



# Statistical tools for data analysis

How to discover a particle for fun and profit  
with the CMS public data

Mario Pelliccioni  
INFN Torino

INFN School of Statistics – Paestum 2024

# Let's do this!

A 4.5 hours class

Cover few relevant cases for statistical analysis in HEP

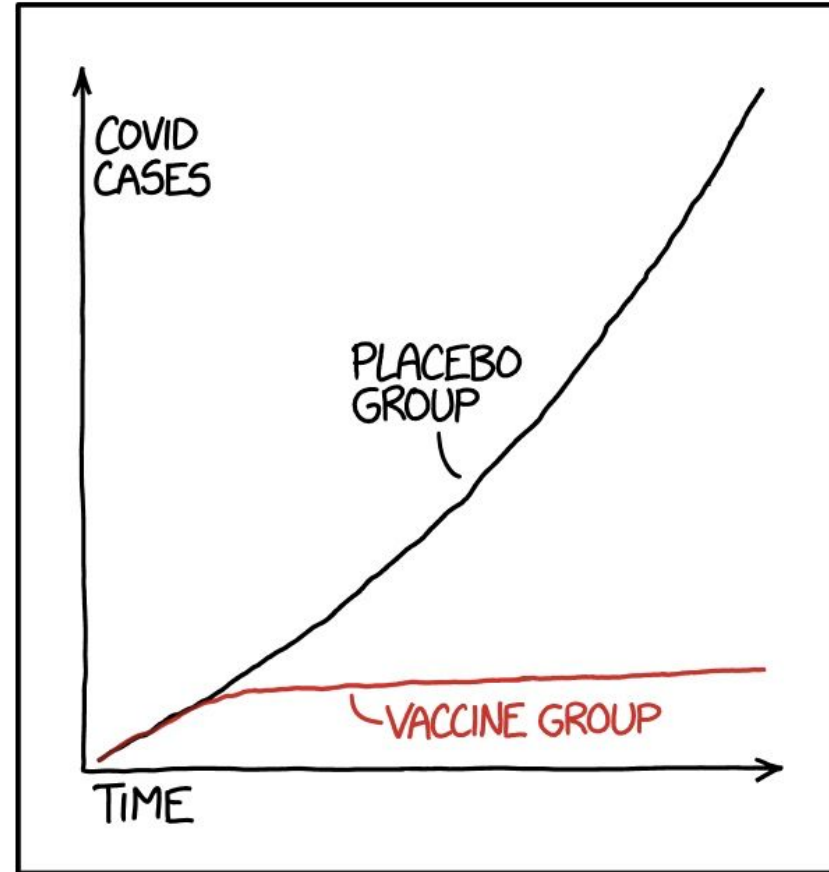
→ Using RooFit and RooStats as main tools

You can use your laptop for this (provided you installed ROOT and python)

→ Exercises will be in pyROOT

CERN/other labs central clusters usually work too

I will flash a few introductory slides for each topic



STATISTICS TIP: ALWAYS TRY TO GET DATA THAT'S GOOD ENOUGH THAT YOU DON'T NEED TO DO STATISTICS ON IT

# Disclaimer

---

The point of this class is to introduce you to some libraries that let you use different statistical tools

I'll try to present as many different approaches as I can

These are not the best (or most appropriate) ways to approach **any** statistical problem

It's your responsibility to find (or build) the best tool for the job!

I will also do some simplifications so that our programs produce a result in a timely manner for this class

# RooFit, RooStats and friends

---

**RooFit:** a ROOT library containing classes that allow to perform multi-dimensional (un)binned maximum likelihood/chi2 fits, toy-MC generation, plotting, etc

**RooStats:** a ROOT library that uses RooFit and provides classes to perform statistical interpretation of your results

**Combine:** an interface to RooFit+RooStats (with some very nifty tools!) created by and for the ATLAS and CMS collaborations

# Documentation

---

For most of what I do, I refer to the ROOT reference guide:

<https://root.cern.ch/doc/master/classes.html>

This includes RooFit and RooStats reference

RooFit manual (a bit outdated):

[https://root.cern.ch/download/doc/RooFit\\_Users\\_Manual\\_2.91-33.pdf](https://root.cern.ch/download/doc/RooFit_Users_Manual_2.91-33.pdf)

RooStats documentation

<https://twiki.cern.ch/twiki/bin/view/RooStats/WebHome>

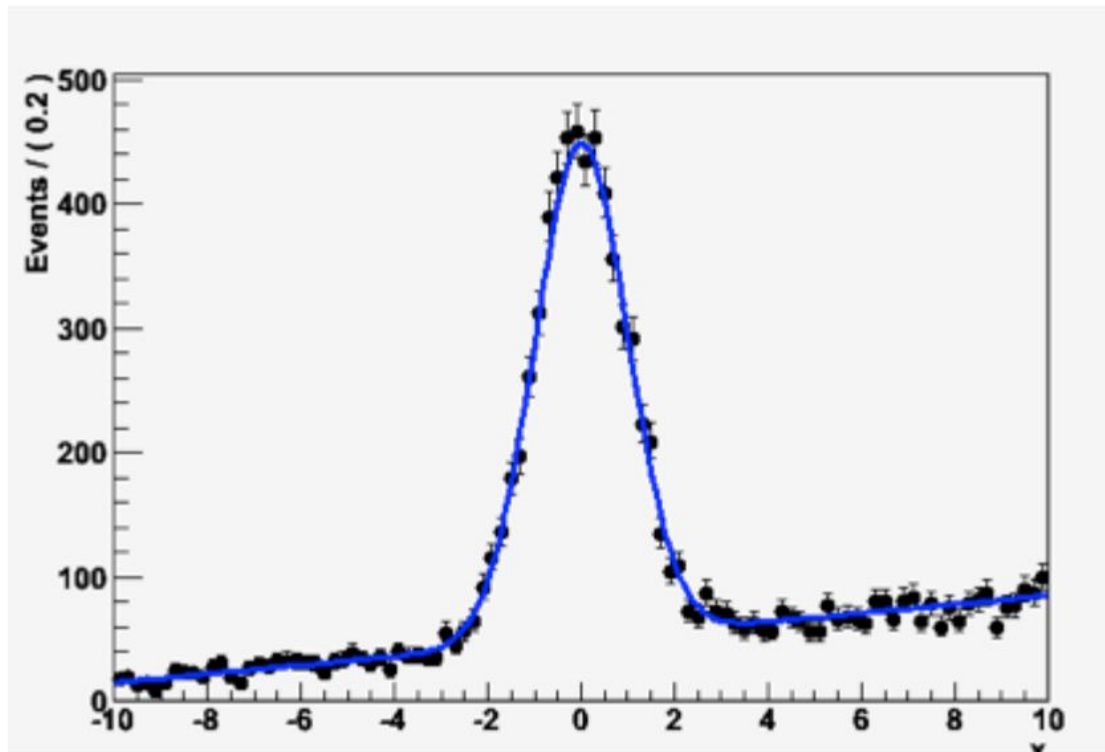
More RooFit/RooStats examples

[https://github.com/pellicci/UserCode/tree/master/RooFitStat\\_class](https://github.com/pellicci/UserCode/tree/master/RooFitStat_class) (C++ based)

