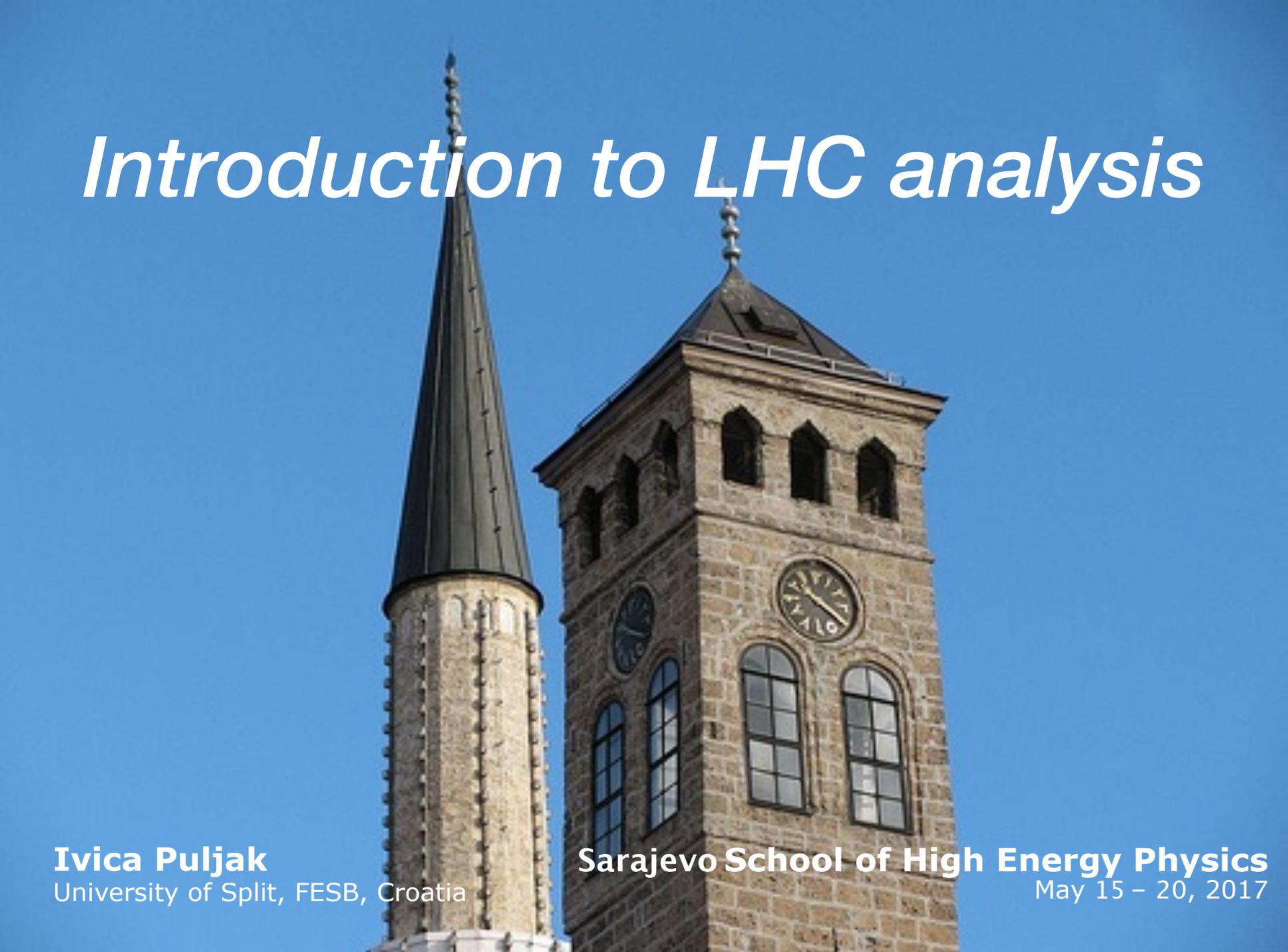


# *Introduction to LHC analysis*



**Ivica Puljak**  
University of Split, FESB, Croatia

**Sarajevo School of High Energy Physics**  
May 15 – 20, 2017

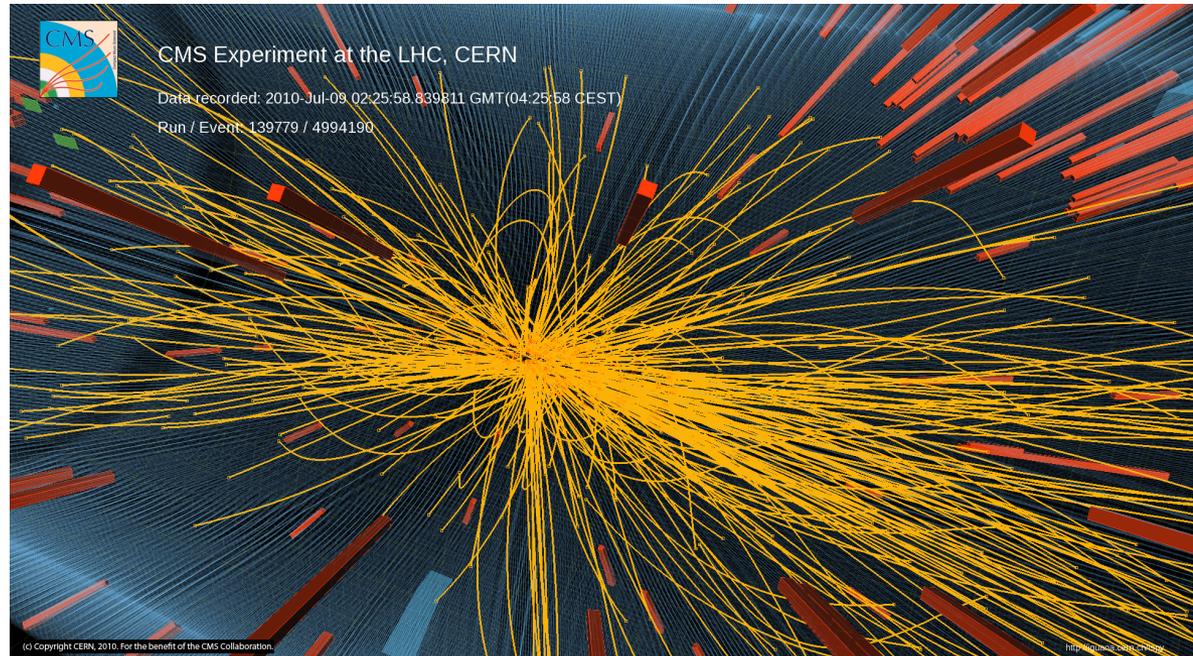
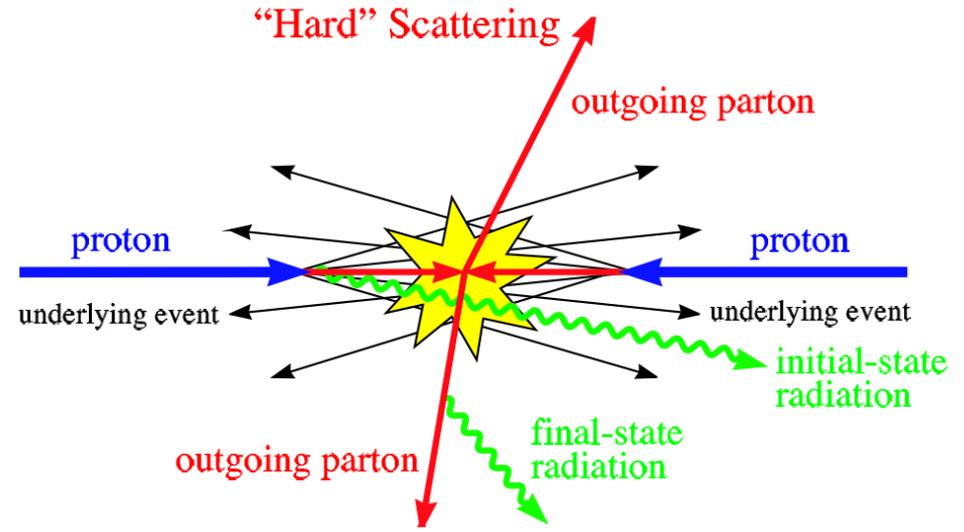
# This lecture ...

---

- ▶ Based on lectures given at the "CERN School of Computing"
  - We will go quickly through these slides
  - They are here mainly to refresh your knowledge in data analysis
    - So that you can follow the other lectures

# Proton – proton collisions

- ▶ Scattering processes at hadron colliders
  - Hard scattering
  - Soft scattering
    - QCD is underlying theory for all of them
- ▶ **Hard scattering** processes
  - Example: W, Z, Higgs, high- $p_T$  jets production
  - Rates and event properties can be predicted with good precision using perturbative QCD (pQCD)
- ▶ **Soft scattering** processes
  - Example: total cross section, underlying event
  - Dominated by less understood non-perturbative QCD effects



# Some of the physicists' jargon

## ▶ **Cross section ( $\sigma$ )**

- A measure of a 'frequency' of the physical process
- Units: barns ( $10^{-28} \text{ m}^2$ )
  - **Typical values:** femtobarns (fb), picobarns (pb)

## ▶ **Luminosity (L)**

- Or *instantaneous luminosity*
- A measure of collisions 'frequency'
  - **Typical at LHC:**  $L = 10^{34} \text{ cm}^{-2}\text{s}^{-1}$

## ▶ **Integrated luminosity ( $\mathcal{L} = \int L dt$ )**

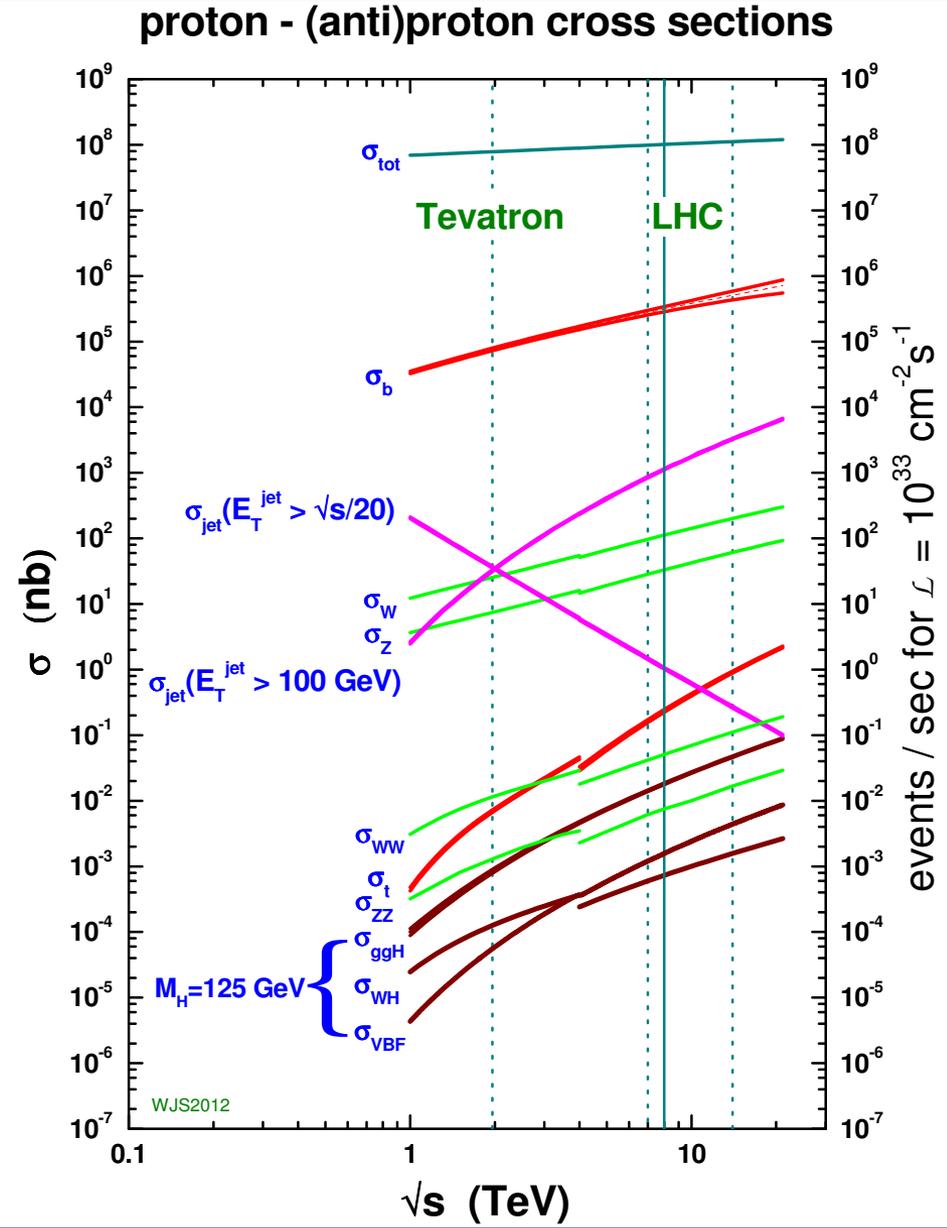
- A measure of number of accumulated collisions after a certain time period
- Units: (cross section) $^{-1}$  .... E.g.  $1 \text{ fb}^{-1} = 1000 \text{ pb}^{-1}$ 
  - **Typical at LHC:** few  $\text{fb}^{-1}$

## ▶ **Number of events (N)**

- Number of (expected) events (N) after a certain time of running

$$N = \sigma \cdot \mathcal{L}$$

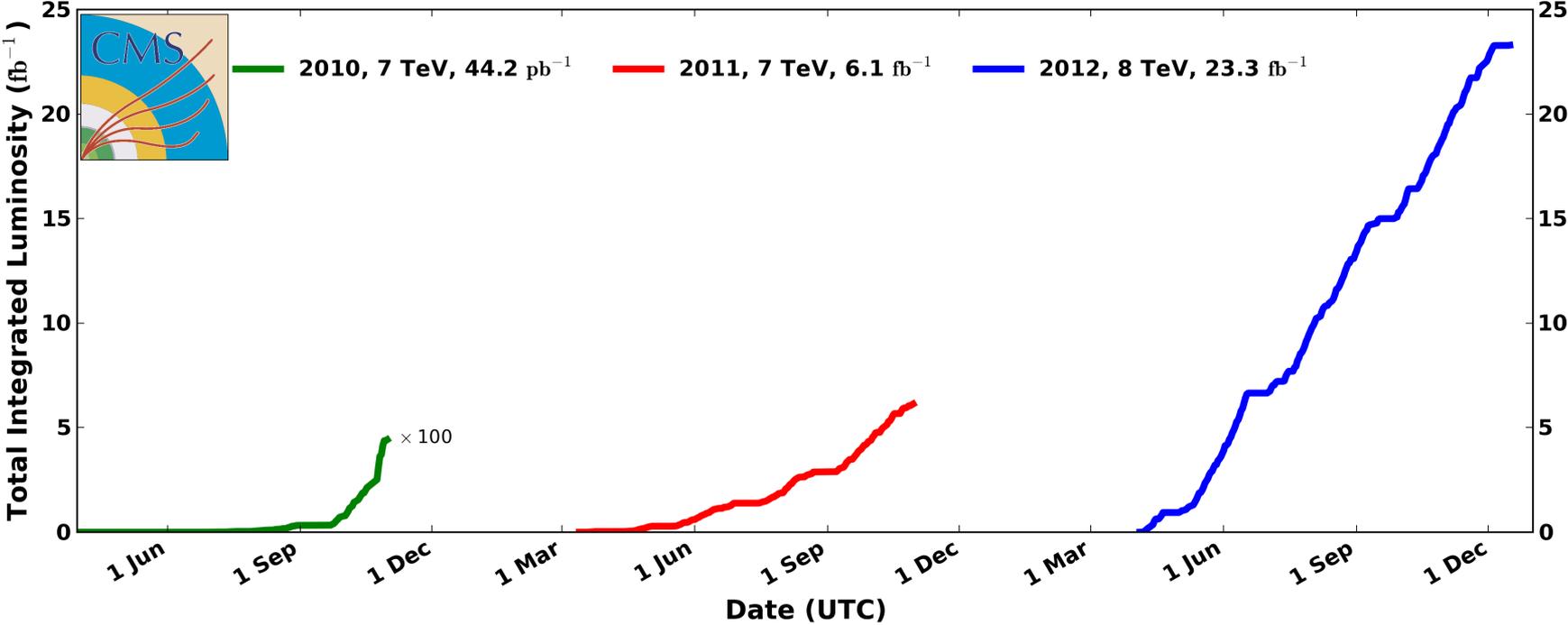
# Cross sections at LHC (and Tevatron)



# Data collected by CMS or ATLAS

### CMS Integrated Luminosity, pp

Data included from 2010-03-30 11:21 to 2012-12-16 20:49 UTC



# What we (will) measure at LHC?

## Something we already know

- At the very beginning of the LHC operation
- For example: production of W and Z bosons

## Something that (probably) exists but wasn't measured yet

- Simply because we are exploring new energy domain
- Standard Model processes
- But surprises are always possible

## Hopefully something new but reasonably expected

- Although "reasonably" is not very well defined 😊
- For example we all expected to find the Higgs boson → and we did find it!
- Heavy neutrinos?

## Maybe something new but less likely

- New heavy bosons ( $Z'$ ,  $W'$ )
- Micro black holes
- Extra dimensions

## Something completely unexpected

- Well, it's hard to look for unexpected 😊



## Nobel Prizes and Laureates

Physics Prizes - < 2013 >

### ▼ About the Nobel Prize in Physics 2013

- Summary
- [Prize Announcement](#)
- [Press Release](#)
- [Advanced Information](#)
- [Popular Information](#)
- [Greetings](#)
- [Award Ceremony Video](#)
- [Award Ceremony Speech](#)
- ▶ [François Englert](#)
- ▶ [Peter Higgs](#)

[All Nobel Prizes in Physics](#)  
[All Nobel Prizes in 2013](#)



The Nobel Prize in Physics 2013  
François Englert, Peter Higgs

Share this:     1.9K 

# The Nobel Prize in Physics 2013



Photo: A. Mahmoud  
**François Englert**  
Prize share: 1/2



Photo: A. Mahmoud  
**Peter W. Higgs**  
Prize share: 1/2

The Nobel Prize in Physics 2013 was awarded jointly to François Englert and Peter W. Higgs *"for the theoretical discovery of a mechanism that contributes to our understanding of the origin of mass of subatomic particles, and which recently was confirmed through the discovery of the predicted fundamental particle, by the ATLAS and CMS experiments at CERN's Large Hadron Collider"*

31 Jul 2012

arXiv:1207.7235v1 [hep-ex]

# Observation of a new boson at a mass of 125 GeV with the CMS experiment at the LHC

The CMS Collaboration\*

## Abstract

Results are presented from searches for the standard model Higgs boson in proton-proton collisions at  $\sqrt{s} = 7$  and 8 TeV in the CMS experiment at the LHC, using data samples corresponding to integrated luminosities of up to  $5.1 \text{ fb}^{-1}$  at 7 TeV and  $5.3 \text{ fb}^{-1}$  at 8 TeV. The search is performed in five decay modes:  $\gamma\gamma$ ,  $ZZ$ ,  $WW$ ,  $\tau^+\tau^-$ , and  $bb$ . An excess of events is observed above the expected background, a local significance of 5.0 standard deviations, at a mass near 125 GeV, signalling the production of a new particle. The expected significance for a standard model Higgs boson of that mass is 5.8 standard deviations. The excess is most significant in the two decay modes with the best mass resolution,  $\gamma\gamma$  and  $ZZ$ ; a fit to these signals gives a mass of  $125.3 \pm 0.4$  (stat.)  $\pm 0.5$  (syst.) GeV. The decay to two photons indicates that the new particle is a boson with spin different from one.

*This paper is dedicated to the memory of our colleagues who worked on CMS but have since passed away.*

*In recognition of their many contributions to the achievement of this observation.*

*Submitted to Physics Letters B*

31 Jul 2012

arXiv:1207.7214v1 [hep-ex]

# Observation of a New Particle in the Search for the Standard Model Higgs Boson with the ATLAS Detector at the LHC

The ATLAS Collaboration

## Abstract

A search for the Standard Model Higgs boson in proton-proton collisions with the ATLAS detector at the LHC is presented. The datasets used correspond to integrated luminosities of approximately  $4.8 \text{ fb}^{-1}$  collected at  $\sqrt{s} = 7 \text{ TeV}$  in 2011 and  $5.8 \text{ fb}^{-1}$  at  $\sqrt{s} = 8 \text{ TeV}$  in 2012. Individual searches in the channels  $H \rightarrow ZZ^{(*)} \rightarrow 4\ell$ ,  $H \rightarrow \gamma\gamma$  and  $H \rightarrow WW^{(*)} \rightarrow e\nu\mu\nu$  in the 8 TeV data are combined with previously published results of searches for  $H \rightarrow ZZ^{(*)}$ ,  $WW^{(*)}$ ,  $bb$  and  $\tau^+\tau^-$  in the 7 TeV data and results from improved analyses of the  $H \rightarrow ZZ^{(*)} \rightarrow 4\ell$  and  $H \rightarrow \gamma\gamma$  channels in the 7 TeV data. Clear evidence for the production of a neutral boson with a measured mass of  $126.0 \pm 0.4 \text{ (stat)} \pm 0.4 \text{ (sys)} \text{ GeV}$  is presented.

This observation, which has a significance of 5.9 standard deviations, corresponding to a background fluctuation probability of  $1.7 \times 10^{-9}$ , is compatible with the production and decay of the Standard Model Higgs boson.

# Expectations vs measurements

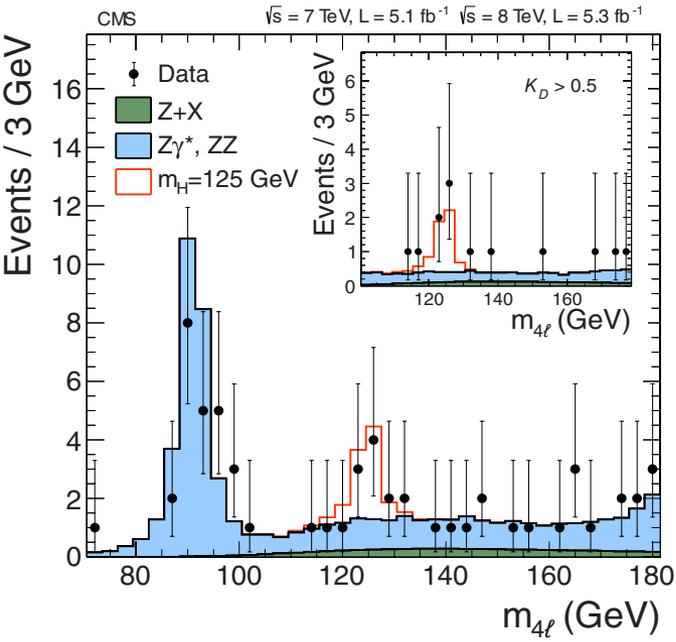


Figure 4: Distribution of the four-lepton invariant mass for the  $ZZ \rightarrow 4\ell$  analysis. The points represent the data, the filled histograms represent the background, and the open histogram shows the signal expectation for a Higgs boson of mass  $m_H = 125$  GeV, added to the background expectation. The inset shows the  $m_{4\ell}$  distribution after selection of events with  $K_D > 0.5$ , as described in the text.

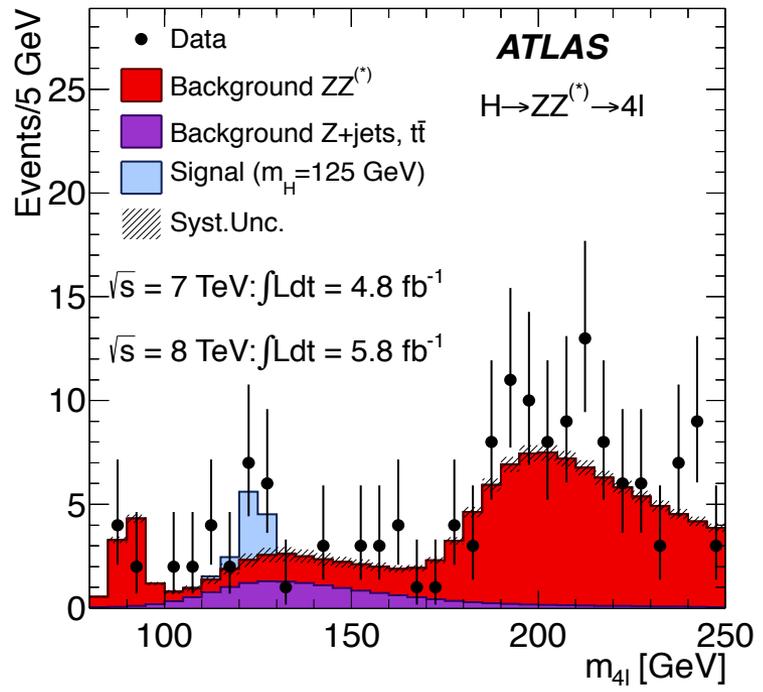
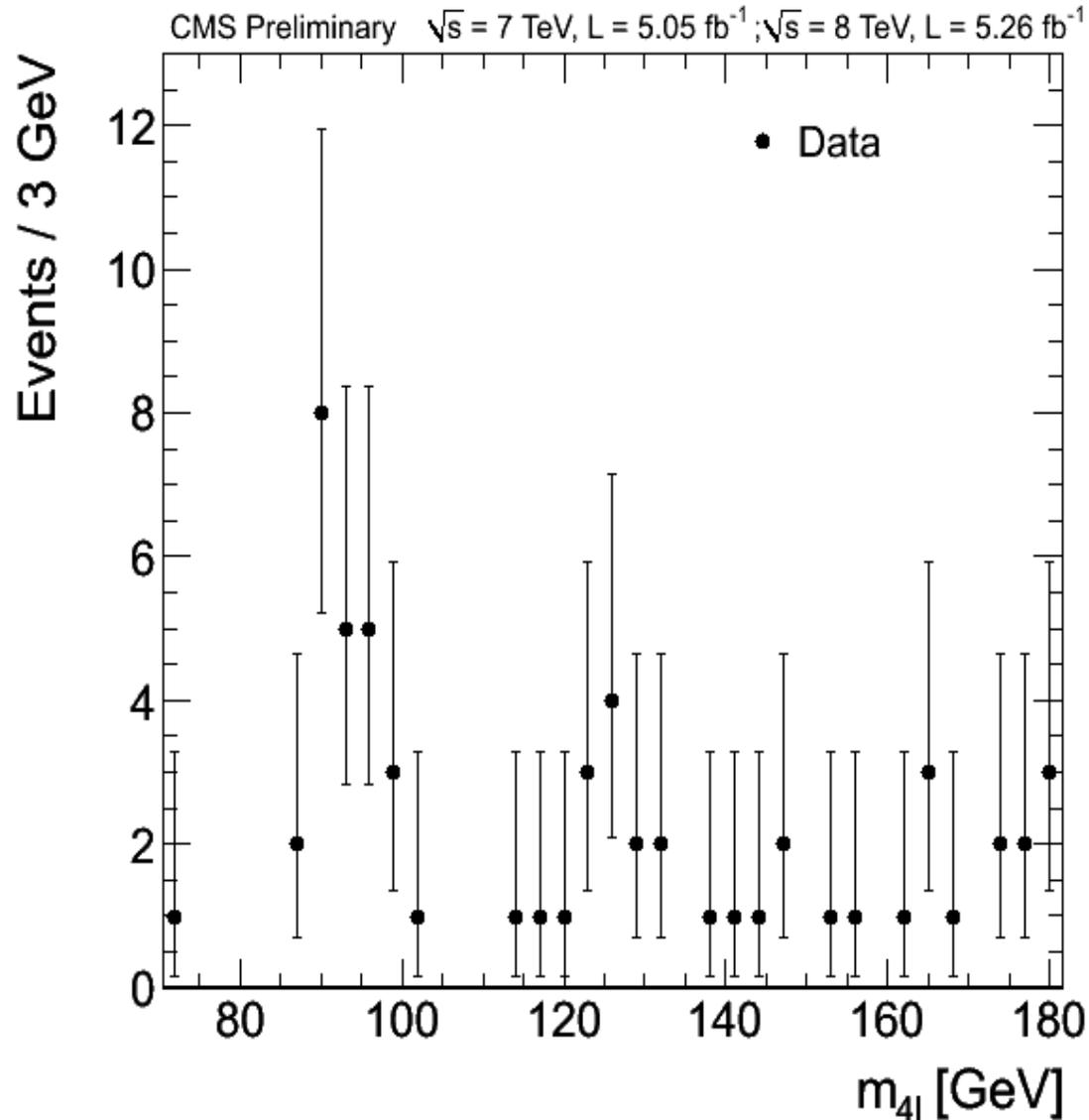
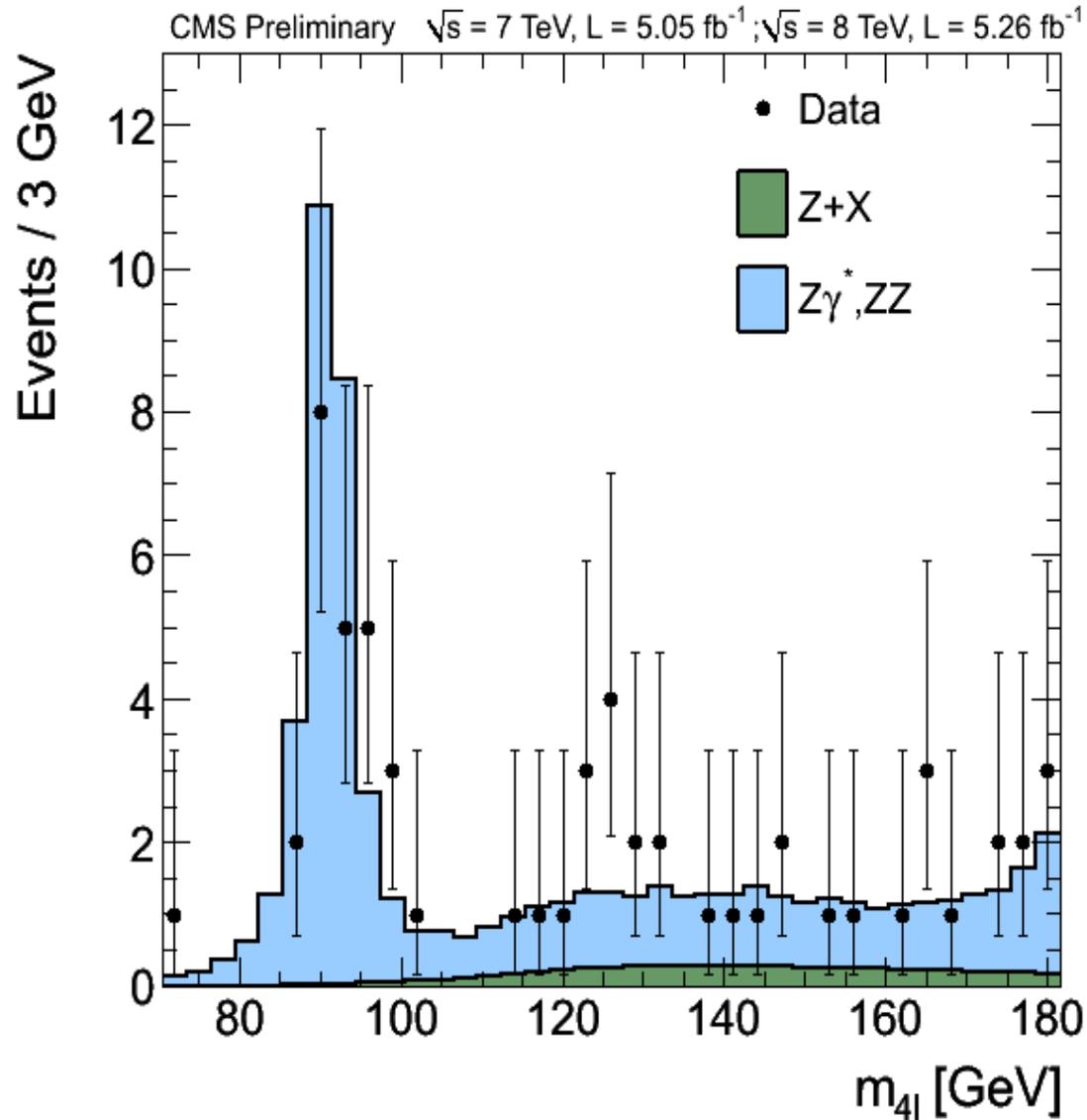


Figure 2: The distribution of the four-lepton invariant mass,  $m_{4\ell}$ , for the selected candidates, compared to the background expectation in the 80–250 GeV mass range, for the combination of the  $\sqrt{s} = 7$  TeV and  $\sqrt{s} = 8$  TeV data. The signal expectation for a SM Higgs with  $m_H = 125$  GeV is also shown.

