Statistics Tools for Particle Physics

An introduction...



Aristotle University of Thessaloniki



European Union European Social Fund Co- financed by Greece and the European Union This research has been co-financed by the European Union (European Social Fund – ESF) and Greek national funds through the Operational Program "Education and Lifelong Learning" of the National Strategic Reference Framework (NSRF) - Research Funding Program: Heracleitus II. Investing in knowledge society through the European Social Fund.

Introduction



For example, Statistics in Particle Physics involves:

- Iooking at histograms, probability distributions
- fits for parameter estimation
- data unfolding, cross-section determination
- setting limits in absence of signal

Introduction

Why become a Statistics expert

Necessary for Particle Physics

- See Higgs boson discovery
- Applications in finance and risk management
 - See job offers in consulting companies... Ph.D. in physics a common requirement!

In this lecture:

- Basics of Statistics
- Introduction to Roofit / Roostats
- Some code examples





Probability

- Random variable: Its value cannot be predicted, instead the probability for obtaining a value can be expressed
- **Probability** (frequentistic): $P_i = \lim_{N \to \infty} \frac{n_i}{N} \simeq \frac{n_i}{N} = P(n_i)$



Rolling dices: a simple discrete-variable experiment

The set of $P(n_1)$, $P(n_2)$, ... $P(n_i)$ defines a **probability distribution**

Probability density function

Probability density function (pdf): defined for continuous variables

f(x) the pdf, not a probability $P(x \in [x_i, x_i + dx]) = f(x_i)dx$ the probability $\int f(x)dx = 1$ pdfs are always normalized to unity

f(x| heta) a parametrized pdf

- x a set of random variables
- heta a set of parameters
- The random variable (*x*) can be independent observations or repetitions of a given experiment
- In HEP it usually is the number of events found in a given experiment, where event refers to the measurement of a given quantity (i.e. mass, cross-section)

Likelihood function

The probability of obtaining a set of observations / measurements derives from the product of the pdf for each individual measurement (x_i)

$$P(x|\theta) = \prod_{i=1}^{N} f(x_i|\theta)$$

- For a given set of observations $x \rightarrow x^{\text{observed}}$, $P(x^{\text{observed}} | \theta)$ is not a pdf and is a function of θ only
- Likelihood function: $L(heta) = P(x^{obs}| heta)$

Binned and unbinned likelihood

N



➡ much faster, but result depends on the binning

 \Rightarrow ideal for large N

х

Likelihood vs x² method

- Likelihood function can be easily interpreted and applied on event-by-event level. The χ² method cannot be applied on event level... (or at least it is not straightforward)
- It can be shown that ML converges faster and with better efficiency. Unless there are special conditions (i.e. limitations in computing power), experts recommend to use ML method.
- ML does not have a problem when there are bins with zero or few events. The χ^2 method is wrong when bin has less than 4 events!

Terminology

- observable: quantity directly measured by an experiment and present in a data set.
- model: a pdf that describes/predicts the distribution of certain observables, normalized to unity.
- parameter of interest: the model parameter that we want to estimate / set limits on (eg. mass, x-section)
- nuisance parameters: the model parameters that are known with limited precision and are not 'interesting' in the above sense (eg. background normalization, shape parameters, systematic unc.)

- **fitting:** the set of statistical tests we run to estimate a parameter θ given some data x and a model $P(x | \theta)$
- Maximum likelihood estimate: the result of the ML fit

ROOT, RooFit, RooStats

- ROOT is the standard analysis package used in HEP
 - Can be downloaded from the site: <u>http://root.cern.ch/drupal/</u>
 - Available binaries and source code to compile in your system
- RooFit is a ROOT library providing a complete toolkit for data analysis
 - Comes with ROOT; if compiling ROOT locally, when running configure use flag: --enable-roofit
 - To use it, just add in your macros: using namespace RooFit;
 - Documentation: <u>http://root.cern.ch/drupal/content/roofit</u>
- RooStats is a set of statistical tools built on top of RooFit, distributed with ROOT
 - To use it: using namespace RooFit; using namespace RooStats;
- Apart from online documentation, plenty of helpful examples in ROOT: \$ROOTSYS/tutorials/roofit, \$ROOTSYS/tutorials/roostats