[https://github.com/pellicci/UserCode/tree/master/RooFitStat\\_class\\_python](https://github.com/pellicci/UserCode/tree/master/RooFitStat_class_python)

# Why do we need RooFit?

- Focus on one practical aspect of many data analysis in HEP: **How do you formulate your p.d.f. in ROOT**
  - For 'simple' problems (gauss, polynomial) this is easy



- But if you want to do unbinned ML fits, use non-trivial functions, or work with multidimensional functions you quickly find that you need some tools to help you

# The origins

---

- **BaBar experiment at SLAC:** Extract  $\sin(2\beta)$  from time-dependent CP violation of B decay:  $e^+e^- \rightarrow Y(4s) \rightarrow BB$ 
  - Reconstruct both Bs, measure decay time difference
  - Physics of interest is in decay time dependent oscillation

$$f_{sig} \cdot \left[ \text{SigSel}(m; \bar{p}_{sig}) \cdot \left( \text{SigDecay}(t; q_{sig}, \sin(2\beta)) \otimes \text{SigResol}(t \mid dt; r_{sig}) \right) \right] + (1 - f_{sig}) \left[ \text{BkgSel}(m; \bar{p}_{bkg}) \cdot \left( \text{BkgDecay}(t; q_{bkg}) \otimes \text{BkgResol}(t \mid dt; r_{bkg}) \right) \right]$$

- Many issues arise
  - Standard ROOT function framework clearly insufficient to handle such complicated functions  $\rightarrow$  **must develop new framework**
  - **Normalization of p.d.f. not always trivial to calculate**  $\rightarrow$  may need numeric integration techniques
  - Unbinned fit, >2 dimensions, many events  $\rightarrow$  computation performance important  $\rightarrow$  **must try optimize code** for acceptable performance
  - Simultaneous fit to control samples to account for detector performance

# “Dictionary”

- Mathematical objects are represented as C++ objects

Mathematical concept			RooFit class
variable	$x$	➔	<b>RooRealVar</b>
function	$f(x)$	➔	<b>RooAbsReal</b>
PDF	$f(x)$	➔	<b>RooAbsPdf</b>
space point	$\vec{x}$	➔	<b>RooArgSet</b>
integral	$\int_{x_{\min}}^{x_{\max}} f(x) dx$	➔	<b>RooRealIntegral</b>
list of space points		➔	<b>RooAbsData</b>

RooFit uses MINUIT for most of its work, it just provides an easy to use interface and optimizations



# Variables

---

All variables ([observables](#) or [parameters](#)) are defined as **RooRealVar**

RooFit needs to be told which one is which

Several constructors available, depending on the needs:

```
var1 = ROOT.RooRealVar("var1", "My first var", 4.15)      #constant variable
var2 = ROOT.RooRealVar("var2", "My second var", 1., 10.); #range, no initial value
var3 = ROOT.RooRealVar("var3", "My third var", 3., 1., 10.); #valid range, initial value
```

You can also specify the unit (mostly for plotting purposes)

```
time = ROOT.RooRealVar("time", "Decay time", 0., 100., "[ps]");
```

You can change the properties of your RooRealVar later (setRange, setBins, etc.)

If you want to be 100% sure a variable will stay constant, use RooConstVar

For discrete variables, use RooCategory

# Probability Density Functions

---

Each PDF in RooFit must inherit from RooAbsPdf

RooAbsPdf provides methods for numerical integration, events generation (hit & miss), fitting methods, etc.

RooFit provides extensive list of predefined functions (RooGaussian, RooPolynomial, RooCBSShape, RooExponential, RooLandau, etc...)

If possible, use a predefined function (if analytical integration or inversion method for generation available, will speed your computation)

You can define a custom function using RooGenericPdf

# Data Handling

---

Two basic classes to handle data in RooFit:

- **RooDataSet**: an unbinned dataset (think of it as a TTree). An ntuple of data
- **RooDataHist**: a binned dataset (think of it as a TH1F)

Both types of data handlers can have multiple dimensions, contain discrete variables, weights, etc.

# The perfect container

---

In order to “move” information among different RooFit/RooStats programs, one can use the RooWorkspace class

A **RooWorkspace** can contain:

- Variables
- PDFs
- DataSets

A RooWorkspace can be saved into a ROOT file

We'll see how to use it

# The problem at hand

We'll be analyzing a sample from the Run-1 CMS dataset

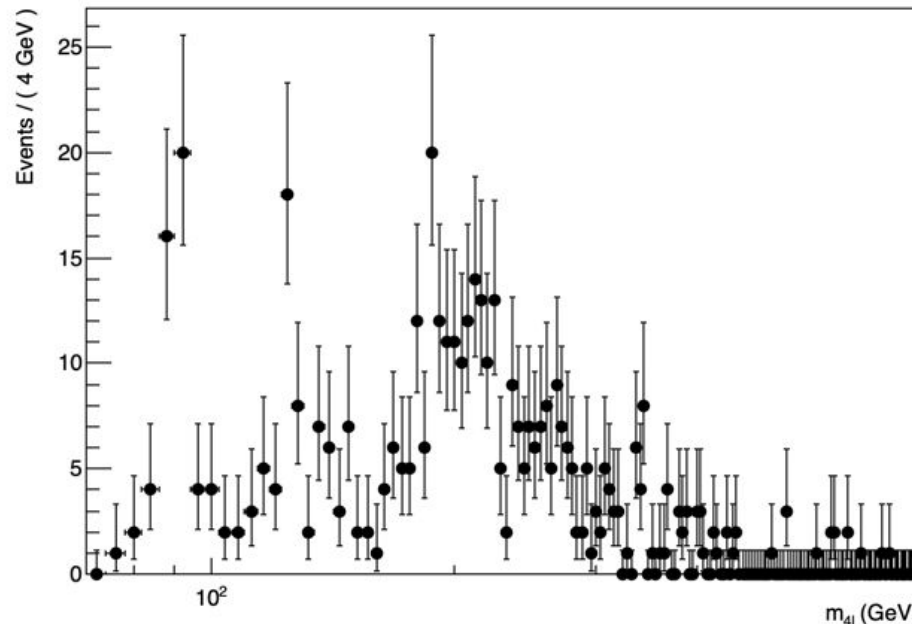
All CMS data from Run1-2 is public → [opendata.cern.ch](http://opendata.cern.ch)