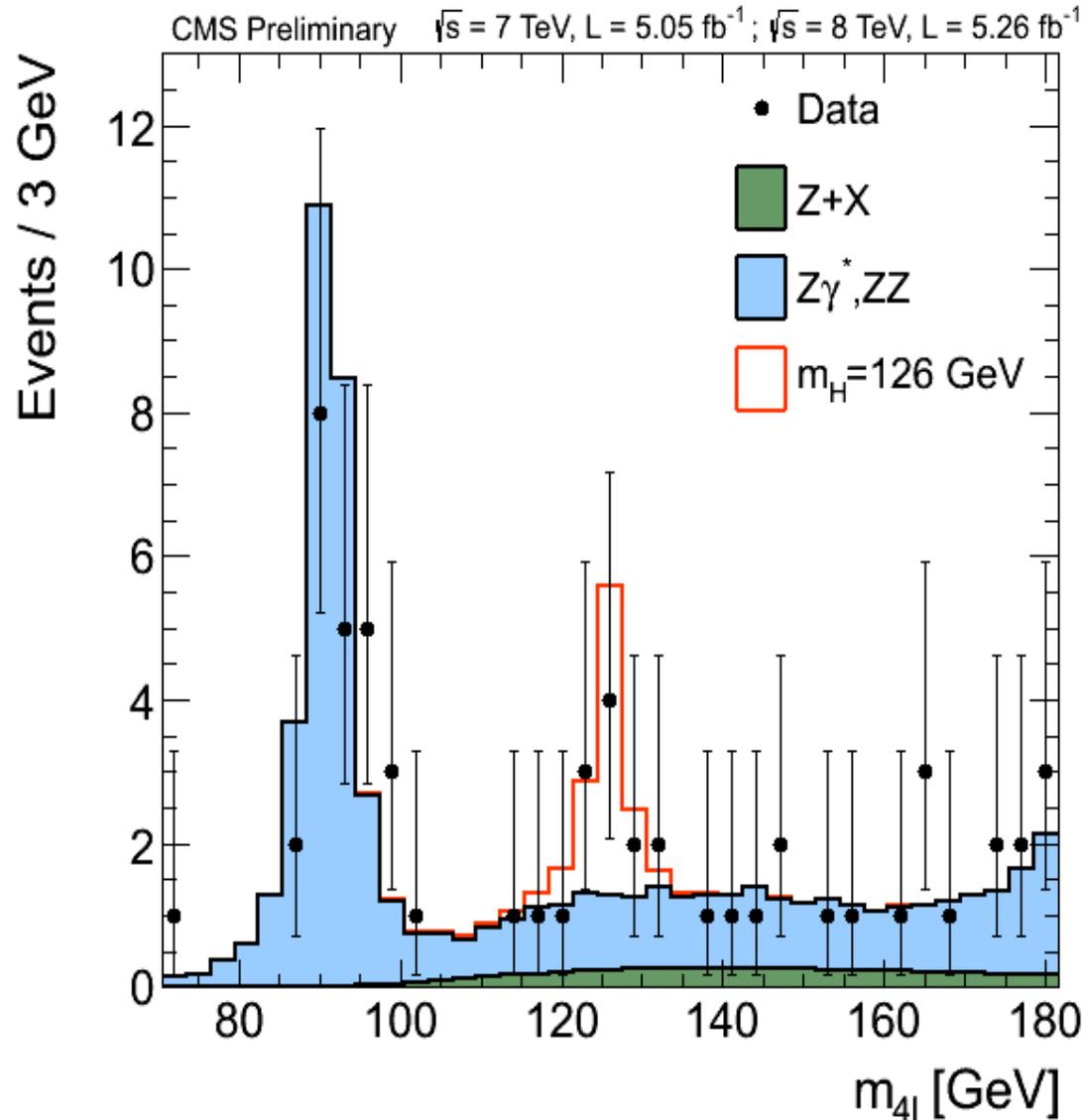
# $H \rightarrow ZZ \rightarrow l-l^+l-l^+$ events distribution



# $H \rightarrow ZZ \rightarrow l-l^+l-l^+$ events distribution



# $H \rightarrow ZZ \rightarrow l-l^+l-l^+$ events distribution



# $H \rightarrow \gamma\gamma$ : Example of fitting

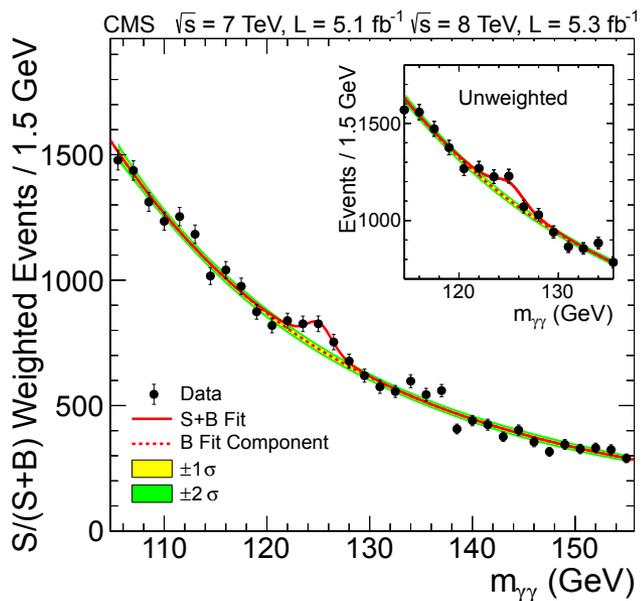


Figure 3: The diphoton invariant mass distribution with each event weighted by the  $S/(S+B)$  value of its category. The lines represent the **fitted background and signal**, and the coloured bands represent the  $\pm 1$  and  $\pm 2$  standard deviation uncertainties on the background estimate. The inset shows the central part of the unweighted invariant mass distribution.

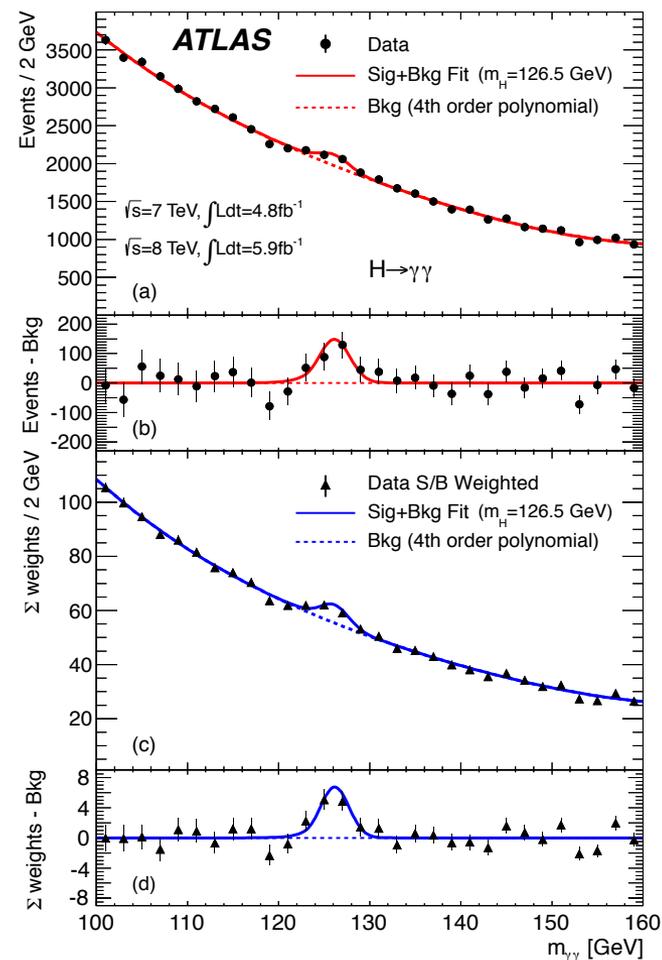


Figure 4: The distributions of the invariant mass of diphoton candidates after all selections for the combined 7 TeV and 8 TeV data sample. The inclusive sample is shown in a) and a weighted version of the same sample in c); the weights are explained in the text. The result of a fit to the data of the sum of a signal component fixed to  $m_H = 126.5$  GeV and a background component described by a fourth-order Bernstein polynomial is superimposed. The residuals of the data and weighted data with respect to the respective fitted background component are displayed in b) and d).

# H→bb: example of Multivariate analysis (MVA)

For the multivariate analysis, a boosted decision tree (BDT) [115, 116] is trained to give a high output value (score) for signal-like events and for events with good diphoton invariant mass resolution, based on the following observables: (i) the photon quality determined from electromagnetic shower shape and isolation variables; (ii) the expected mass resolution; (iii) the per-event estimate of the probability of locating the diphoton vertex within 10 mm of its true location along the beam direction; and (iv) kinematic characteristics of the photons and the diphoton system. The kinematic variables are constructed so as to contain no information about the invariant mass of the diphoton system. The diphoton events not satisfying the dijet selec-

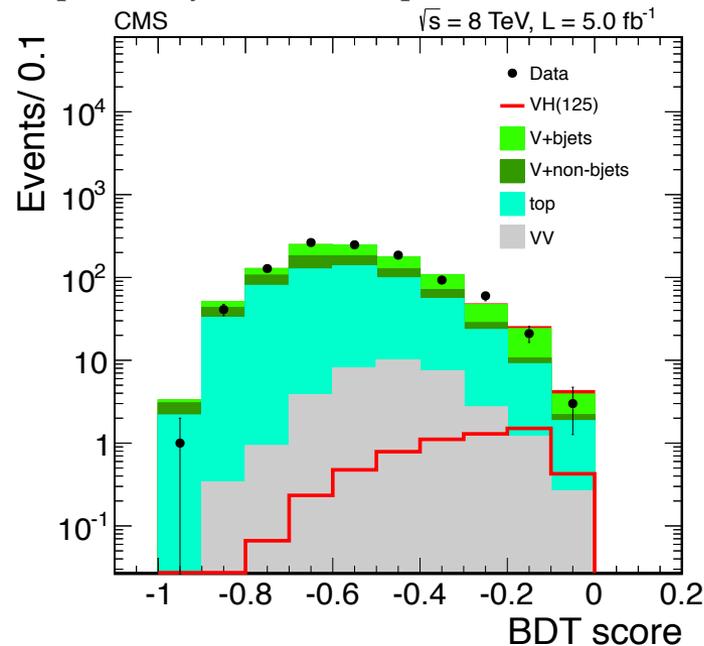


Figure 11: Distribution of BDT scores for the high- $p_T$  subchannel of the  $Z(\nu\nu)H(bb)$  search in the 8 TeV data set after all selection criteria have been applied. The signal expected from a Higgs boson ( $m_H = 125$  GeV), including  $W(\ell\nu)H$  events where the charged lepton is not reconstructed, is shown added to the background and also overlaid for comparison with the diboson background.

# Example of limits

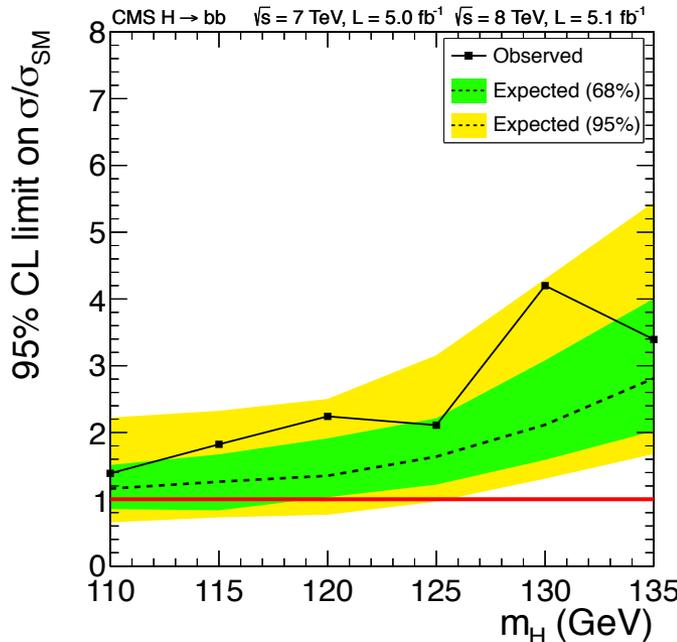


Figure 12: The 95% CL limit on the signal strength  $\sigma/\sigma_{SM}$  for a Higgs boson decaying to two b quarks, for the combined 7 and 8 TeV data sets. The symbol  $\sigma/\sigma_{SM}$  denotes the production cross section times the relevant branching fractions, relative to the SM expectation. The background-only expectations are represented by their median (dashed line) and by the 68% and 95% CL bands.

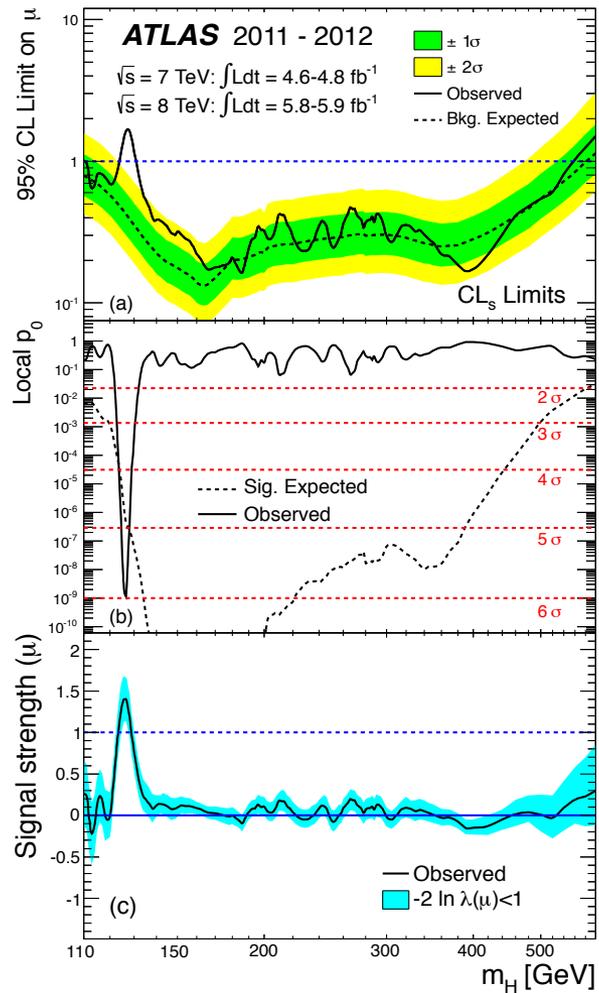


Figure 7: Combined search results: (a) The observed (solid) 95% CL limits on the signal strength as a function of  $m_H$  and the expectation (dashed) under the background-only hypothesis. The dark and light shaded bands show the  $\pm 1\sigma$  and  $\pm 2\sigma$  uncertainties on the background-only expectation. (b) The observed (solid) local  $p_0$  as a function of  $m_H$  and the expectation (dashed) for a SM Higgs boson signal hypothesis ( $\mu = 1$ ) at the given mass. (c) The best-fit signal strength  $\hat{\mu}$  as a function of  $m_H$ . The band indicates the approximate 68% CL interval around the fitted value.

# p-value and hypothesis testing

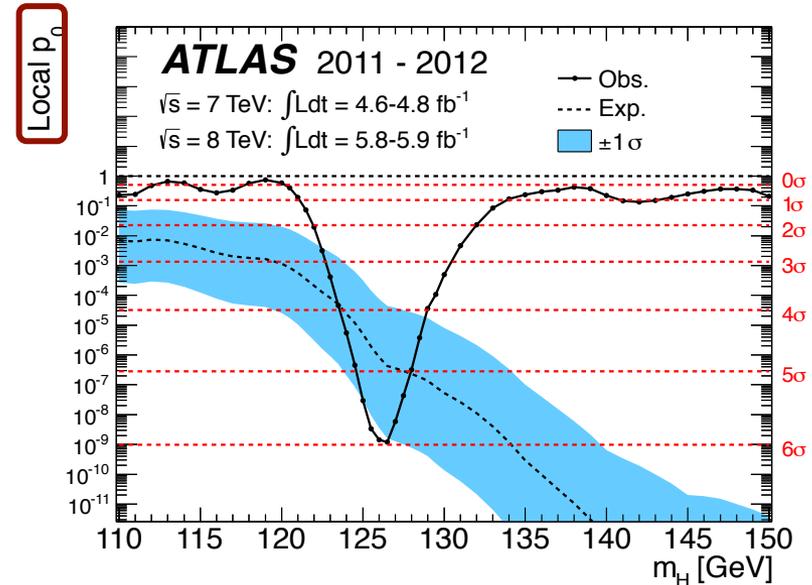
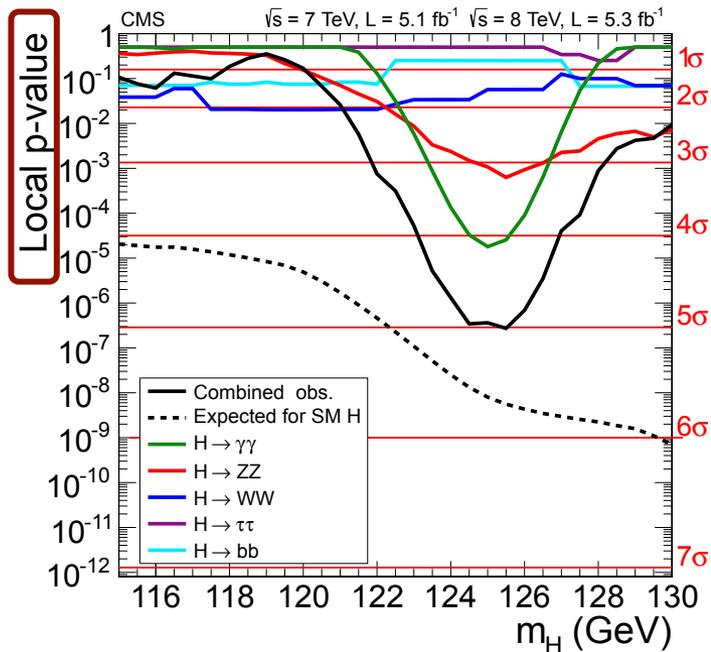


Figure 9: The observed (solid) local  $p_0$  as a function of  $m_H$  in the low mass range. The dashed curve shows the expected local  $p_0$  under the hypothesis of a SM Higgs boson signal at that mass with its  $\pm 1\sigma$  band. The horizontal dashed lines indicate the  $p$ -values corresponding to significances of 1 to 6  $\sigma$ .

Figure 15: The observed local  $p$ -value for the five decay modes and the overall combination as a function of the SM Higgs boson mass. The dashed line shows the expected local  $p$ -values for a SM Higgs boson with a mass  $m_H$ .

# Measuring properties

Asymptotically, **the test statistic**  $-2 \ln \lambda(\mu, m_H)$  is distributed as a  $\chi^2$  distribution with two degrees of freedom. The resulting 68% and 95% CL contours for the  $H \rightarrow \gamma\gamma$  and  $H \rightarrow WW^{(*)} \rightarrow \ell\nu\ell\nu$  channels are shown in

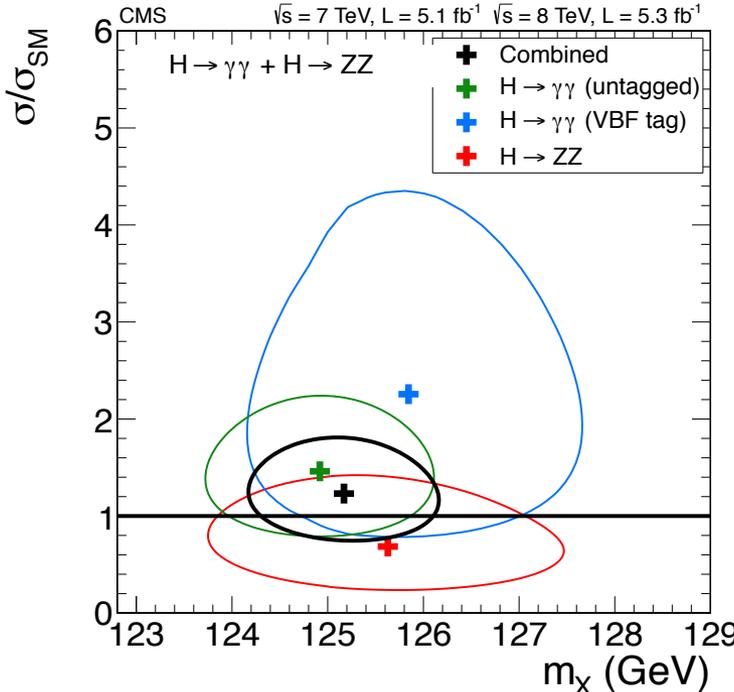


Figure 17: **The 68% CL contours** for the signal strength  $\sigma/\sigma_{SM}$  versus the boson mass  $m_\chi$  for the untagged  $\gamma\gamma$ ,  $\gamma\gamma$  with VBF-like dijet,  $4\ell$ , and their combination. The symbol  $\sigma/\sigma_{SM}$  denotes the production cross section times the relevant branching fractions, relative to the SM expectation. In this combination, the relative signal strengths for the three decay modes are constrained by the expectations for the SM Higgs boson.

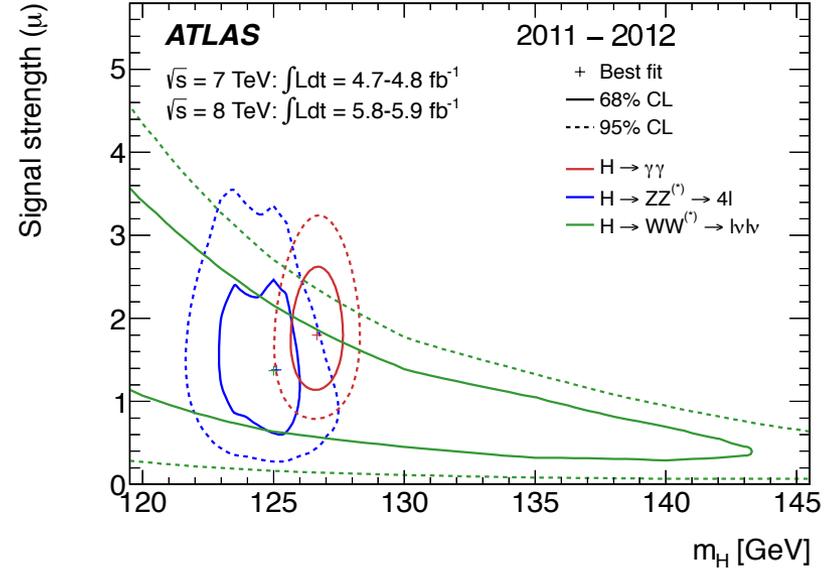
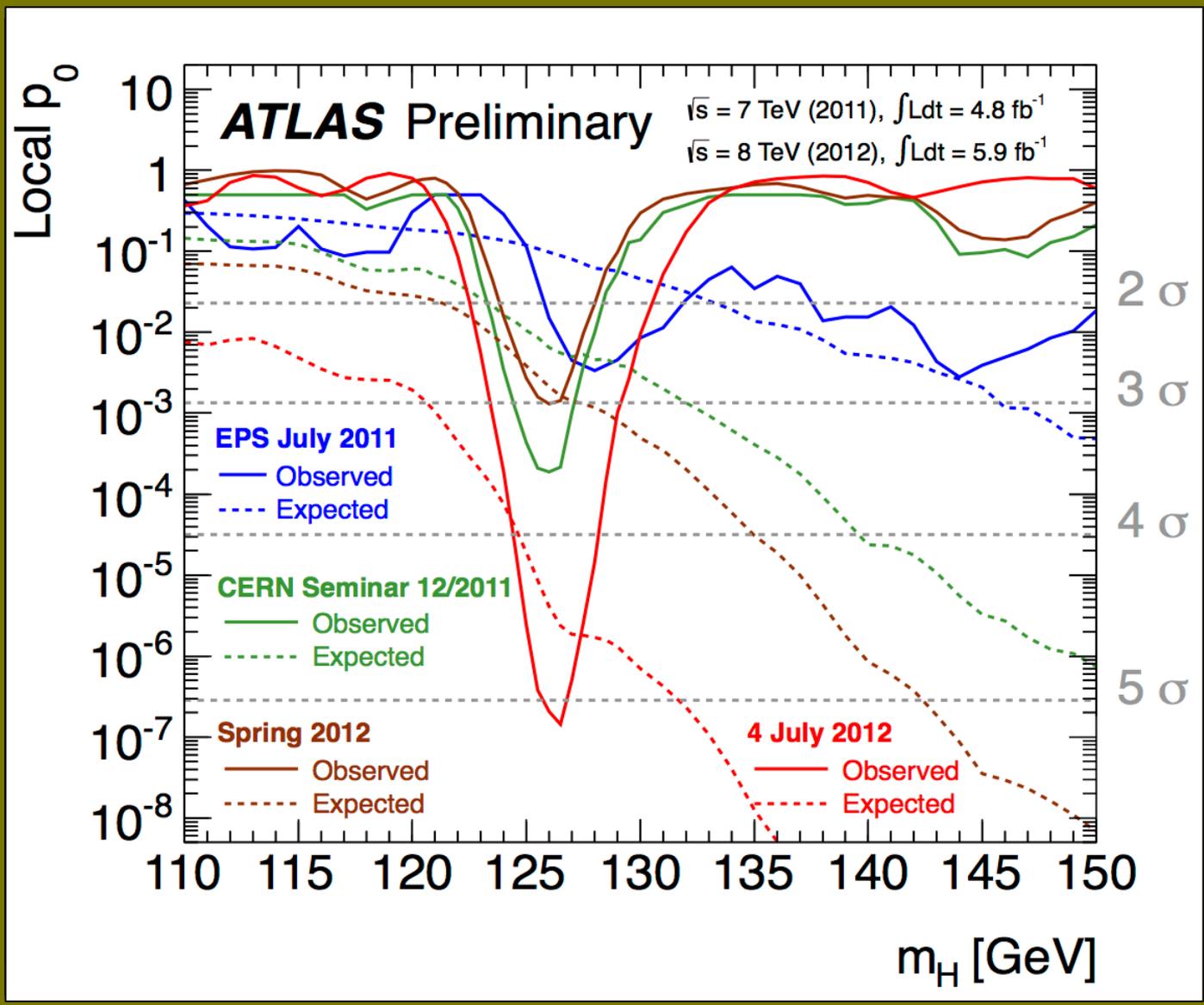


Figure 11: **Confidence intervals** in the  $(\mu, m_H)$  plane for the  $H \rightarrow ZZ^{(*)} \rightarrow 4\ell$ ,  $H \rightarrow \gamma\gamma$ , and  $H \rightarrow WW^{(*)} \rightarrow \ell\nu\ell\nu$  channels, including all systematic uncertainties. The markers indicate the maximum likelihood estimates ( $\hat{\mu}, \hat{m}_H$ ) in the corresponding channels (the maximum likelihood estimates for  $H \rightarrow ZZ^{(*)} \rightarrow 4\ell$  and  $H \rightarrow WW^{(*)} \rightarrow \ell\nu\ell\nu$  coincide).

