

COURSE OF STUDY *Physics (LM-17)*
ACADEMIC YEAR 2024-2025

ACADEMIC SUBJECT *Scientific Data Analysis Laboratory*

General information	
Year of the course	2nd
Academic calendar (starting and ending date)	1 st semester: September – December 2024
Credits (CFU/ECTS):	6
SSD	FIS/01 – Experimental Physics
Language	English
Mode of attendance	Compulsory

Professor/ Lecturer	
Name and Surname	Alexis Pompili
E-mail	alexis.pompili@ba.infn.it , alexis.pompili@uniba.it
Telephone	+39 080 5442436
Department and address	Dipartimento Interateneo di Fisica (Office R15) - Via G.Amendola 173, 70125 Bari
Virtual room	Personal Zoom room
Office Hours (and modalities: e.g., by appointment, online, etc.)	Receiving hours for students: - Thursdays: 11-13 and 15-17 (@ office or online/zoom) - by appointment (@ office or online/zoom)

Work schedule			
Hours			
Total	Lectures	Hands-on (laboratory, workshops, working groups, seminars, field trips)	Out-of-class study hours/ Self-study hours
150	16	60	74
CFU/ECTS			
6	2	4	

Learning Objectives	Learners are expected to achieve, by the end of the course, a good knowledge about how to apply – critically – advanced statistics concepts in problems borrowed from real life research, by using proper statistical methods and techniques used in the field of Sub-nuclear and Nuclear Physics. Computational implementation of these methods and techniques is meant to be carried out within the established HEP toos/frameworks by using suitable code in C++ and/or python.
Course prerequisites	- Statistical Data Analysis Course (1 st year Master degree course) - courses dealing with Programming languages during the Bachelor degree course (for instance Informatica).

Teaching strategies	Leveraging the theoretical statistics concepts learned in the Statistical Data Analysis Course, the students are guided, through a fully hands-on approach, to handle from basic to advanced methods and techniques of data analysis, with applications and examples borrowed from the High Energy Physics field. Students
----------------------------	--

	<p>work, along the whole course, with their own personal account on a dedicated Linux virtual machine hosted by ReCas-Bari Data Center, properly configured and accessible also from home to allow extended individual exercise.</p> <p>In-class examples/exercises are fully supported by code files (in C/C++ and /or python) and complemented by slides-based presentations. All the needed files are provided; often some additional homework is proposed as further development of the exercise started in the laboratory lesson. A webpage is provided as additional support, with useful links to organize the given material and also some additional selected documentation, online manuals etc... Supplementary worksheet can be carried out individually but working in small teams is encouraged.</p>
<p>Expected learning outcomes in terms of</p>	
<p>Knowledge and understanding on:</p>	<ul style="list-style-type: none"> ● Understanding the scientific method, the nature, and the methods of research in Physics ● Knowledge of advanced computational techniques ● Knowledge of advanced mathematical tools commonly used in basic and applied research fields ● knowledge to configure an analysis of given (real or simulated) data and understanding the obtained results concerning the extraction or characterization of a physical signal or a physical trend, by means of a correct statistical treatment; ● knowledge of several statistical methods and techniques to handle data and understanding the context in which they can be properly applied and the approximations and uncertainties involved.
<p>Applying knowledge and understanding on:</p>	<ul style="list-style-type: none"> ● Ability to identify the essential elements of a phenomenon ● Ability to use analogy to apply known solutions to new problems (problem solving) ● Ability to use analytical and numerical mathematical computation tools ● Ability to use electronic and computer technologies and their application to experimental data acquisition ● knowledge and ability to set up a statistically proper and computationally reliable data analysis task, by implementing suitable code (in C++ and/or python) within the established HEP tools/frameworks; ● Ability to build and test models and hypotheses.
<p>Soft skills</p>	<ul style="list-style-type: none"> ● <i>Making informed judgments and choices</i> <ul style="list-style-type: none"> ○ Ability to work with increasing levels of autonomy, including taking responsibility in project planning and managing facilities. ○ Awareness of safety issues in laboratory activities ○ Develop the ability to autonomously design and set up a correct data analysis task and to evaluate the obtained results, also understanding how they can be eventually improved. ● <i>Communicating knowledge and understanding</i>

	<ul style="list-style-type: none"> o Competence in communication in Italian and English in advanced fields of Physics o Develop computer skills related to experimental data processing/analysis and to how-to support their exposition/presentation, using specific and technical terminology; o communication and presentation skills - in English language - based on the specific terminology used in the field of experimental data analysis; moreover, dissemination of knowledge with appropriate scientific language; o general ability to work in a group, and to be inserted quickly and effectively in a workplace. <ul style="list-style-type: none"> ● Capacities to continue learning <ul style="list-style-type: none"> o Acquisition of basic knowledge tools for continuous learning and knowledge updates o Ability to learn and to transfer new techniques and methods in experimental data analyses; o Ability to handle and configure a data analysis - given a dataset of a wide variety in nature - by building a model and verifying/testing its validity and by extracting an underlying signature/feature and evaluating its significance in statistical terms.
Syllabus	