- Events with two 4 leptons (=electron, muon)
- Applied a *somewhat* similar “historical” selection for  $H \rightarrow ZZ$
- Calculated the invariant mass of the system
- Saved it into a RooDataSet (a 1D ntuple containing “m4l” variable)

Also extracted the signal and background distributions from MC

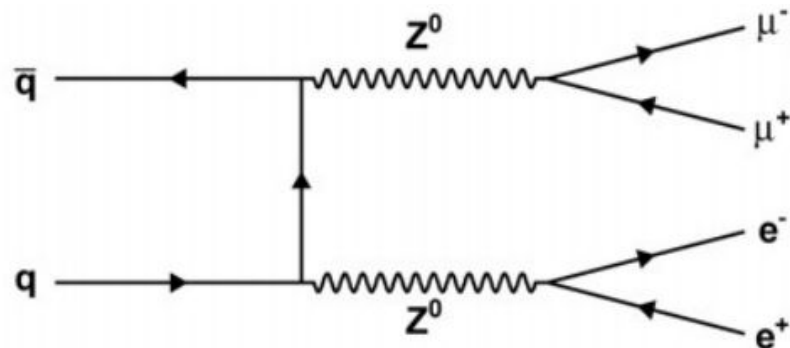
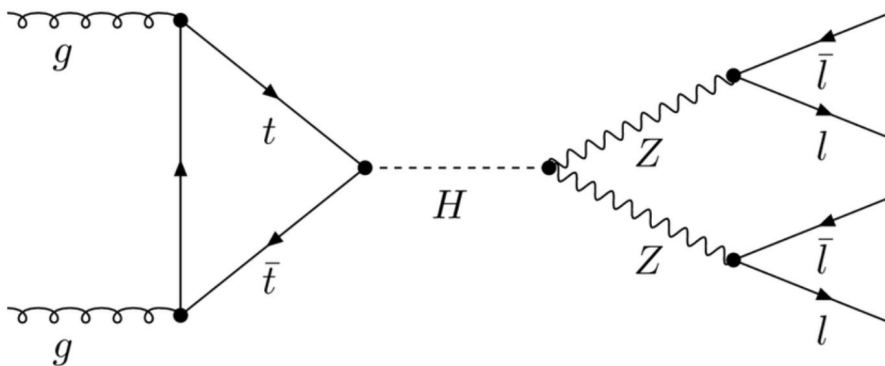
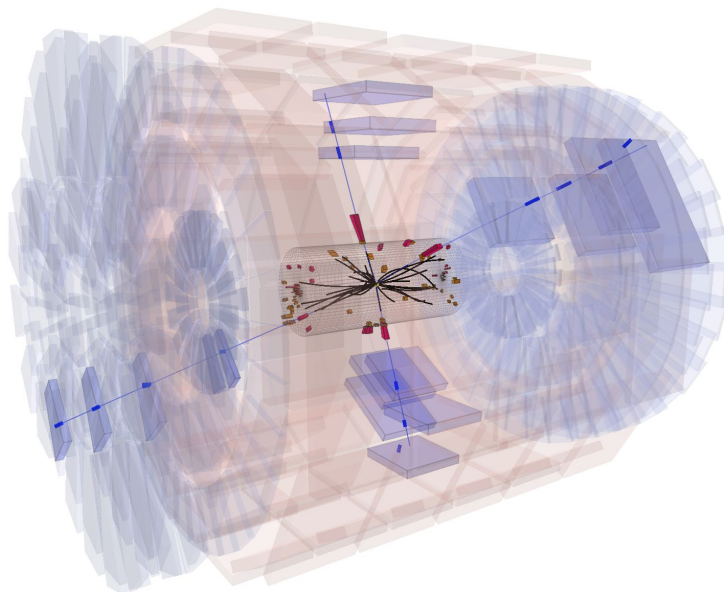
This will be our input distribution

A RooPlot of "m<sub>4l</sub>"



# Context: $H \rightarrow ZZ \rightarrow 4l$

- Very clean final state: four high momentum leptons
- Small number of events:  $BR(H \rightarrow ZZ) \sim 3\%$ 
  - Low statistic, high  $p_T$
- Final states considered:
  - $H \rightarrow \mu^+ \mu^- \mu^+ \mu^-$
  - $H \rightarrow e^+ e^- e^+ e^-$  **Merged in our sample**
  - $H \rightarrow \mu^+ \mu^- e^+ e^-$
- Three main backgrounds
  - $qq \rightarrow ZZ \rightarrow 4l$
  - $gg \rightarrow ZZ \rightarrow 4l$
  - $Z+X$  (jets mis-IDed as leptons)



# Exercise #0

---

The first exercise involves RooFit only

- Construct a signal + backgrounds PDF
  - We will use the MC histograms describing the expected distributions for this first fit
- For now, we will fit for the signal number of events
- Fit it, plot it, save it

We are going to use this program all the way through the exercises

# Intermezzo

DID THE SUN JUST EXPLODE?  
(IT'S NIGHT, SO WE'RE NOT SURE.)