# Evolution of the excess with time



Energy-scale systematics not included

# Conclusions of papers - ATLAS

## 10. Conclusion

Searches for the Standard Model Higgs boson have been performed in the  $H \rightarrow ZZ^{(*)} \rightarrow 4\ell$ ,  $H \rightarrow \gamma\gamma$  and  $H \rightarrow WW^{(*)} \rightarrow e\nu\mu\nu$  channels with the ATLAS experiment at the LHC using 5.8–5.9 fb<sup>-1</sup> of  $pp$  collision data recorded during April to June 2012 at a centre-of-mass energy of 8 TeV. These results are combined with earlier results [17], which are based on an integrated luminosity of 4.6–4.8 fb<sup>-1</sup> recorded in 2011 at a centre-of-mass energy of 7 TeV, except for the  $H \rightarrow ZZ^{(*)} \rightarrow 4\ell$  and  $H \rightarrow \gamma\gamma$  channels, which have been updated with the improved analyses presented here.

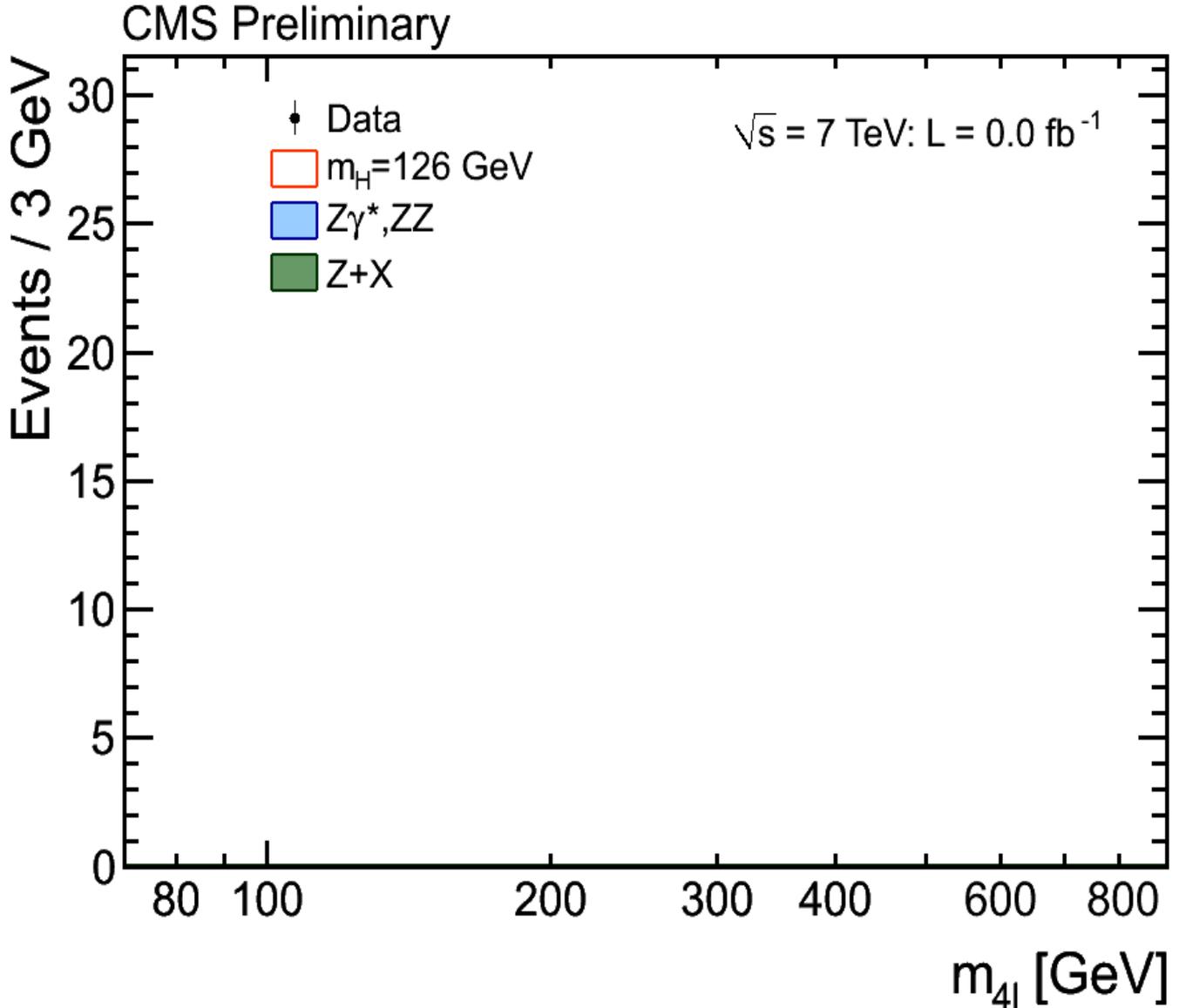
The Standard Model Higgs boson is excluded at 95% CL in the mass range 111–559 GeV, except for the narrow region 122–131 GeV. In this region, an excess of events with significance  $5.9\sigma$ , corresponding to  $p_0 = 1.7 \times 10^{-9}$ , is observed. The excess is driven by the two channels with the highest mass resolution,  $H \rightarrow ZZ^{(*)} \rightarrow 4\ell$  and  $H \rightarrow \gamma\gamma$ , and the equally sensitive but low-resolution  $H \rightarrow WW^{(*)} \rightarrow \ell\nu\ell\nu$  channel. Taking into account the entire mass range of the search, 110–600 GeV, the global significance of the excess is  $5.1\sigma$ , which corresponds to  $p_0 = 1.7 \times 10^{-7}$ .

These results provide conclusive evidence for the discovery of a new particle with mass  $126.0 \pm 0.4$  (stat)  $\pm 0.4$  (sys) GeV. The signal strength parameter  $\mu$  has the value  $1.4 \pm 0.3$  at the fitted mass, which is consistent with the SM Higgs boson hypothesis  $\mu = 1$ . The decays to pairs of vector bosons whose net electric charge is zero identify the new particle as a neutral boson. The observation in the diphoton channel disfavors the spin-1 hypothesis [140, 141]. Although these results are compatible with the hypothesis that the new particle is the Standard Model Higgs boson, more data are needed to assess its nature in detail.

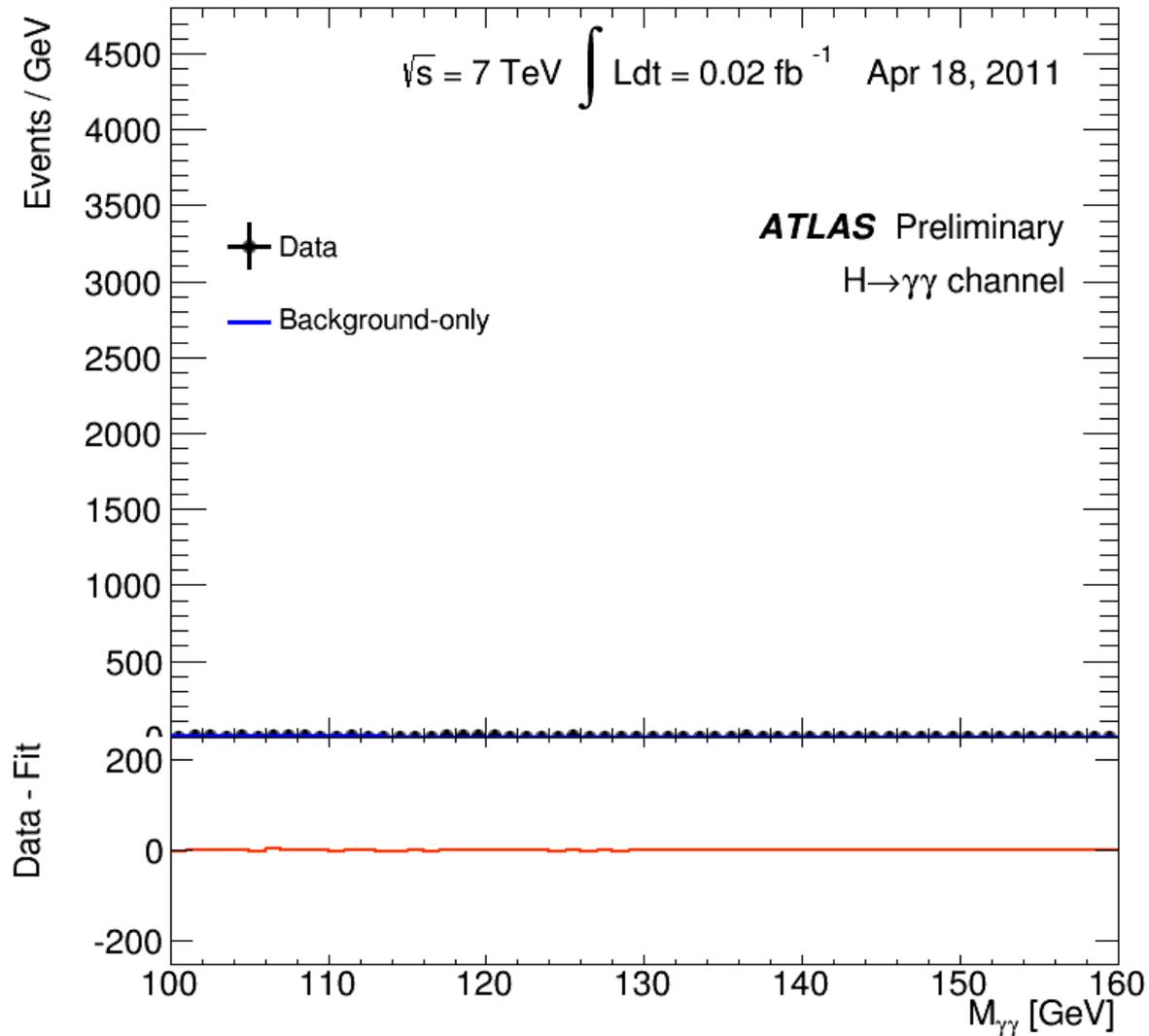
# Conclusions of papers - CMS

Results are presented from searches for the standard model Higgs boson in proton-proton collisions at  $\sqrt{s} = 7$  and 8 TeV in the CMS experiment at the LHC, using data samples corresponding to integrated luminosities of up to  $5.1 \text{ fb}^{-1}$  at 7 TeV and  $5.3 \text{ fb}^{-1}$  at 8 TeV. The search is performed in five decay modes:  $\gamma\gamma$ ,  $ZZ$ ,  $W^+W^-$ ,  $\tau^+\tau^-$ , and  $b\bar{b}$ . An excess of events is observed above the expected background, with a local significance of  $5.0\sigma$ , at a mass near 125 GeV, signalling the production of a new particle. The expected local significance for a standard model Higgs boson of that mass is  $5.8\sigma$ . The global  $p$ -value in the search range of 115–130 (110–145) GeV corresponds to  $4.6\sigma$  ( $4.5\sigma$ ). The excess is most significant in the two decay modes with the best mass resolution,  $\gamma\gamma$  and  $ZZ$ , and a fit to these signals gives a mass of  $125.3 \pm 0.4$  (stat.)  $\pm 0.5$  (syst.) GeV. The decay to two photons indicates that the new particle is a boson with spin different from one. The results presented here are consistent, within uncertainties, with expectations for a standard model Higgs boson. The collection of further data will enable a more rigorous test of this conclusion and an investigation of whether the properties of the new particle imply physics beyond the standard model.

# Evolution of excess: CMS $H \rightarrow ZZ \rightarrow 4l$



# Evolution of excess: ATLAS $H \rightarrow \gamma\gamma$



# Evolution of language

## ► February 2012

- Combined results of **searches for the standard model Higgs boson** in pp collisions at  $\sqrt{s} = 7$  TeV
- By CMS Collaboration, *Phys. Lett. B* 710 (2012) 26-48

## ► July 2012

- **Observation** of a **new boson** with a mass of 125 GeV with the CMS experiment at the LHC
- By CMS Collaboration, *Phys. Lett. B* 716 (2012) 30-61

## ► December 2012

- Study of the Mass and Spin-Parity of the **Higgs Boson Candidate** Via Its Decays to Z Boson Pairs
- By CMS Collaboration, *Phys. Rev. Lett.* 110 (2013) 081803

## ► July 2013

- Measurements of **Higgs boson** production and couplings in diboson final states with the ATLAS detector at the LHC
- By ATLAS Collaboration, *Phys. Lett. B* 726 (2013) 88

## ► March 2015

- Combined Measurement of **the Higgs Boson** Mass in pp Collisions at  $\sqrt{s} = 7$  and 8 TeV with the ATLAS and CMS Experiments
- By ATLAS and CMS Collaborations, *Phys.Rev.Lett.* 114 (2015) 191803

# Data analysis - general picture



1

2 Sampling a reality  
**Experiment**

**Physical phenomena**  
*Described by a theory*

$$e(W_\mu^- W_\mu^+ - W_\mu^+ W_\mu^-)|^2 -$$

$$- W_\mu^+ A_\mu) + i g' c_w (W_\mu^+ Z_\nu -$$

$$- g_\nu Z_\mu + i g' c_w (W_\mu^+ W_\nu - W$$

Described by PDFs,  
depending on  $p$  **unknown**  
parameters with **true values**

$$\theta^{true} = (\theta_1^{true}, \theta_2^{true}, \dots, \theta_p^{true})$$

For example:

$$\theta^{true} = (m_H^{true}, \Delta m_s^{true}, \dots, \sigma_{tot}^{true})$$

3

**Data sample**  
 $x = (x_1, x_2, \dots, x_N)$

For example:  
 $x = (event_1, \dots, event_N)$

In statistics  $x$  is a multivariate **random variable** (each event has many properties, all potential variables)

4

**Data analysis**

5

**Results**

- parameter estimates
- confidence limits
- hypothesis tests

# Data analysis – general picture

The main goal:  
learn more about NATURE

For example, let's suppose the TRUE state of nature is:

Higgs boson exists with the mass of  $m_H(\text{true}) = 134.26 \text{ GeV}$

Make an experiment and obtain a DATA SAMPLE

Use data sample to examine this!

Events collected after some time of LHC running
Event 1
Event 2
...
Event N

Event 1
Object 1
Object 2
...
Object k

If Object 1 == electron
$p_x$
$p_y$
$p_z$
E
...

$N \sim 100/\text{s} \times 10^7 \text{ s/year}$   
 $N \sim 10^9 \text{ events per year}$

Objects  $\equiv$  reconstructed objects  
i. e. electrons, photons, jets, muons ...

# Data analysis, statistics and probability

- ▶ **Data analysis** is the process of transforming raw data into usable information



- ▶ Data analysis uses **statistics** for presentation and interpretation (explanation) of data
  - *Descriptive statistics*
    - Describes the main features of a collection of data in quantitative terms
  - *Inductive statistics*
    - Makes *inference* about a random process from its observed behavior during a finite period of time
- ▶ A mathematical foundation for statistics is the **probability theory**

# Confirmatory and exploratory data analysis

## ► Confirmatory data analysis = Statistical hypothesis testing

- A method of making statistical decisions using experimental data
- Two main methods
  - **Frequentist** hypothesis testing
    - Hypothesis is either true or not
  - **Bayesian** inference
    - Introduces a “degree of belief”

## ► Exploratory data analysis

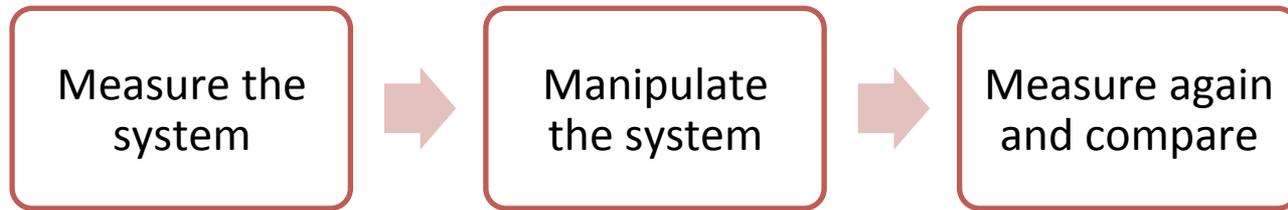
- Uses data to suggest hypothesis to test
- Complements confirmatory data analysis
- Main objectives:
  - Suggest hypothesis about the causes of observed phenomena
  - Asses assumptions on which statistical inference will be based
  - Select appropriate statistical tools and techniques
  - Eventually suggest further data collection

# Quantitative vs graphical techniques

- ▶ **Quantitative techniques** yield numeric or tabular output
  - Hypothesis testing
  - Analysis of variance
  - Point estimation
  - Interval estimation
- ▶ **Graphical techniques**
  - Used for gaining insight into data sets in terms of testing assumptions, model selection, estimator selection ...
  - Provide a convincing mean of presenting results
  - Includes: graphs, histograms, scatter plots, probability plots, residual plots, box plots, block plots, biplots
  - Four main objectives:
    - Exploring the **content** of a data set
    - Finding **structure** in data
    - Checking **assumptions** in statistical models
    - **Communicate** the results of an analysis

# Experimental vs observational studies

## ▶ Experimental studies



- Example: Study of whether and how much a free coffee would improve working performance of scientists in Building 40 at CERN

## ▶ Observational studies

- No experimental manipulation
- Data are gathered and analysed
- Example:
  - Study of correlation between number of beers drunk in a pub on Wednesday evening on performance on the exam the day after

# Probability – basic concepts

## ► Definitions of probability

- Mathematical probability
  - Probability is a basic and an abstract concept
- Frequentist probability
  - Using only measured frequencies
- Bayesian probability
  - Based on a *degree of belief*

# Mathematical probability

- Developed in 1933 by Kolmogorov in his “*Foundations of the Theory of Probability*”
- Define  $\Omega$  as 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
- We define the probability of the occurrence of  $x_i$ ,  $P(x_i)$  to obey the **Kolmogorov axioms**:

$$(a) P(x_i) \geq 0 \quad \text{for all } i$$

$$(b) P(x_i \text{ or } x_j) = P(x_i) + P(x_j)$$

$$(c) \sum_{\Omega} P(x_i) = 1$$

- 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

# Frequentist probability

- ▶ Experiment:
  - $N$  events observed
  - Out of them  $n$  is of type  $x$
- ▶ **Frequentist probability** that any single event will be of type  $x$

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

- ▶ Important restriction: such a probability can only be applied to repeatable experiments
  - For example one can't define a probability that it'll snow tomorrow
  - Although this seems to be a serious problem, a job of scientist is to try to get as close as possible to repeatable experiments and produce reproducible results
- ▶ Frequentist statistics is often associated with the names of *Jerzy Neyman* and *Egon Pearson*

# Bayesian probability

- ▶ Based on a concept of “degree of belief”
- ▶ An operational definition of belief is based on coherent bet by Finneti
  - What’s amount of money one ‘s willing to bet based on her/his belief on the future occurrence of the event

- ▶ Bayesian inference uses Bayes’ formula for conditional probability:

$$P(H | D) = \frac{P(D | H)P(H)}{P(D)}$$

- ▶  $H$  is a **hypothesis**, and  $D$  is the **data**.
- ▶  $P(H)$  is the **prior probability** of  $H$ : the probability that  $H$  is correct before the data  $D$  was seen.
- ▶  $P(D|H)$  is the **conditional probability** of seeing the data  $D$  given that the hypothesis  $H$  is true.  $P(D|H)$  is called the **likelihood**.
- ▶  $P(D)$  is the **marginal probability** of  $D$ .
  - $P(D)$  is the prior probability of witnessing the data  $D$  under all possible hypotheses
- ▶  $P(H|D)$  is the **posterior probability**: the probability that the hypothesis is true, given the data and the previous state of belief about the hypothesis

# Properties of distributions

▶ **Probability density function** (PDF) =  $f(x)$

▶ **Expectation**

- Expectation of any random function  $g(x)$ :

$$E(g) = \int g(x) f(x) dX$$

- Expectation of  $x \equiv$  **mean** of the  $f(x) \equiv$  **expected value** of  $x$  :

$$E(x) = \mu = \bar{x} = \langle x \rangle = \int x f(x) dx$$

▶ **Variance**

$$V(x) = \sigma^2 = E[(x - \mu)^2] = E(x^2) - \mu^2 = \int (x - \mu)^2 f(x) dx$$

- $\sigma$  is called the **standard deviation**

▶  $E(x)$  is a measure of the **location** of the distribution

▶  $V(x)$  is a measure of the **spread** of the distribution

# Moments

$\mu_n = E(x^n)$  is the  $n^{\text{th}}$  algebraic moment

$V_n = E\{[x^n - E(x)]^n\}$  is the  $n^{\text{th}}$  central moment

$\mu'_n = E(|x^n|)$  is the  $n^{\text{th}}$  absolute moment

$V'_n = E\{|x^n - E(x)|^n\}$  is the  $n^{\text{th}}$  absolute central moment

► **The coefficient of skewness**

*A measure of the skewness of the distribution*

$$\gamma_1 = \frac{V_3}{V_2^{3/2}}$$

► **The coefficient of kurtosis**

*A measure of the "peakedness" of the distribution*

$$\gamma_2 = \frac{V_4}{V_2^2} - 3$$

# Poisson distribution

Variable	$r$ , positive integer	 <p>Siméon-Denis Poisson (1781-1840)</p>
Parameters	$\mu$ , positive real number	
Probability function	$P(r; \mu) = \frac{\mu^r e^{-\mu}}{r!}$	
Mean	$E(r) = \mu$	
Variance	$V(r) = \mu$	
Usage example	Number of events $r$ collected after integrated luminosity $\int \mathcal{L} dt$ . Expected number of events is $\mu = \sigma \int \mathcal{L} dt$ . $\sigma$ is the cross section.	
Comments	<ul style="list-style-type: none"> <li>• <math>P(r; \mu)</math> expresses the probability of a number of events occurring in a fixed period of time if these events occur with a known average rate and independently of the time since the last event</li> <li>• <math>\mu</math> represents expected number of events in a given time interval</li> <li>• Time between two successive events is exponentially distributed</li> <li>• Poisson distribution is also called Poissonian</li> </ul>	

# Poisson distribution

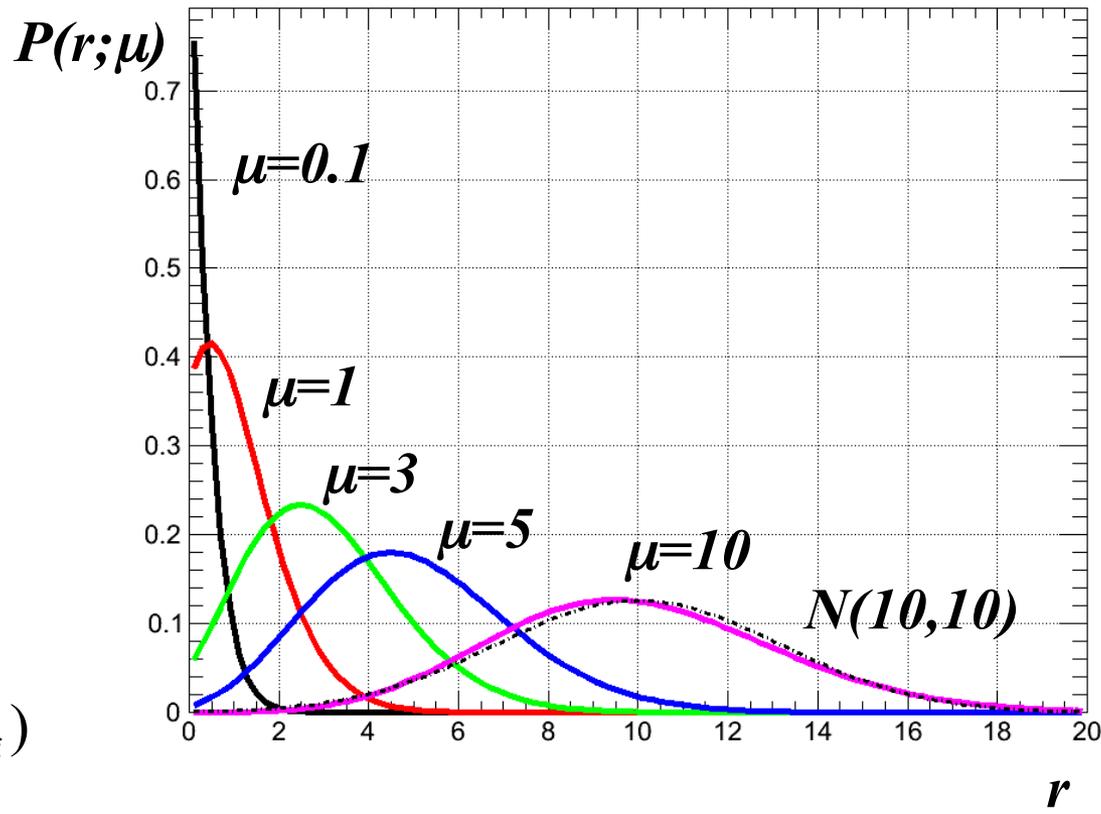
► For a large  $\mu$  Poisson distribution converges towards a Gaussian distribution

$$Pois(r; \mu) \xrightarrow{N \gg} Gauss(r; \mu, \sigma^2 = \mu)$$

► Sum of Poisson distributed random variables also follows a Poisson distribution whose parameter is sum of the component parameters

$$X_i \sim Pois(r; \mu_i)$$

$$Y = \sum_i X_i \sim Pois(r; \sum_i \mu_i)$$



▪ F.g. When combining signal (s) and background (b)

$$P(r; s, b) \sim Pois(r; s+b)$$

# Covariances and correlations

- ▶ Joint PDF for two random variables =  $f(x,y)$
- ▶ The **mean** and the **variance** of  $x$  and  $y$ :

$$\mu_x = E(x) = \iint xf(x,y)dxdy \quad \mu_y = E(y) = \iint yf(x,y)dxdy$$
$$\sigma_x^2 = E[(x - \mu_x)^2] \quad \sigma_y^2 = E[(y - \mu_y)^2]$$