<p>Content knowledge</p>	<p>The course contents and structure can be summarized as follow.</p> <p>The course aims to provide, at different levels of approach complexity, the knowledge and the related skills needed to configure, set up and handle, in a statistical and computational proper way, the analysis of a large variety of data, either real or simulated according to a certain model, and the evaluation of the results of this analysis, included the awareness of the implicit approximations and the comprehension of the statistical/systematic uncertainties and correlations involved.</p> <p>An overview of the suitable software tools and frameworks commonly established in the HEP sector are given, together with an extended use of them in a large set of examples of application and exercises.</p> <p>The course is designed into 3 modules (A, B, C) of increasing complexity. The hours can be divided among these 3 modules as follows:</p> <p>A : 14 (2 class + 12 lab) B : 46 (10 class + 36 lab) C : 16 (4 class + 12 lab)</p> <p>The usage of data analysis software tools is introduced gradually through the modules' sequence: ROOT in (A), Roofit package in (B), Jupyter, Pandas "ecosystem", and RDataFrame in (C).</p> <p>The detailed course program follows below.</p> <p><u>Mod. A</u></p> <p>Review of main useful Unix/Linux commands. Review of the main useful commands of the editors vi and emacs. Introduction to the ROOT data analysis software tool.</p> <p>Handling and representation of mathematical functions.</p> <p>Simple applications of representation of spreadsheet data using ROOT (TGraphErrors).</p> <p>Simple <i>random generation</i> of variables starting from a sampling distribution, storage into structures (histograms, <i>trees</i>), handling and visualization of generated data.</p> <p>Histograms' handling (operations with histograms), comparison between two or more histograms (<i>relative</i> normalization). Example of data-<i>Monte Carlo</i> (MC) comparison: between a data distribution and a set of different simulated components (to be composed into a <i>stacked</i> histogram); <i>absolute</i> normalization. Data-MC ratio; <i>rebinning</i> in presence of relevant fluctuations due to low statistics. Histogram of an efficiency and binomial bin error.</p> <p>Example of a simple binary classifier to extract a signal from a set of backgrounds: selection variables and evaluation of their background rejection power by comparing their ROC curves. Choice of the Working Point.</p> <p><u>Mod. B</u></p>
---------------------------------	--

	<p>Introduction to the ROOT package RooFit.</p> <p>Maximum Likelihood interpolation method; implementation of a complex Probability Density Function and components of a theoretical fit model. <i>Extended</i> maximum likelihood. <i>Binned</i> and <i>unbinned</i> fits with examples. Introduction to the support of the modeling within RooFit; libraries, <i>workspace</i> and <i>factory</i> in RooFit.</p> <p>Background modelling: standard polynomials, <i>Chebyshev</i> and <i>Bernstein</i> polynomials, <i>Argus</i> function, polynomials functions with thresholds.</p> <p>Sigmoidal functions; example of fitting an efficiency curve with an <i>error</i> function.</p> <p><i>Goodness-of-fit</i>: normalized chi-squared, fit probability and bin-by-bin pulls. Pulls as tool to check for biases and under/over-estimation of uncertainties.</p> <p>Examples of fits to invariant mass distributions; case of a particle /resonance with an intrinsic width much smaller than experimental mass resolution (one or more gaussians fit) and case of a particle/resonance with an intrinsic width of the same order of magnitude of the experimental resolution (<i>Voigtian</i> function, using a non-relativistic <i>Breit-Wigner</i> function). Fit with an explicit numerical convolution of a relativistic Breit-Wigner and gaussian resolution function.</p> <p>Empirical functions: <i>Crystal Ball</i> function (single- and double-sided) in case of radiative tails, <i>Johnson</i> function (as better alternative to more than two gaussians), <i>Landau</i> function.</p> <p>Exponential convoluted with a time resolution function (RooDecay model); inclusion of an event-by-event time error.</p> <p>Covariance/correlation matrix and estimation of the parameters' uncertainties with Hesse (parabolic approximation) and with Minos (possibly asymmetric errors) methods.</p> <p>Likelihood <i>scan</i> and <i>profile likelihood</i> with examples; verification of the equivalence between the method of likelihood <i>profiling</i> and Minos.</p> <p>Computation of the effective resolution in case of more than one gaussian in the fit, with estimation of its uncertainty by error propagation taking properly into account the correlations among fit parameters (<i>getPropagatedError</i> function).</p> <p>Simultaneous (multidimensional) fits in RooFit. Example: 2D fit of in-variant mass and proper time distributions (mass-lifetime fit of a particle).</p> <p>Non-parametric interpolation technique of the distribution of one (or more) random variable(s). Examples of application of the <i>kernel density estimation</i> (KDE) in one or two dimensions (<i>RooKeysPdf</i> and <i>Roo2DKeysPdf</i>).</p> <p>RooFit as a tool of distributions' generation according to chosen theoretical models (<i>MC toys</i>). Examples of generation of a distribution and its following fit; how to settle the <i>seed</i> of the random generator.</p> <p><i>Closure test</i> of a fitting task. <i>Timing</i> evaluation of a <i>fitting task</i>.</p> <p>Kolmogorov-Smirnov test.</p> <p>Statistical significance of a signal by means of a likelihood ratio in the context of the Wilks theorem: example of approximate estimation. Brief mention to the statistical significance of a signal in the frequentist context (by using pseudo-experiments) resumed in the following module C.</p>
--	---