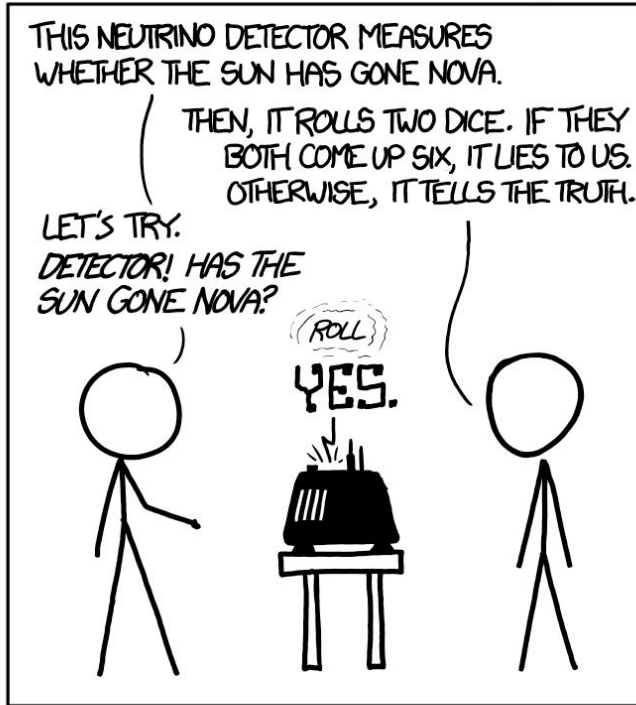
THIS NEUTRINO DETECTOR MEASURES  
WHETHER THE SUN HAS GONE NOVA.

THEN, IT ROLLS TWO DICE. IF THEY  
BOTH COME UP SIX, IT LIES TO US.  
OTHERWISE, IT TELLS THE TRUTH.

LET'S TRY.

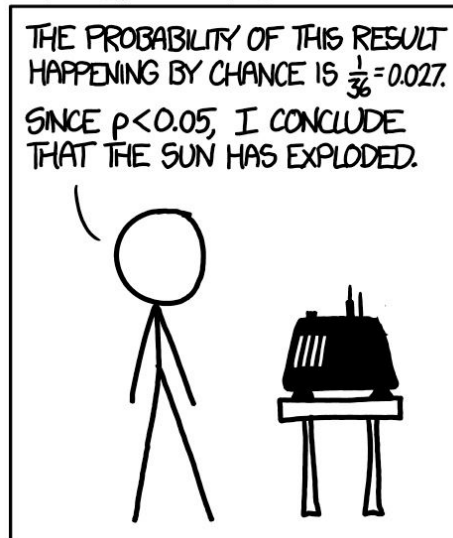
DETECTOR! HAS THE  
SUN GONE NOVA?

ROLL  
YES.



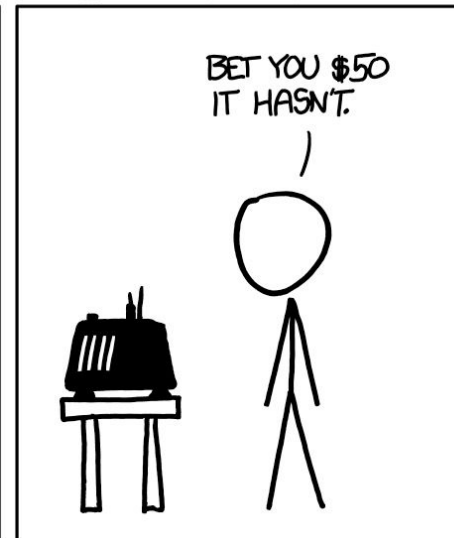
FREQUENTIST STATISTICIAN:

THE PROBABILITY OF THIS RESULT  
HAPPENING BY CHANCE IS  $\frac{1}{36} = 0.027$ .  
SINCE  $p < 0.05$ , I CONCLUDE  
THAT THE SUN HAS EXPLODED.



BAYESIAN STATISTICIAN:

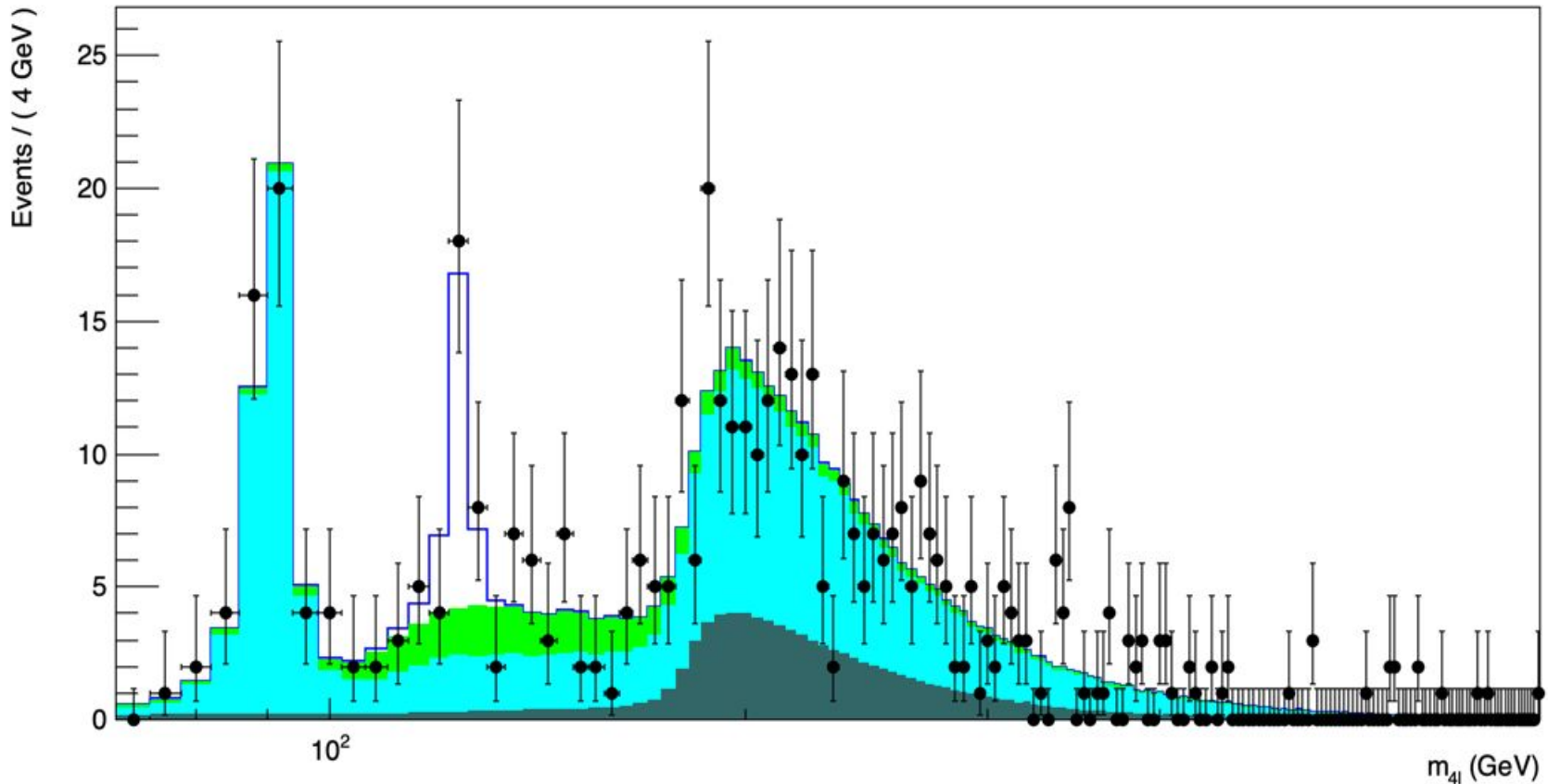
BET YOU \$50  
IT HASN'T.