- ▶ **Covariance**

$$\text{cov}(x, y) = E[(x - \mu_x)(y - \mu_y)] = E(xy) - E(x)E(y)$$

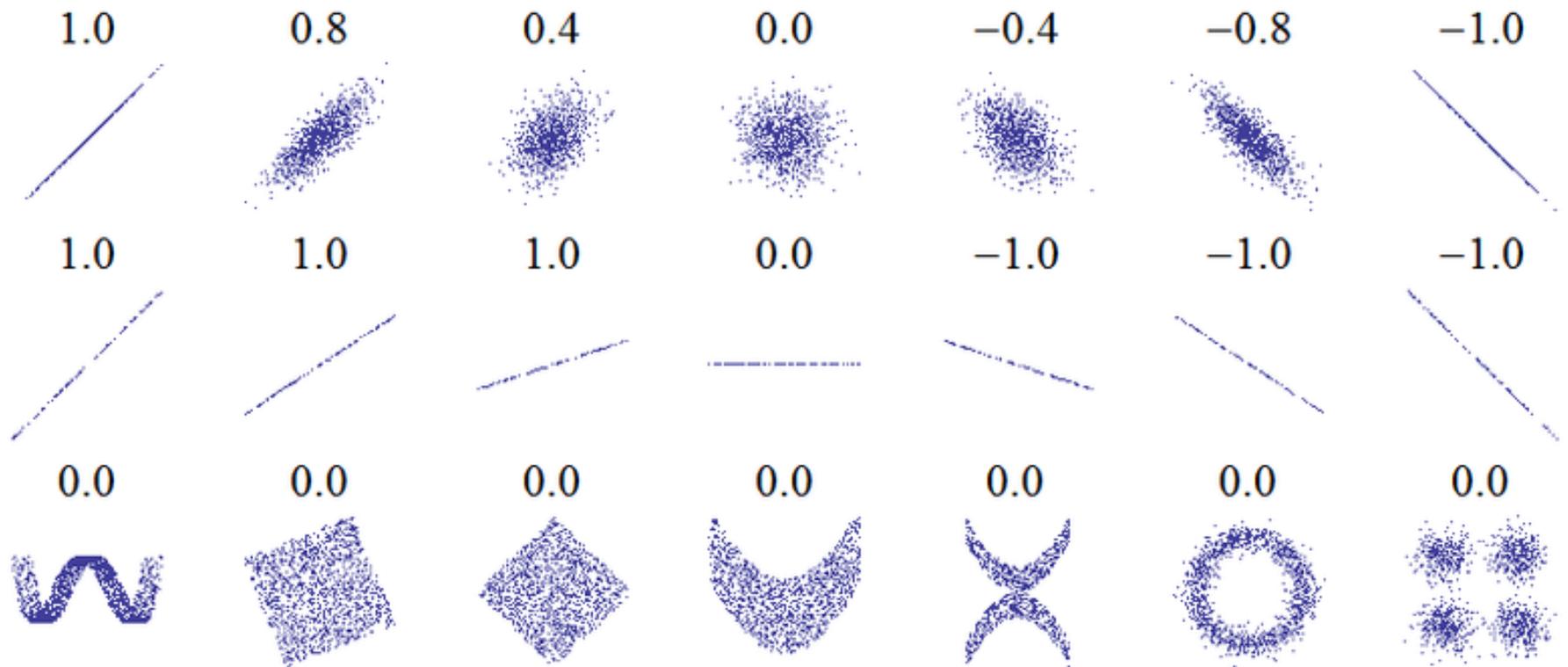
- ▶ **Correlation coefficient**

$$\text{corr}(x, y) = \rho(x, y) = \rho_{xy} = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

- ▶ **Covariance/Variance/Error matrix:**

$$V = \begin{bmatrix} \text{cov}(x, x) & \text{cov}(x, y) \\ \text{cov}(x, y) & \text{cov}(y, y) \end{bmatrix} = \begin{bmatrix} \sigma_x^2 & \rho_{xy} \sigma_x \sigma_y \\ \rho_{xy} \sigma_x \sigma_y & \sigma_y^2 \end{bmatrix}$$

# Correlations - illustration



# General picture



1

Sampling a reality  
**Experiment**

**Physical phenomena**  
*Described by a theory*

$$e(W_\mu^- W_\mu^+ - W_\mu^+ W_\mu^-)|^2 -$$

$$- W_\mu^+ A_\mu + i g' c_w (W_\mu^+ Z_\nu -$$

$$- g_\nu Z_\mu + i g' c_w (W_\mu^+ W_\nu - W$$

Described by PDFs,  
depending on  $p$  **unknown**  
parameters with **true values**

$$\theta^{true} = (\theta_1^{true}, \theta_2^{true}, \dots, \theta_p^{true})$$

For example:

$$\theta^{true} = (m_H^{true}, \Delta m_s^{true}, \dots, \sigma_{tot}^{true})$$

3

**Data sample**  
 $x = (x_1, x_2, \dots, x_N)$

For example:  
 $x = (event_1, \dots, event_N)$

In statistics  $x$  is a multivariate **random variable** (each event has many properties, all potential variables)

4

**Data analysis**

5

**Results**

- parameter estimates
- confidence limits
- hypothesis tests

# Physicists and statisticians

## ● Example: histogram fitting

### Physicists

1. Determining the “best fit” parameters of a curve



2. Determining the errors on the parameters



3. Judging the goodness of a fit

### Statisticians

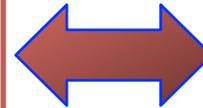
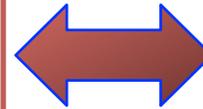
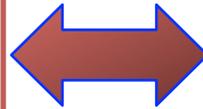
1. Point estimation



2. Confidence interval estimation



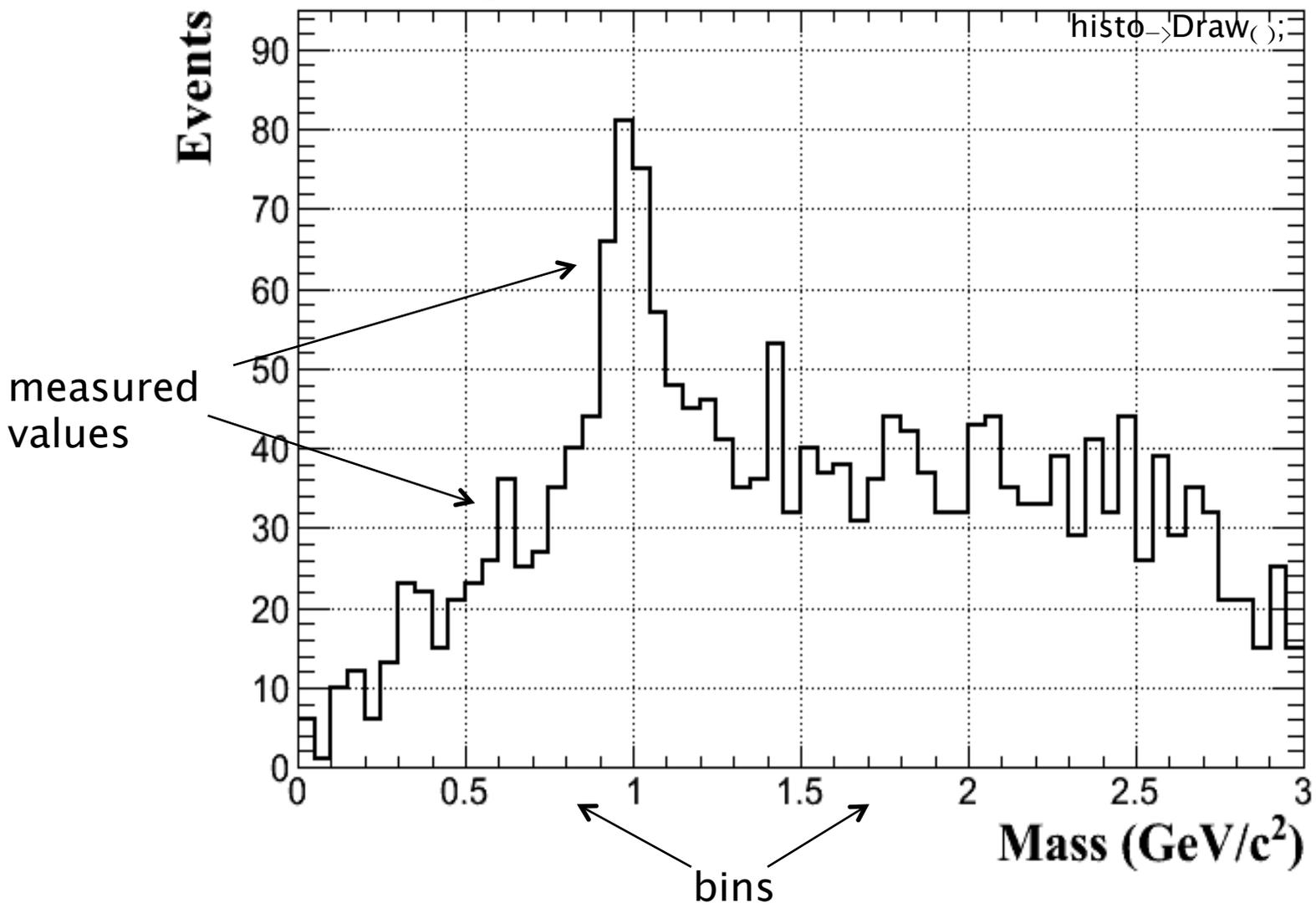
3. Goodness-of-fit (Hypothesis) testing



Adopted from [Baker, Cousins, 1984]

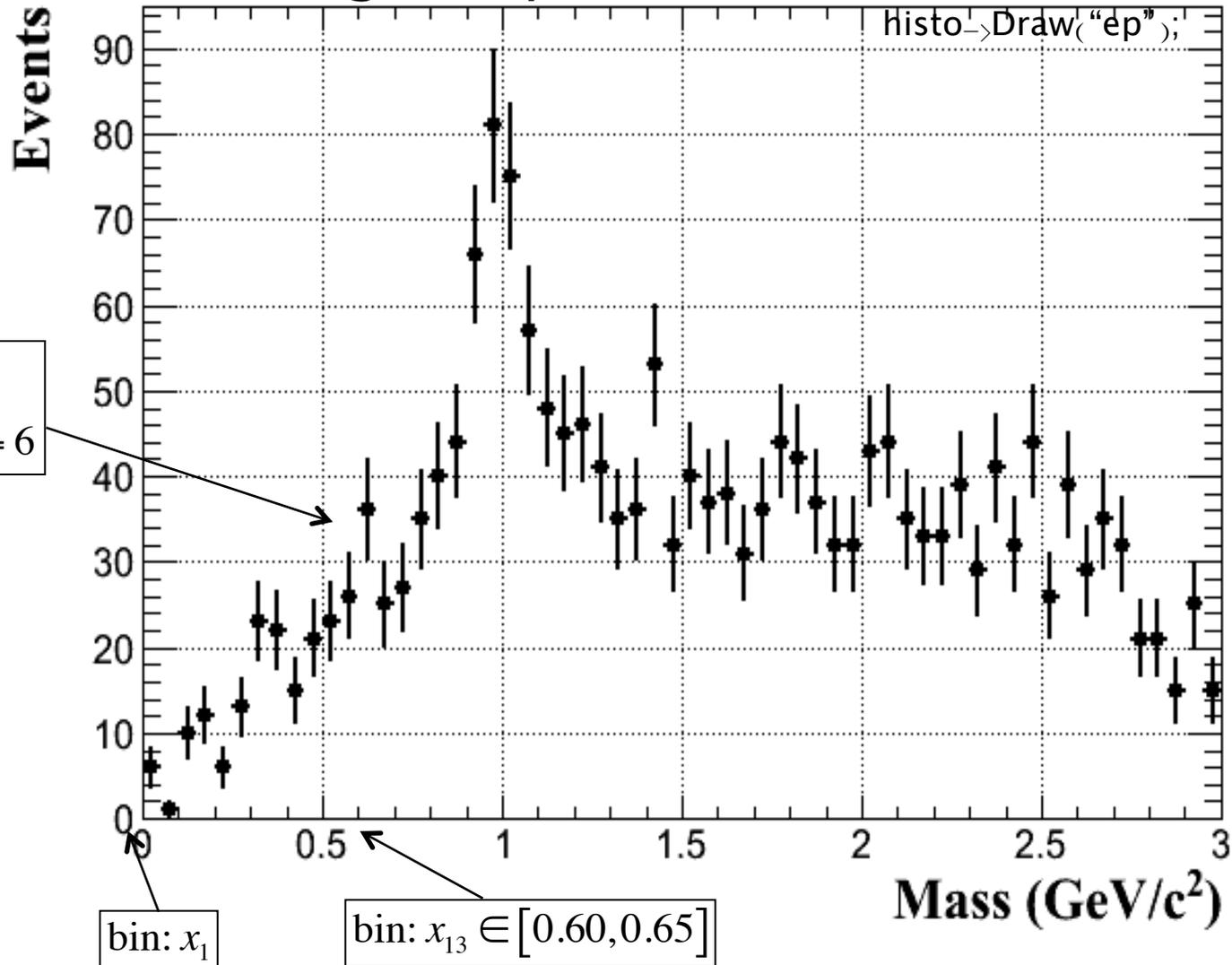
# Example: mass measurement

## Histogram: default representation



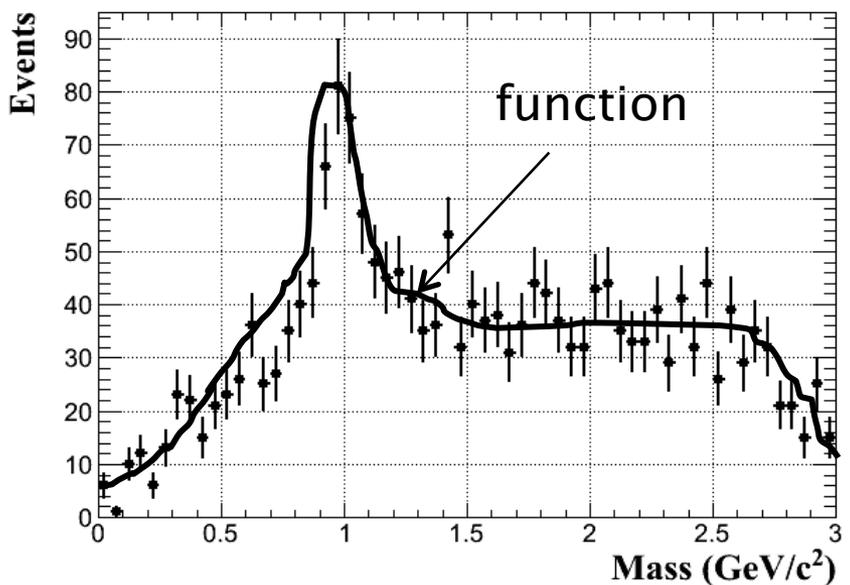
# Example: mass measurement

## Histogram: points with errors

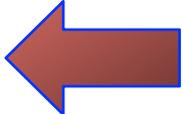


# Example: mass measurements

- Therefore we have
  - a set of precisely known values  $x = (x_1, \dots, x_N)$  – histograms bins
  - At each  $x_i$ 
    - a measured value  $y_i$  – number of events in a given bin
    - a corresponding error on measured value  $\sigma_i$



1 Physical phenomena (theory)



Described by a function, depending on  $p$  **unknown** parameters with **true values**

$$\theta^{true} = (\theta_1^{true}, \theta_2^{true}, \dots, \theta_p^{true})$$

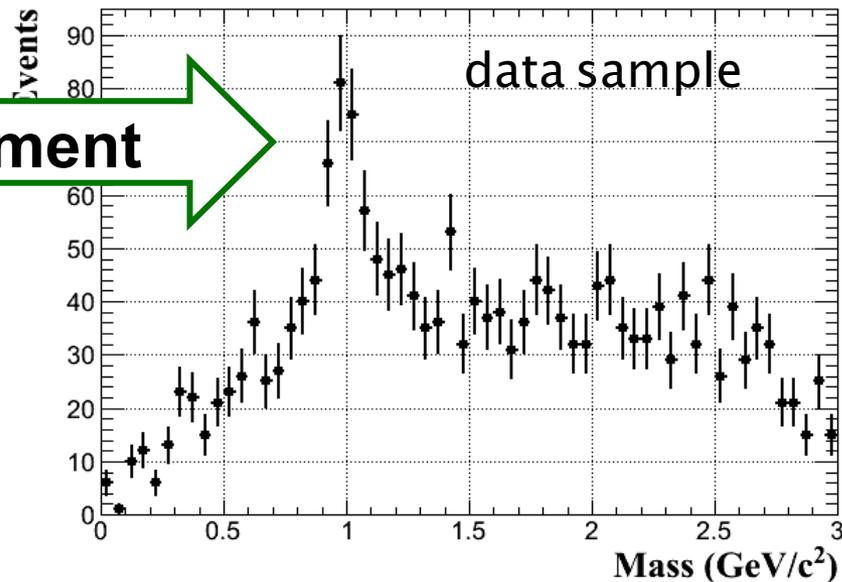
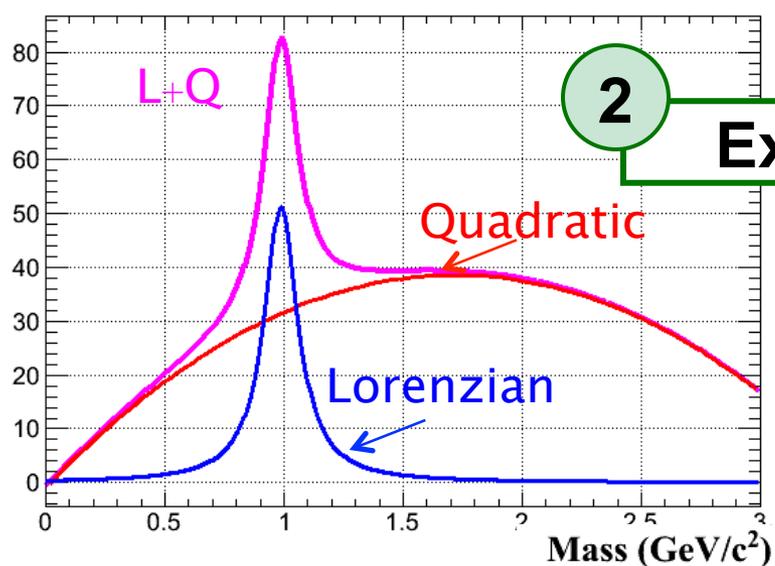
- What we want is to estimate the values of
- This is what we call the parameter ESTIMATOR:

$$\theta_i^{true}$$

$$\hat{\theta}_i$$

# Experiment: sampling the reality

## Physical phenomena

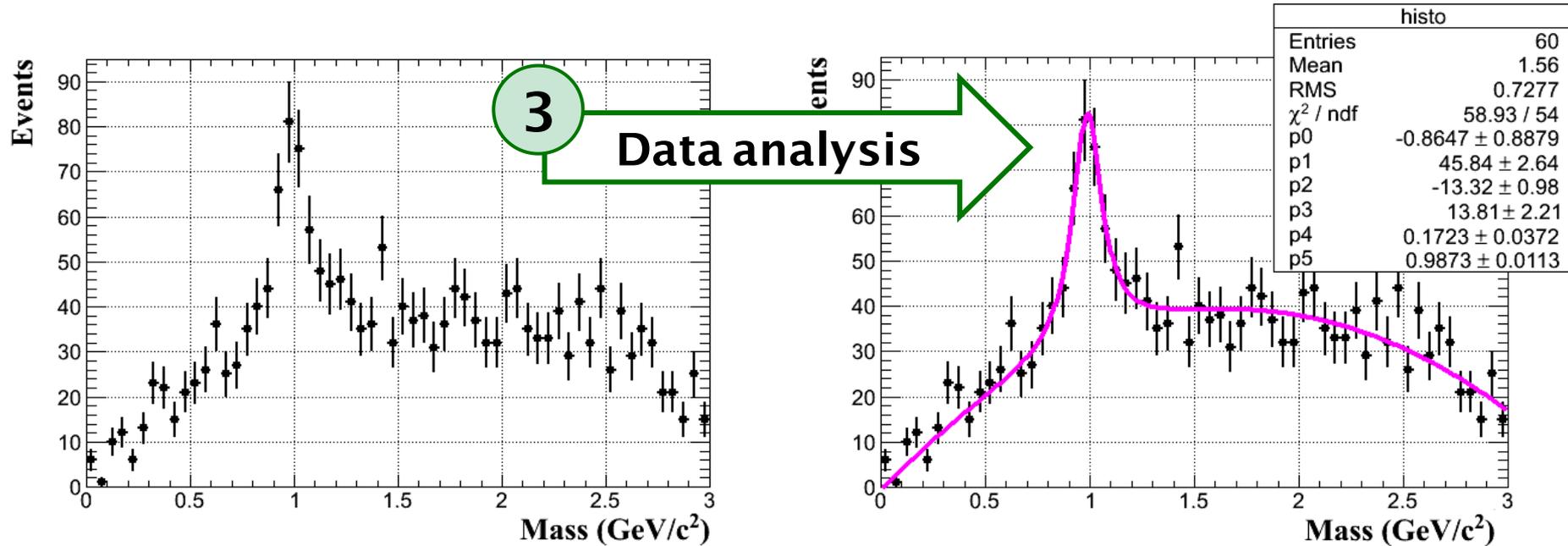


$$\text{Lorentzian} = L(x; D, \Gamma, M) \sim \frac{D\Gamma}{(x^2 - M^2)^2 + 0.25\Gamma^2} \quad \text{Quadratic} = Q(x; A, B, C) \sim A + Bx + Cx^2$$

## Underlying phenomena depends on 6 unknown parameters:

$$F(x; D, \Gamma, M, A, B, C) = L(x; D, \Gamma, M) + Q(x; A, B, C) = F(x; \theta)$$

# Data analysis: estimating parameters



- From data sample we are looking for the **function** that describes the measurements the best
- The parameters of that **function** are estimators of unknown parameters

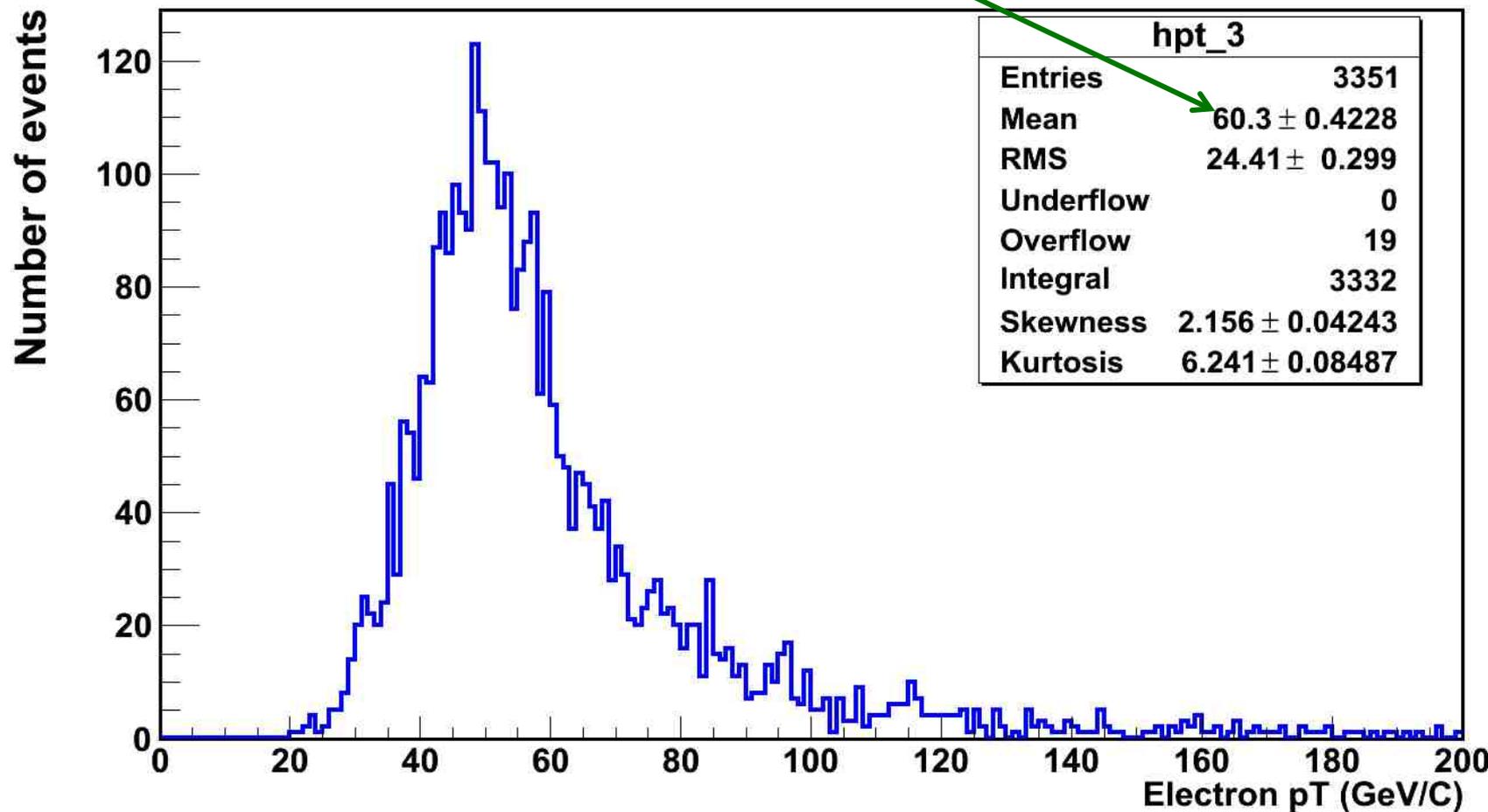
$$F(x; \hat{D}, \hat{\Gamma}, \hat{M}, \hat{A}, \hat{B}, \hat{C}) = L(x; \hat{D}, \hat{\Gamma}, \hat{M}) + Q(x; \hat{A}, \hat{B}, \hat{C}) = F(x; \hat{\theta})$$

# Statistic

- ▶ Be careful: **statistic** is not statistics!
- ▶ Any new random variable (f.g.  $T$ ), defined as a function of a measured sample  $x$  is called a **statistic**  $T = T(x_1, \dots, x_N)$
- ▶ For example, the sample mean  $\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$  is a statistic!
- ▶ A statistic used to estimate a parameter is called an **estimator**
  - For instance, the **sample mean** is a statistic and an estimator for the **population mean**, which is an unknown parameter
  - **Estimator** is a function of the data
  - **Estimate**, a value of estimator, is our “best” guess for the true value of parameter
- ▶ Some other examples of statistics: sample median, variance, standard deviation, quartiles, percentiles, t-statistic, chi-square statistic, kurtosis, skewness etc.