	<p>Review of the background subtraction methods: <i>sidebands</i> technique, <i>bin-wise method</i> and sPlot (by using the ROOT/TMVA package), with practical applications. Rejection of backgrounds and extraction of a signal. Example/exercise of a multivariate analysis by using a <i>boosted decision tree</i> (BDT).</p> <p>Hypothesis testing and presence of a signal beyond background events. Computation of the <i>p-value</i>, and thus of the statistical significance, for the observation of a signal; definition of two models (null model for background-only hypothesis and alternative model for signal+background hypothesis) and, given the observed data, two calculation approaches:</p> <ol style="list-style-type: none"> 1) <i>frequentist</i>, with a test statistics of <i>one-side Profile likelihood</i> type, by means of pseudo-experiments generation (to be executed in out-of-class hours, being computationally intensive), and 2) <i>asymptotic</i> (AsymptoticCalculator) with the same test statistics. <p>Additional/optional exercise foresees the <i>p-value</i> computation as a function of a signal feature (for instance the signal mass), by using the first approach.</p> <p>Mod. C</p> <p>Introduction to Python scripting language: basic features, flux control, functions. Libraries for scientific computation: numPy, sciPy, matplotlib, uproot (Pandas “ecosystem”). Introduction to the Jupyter framework and examples of modern data analysis with signal extraction in the HEP field (extraction of a signal associated to a decaying particle by means of a full cut-based selection configured within the Jupyter framework).</p>
<p>Texts and readings</p>	<p>The material provided by the teacher covers the full course; it consists of a set of several pdf files inserted in a webpage dedicated to the course. This webpage contains also useful links on valuable online documentation and additional tutorials selected by the teacher and freely available on the internet.</p> <p>If needed by the topic, these files contain a brief and specific theoretical introduction/recall. As a further reference the student can consider the teaching material provided and the reference books suggested in the <i>Statistical Data Analysis</i> course, and in particular the well-known G. Cowan’s (*) and L. Lista’s (**) textbooks.</p> <p>Suggested reference for additionally exploring basic concepts of Machine Learning, at the end of module C, is the textbook by I. Goodfellow <i>et al.</i> (***).</p> <p>The examples and exercises carried out in the lab hours are run on a ReCas virtual machine dedicated to the course, where the students are given a personal account. The files to be used for the examples or to inspire (further) exercises are illustrated in the provided teaching material. Open data from the CMS experiment are used.</p> <p>(*) G. Cowan, <i>Statistical Data Analysis</i>, 1998, Clarendon Press. (**) L. Lista, <i>Statistical method for Data Analysis in Particle Physics</i>, 2019 (3rd edition), Springer Verlag. (***) I. Goodfellow <i>et al.</i>, <i>Deep Learning</i>, 2016, MIT Press.</p>

Notes, additional materials	The proposed books are far richer than the conceptual content of the course. They can be used as a reference guideline concerning the basics concepts, methods and techniques used in the course.
Repository	https://web2.ba.infn.it/~pompili/teaching.html
Assessment	
Assessment methods	- Laboratory reports (20%), - Laboratory exam (80%)
Assessment criteria	The student is expected to have learned: <ul style="list-style-type: none"> - the conceptual and practical/technical knowledge of methods and techniques of data analysis, and the awareness of the proper statistics treatment and computational issues involved in the analysis; - the ability to design, configure and set up a data analysis task; - the knowledge and understanding of the statistically appropriate interpretation of the analysis results, including systematic uncertainties, eventual approximations, possible improvements.
Final exam and grading criteria	Evaluation/grading criteria for the final exam: <ul style="list-style-type: none"> • <i>Knowledge and understanding</i>: 30% • <i>Applying knowledge and understanding</i>: 45% • <i>Autonomy of judgment</i>: 10% • <i>Communicating knowledge and understanding / Communication skills</i>: 10% • <i>Capacities to continue learning</i>: 5%
Further information	
	<p>The three modules are focused on laboratory applications, with examples and exercises, of statistical concepts introduced and studied in depth in the Statistical Data Analysis course. The in-class hours will be devoted to recall/refresh this knowledge but with concrete reference to the relevant functions, methods and algorithms that are involved and are available in the software tools/frameworks to be used.</p> <p>The very end of the third module foresees in-class hours to introduce a coherent and focused/guided overview of machine learning concepts and methods (slightly mentioned in the Statistical Data Analysis course) in the specific HEP context. It is implicitly assumed that a wider theoretical introduction and a wider range of applications in contexts other than HEP are provided in a dedicated Machine Learning course among those left as free students' choice.</p>