# Result of exercise #0

A RooPlot of " $m_{4l}$ "



Correlation with the Z+X background?

# Parameter of interest

---

→ a variable you want to know to the best precision and accuracy possible.  
Depends on problem

Number of h125 could be considered the POI

In reality, probably more interested in Higgs production cross section → real connection with theory

$$\sigma(pp \rightarrow H + X) \cdot BR(H \rightarrow ZZ) \cdot BR(Z \rightarrow \ell^+ \ell^-)^2 = \frac{N_{h125}}{\epsilon_{4\ell} \cdot \mathcal{L}}$$

How do we express our problem in this way?

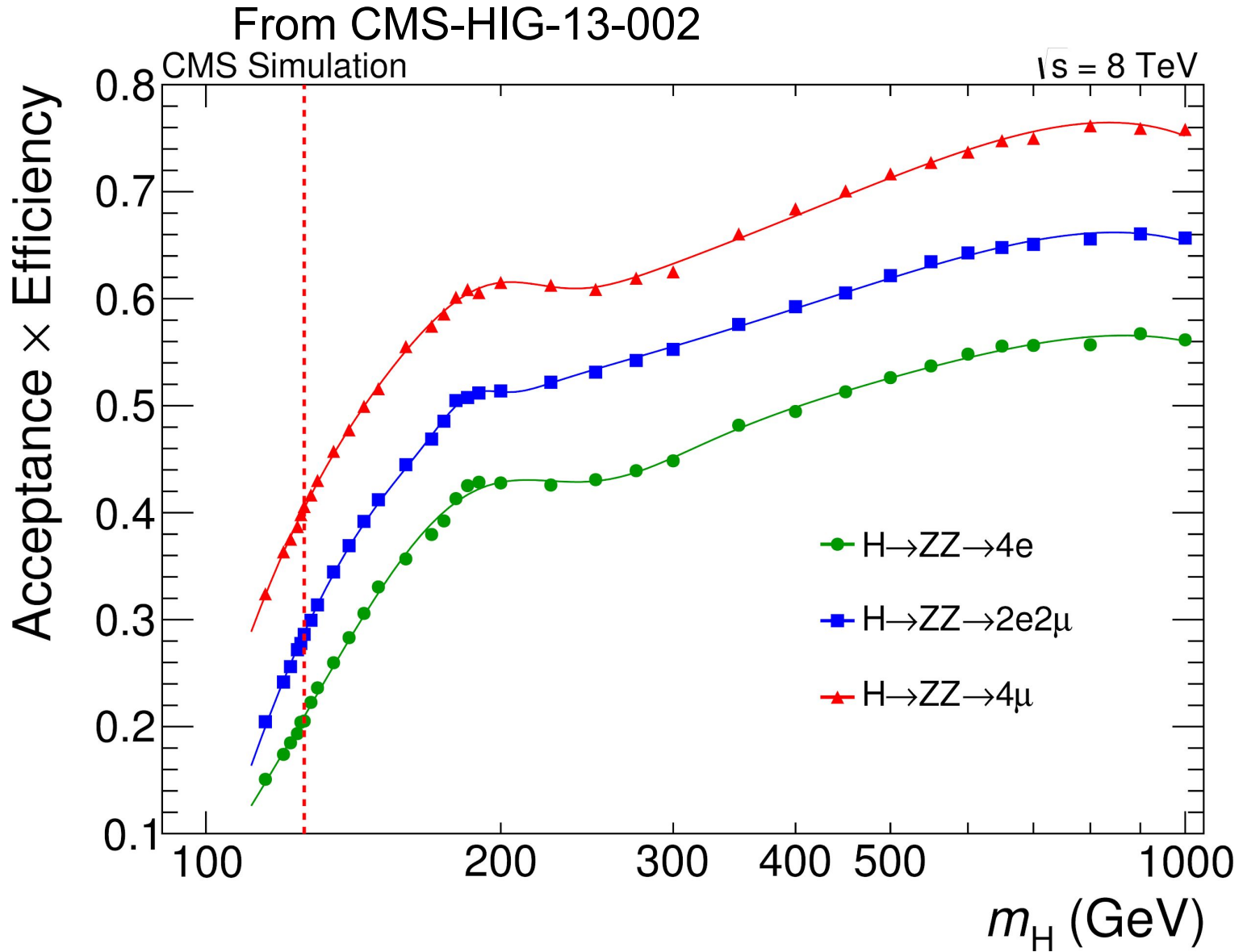
We'll assume:

35% total efficiency

A luminosity of  $24.8 \text{ fb}^{-1}$

Both efficiency and luminosity uncertainties are negligible (for now!)

# Signal efficiency



# Intermezzo

---

```
int getRandomNumber()  
{  
    return 4; // chosen by fair dice roll.  
              // guaranteed to be random.  
}
```

# RooStats

---

Set of libraries for statistical interpretation of your results  
→ communicates with RooFit via RooWorkspace

RooStats does essentially two things:



Interval calculation

Hypothesis testing

To do this, it uses “calculators”





# Main RooStats calculators

---

## ProfileLikelihood calculator

- interval estimation using asymptotic properties of the likelihood function

## Bayesian calculators

- interval estimation using Bayes theorem

**BayesianCalculator** (analytical or adaptive numerical integration)

**MCMCCalculator** (Markov-Chain Monte Carlo)

## HybridCalculator, FrequentistCalculator

- frequentist hypothesis test calculators using toy data (difference in treatment of nuisance parameters)

## AsymptoticCalculator

- hypothesis tests using asymptotic properties of likelihood function

## HypoTestInverter

- invert hypothesis test results (from Asymptotic, Hybrid or FrequentistCalculator) to estimate an interval
- main tools used for limits at LHC (limits using CLs procedure)

## NeymanConstruction and FeldmanCousins

- frequentist interval calculators

# Exercise #1: significance

---

From exercise#0, we can clearly see a peak at 125 GeV

Is this actually clear? How do we quantify?