# Estimators in ROOT - example

Notice an influence of the tail on the mean value



# How to find a good estimator?

## The Method of Moments

- Giving consistent and asymptotically unbiased estimators
- But are not as efficient as the maximum likelihood estimates
- Not covered in this lecture

## The Maximum Likelihood Method

- Also giving consistent and asymptotically unbiased estimators
- Widely used in practice

## The Least Squares Method (Chi-Square)

- Giving consistent estimator
- Linear chi-square estimator is unbiased
- Frequently used in histogram fitting

# Likelihood function

- ▶ Assume that observations (events) are independent
  - With the PDF depending on parameters  $\theta$  :  $f(x_i; \theta)$
- ▶ **The probability that all  $N$  events will happen**, i.e. the PDF of  $x$  is, by independence, a product of all single events PDFs

$$P(\mathbf{x}; \theta) = P(x_1, \dots, x_N; \theta) = \prod_{i=1}^N f(x_i; \theta)$$

- ▶ When the variable  $x$  is replaced by the observed data  $\mathbf{x}^0$ , then  $P$  is no longer a PDF
- ▶ It is usual to denote it by  $L$  and call  $L(\mathbf{X}^0; \theta)$  the **likelihood function**
  - Which is now a function of  $\theta$  only

$$L(\theta) = P(\mathbf{X}^0; \theta)$$

- ▶ Often in the literature, and through this lectures, it's convenient to keep  $X$  as a variable and continue to use notation  $L(\mathbf{X}; \theta)$

# Maximum likelihood method

- ▶ Reminder: the probability that all  $N$  independent events will happen is given by the **likelihood function**  $L(\mathbf{x}; \boldsymbol{\theta}) = \prod_{i=1}^N f(x_i; \boldsymbol{\theta})$

- ▶ The principle of maximum likelihood (ML) says:

**The maximum likelihood estimator  $\hat{\boldsymbol{\theta}}$  is the value of  $\boldsymbol{\theta}$  for which the likelihood is a maximum!**

- In words of R. J. Barlow: "You determine the value of  $\theta$  that makes the probability of the actual results obtained,  $\{x_1, \dots, x_N\}$ , as large as it can possible be."
- ▶ In practice it's easier to maximize the **log-likelihood function**

$$\ln L(\mathbf{x}; \boldsymbol{\theta}) = \sum_{i=1}^n \ln f(x_i; \boldsymbol{\theta})$$

- ▶ For  $p$  parameters we get a set of  $p$  likelihood equations

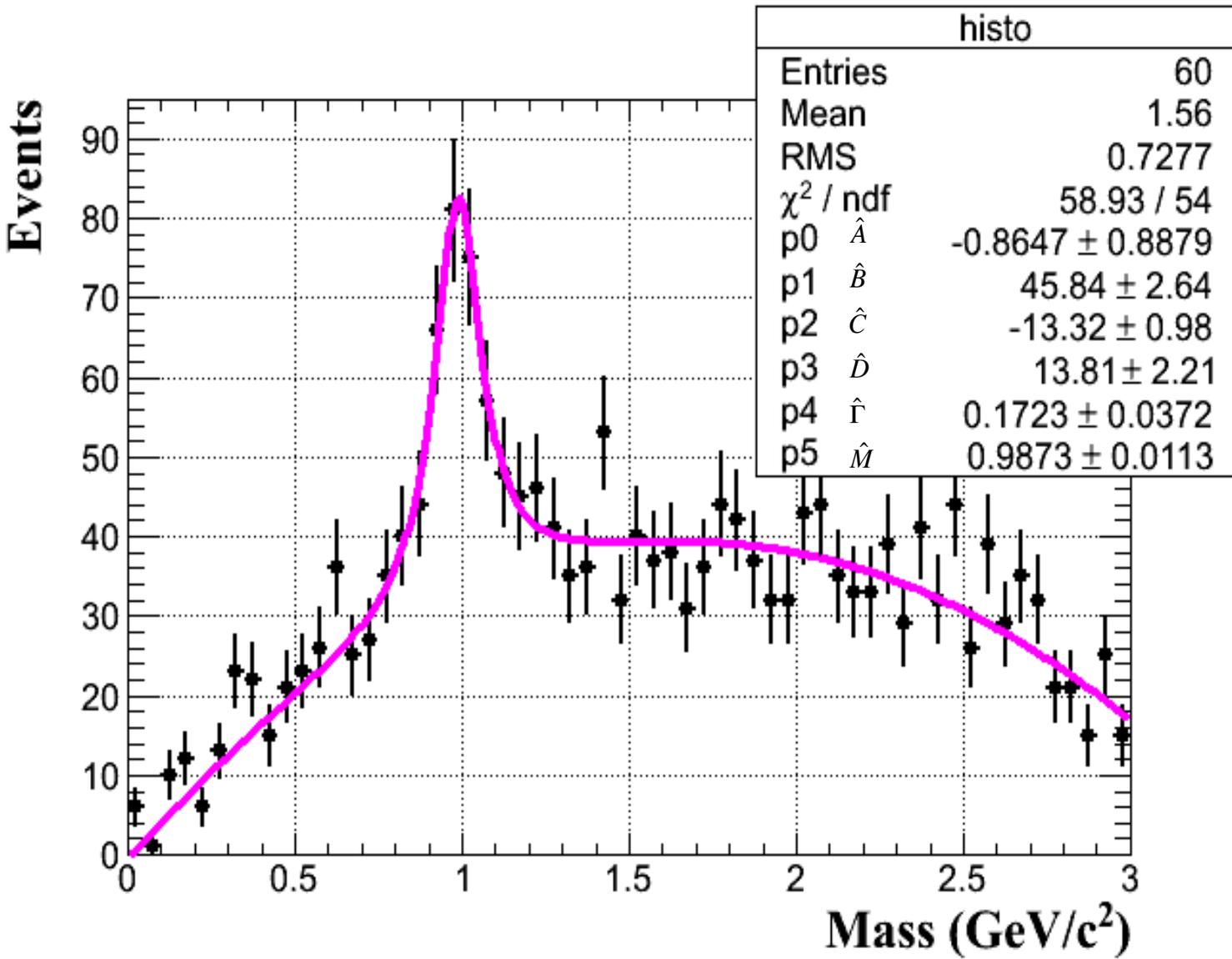
$$\frac{\partial \ln L(\mathbf{x}; \boldsymbol{\theta})}{\partial \theta_j} = 0, \quad j = 1, 2, \dots, p$$

- ▶ It is often more convenient the **minimize  $-\ln L$  Or  $-2\ln L$** 
  - Minimization with MINUIT/MIGRAD or FUMILI in ROOT

# Maximum Likelihood - comments

- ▶ ML estimator is **consistent**
- ▶ ML estimate is approximately **unbiased** and **efficient** for large samples
  - Still useful for small samples, but with extra care!
- ▶ ML estimate is **invariant**
  - A transformation of parameter won't change the answer
- ▶ **ML estimate is not the most likely value of parameter; it is the estimate that makes your data the most likely!**
- ▶ What was presented up to now is sometimes called **unbinned maximum likelihood**
- ▶ **Binned maximum likelihood**: when data are organized in bins
- ▶ Extra care to be taken when the best value of parameters are near imposed limits
- ▶ ML has many advantages, but a few drawbacks too
  - F.g. goodness-of-fit for ML is non-trivial issue, still open and debated

# Example: results of the fit



# Least squares method

- ▶ Suppose we have
  - A set of precisely known values  $\mathbf{x} = (x_1, \dots, x_N)$ 
    - For example histograms bins
  - At each  $x_i$ 
    - a measured value  $y_i$ 
      - For example number of events in the given histogram bin
    - corresponding error on measured value  $\sigma_i$
    - predicted value of measurement that depends on parameters  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)$  we want to estimate:  $F(x_i; \boldsymbol{\theta})$
  - Suppose that measurements are independent
- ▶ To find best estimate of  $\boldsymbol{\theta}$  we minimize the suitably weighted summ of squared differences between measured and predicted values → so called “**least squares**” or “**chi-square**”

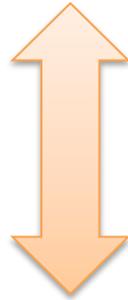
$$\chi^2(\boldsymbol{\theta}) = \sum_{i=1}^N \frac{(y_i - F(x_i; \boldsymbol{\theta}))^2}{\sigma_i^2}$$

# Least squares method

- ▶ If  $y_i$  are Gaussian distributed with variances  $\sigma_i$

$$\chi^2(\boldsymbol{\theta}) = \sum_{i=1}^N \frac{(y_i - F(x_i; \boldsymbol{\theta}))^2}{\sigma_i^2} = -2 \ln L(\boldsymbol{\theta}) + \text{constant}$$

Minimizing chi-square  $\chi^2$



Maximizing log-likelihood  $\ln L$

*or minimizing  $-2\ln L$*

# Pearson's vs Neyman's chi-square

- ▶ If  $y_i$  are Poissonian distributed, there are two choices
  - Reminder first: for Poissonian **variance = mean value** ( $\sigma^2 = \mu$ )
  - So called **Pearson's chi-square** (or "chi-square")

$$\chi^2(\boldsymbol{\theta}) = \sum_{i=1}^N \frac{(y_i - F(x_i; \boldsymbol{\theta}))^2}{F(x_i; \boldsymbol{\theta})}$$

- But now  $\sigma_i$  depends on  $\boldsymbol{\theta}$  which complicates the minimization

- So called **Neyman's chi-square** (or "modified chi-square")
- Minimization simpler
- Easier to combine data with different basic accuracies
- Problem with  $y_i = 0$

$$\chi^2(\boldsymbol{\theta}) = \sum_{i=1}^N \frac{(y_i - F(x_i; \boldsymbol{\theta}))^2}{y_i}$$

- For example in ROOT this bin ignored
- For small samples better use ML

- The best values of parameters  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)$  are found by solving  $p$  equations

$$\frac{\partial \chi^2(\boldsymbol{\theta})}{\partial \theta_i} = 0, \quad i = 1, \dots, p$$

# Reminder

## ● Example: histogram fitting

### Physicists

1. Determining the “best fit” parameters of a curve

2. Determining the errors on the parameters

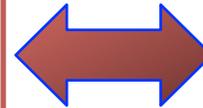
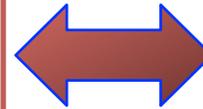
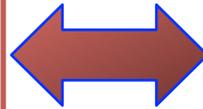
3. Judging the goodness of a fit

### Statisticians

1. Point estimation

2. Confidence interval estimation

3. Goodness-of-fit (Hypothesis) testing



Adopted from [Baker, Cousins, 1984]

# Confidence intervals<sup>7</sup>

For a Gaussian estimator the result of an experiment is usually expressed by

- The parameter's estimated value, plus/minus an estimate of the **standard deviation**,  $\hat{\theta} \pm \sigma_{\hat{\theta}}$

If the pdf is not Gaussian, or in the presence of physical boundaries

- One usually quotes instead an **interval**.

The quoted interval or limit should:

- Objectively communicate the result of the experiment,
- **Communicate incorporated prior beliefs** and relevant assumptions,
- Provide interval that covers the true value of the  $\theta$  with specified probability,
- Make possible to draw conclusions about the parameter.

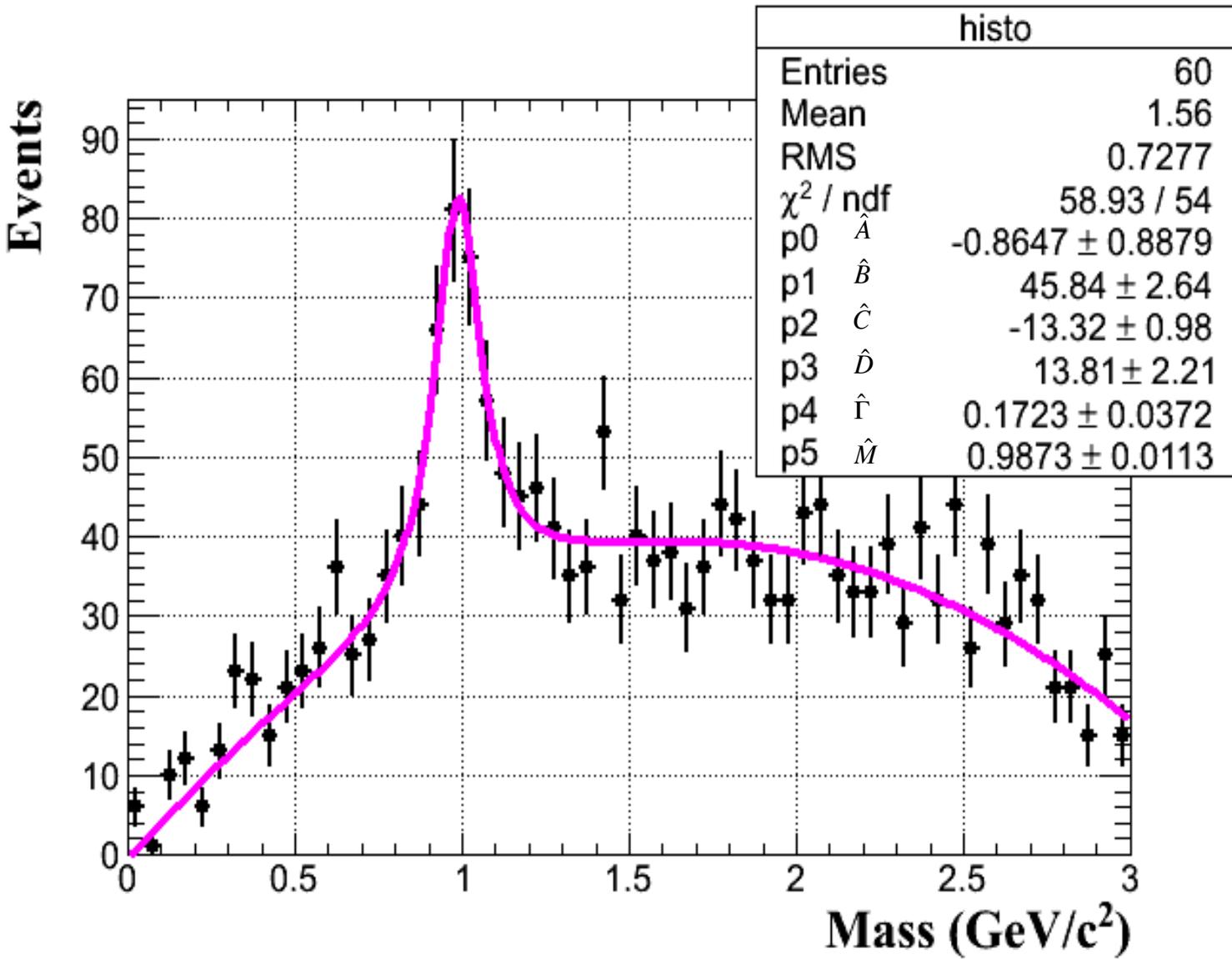
These goals are satisfied in **case of large data sample** by  $\hat{\theta} \pm \sigma_{\hat{\theta}}$ , and in the multi-parameter case by

- The parameter estimates and covariance matrix.

---

<sup>7</sup>Adapted from [Particle Data Group](#).

# Example: results of the fit



# Errors on the ML estimates (1/4)

► How to obtain errors on the parameters estimated by the ML?

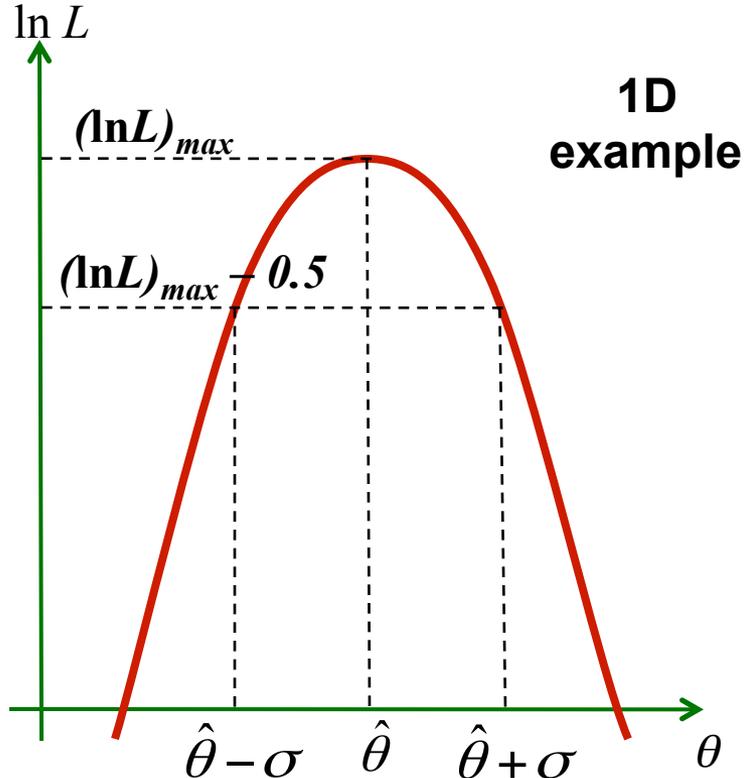
► Option 1: **Matrix inversion**

- Covariance matrix is minus the inverse of the matrix of second derivatives
- Done with MINUIT/HESSE in ROOT

$$\text{cov}^{-1}(\theta_i, \theta_j) = - \left. \frac{\partial^2 L}{\partial \theta_i \partial \theta_j} \right|_{\theta = \hat{\theta}}$$

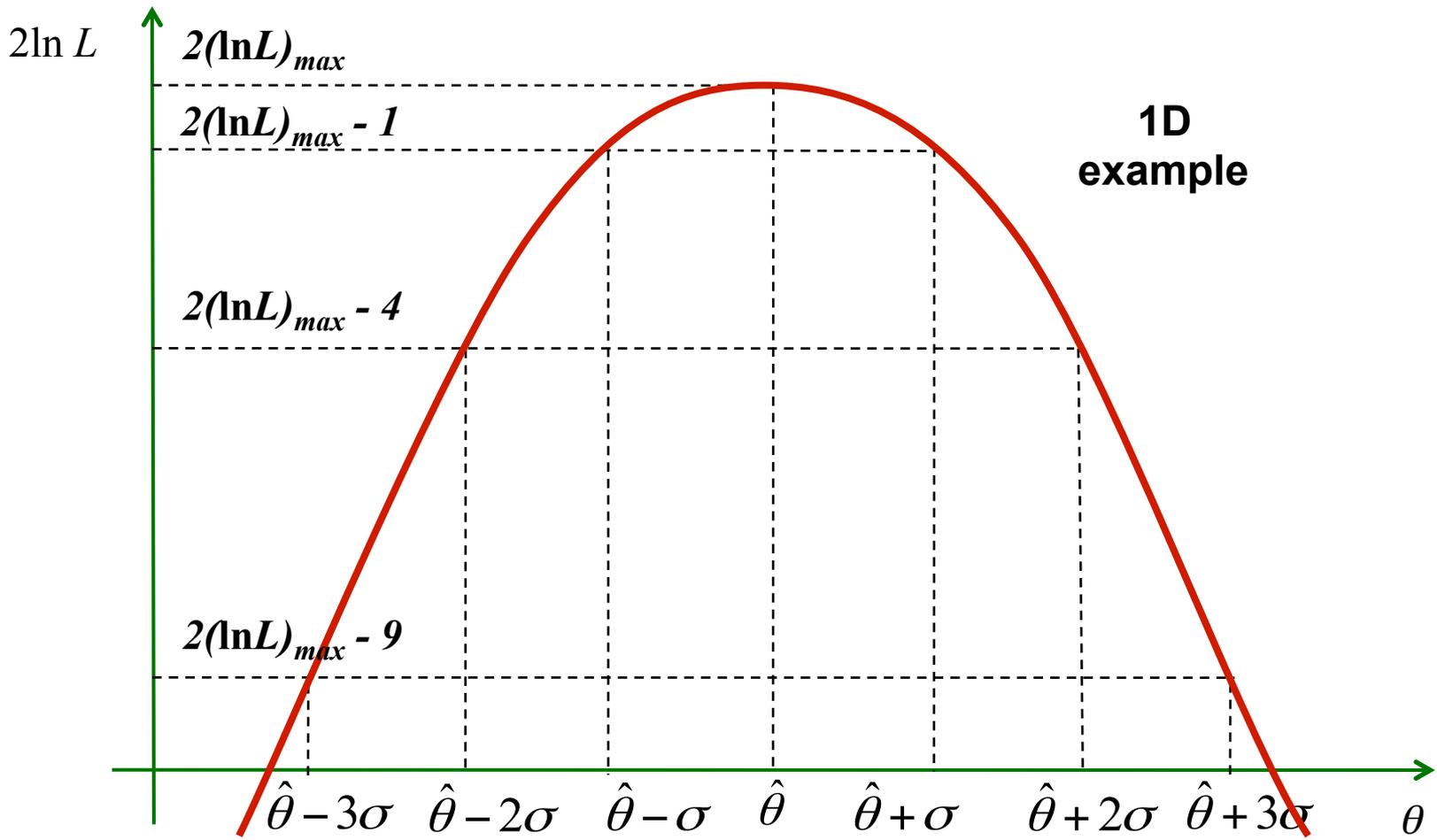
► Option 2: **Log – likelihood curve**

- In the large N limits the likelihood function is Gaussian and the log-likelihood is parabola
- By definition  $(\ln L)_{max} = \ln L(\hat{\theta})$
- $\pm 1\sigma$  limits on  $\theta$  are those values of  $\theta$  for which  $\ln L$  falls by 0.5 from its maximum value  $L_{max}$
- For  $\pm 2\sigma$  ( $\pm 3\sigma$ ) limits  $\ln L$  falls by 2 (4.5)
- Done with MINUIT/MINOS in ROOT



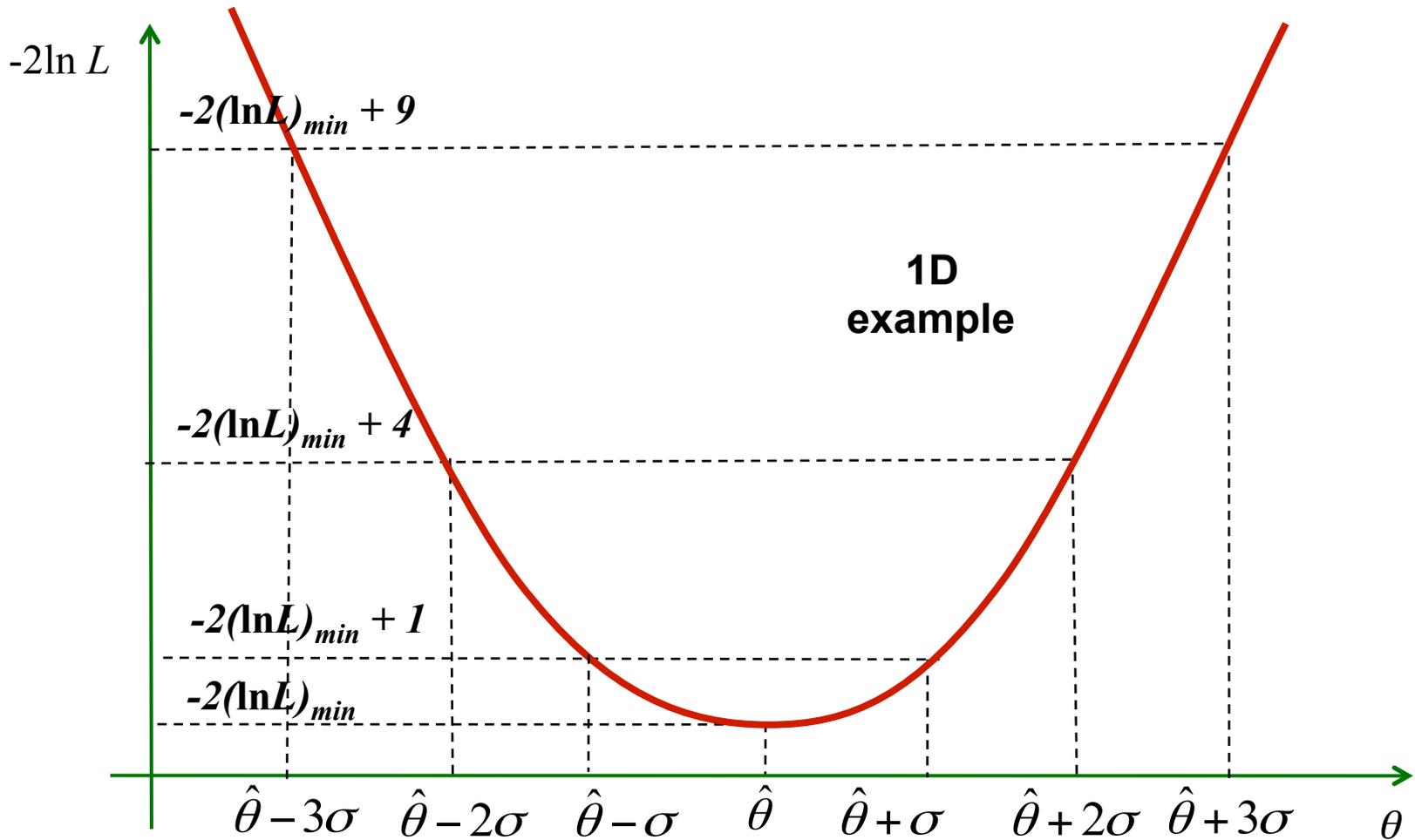
# Errors on the ML estimates (2/4)

► The same, but now maximizing  $2\ln L$



# Maximisation → Minimisation

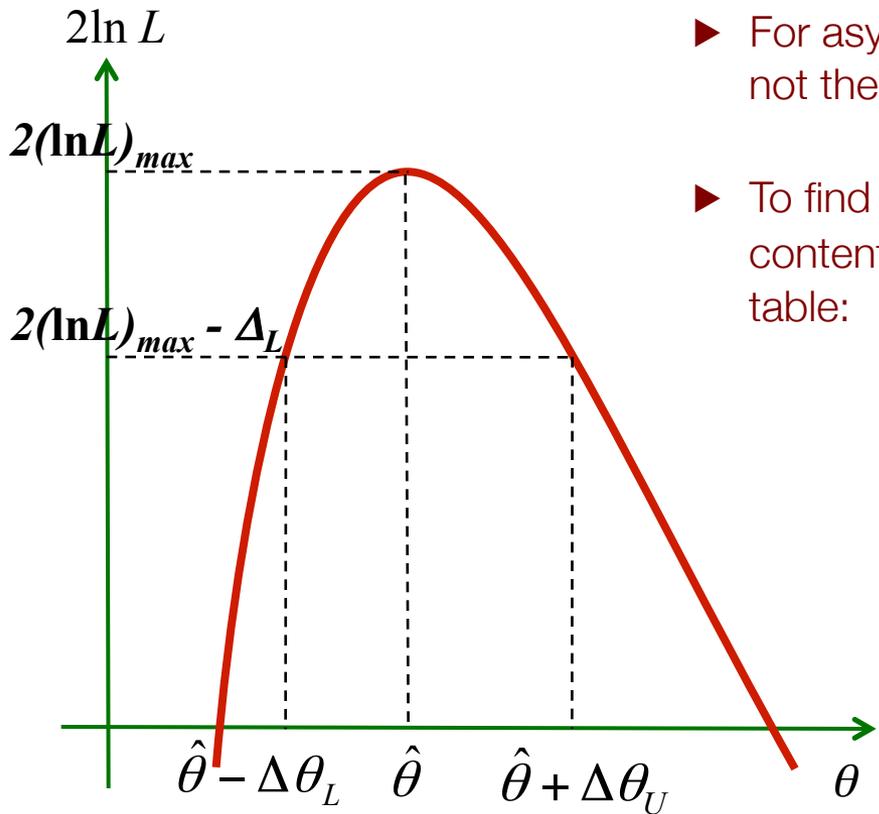
► In reality one is minimizing  $-2\ln L$



# Errors on the ML estimates (3/4)

## ► Asymmetric example

- For finite samples and/or non-linear problems  $\ln L$  is not necessarily parabolic nor symmetric
- Confidence intervals can still be extracted from the  $\ln L$  curve



► For asymmetric  $\ln L$  curve **upper** and **lower** limits on  $\theta$  are not the same

$$\theta = \hat{\theta}^{+\Delta\theta_U}_{-\Delta\theta_L}$$

► To find upper and lower limits with a certain probability content ( $\beta$ ) of the confidence region  $\rightarrow$  use  $\Delta_L$  from the table:

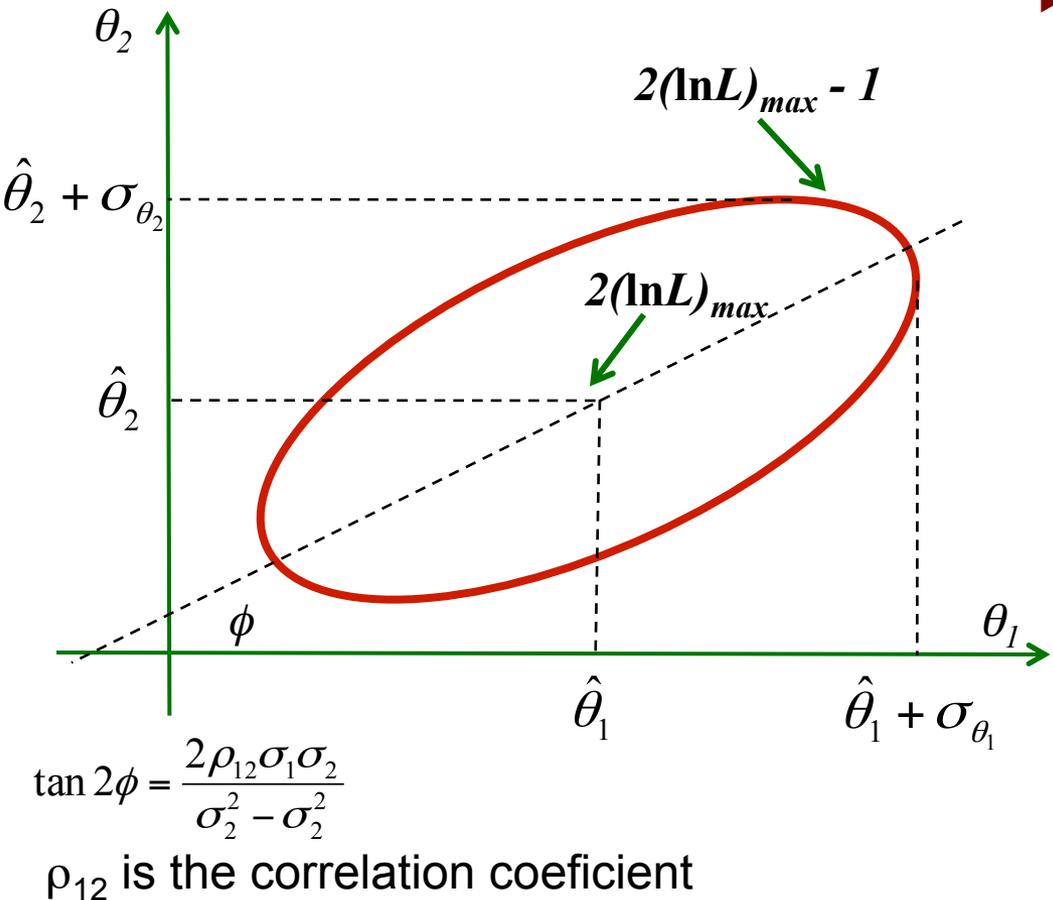
$\Delta_L$	$\beta$ (%)
1	68.27
4	95.45
9	99.73

► ROOT uses Minuit/MINOS to extract limits (errors) in this way

1D example

# Errors on the ML estimates (4/4)

- ▶ 2D example: Standard error ellipse
  - For more information see f.g. PDG
- ▶ This is so called the **plane tangent method**



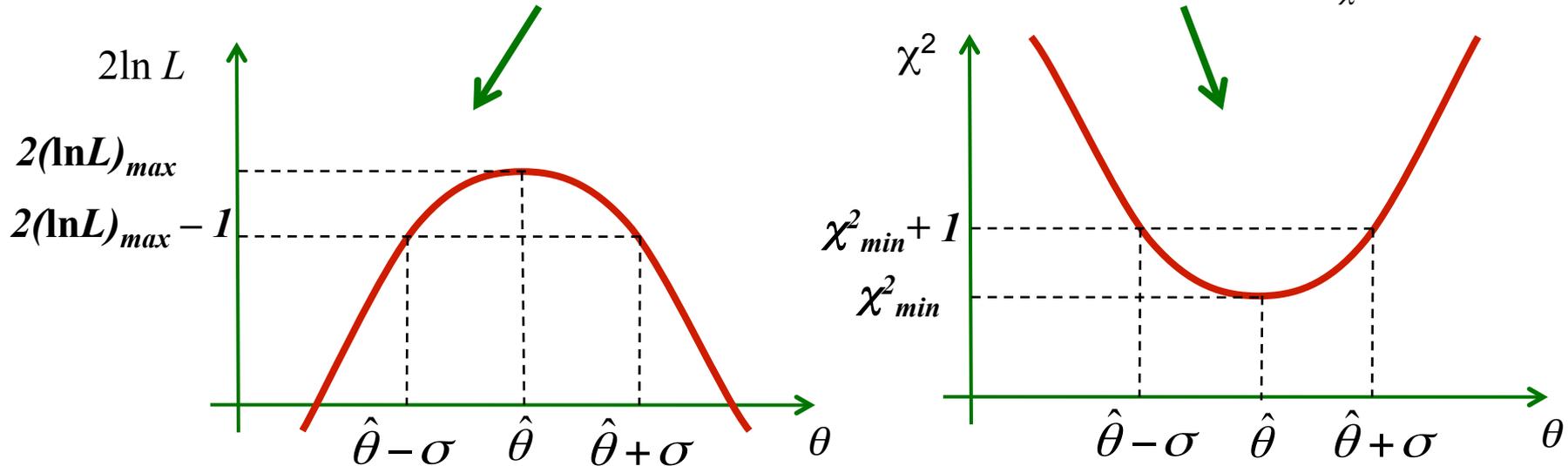
- ▶ ROOT uses **Minuit/MINOS**
  - Works well also with non-regular iso-probability curves
  - Upper and lower limits for parameter  $\theta_i$  are those values of  $\theta_i$  for which
 
$$\max_{\theta_j, j \neq i} [2 \ln L] = 2(\ln L)_{\max} - \Delta_L$$
 with  $\Delta_L$  from the table on the slide before
  - This is OK when interested in errors for only **one** parameter, regardless all others
  - Case of **simultaneous errors** estimate for more parameters → later in this lecture

# Chi-square: Finding errors

► **Errors** (or limits) on parameters are found in the equivalent way as for the ML method

- Matrix inversion
- Shape of  $\chi^2$  around it's minimum value

$$\text{Prob}(2 \ln L \geq 2 \ln L_{\max} - \Delta_L) \Leftrightarrow \text{Prob}(\chi^2) \leq \chi_{\min}^2 + \Delta_{\chi^2}$$



# Multiparameters errors

▶ When interested in simultaneous error estimation on more than one parameter, then the probability content (coverage probability) of the constant  $-2\ln L$  or  $\chi^2$  contours is much smaller than in 1D case

▶ Example (recall 2D Gaussians probabilities):

	$\Delta_L / \Delta_{\chi^2}$	$P_{1D}$	$P_{2D}$
$1\sigma$	1	0.68	0.39
$2\sigma$	4	0.96	0.86

▶ Therefore, to increase the coverage probability we have to increase  $\Delta_L$  or  $\Delta_{\chi^2}$   
 → see the values in the table (from PDG)

**Table 32.2:**  $\Delta\chi^2$  or  $2\Delta\ln L$  corresponding to a coverage probability  $1 - \alpha$  in the large data sample limit, for joint estimation of  $m$  parameters.