Simple pdf

Create a simple pdf, plot it and project with pseudo-data

```
using namespace RooFit;
RooRealVar x("x","some random variable",-1,1);
RooRealVar mu("mu","gaussian mean value",0,-1,1);
RooRealVar sigma("sigma","gaussian width",0.2,0.1,0.3);
RooGaussian myPdf("myPdf","gaussian pdf",x,mu,sigma);
```



Adding pdfs



Signal+background pdf

Fit a signal+background mass distribution



RooAddPdf model("model","",RooArgList(signal_model,bkg_model),fsig);

Extended likelihood



RooAddPdf model("model","",RooArgList(signal_model,bkg_model),RooArgList
(nsig,nbkg));



Examples

A two-dimensional fit on mass and lifetime





 (μ, σ)



ood ra-

Signal+background fit in RooStats

Simple mass fit model with RooStats formalism

```
using namespace RooFit;
using namespace RooStats;
RooWorkspace* wspace = new RooWorkspace("wspace");
RooRealVar *m = wspace->factory("m[5000,5700]");
...
wspace->factory("Gaussian:gauss1(m,mean[5278,5275,5285],sigma1[20,0,100])");
wspace->factory("Gaussian:gauss2(m,mean,sigma2[80,40,200])");
wspace->factory("SUM:signal_model(a[0.3,0,1]*gauss1,gauss2)");
wspace->factory("Polynomia1:bkg_model(m,b[0.0005,-0.001,0.001])");
RooAbsPdf *model = wspace->factory(
"SUM:model(nsig[50000,0,500000]*signal_model,nbkg[10000,0,1000000]*bkg_model)");
```

- Based on the concept of RooWorkspace: the technology to record and share arbitrarily complex models in ".root" format ('*digital publishing'*)
- Combining models from various workspaces one can perform combinations across channels, experiments, etc.

Confidence interval

Confidence interval for simple counting experiment

Observed 10 events, expect **exactly 2** background events.

95% confidence interval for signal?

 $P(N_{obs}|s+b)$



```
wspace->factory("obs[10]");
wspace->factory("Poisson::countingModel(obs,sum(s[8,0,20],b[2]))");
ModelConfig modelConfig(new RooWorkspace());
modelConfig.SetWorkspace(*wspace);
modelConfig.SetPdf(*wspace->pdf("countingModel"));
modelConfig.SetPdf(*wspace->pdf("countingModel"));
modelConfig.SetParametersOfInterest(*wspace->var("s"));
ProfileLikelihoodCalculator plc(*data, modelConfig);
plc.SetConfidenceLevel( 0.95 );
LikelihoodInterval* plInt = plc.GetInterval();
```

Including systematics



Observed 10 events, expect **2** ± **1** background events.

95% confidence interval for signal?

 $P(N_{obs}|s+b) \cdot \pi(b)$



```
wspace->factory("obs[10]");
wspace->factory("Poisson::countingModel(obs, sum(s[8,0,20],b[2,1,3]))");
wspace->factory("Gaussian::bkgConstraint(b,meanb[2],sigmab[1])");
wspace->factory("PROD::modelWithConstraints(countingModel,bkgConstraint)");
ModelConfig modelConfig(new RooWorkspace());
modelConfig.SetWorkspace(*wspace);
modelConfig.SetPdf(*wspace->pdf("modelWithConstraints"));
modelConfig.SetPdf(*wspace->pdf("modelWithConstraints"));
modelConfig.SetParametersOfInterest(*wspace->var("s"));
ProfileLikelihoodCalculator plc(*data, modelConfig);
plc.SetConfidenceLevel( 0.95 );
LikelihoodInterval* plInt = plc.GetInterval();
```

Profile likelihood

Cross-section estimation with profile likelihood

We want to measure the cross-section given the number of observed signal

But the observed signal is a function of the cross-section!

$$\sigma = \frac{s}{L \cdot \varepsilon}$$

 $s = \sigma \cdot L \cdot \varepsilon$

$$\mathcal{L}(\sigma) = P(N_{obs}|s+b) \cdot G(b,\delta b) \cdot G(L,\delta L) \cdot G(\varepsilon,\delta\varepsilon)$$