Let's use the likelihood ratio!



# Intermezzo

$$P\left(\begin{array}{l} \text{I'M NEAR} \\ \text{THE OCEAN} \end{array} \middle| \begin{array}{l} \text{I PICKED UP} \\ \text{A SEASHELL} \end{array}\right) =$$

$$\frac{P\left(\begin{array}{l} \text{I PICKED UP} \\ \text{A SEASHELL} \end{array} \middle| \begin{array}{l} \text{I'M NEAR} \\ \text{THE OCEAN} \end{array}\right) P\left(\begin{array}{l} \text{I'M NEAR} \\ \text{THE OCEAN} \end{array}\right)}{P\left(\begin{array}{l} \text{I PICKED UP} \\ \text{A SEASHELL} \end{array}\right)}$$

$$P\left(\begin{array}{l} \text{I PICKED UP} \\ \text{A SEASHELL} \end{array}\right)$$



STATISTICALLY SPEAKING, IF YOU PICK UP A SEASHELL AND DON'T HOLD IT TO YOUR EAR, YOU CAN PROBABLY HEAR THE OCEAN.

# Exercise #2

---

One important parameter is the mass of the Higgs boson  
→ good connection with theory!

So let's try to obtain confidence and probability intervals on this parameter

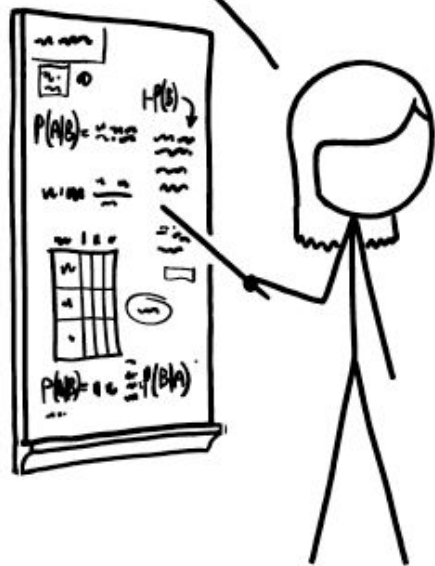
First, we need to modify exercise#0 so our model depends on this variable

Then we'll use a frequentist and a bayesian approach to determine the interval

# Intermezzo

GIVEN THESE PREVALENCES,  
IS IT LIKELY THAT THE TEST  
RESULT IS A FALSE POSITIVE?

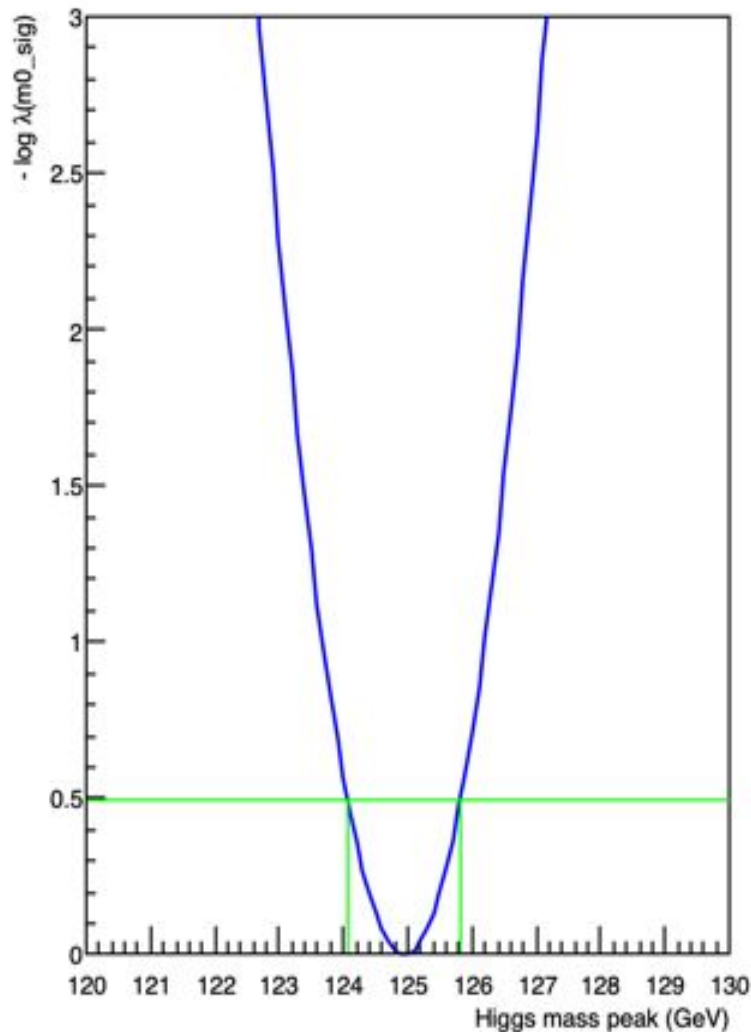
WELL, THIS CHAPTER IS ON  
BAYES' THEOREM, SO YES.



