-1) and P(honest) = P(honest | wins -1), where wins 1 indicates the probability assigned after the previous win
- P(wins=1 | cheat) = P(win | cheat) = 1 and P(wins = 1 | honest) = 0.5 Iterative application of the Bayes' formula:

$$P(cheat | wins) = \frac{P(win P(win P($$

n | cheat) P(cheat | wins - 1)vins - 1) + P(win | honest)P(honest | wins - 1)

P(cheat | wins - 1) $P(cheat | wins) = \frac{1}{P(cheat | wins - 1) + 0.5 * P(honest | wins - 1)}$

P(cheat)	P(cheat I wins)		
%	wins=5	wins=10	N
1	24%	91%	
5	63%	98%	g
50	97%	99.9%	g

vins=15 99.7%

99.94%

99.99%

When you learn from the experience, your conclusion does not longer depend on the initial assumptions!





RANDOM VARIABLES

- Outcome not predictable, only the probabilities known
- **Random event** is an event having more than one possible outcome • Each outcome may have associated probability
- \bullet Different possible outcomes may take different possible numerical values x_1 , X₂, ...
- The corresponding probabilities $P(x_1)$, $P(x_2)$, ... form a **probability** distribution
- If observations are independent the distribution of each random variable is unaffected by knowledge of any other observation • When an experiment consists of N repeated observations of the same random variable x, this can be considered as the single observation of a random vector
 - **x**, with components x_1, x_2, \ldots, x_N





DISCRETE RANDOM VARIABLES

- Rolling a die:
 - Sample space = $\{1, 2, 3, 4, 5, 6\}$
 - Random variable x is the number rolled
- Discrete probability distribution:









CONTINUOUS RANDOM VARIABLES

• A spinner:

- Can choose a real number from [0,2n]
- All values equally likely
- x =the number spun
- Probability to select any real number = 0
- Probability to select any range of values > 0• Probability to choose a number in [0,n] = 1/2
- Probability to select a number from any range Δx is $\Delta x/2n$
- Now we say that **probability density** p(x) of x is 1/2n

• More general: $P(A < x < B) = \int_{A}^{B} p(x) dx$