$(1 - \alpha)$ (%)	$m = 1$	$m = 2$	$m = 3$
68.27	1.00	2.30	3.53
90.	2.71	4.61	6.25
95.	3.84	5.99	7.82
95.45	4.00	6.18	8.03
99.	6.63	9.21	11.34
99.73	9.00	11.83	14.16

▶ ROOT **Tminuit::Contour** draws contours of constant  $-2\ln L$  or  $\chi^2$  with a given probability coverage

# Bayesian Confidence Intervals

In Bayesian statistics, all knowledge about parameter  $\theta$  is summarized by the posterior pdf  $p(\theta|\mathbf{x})$ ,

$$p(\theta|\mathbf{x}) = \frac{L(\mathbf{x}|\theta)\pi(\theta)}{\int L(\mathbf{x}|\theta')\pi(\theta')d\theta'}$$

which gives the degree of belief for  $\theta$  to have values in a certain region given the data  $\mathbf{x}$ .

- $\pi(\theta)$  is the prior pdf for  $\theta$ , reflecting experimenter's subjective degree of belief about  $\theta$  before the measurement.
- $L(\mathbf{x}|\theta)$  is the likelihood function, i.e. the joint pdf for the data given a certain value of  $\theta$ .
  - $L(\mathbf{x}|\theta)$  **should be published whenever possible**, to enable readers to calculate their own posterior pdf.
- The denominator simply normalizes the posterior pdf to unity.

# Uncertainty in physics

The sources of uncertainty in measurement<sup>9</sup>:

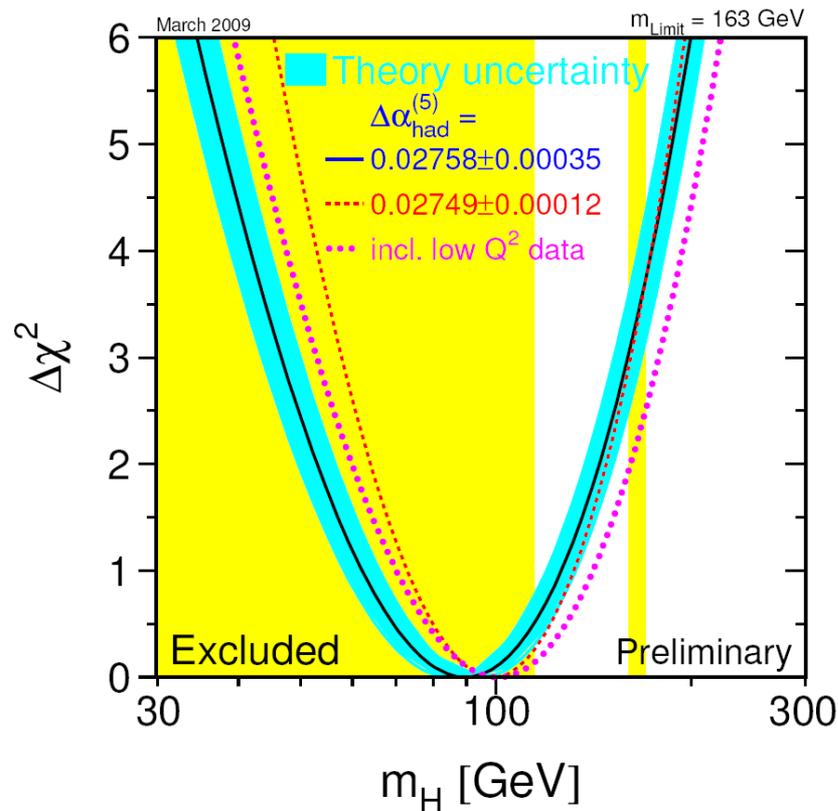
- **Incomplete definition** of the measurand; or its imperfect realization
- **Non-representative sampling**
- inadequate knowledge of the effects of environmental conditions; or imperfect measurements of these conditions
- **Personal bias** in reading instruments
- **Finite instrument resolution**
- Inexact values of measurement standards and reference materials
- **Inexact values of constants** and other parameters obtained from external sources and used in the data-reduction algorithm
- **Approximations and assumptions** incorporated in the measurement procedure
- **Variations of repeated observations** of the measurand under apparently identical conditions

---

<sup>9</sup>Adapted from the The International Organization for Standardization (ISO) Guide to the Expression of Uncertainty in Measurement.

# Example

## Higgs boson mass constraints from Electroweak precision tests

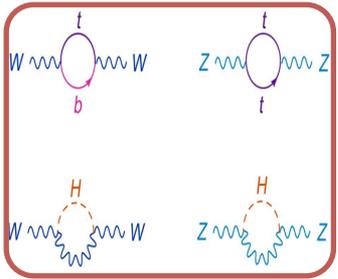


# Method



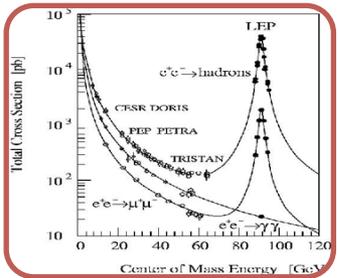
## Step 1 – Very precise measurements of SM

- Measure SM parameters extremely well
- $\alpha$ ,  $M_Z$ ,  $G_F$
- $\mu$  lifetime,  $(g-2)_e$ , LEP ...



## Step 2 – Predictions (assuming Higgs boson)

- Calculate quantum corrections to other observables
  - $m_W$ ,  $A_{LR}$ ,  $\sin^2\theta_W$  ...
- Depending on  $\alpha$ ,  $M_Z$ ,  $G_F$ , but also on  $m_t$ ,  $m_H$  ...



## Step 3 – Precise electroweak measurements

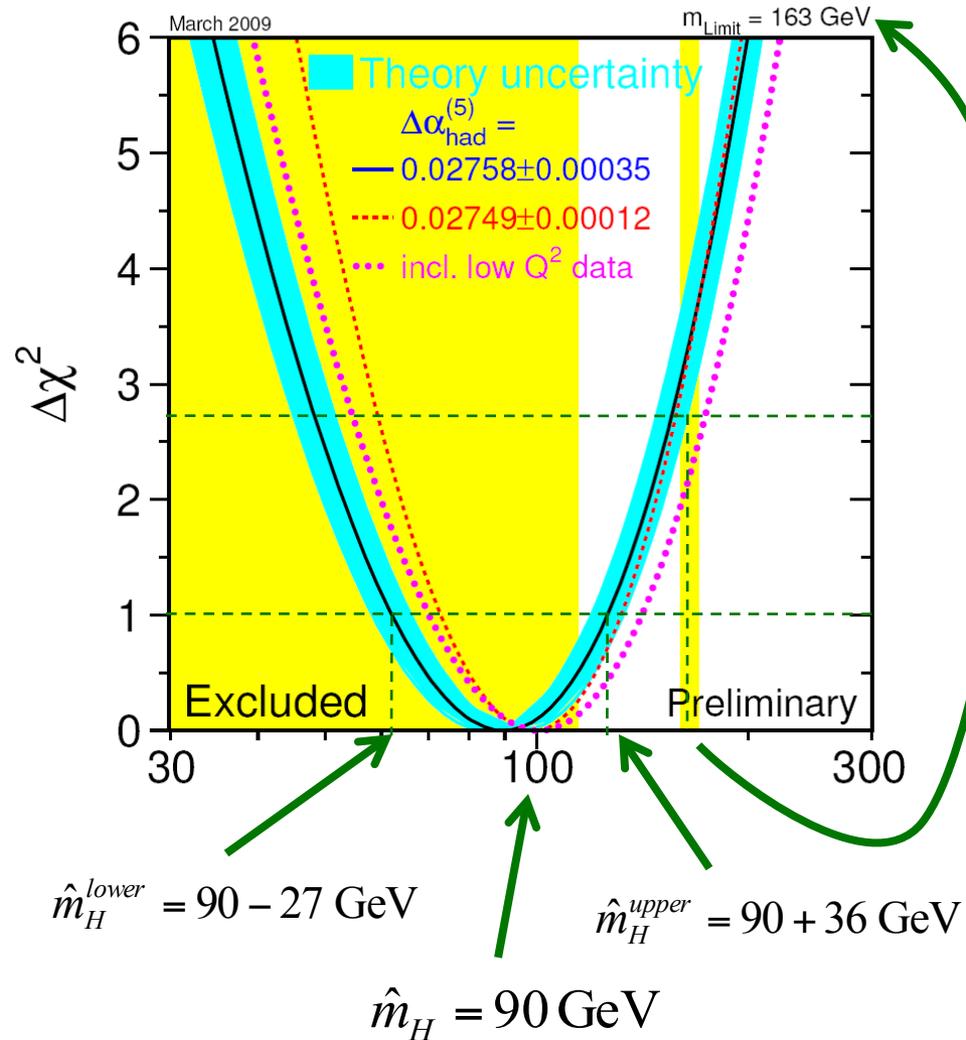
- Measure very precisely observables from Step 2
- @ SLC, LEP, Tevatron ...

# Results from step 2 and 3



March 2009

# The best fit



- From the LEP Electroweak Working group:

- “The preferred value for its mass, corresponding to the minimum of the curve, is at 90 GeV, with an experimental uncertainty of +36 and -27 GeV (at 68 percent confidence level derived from  $\Delta\chi^2 = 1$  for the black line, thus not taking the theoretical uncertainty shown as the blue band into account).”
- “The precision electroweak measurements tell us that the mass of the Standard-Model Higgs boson is lower than about 163 GeV (one-sided 95 percent confidence level upper limit derived from  $\Delta\chi^2 = 2.7$  for the blue band, thus including both the experimental and the theoretical uncertainty).”

# Mass of a new boson

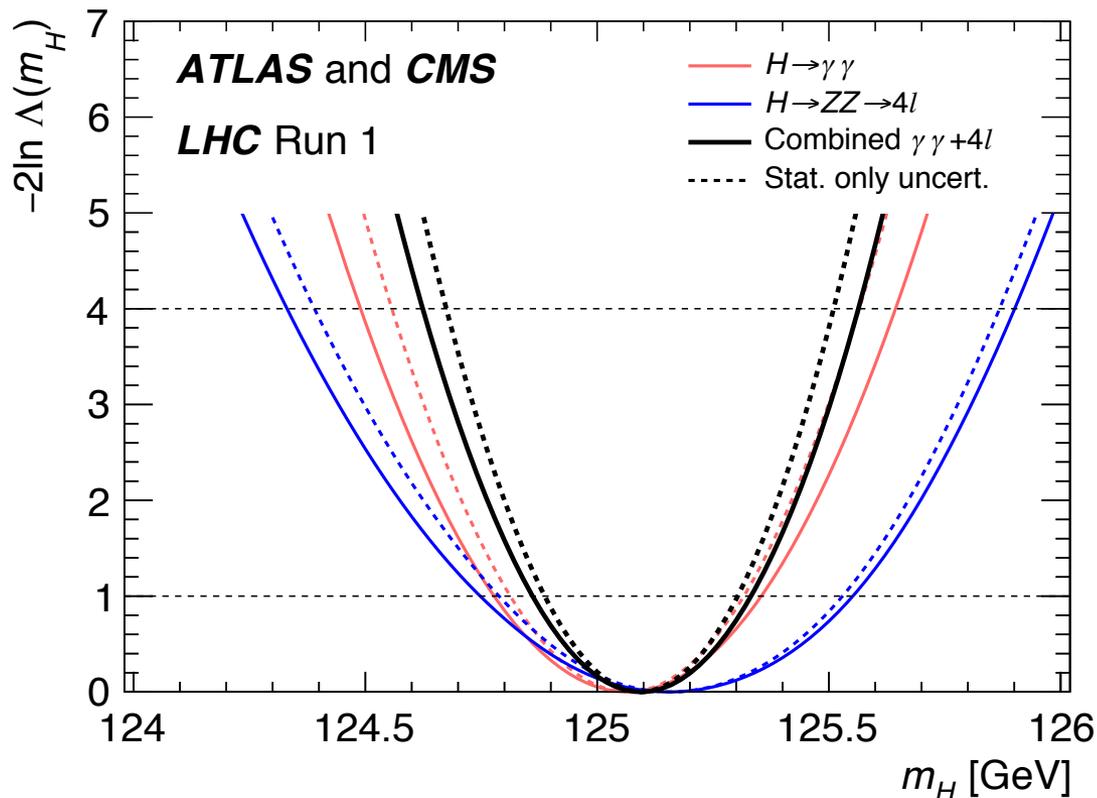


Figure 1: Scans of twice the negative log-likelihood ratio  $-2 \ln \Lambda(m_H)$  as functions of the Higgs boson mass  $m_H$  for the ATLAS and CMS combination of the  $H \rightarrow \gamma\gamma$  (red),  $H \rightarrow ZZ \rightarrow 4\ell$  (blue), and combined (black) channels. The dashed curves show the results accounting for statistical uncertainties only, with all nuisance parameters associated with systematic uncertainties fixed to their best-fit values. The 1 and 2 standard deviation limits are indicated by the intersections of the horizontal lines at 1 and 4, respectively, with the log-likelihood scan curves.

$$\begin{aligned} m_H &= 125.09 \pm 0.24 \text{ GeV} \\ &= 125.09 \pm 0.21 \text{ (stat.)} \pm 0.11 \text{ (syst.) GeV,} \end{aligned}$$

# Reminder

## ● Example: histogram fitting

### Physicists

1. Determining the “best fit” parameters of a curve

2. Determining the errors on the parameters

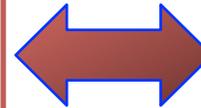
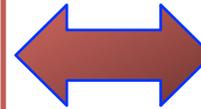
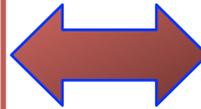
3. Judging the goodness of a fit

### Statisticians

1. Point estimation

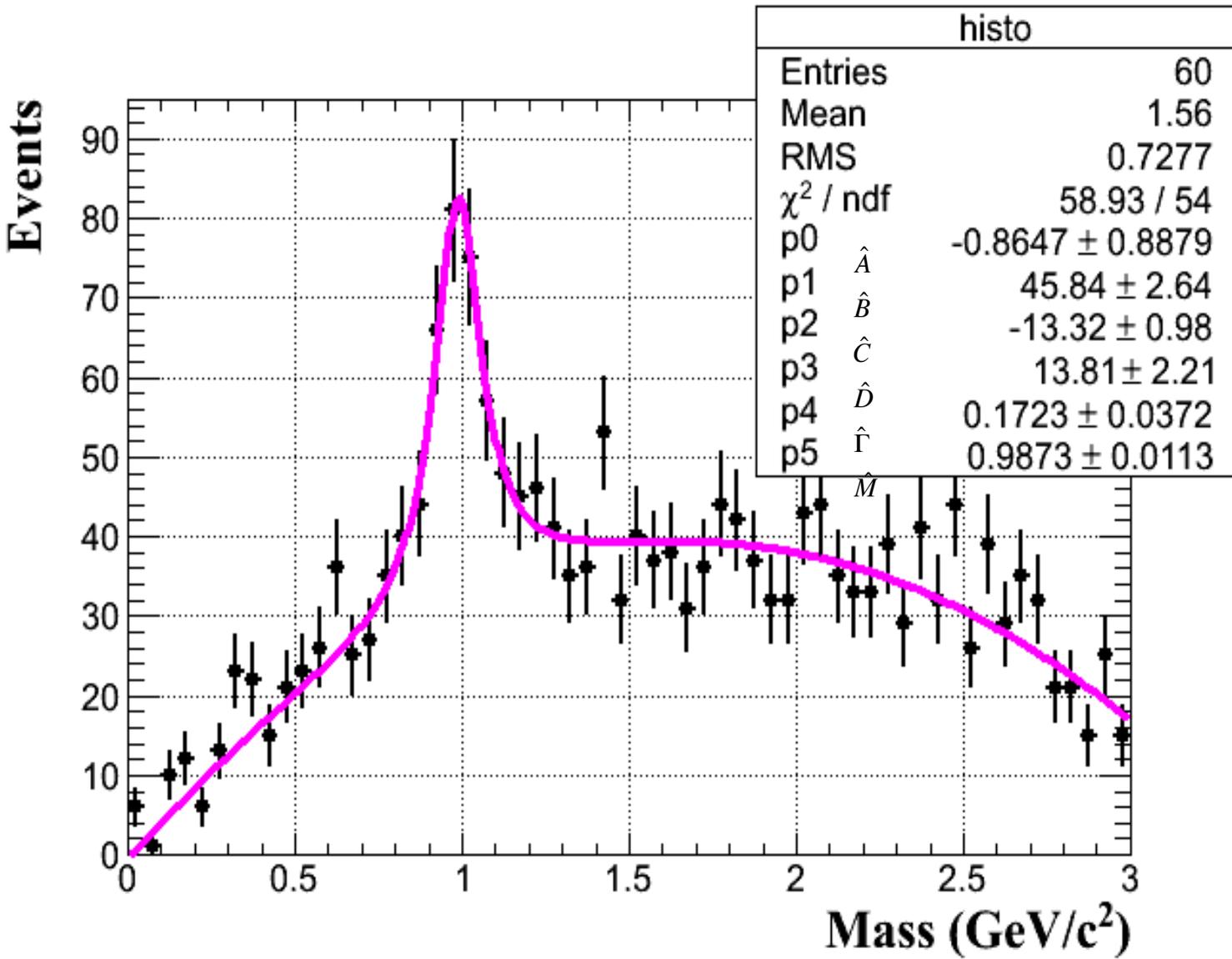
2. Confidence interval estimation

3. Goodness-of-fit testing



Adopted from [Baker, Cousins, 1984]

# Example: results of the fit



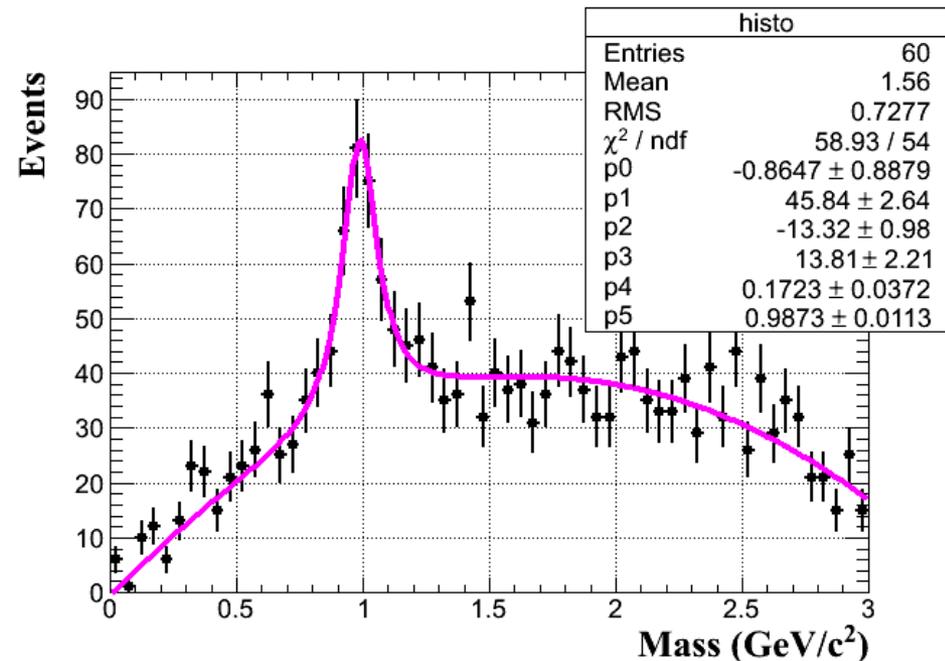
# Goodnes-of-fit tests

## ► We are now interested in this kind of questions

- Is the fit good or not?
- How **significant** is discrepancy between data and obtained functional form?
- How well does the vector of measurements in the histogram  $\mathbf{n} = (n_1, \dots, n_k)$  compare with predicted values  $\mathbf{v} = E[\mathbf{n}] = (v_1, \dots, v_k)$ ?

## ► These questions can be answered with a **goodnes-of-fit test**

- Which is itself a part of a so called HYPOTHESIS TESTING

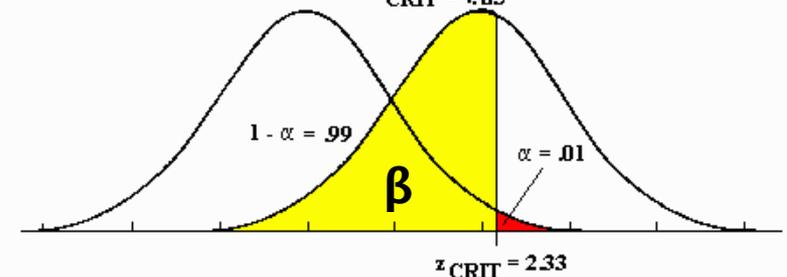
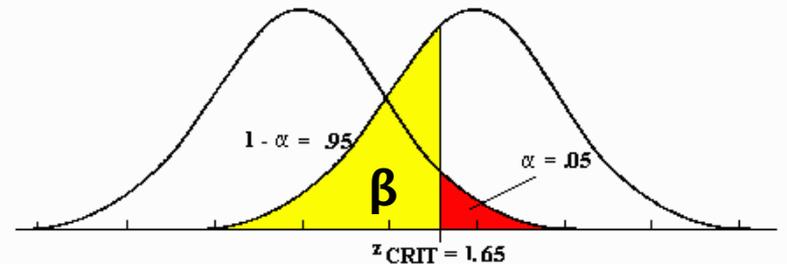
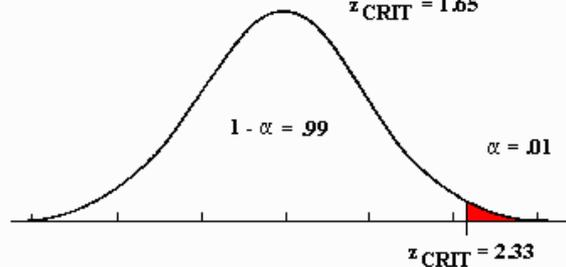
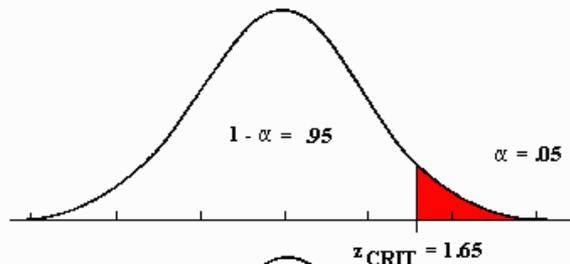