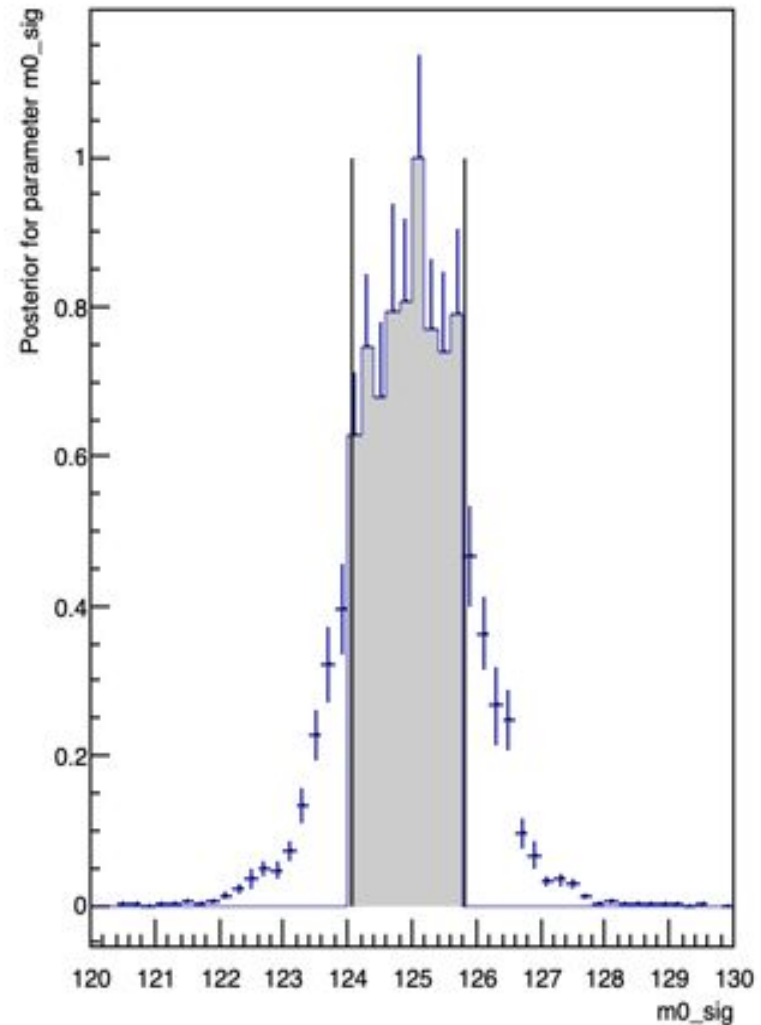
SOMETIMES, IF YOU UNDERSTAND  
BAYES' THEOREM WELL ENOUGH,  
YOU DON'T NEED IT.

# Result of exercise #2

Profile Likelihood Ratio



Bayesian probability interval (Markov Chain)



# Exercise #3: upper limit

---

While the 125 GeV excess is convincing, there's also a possible excess around 145 GeV

Did we miss another boson?

We will modify exercise#0 to include an additional resonance, and check its significance

We will then set an upper limit on this possible resonance contribution

For the frequentist method, we will use  $CL_s$ ...

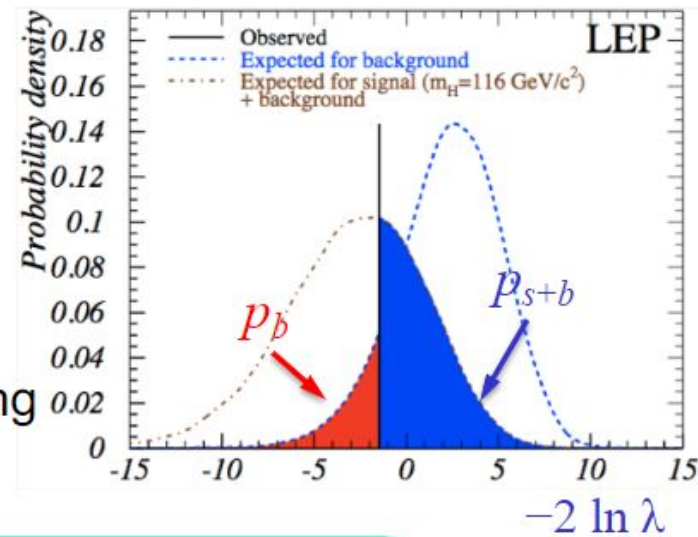
# Understanding $CL_s$

- A **modified approach** was proposed for the first time when combining the limits on the Higgs boson search from the four LEP experiments, ALEPH, DELPHI, L3 and OPAL
- Given a test statistic  $\lambda(x)$ , determine its distribution for the two hypotheses  $H_1(s + b)$  and  $H_0(b)$ , and compute:

$$\left\{ \begin{array}{l} p_{s+b} = P(\lambda(x|H_1) \leq \lambda^{\text{obs}}) \\ p_b = P(\lambda(x|H_0) \geq \lambda^{\text{obs}}) \end{array} \right.$$

- The upper limit is computed, instead of requiring  $p_{s+b} \leq \alpha$ , on the modified statistic  $CL_s \leq \alpha$ :

- Since  $1 - p_b \leq 1$ ,  $CL_s \geq p_{s+b}$ , hence upper limits computed with the  $CL_s$  method are always **conservative**



$$CL_s = \frac{p_{s+b}}{1 - p_b}$$

Note:  $\lambda \leq \lambda^{\text{obs}}$  implies  $-2\ln\lambda \geq \lambda^{\text{obs}}$



# Intermezzo

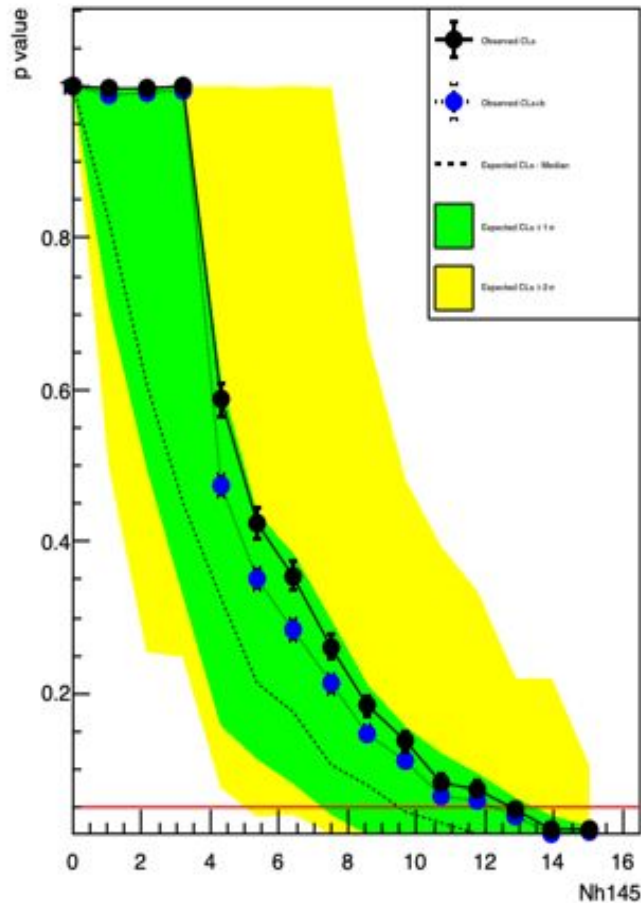
---



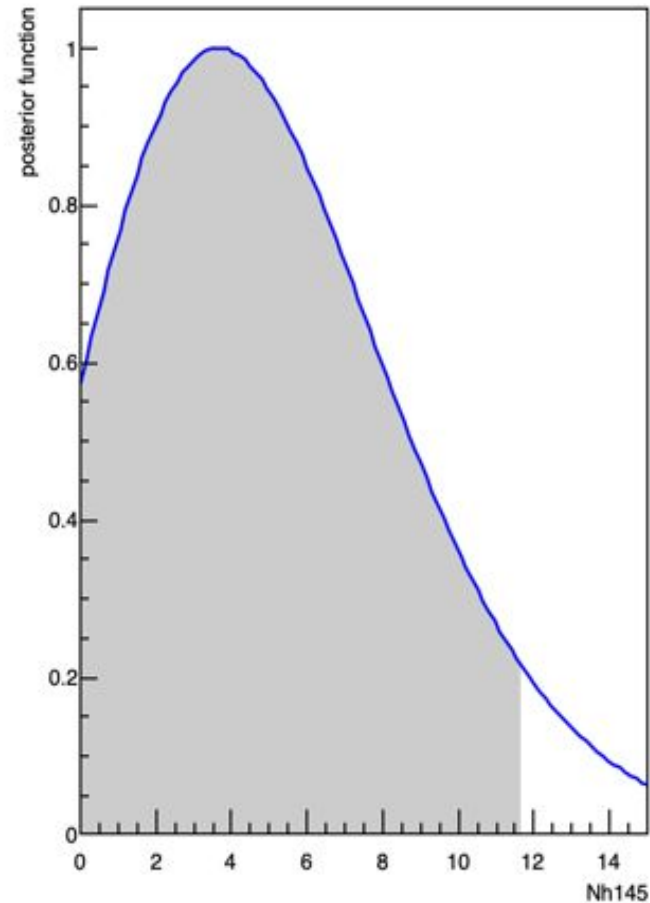
Statisticians have a different version of  
The Boy Who Cried Wolf.

# Results of exercise#3

Frequentist scan result for psi xsec



Posterior probability of parameter "Nh145"



How can we improve the statistics in the belt?



# Exercise #4: test fit with toy-MCs

---

RooFit has tools to test robustness of model via toy-MC generation

Usually healthy to perform, especially on POI

We do this in exercise#4

Approach:

- Use result of fit #0 as *true* model
- Generate 1k experiments, with same stats as CMS, using *true* model
- Fit the 1k experiments with same model
- Compare fit result with the *true* value of the parameters

RooFit can do this pretty easily

# Intermezzo

---

MODIFIED BAYES' THEOREM:

$$P(H|X) = P(H) \times \left( 1 + P(C) \times \left( \frac{P(x|H)}{P(x)} - 1 \right) \right)$$

H: HYPOTHESIS

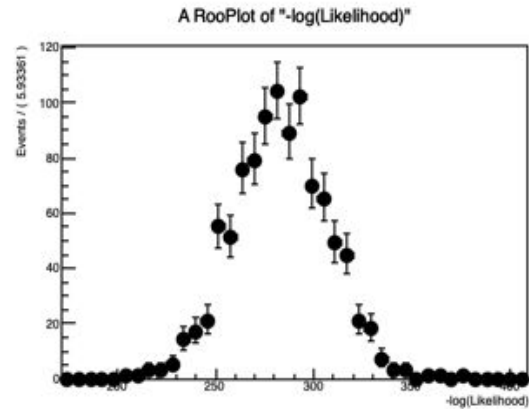
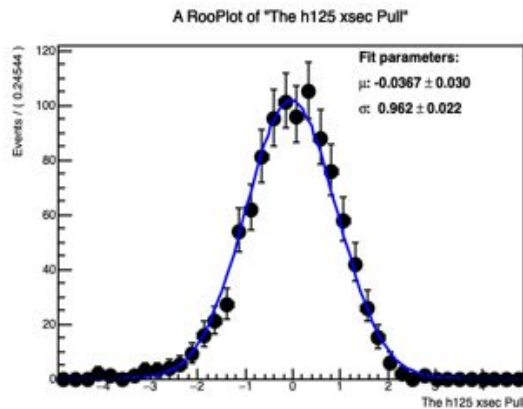
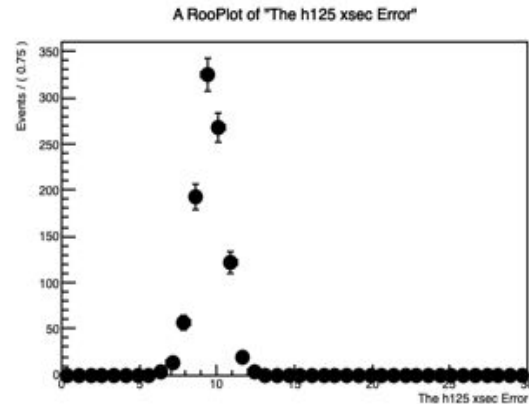
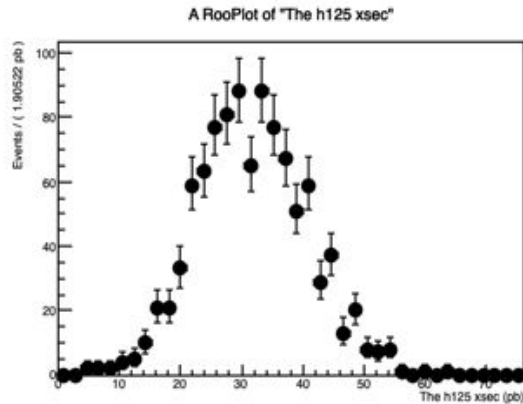
X: OBSERVATION

P(H): PRIOR PROBABILITY THAT H IS TRUE

P(x): PRIOR PROBABILITY OF OBSERVING X

P(C): PROBABILITY THAT YOU'RE USING  
BAYESIAN STATISTICS CORRECTLY

# Results of exercise #4



Additional exercise: what happens if

- you increase number of experiments?
- you increase statistics of each experiment?

→ different effect on what toy-MCs can tell you

# Blinding

---

In general, should not run interpretation on data before analysis strategy is decided

→ Check out the  $\sim 80$  GeV top quark “discovery” for a cautionary tale...

Roofit has tools to ease blinding. For example

```
var1 = ROOT.RooUnblindOffset("var1","blinded var","Daredevil",1.0,my_poi)
```

my\_poi shifted by unknown quantity (seeded by string “Daredevil” of same order)

Can be used to study shifts (systematics!) directly on data while remaining blind

---

**That's all folks!**