# Hypothesis testing: courtroom trial

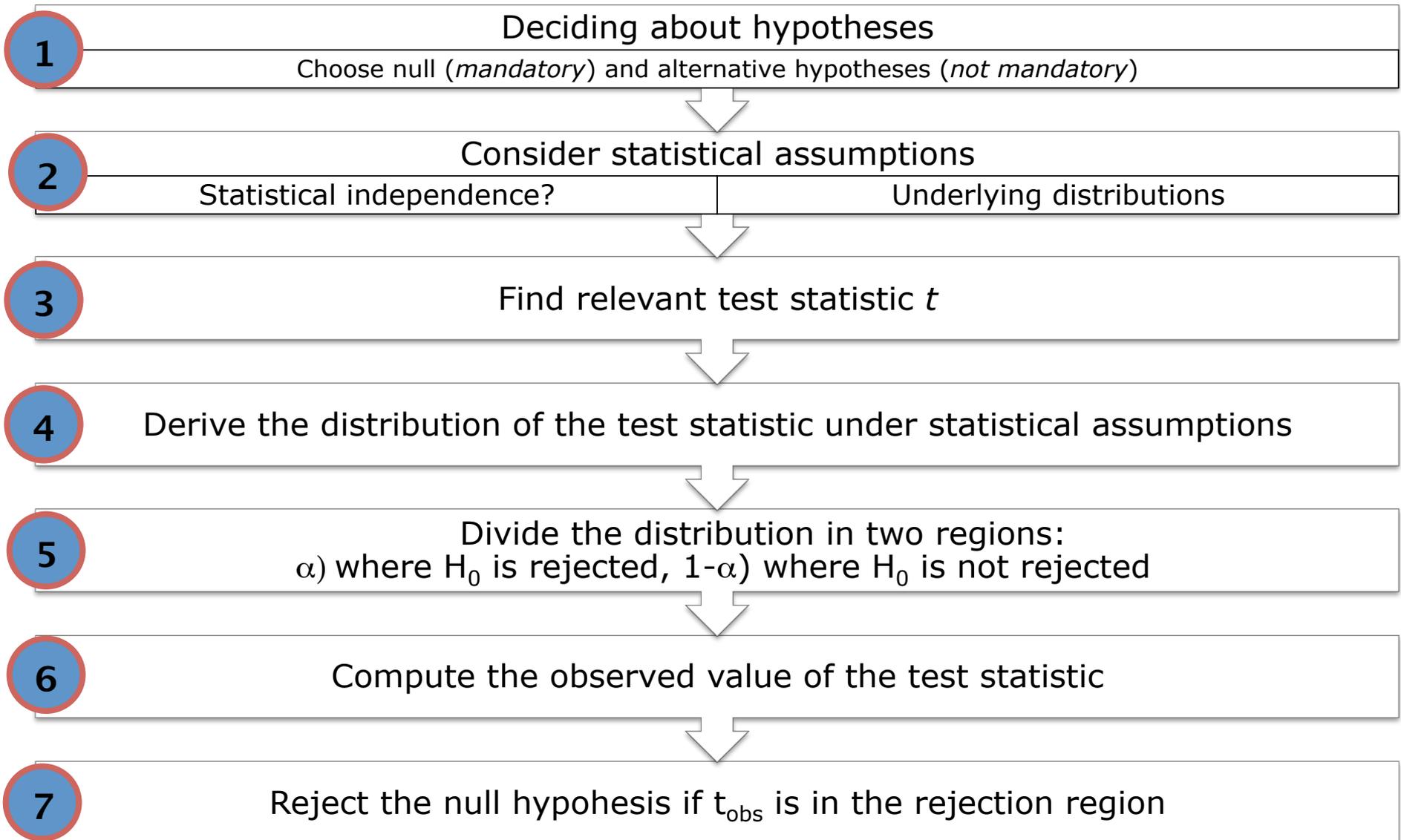
- ▶ Consider a criminal trial
- ▶ There is one simple rule: a defendant is considered **not guilty** as long as his **guilt is not proven**
- ▶ The prosecutor tries to prove the guilt of the defendant
  - Only when there is **enough evidence** the defendant is **convicted**
- ▶ We start with two hypotheses:
  - $H_0$ : the defendant is not guilty → NULL HYPOTHESIS
  - $H_1$ : the defendant is guilty → ALTERNATIVE HYPOTHESIS
- ▶ Null hypothesis is considered accepted for time being
- ▶ Common sense: the hypothesis of innocence is rejected only if the error is very unlikely
  - We don't want to convict innocent person!
  - This is called **Error of the first kind** and we want it to be small
- ▶ **Error of the second kind**: liberating someone who indeed committed the crime
  - This one can be large, but we also want it to be small

# Hypothesis testing: errors

		True state	
		$H_0$ is true (he is not guilty)	$H_1$ is true (he is guilty)
Decision	Accept $H_0$ (acquittal)	Right decision <i>Probability = <math>1-\alpha</math></i> ( <i>significance level</i> )	Wrong decision <b>Type II error</b> Probability = $\beta$ (power)
	Reject $H_0$ (conviction)	Wrong decision <b>Type I error</b> Probability = $\alpha$	Right decision Probability = $1-\beta$



# Testing procedure



# Goodnes-of-fit tests

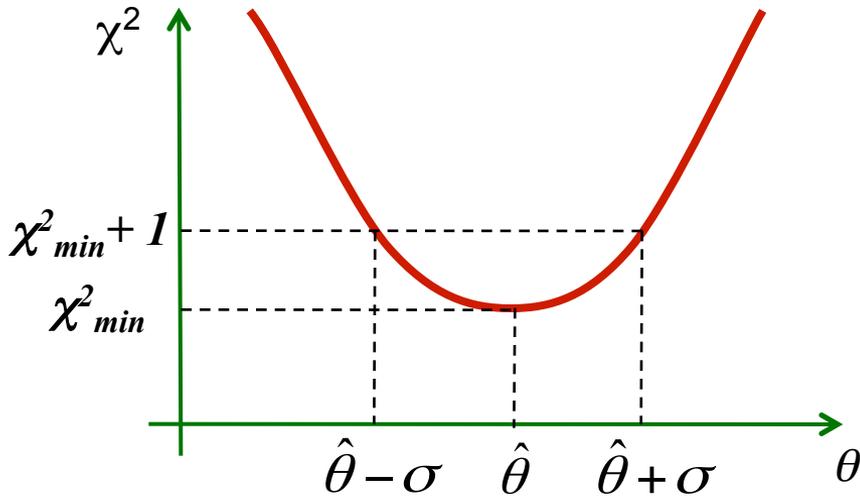
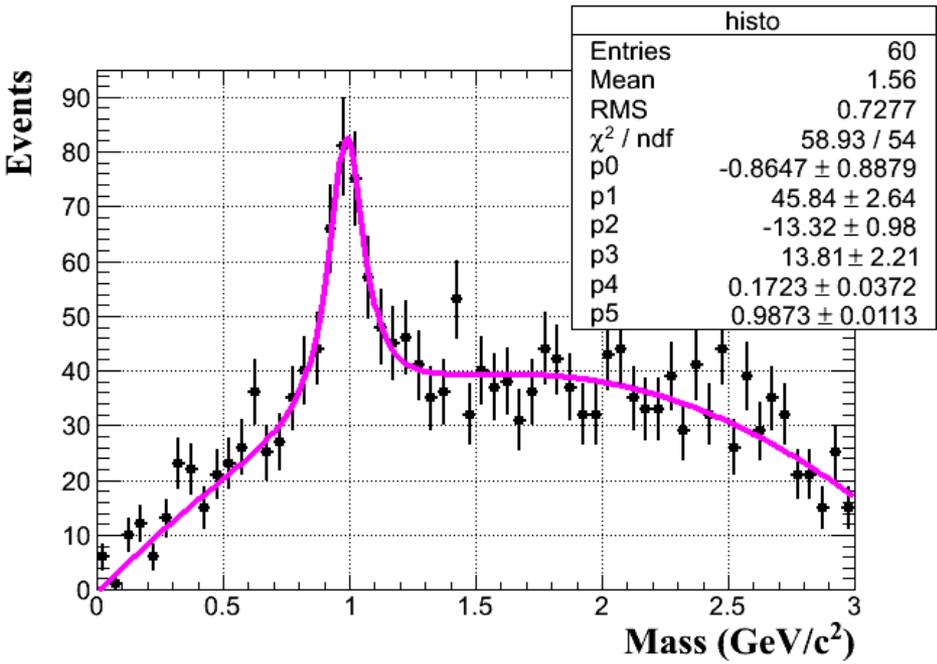
- ▶ In our case **NULL hypothesis  $H_0$**  is: *The **functional form** (or predicted values) describes well our data!*
- ▶ The form (i.e. the parameters that form depends on) is found by one of the methods for parameter estimation (moments, ML, chi-square)
- ▶ We are now looking for a **statistic  $t$**  (usually a single number) whose value reflects an agreement between the data and the hypothesis
  - The most commonly used statistic is the

1

2

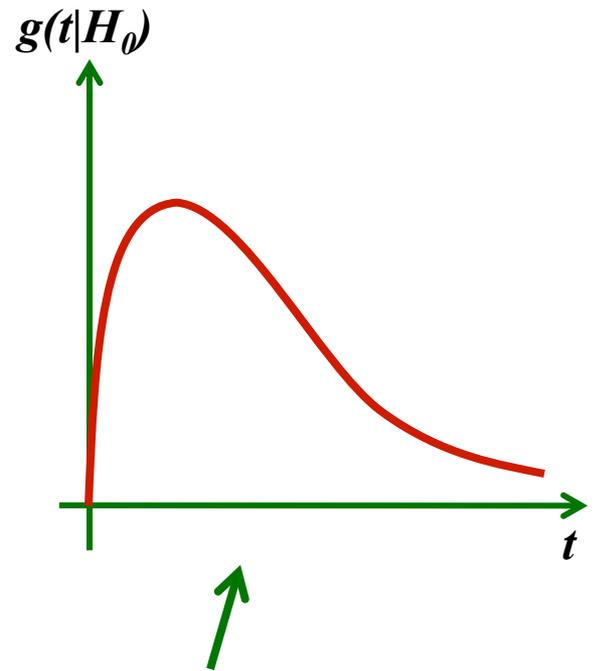
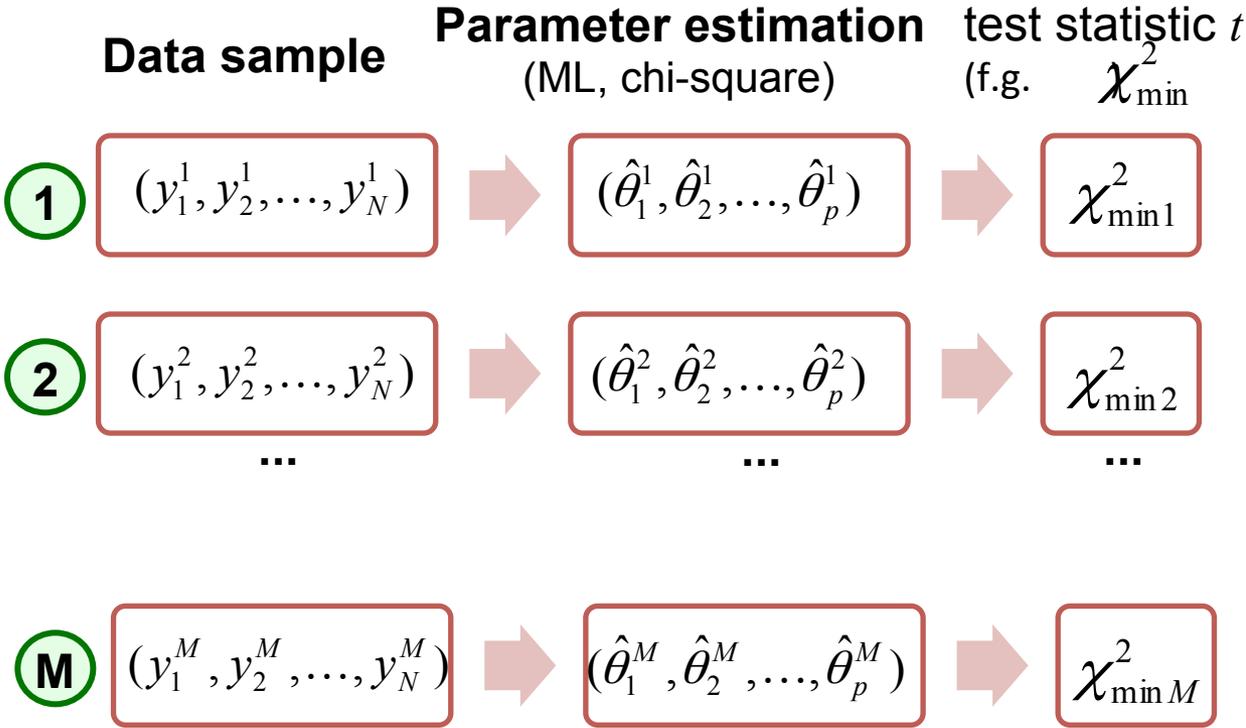
3

$$\chi^2_{\min}$$



# Distribution of the test statistic $t$

► **Imagine** we have many ( $M$ ) experiments (i.e. data samples) trying to test the null hypothesis  $H_0$



► We **would** then obtain a probability distribution function (PDF) of the test statistic, giving the  $H_0$  is true,  $g(t|H_0)$

# Finally: observed test statistic $t_{obs}$

## ► Now divide the distribution in two regions:

- $\alpha$ ) where  $H_0$  is rejected
- $1-\alpha$ ) where  $H_0$  is not rejected

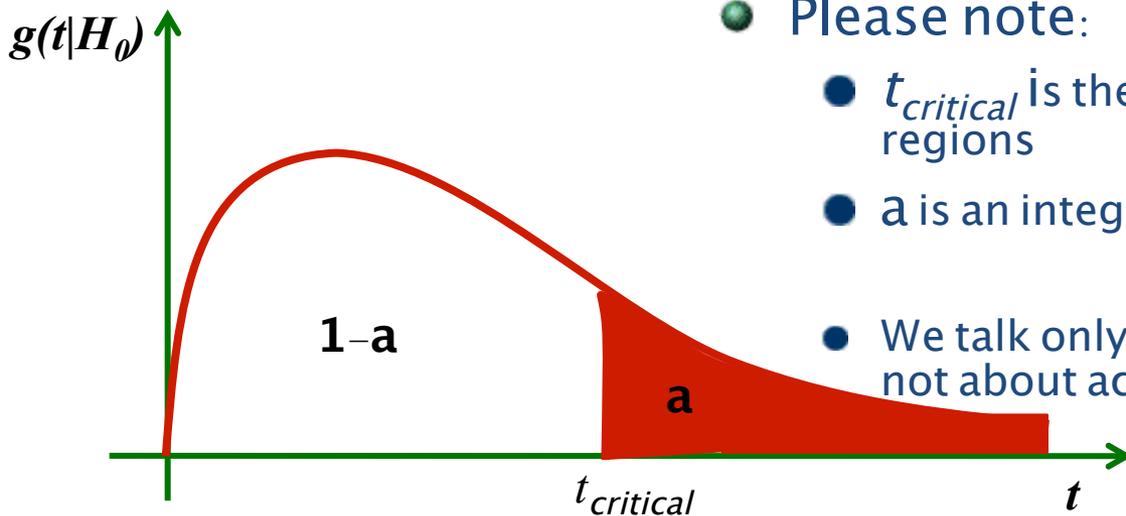
5

### ● Please note:

- $t_{critical}$  is the value of test statistic dividing two regions

- $\alpha$  is an integral: 
$$\alpha = \int_{t_{critical}}^{\infty} g(t | H_0) dt$$

- We talk only about rejecting the null hypothesis  $H_0$ , not about accepting any hypothesis



## ► **Very important!!!**

- We should decide about two regions before looking at the observed value of the test statistics

## ► Now we can calculate the observed test statistic $t_{obs}$

6

## ► At the end

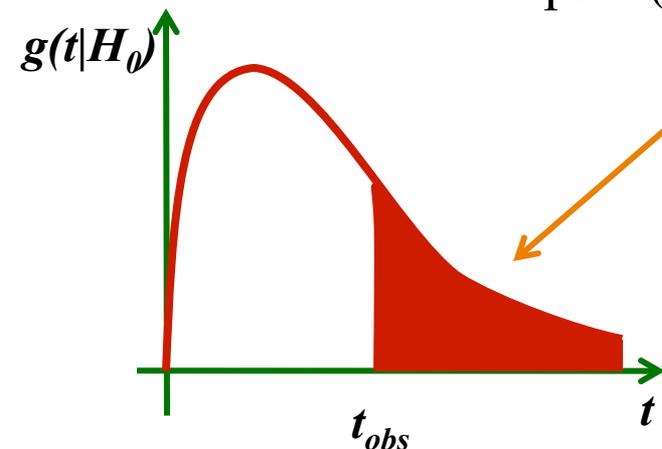
- If  $t_{obs} > t_{critical}$ : reject  $H_0$
- If  $t_{obs} < t_{critical}$ : do not reject  $H_0$

7

# $p$ -value

- ▶ So, our problem is to make the final conclusion based on only one number, the observed test statistic!
- ▶ Let's say the value of test statistic for our experiment is  $t_{obs}$
- ▶ And let's suppose that large value of  $t$  suggest larger discrepancy of the  $H_0$  with observed data (usually the case)
- ▶ Now, having  $g(t|H_0)$  we can for example answer to the question **What is the probability to obtain the value of  $t$  equal or greater than the value  $t_{obs}$  we observed?**
- ▶ The answer is simple an integral of the  $g(t|H_0)$ :

$$\text{prob}(t \geq t_{obs}) = \int_{t_{obs}}^{\infty} g(t | H_0) dt$$

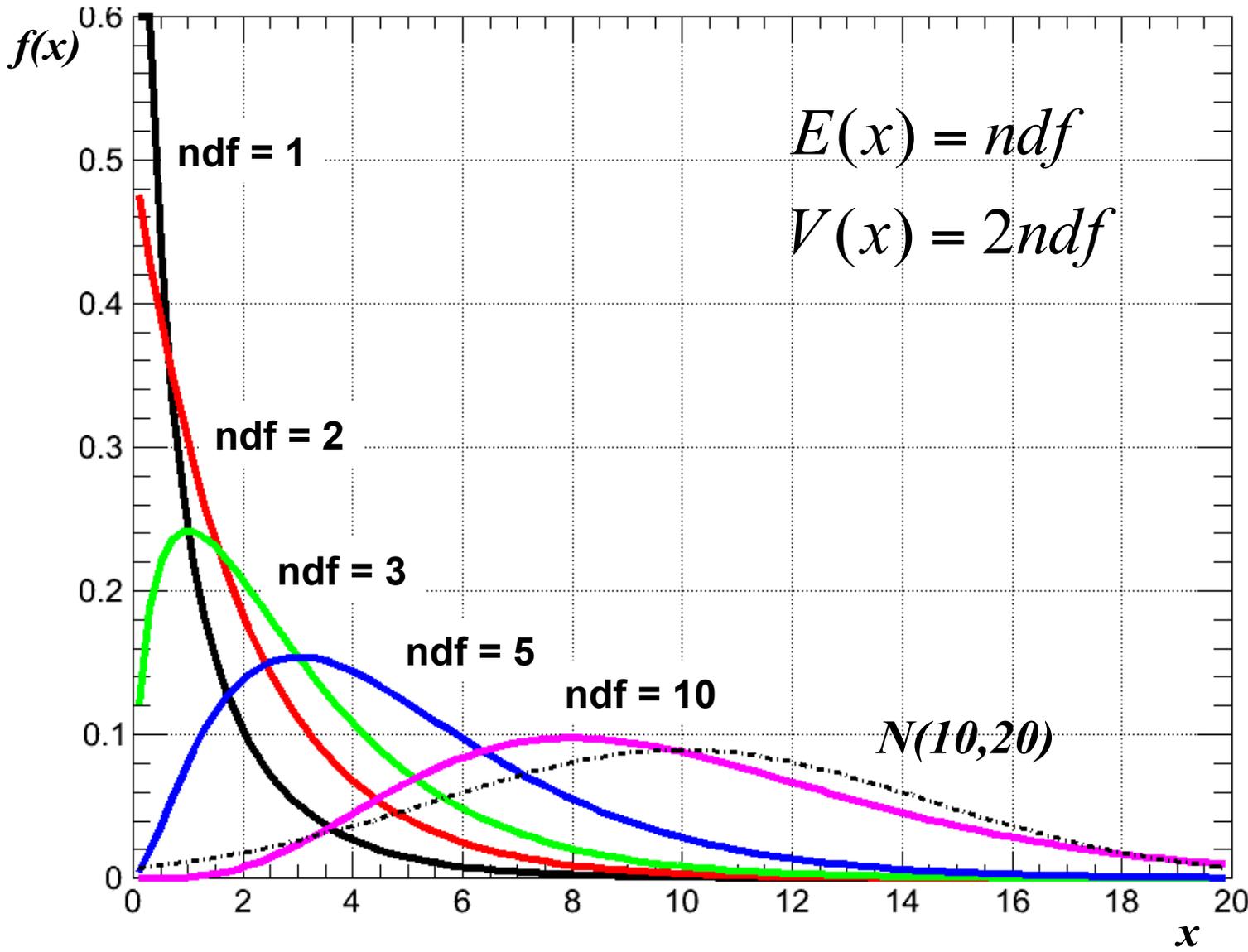


- This probability is so called **p-value**
- From PDG: "...  $p$ -value is defined as the probability to find  $t$  in the region of equal and lesser compatibility with  $H_0$  than the level of compatibility observed with actual data ..."

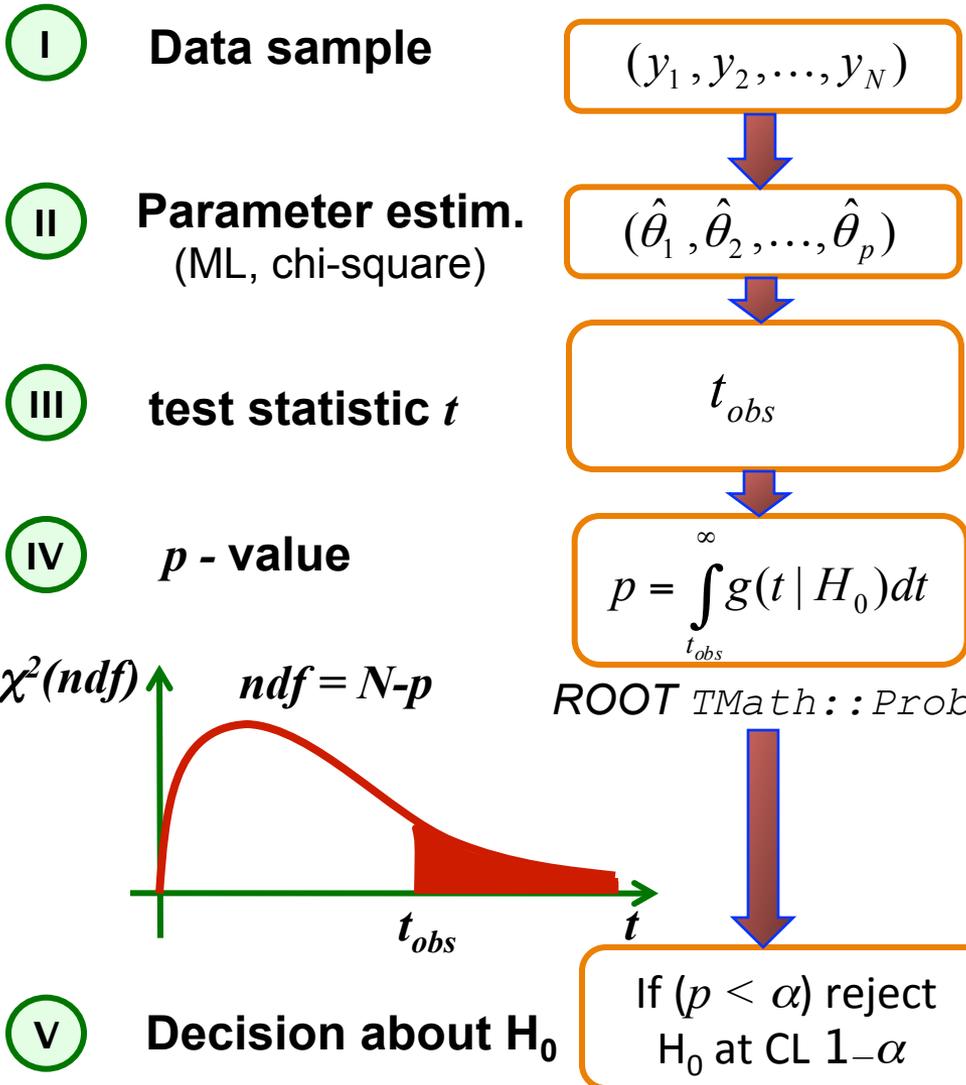
# $\chi^2(\text{ndf})$ distribution

- ▶ Well, this is all nice, but: as we don't have so many experiments, how do we get the PDF for the test statistics,  $g(t|H_0)$ ?
- ▶ For once, it turns out that we are 'lucky': most commonly used statistics for GOF testing are distributed as a  $\chi^2$  distribution!
  - That's actually the reason why they are so often used 😊
  - For example: when fitting histograms with  $N$  bins, with the function depending on  $p$  parameters, then the  $\chi^2_{\min}$  obtained in the fit, is distributed according to the  $\chi^2(N-p)$  function
    - $(N-p)$  is called **number of degrees of freedom (ndf)**
- ▶ If we are not so 'lucky' than we can use so called "**Toy Monte Carlo**" to generate  $g(t|H_0)$  from assumed distribution (describing the null hypothesis)
  - We "just" generate Monte Carlo experiments, find  $t$  for each of them and make a distribution  $g(t|H_0)$
  - We can even directly study the properties of the estimators (like bias, variance) as we can construct their distributions from MC experiments

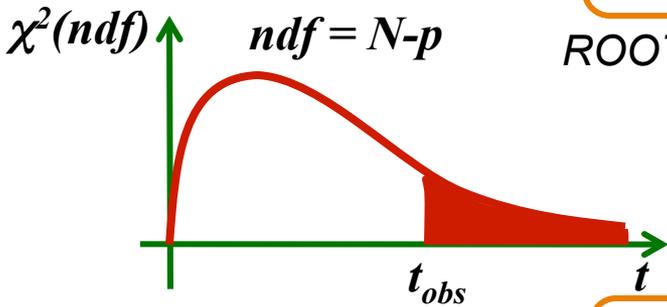
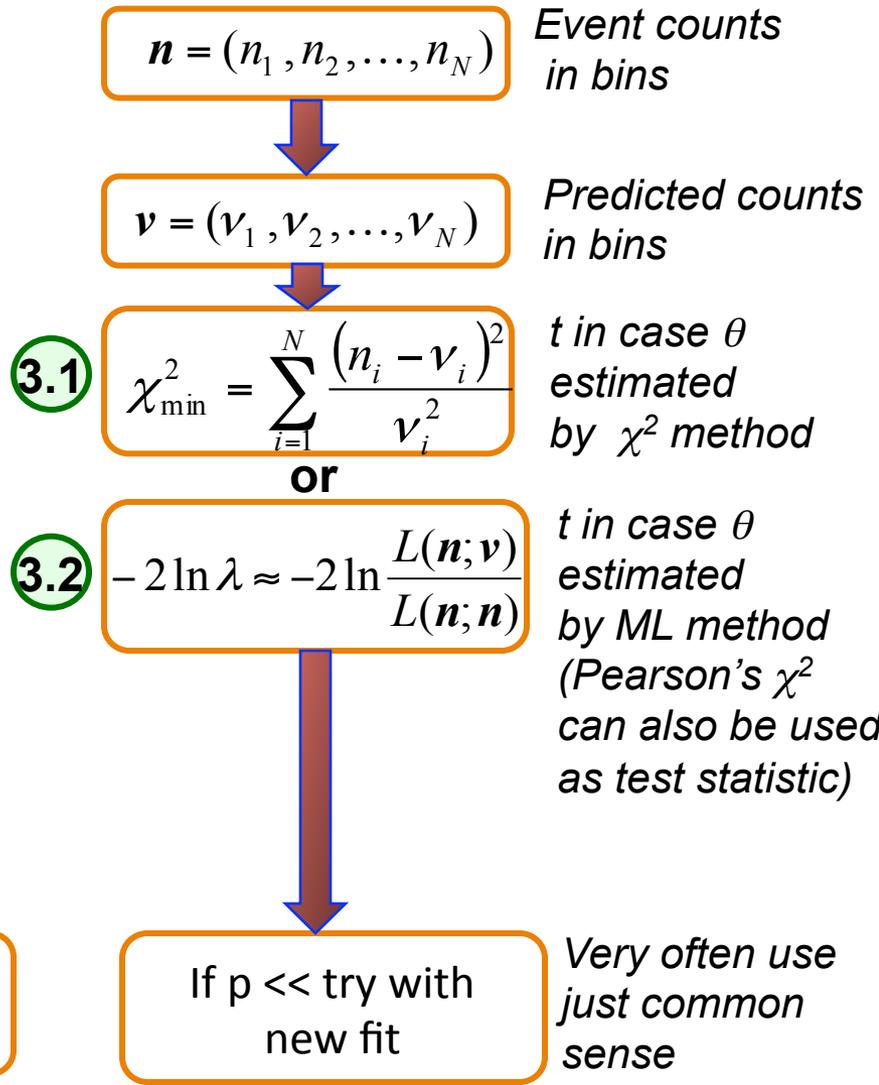
# Reminder - $\chi^2$ distribution



# GOF - overview

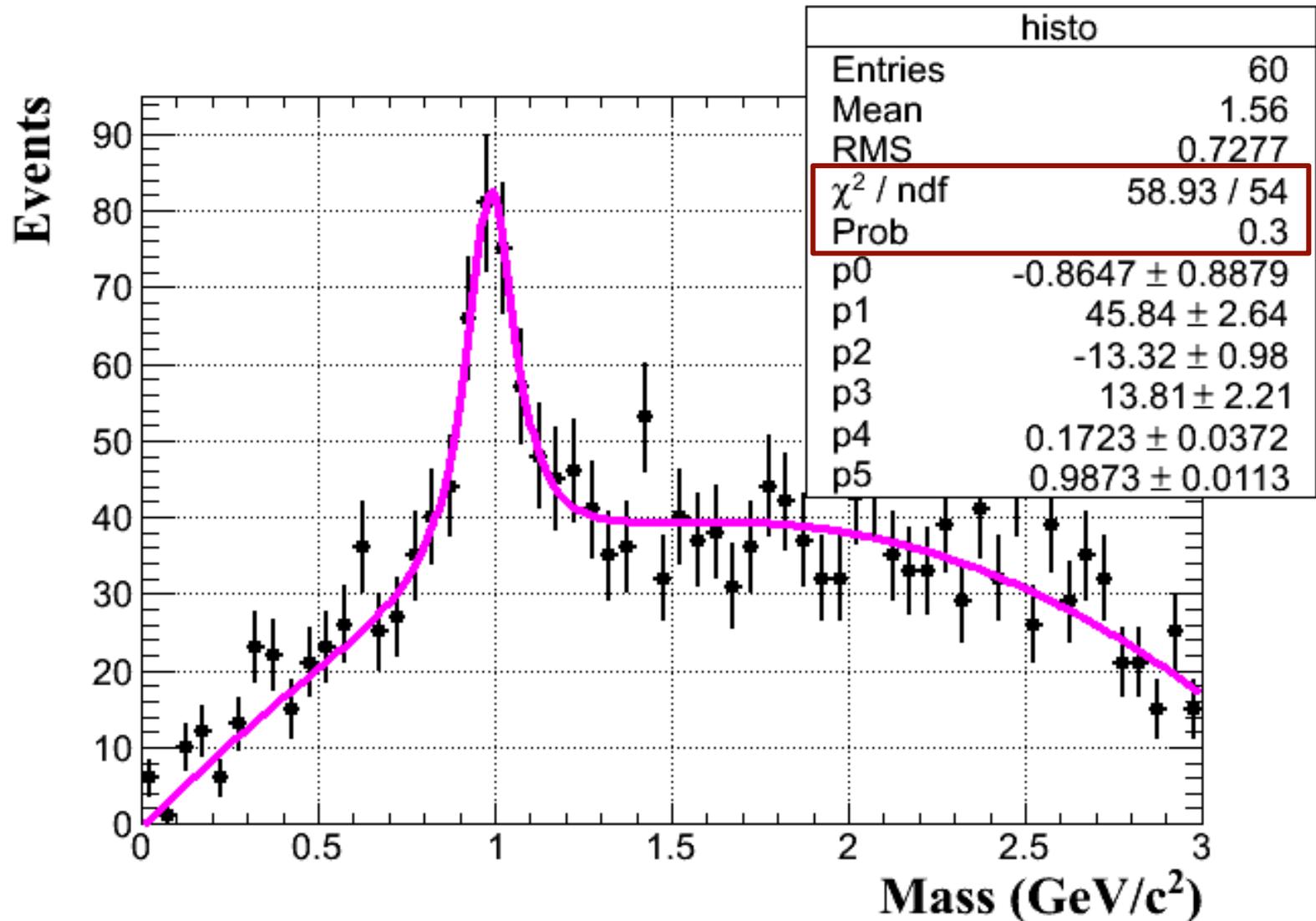


## Example: histogram fitting



ROOT TMath::Prob

# Our example



## Example: p-value in a counting experiment

A theory predicts the decay rate of a radioactive sample to be  $\mu = 17.3 \text{ decays}/h$ , and we measure  $N = 12 \text{ decays}/h$ .

**Q:** Is the measurement compatible with the theory?

**Solution:** We choose the test statistics to be the absolute difference  $t = |n - \mu|$ , so  $t_0 = |N - \mu| = 5.3$ . We use the Poisson distribution to calculate the p-value:

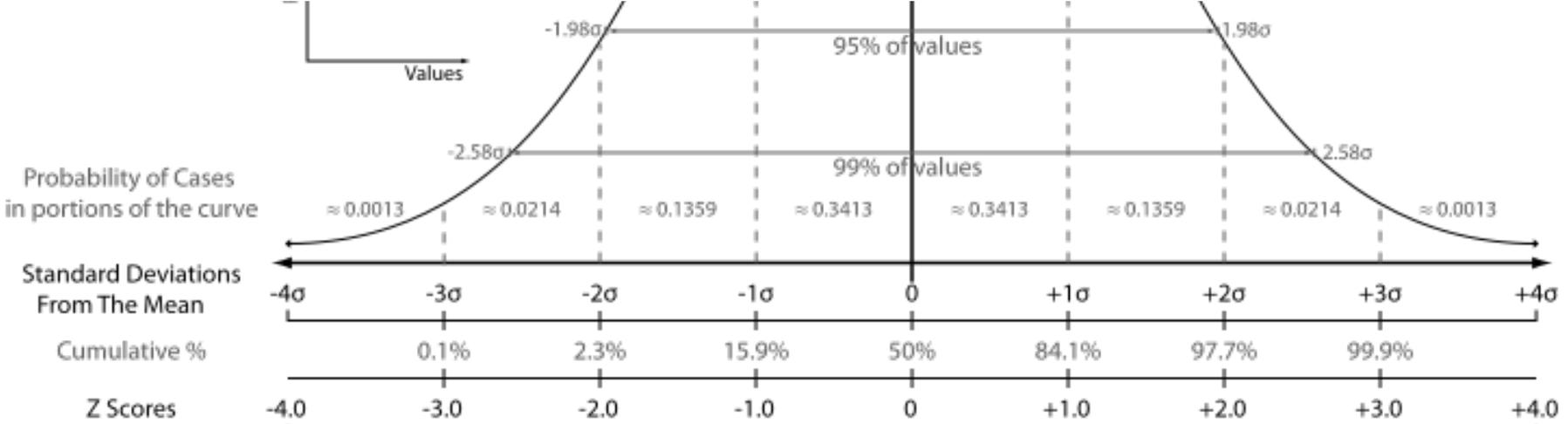
$$p = \sum_{n: |n-\mu| \geq 5.3} \frac{e^{-\mu} \mu^n}{n!} = \sum_{n=1}^{12} \frac{e^{-17.3} 17.3^n}{n!} + \sum_{n=23}^{\infty} \frac{e^{-17.3} 17.3^n}{n!} = 0.23$$

Observation  $N$  is not significantly different from the theoretical prediction  $\mu$ , since we have a 23% probability to measure the decay this far from the expected value.

This result can be confirmed with simple MC simulation that takes samples from the Poisson distribution.

# Converting p-values to significances

p-value	Zscore / significance	Area $\pm n\sigma$	Probability of outcome: 1 in ...
0.159	1	0.68268949	3.15
0.023	2	0.95449974	22.0
1.35E-03	3	0.99730020	370
3.17E-05	4	0.99993666	15,787
2.87E-07	5	0.99999943	1,744,278



# When do we claim a discovery?

- ▶ Claiming discovery is a serious issue
  - It should stay with us for a long long time (if not forever 😊)
- ▶ So, when do we claim a discovery?
  - When we are sure.
  - But we are never sure!
  - That's right, but we can be pretty sure 😊
  - 'Pretty' is not a scientific term!?
  - That's right, therefore we adopted some kind of a convention:
    - Make a hypothesis that the result you obtain is due to the fluctuation of the background (i.e. already know processes)
    - Calculate a probability for that hypothesis
    - Reject the hypothesis if that probability is smaller than 0.000000287 (significance > 5)

# Accepting or rejecting theories?

- ▶ **Imagine we make an experiment and obtain data**
  - Theory 1 agrees with data
  - Theory 2 also agrees
  - Theory 3 also agrees
  - ...
  - Theory n also agrees
  - Than the statement that "*Theory 1 is acceptable*" is not so strong
    - Not wrong neither
- ▶ **But imagine this scenario**
  - Theory 1 gives precise prediction
  - Experiment doesn't quite agree with that prediction
  - Than the statement "*Theory 1 is not acceptable*" is rather strong
  - Therefore we better reject than accept theories

---

# Backup

# Binomial distribution

Variable	$r$ , positive integer $\leq N$
Parameters	$N$ , positive integer; $p$ , $0 \leq p \leq 1$
Probability function	$P(r; N, p) = \binom{N}{r} p^r (1-p)^{N-r}$
Mean	$E(r) = Np$
Variance	$V(r) = Np(1-p)$
Usage example	<p>Example – <math>Z</math> decay:</p> <ul style="list-style-type: none"> <li>- <math>p = BR(Z \rightarrow ee) = 3\%</math></li> <li>- <math>P(5; 80, 0.03) = 6\%</math> probability to find exactly 5 <math>ee</math> events out of 80 <math>Z</math> decays</li> </ul>
Comment	$P(r; N, p)$ is a probability of finding exactly $r$ successes in $N$ trials, when probability of success in each single trial is a constant, $p$

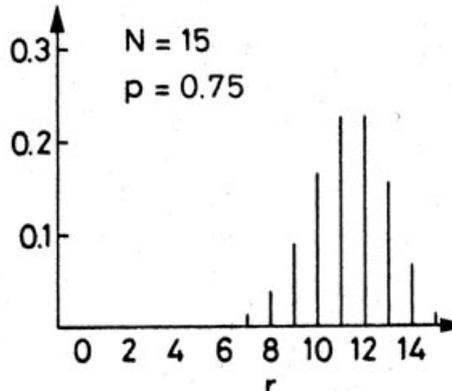
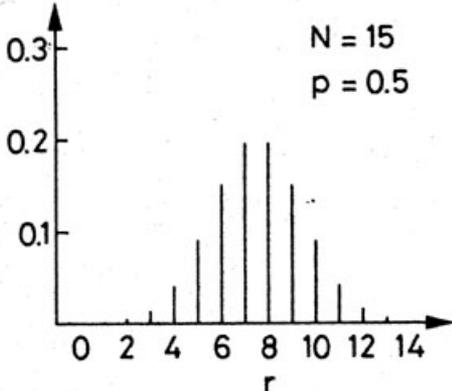
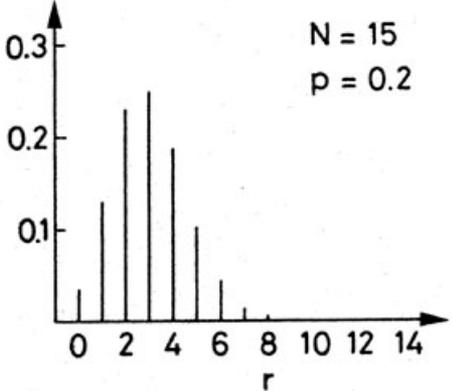


Figure from <http://nedwww.ipac.caltech.edu/level5/Leo/Figures/figure1.jpeg>

# Multinomial distribution

Variable	$r_i, i = 1, \dots, k$ , positive integers $\leq N$
Parameters	$N$ , positive integer $k$ , positive integer $p_i \geq 0, i = 1, \dots, k, \quad \sum_{i=1}^k p_i = 1$
Probability function	$P(r_1, \dots, r_k; N, p_1, \dots, p_k) = \frac{N!}{r_1! \dots r_k!} p_1^{r_1} \dots p_k^{r_k}$
Mean	$E(r_i) = Np_i$
Variance	$V(r_i) = Np_i(1-p_i)$
Usage example	Histogram containing $N$ events distributed in $k$ bins, with $r_i$ events in the $i^{th}$ bin
Comment	<ul style="list-style-type: none"> <li>• Multinomial distribution is the generalization of the binomial distribution to the case of more than two possible outcomes of an experiment</li> <li>• When <math>p_i \ll 1</math> (many bins) <math>V(r_i) \sim Np_i = r_i</math></li> </ul>

# Poisson distribution

Variable	$r$ , positive integer	 <p>Siméon-Denis Poisson (1781-1840)</p>
Parameters	$\mu$ , positive real number	
Probability function	$P(r; \mu) = \frac{\mu^r e^{-\mu}}{r!}$	
Mean	$E(r) = \mu$	
Variance	$V(r) = \mu$	
Usage example	Number of events $r$ collected after integrated luminosity $\int \mathcal{L} dt$ . Expected number of events is $\mu = \sigma \int \mathcal{L} dt$ . $\sigma$ is the cross section.	
Comments	<ul style="list-style-type: none"> <li>• <math>P(r; \mu)</math> expresses the probability of a number of events occurring in a fixed period of time if these events occur with a known average rate and independently of the time since the last event</li> <li>• <math>\mu</math> represents expected number of events in a given time interval</li> <li>• Time between two successive events is exponentially distributed</li> <li>• Poisson distribution is also called Poissonian</li> </ul>	

# Poisson distribution

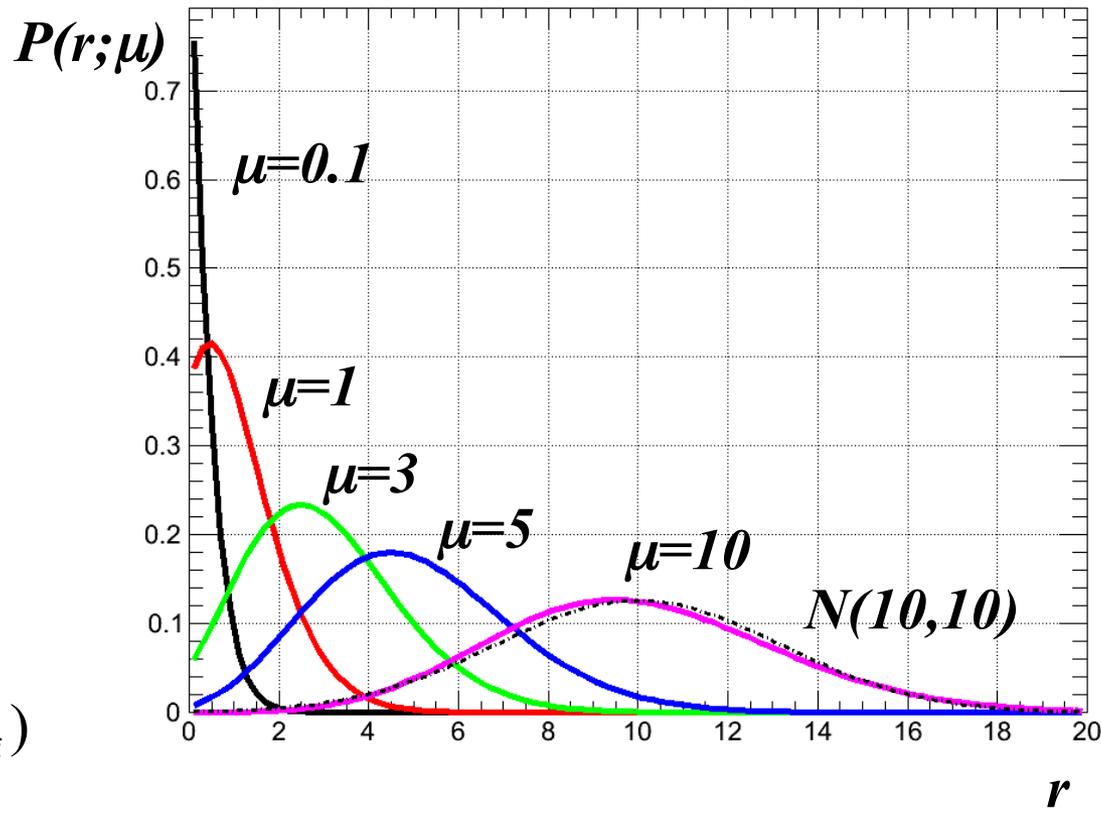
► For a large  $\mu$  Poisson distribution converges towards a Gaussian distribution

$$Pois(r; \mu) \xrightarrow{N \gg} Gauss(r; \mu, \sigma^2 = \mu)$$

► Sum of Poisson distributed random variables also follows a Poisson distribution whose parameter is sum of the component parameters

$$X_i \sim Pois(r; \mu_i)$$

$$Y = \sum_i X_i \sim Pois(r; \sum_i \mu_i)$$



▪ F.g. When combining signal (s) and background (b)

$$P(r; s, b) \sim Pois(r; s+b)$$

# Normal or Gaussian distribution

Variable	$x$ , positive real number	 <p>Carl Friedrich Gauss (1777-1855)</p>
Parameters	$\mu$ , real number $\sigma$ , real number	
Probability density function	$f(x) = N(\mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2}\right]$	
Mean	$E(x) = \mu$	
Variance	$V(x) = \sigma^2$	
Cumulative distribution	$F(x) = \Phi\left(\frac{x - \mu}{\sigma}\right); \quad \Phi(Z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^Z e^{-\frac{1}{2}x^2} dx$	
Comments	<ul style="list-style-type: none"> <li>• The most important distribution in statistics</li> <li>• The half-width at half-height is <math>1.176\sigma</math></li> <li>• <math>N(0, 1)</math> is called <i>standard</i> Normal density</li> <li>• Any linear combination of the <math>x_i</math> is also Normal</li> </ul>	



# Why is Gauss Normal?

## ► **Central limit theorem:**

If we have a set of  $N$  independent variables  $x_i$ , each from a distribution with mean  $\mu_i$  and variance  $\sigma_i^2$ , then the distribution of the sum  $X = \sum x_i$

a) has a mean  $\langle X \rangle = \sum \mu_i$ ,

b) has a variance  $V(X) = \sum \sigma_i^2$ ,

c) becomes Gaussian as  $N \rightarrow \infty$ .

► Therefore, no matter what the distributions of original variables may have been, their sum will be Gaussian in a large  $N$  limit

- Example: measurements errors

► Example (adopted from Barlow):

*"Human heights are well described by a Gaussian distribution, as many other anatomical measurements, as these are due to the combined effects of many genetic and environmental factors."*

# More than two variables

- ▶ Let's say that each event measure three quantities A, B and C
- ▶ We then have three random variables  $x$ ,  $y$  and  $z$
- ▶ Vector of measurements is now a matrix:

Event	A	B	C
1	$x_1$	$y_1$	$z_1$
2	$x_2$	$y_2$	$z_2$
...	...	...	...
N	$x_N$	$y_N$	$z_N$
Mean→	$\mu_x$	$\mu_y$	$\mu_z$

- ▶ Introducing new notation

$$(x, y, z) \rightarrow (x_{(1)}, x_{(2)}, x_{(3)}) = \vec{x} = \mathbf{x}$$

$$(\mu_x, \mu_y, \mu_z) \rightarrow (\mu_{(1)}, \mu_{(2)}, \mu_{(3)}) = \vec{\mu} = \boldsymbol{\mu}$$

- ▶ In case of  $m$  variables
- ▶ Please note: this multivariate vector  $x$  is a vector of  $m$  variables for one event, while in the case of one variable  $x$  is a vector of values of one variable for  $N$  events

$$\mathbf{x} = (x_{(1)}, x_{(2)}, \dots, x_{(m)})$$

# Multivariate Gaussian

- ▶ Multivariate Gaussian for the vector  $\mathbf{x} = (x_{(1)}, x_{(2)}, \dots, x_{(m)})$

$$f(\mathbf{x}; \boldsymbol{\mu}, V) = \frac{1}{(2\pi)^{n/2} |V|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T V^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]$$

- ▶  $x$  and  $\boldsymbol{\mu}$  are column vectors, while  $x^T$  and  $\boldsymbol{\mu}^T$  are row vectors

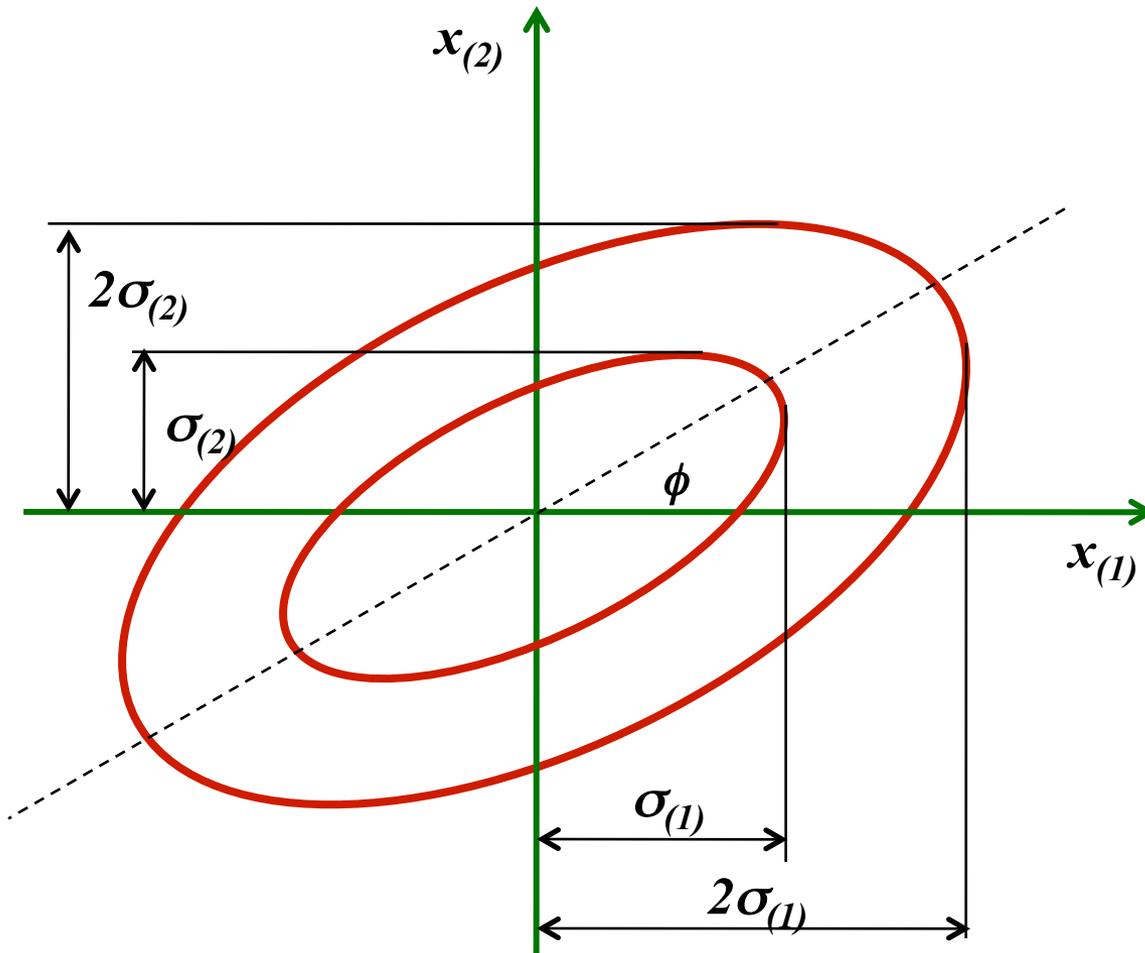
$$\mu_{(i)} = E(x_{(i)}) \quad V_{ij} = \text{cov}[x_{(i)}, x_{(j)}]$$

- ▶ Case of two variables ( $m = 2$ )

$$f(x_{(1)}, x_{(2)}; \mu_{(1)}, \mu_{(2)}, \sigma_{(1)}, \sigma_{(2)}) =$$

$$\frac{1}{2\pi\sigma_{(1)}\sigma_{(2)}\sqrt{1-\rho^2}} \times \exp\left\{-\frac{1}{2}\begin{bmatrix} x_{(1)} - \mu_{(1)} & x_{(2)} - \mu_{(2)} \end{bmatrix} \begin{bmatrix} \sigma_{(1)}^2 & \rho\sigma_{(1)}\sigma_{(2)} \\ \rho\sigma_{(1)}\sigma_{(2)} & \sigma_{(2)}^2 \end{bmatrix}^{-1} \begin{bmatrix} x_{(1)} - \mu_{(1)} \\ x_{(2)} - \mu_{(2)} \end{bmatrix}\right\} =$$
$$\frac{1}{2\pi\sigma_{(1)}\sigma_{(2)}\sqrt{1-\rho^2}} \times \exp\left\{-\frac{1}{2(1-\rho^2)} \left[ \left(\frac{x_{(1)} - \mu_{(1)}}{\sigma_{(1)}}\right)^2 + \left(\frac{x_{(2)} - \mu_{(2)}}{\sigma_{(2)}}\right)^2 - 2\rho \left(\frac{x_{(1)} - \mu_{(1)}}{\sigma_{(1)}}\right) \left(\frac{x_{(2)} - \mu_{(2)}}{\sigma_{(2)}}\right) \right]\right\}$$

# 2D Gaussian: iso-probability curves



	$P_{1D}$	$P_{2D}$
$1\sigma$	0.6827	0.3934
$2\sigma$	0.9545	0.8647
$3\sigma$	0.9973	0.9889
$1.515\sigma$		0.6827
$2.486\sigma$		0.9545
$3.439\sigma$		0.9973

**Remember (roughly) these values, we'll use them later in error estimates!**

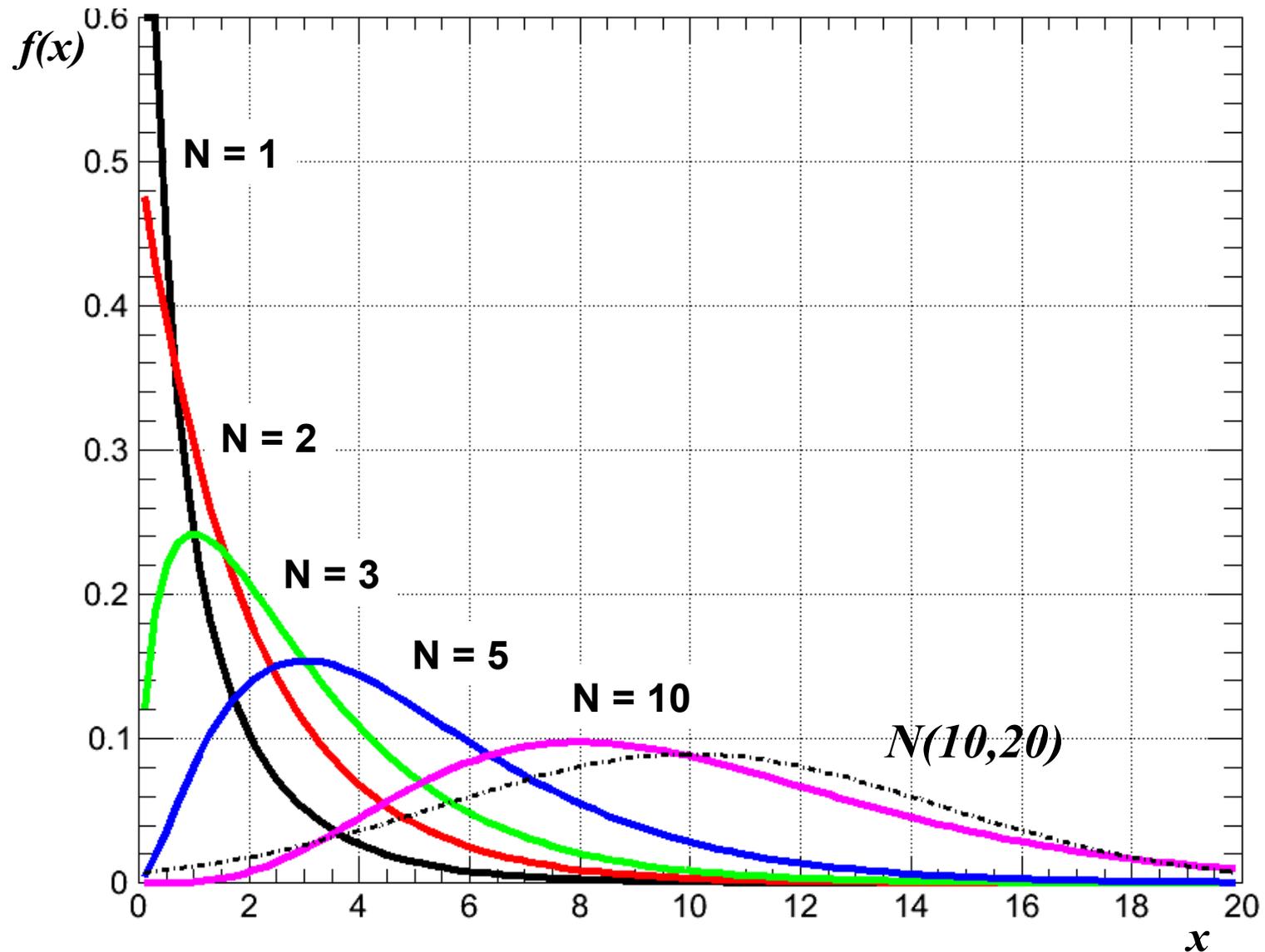
$\phi$  is a measure of the correlation (more details later)

Adopted from L. Lista

# Chi-square distribution

Variable	$x$ , positive real number
Parameters	$N$ , positive integer (number of “degrees of freedom”)
Probability function	$f(x) = \left( \frac{1}{2} \left( \frac{x}{2} \right)^{\frac{N}{2}-1} e^{-\frac{x}{2}} \right) / \Gamma\left(\frac{N}{2}\right)$
Mean	$E(x) = N$
Variance	$V(x) = 2N$
Usage example	Chi-square test for goodness of fit
Comments	<ul style="list-style-type: none"><li>• If <math>x_i</math> are <math>k</math> independent, normally distributed random variables with mean 0 and variance 1, then the random variable <math>Q = \sum x_i^2</math> is distributed according to the chi-square distribution with <math>k</math> degrees of freedom</li><li>• The chi-square distribution is a special case of the gamma distribution.</li></ul>

# Chi-square distribution



# Some other distributions

## ▶ Student's $t$ -distribution

- Used for hypothesis testing
- First published in 1908 by W. S. Gosset, while he worked at a Guinness Brewery, under the pseudonym *Student*



## ▶ Beta distribution

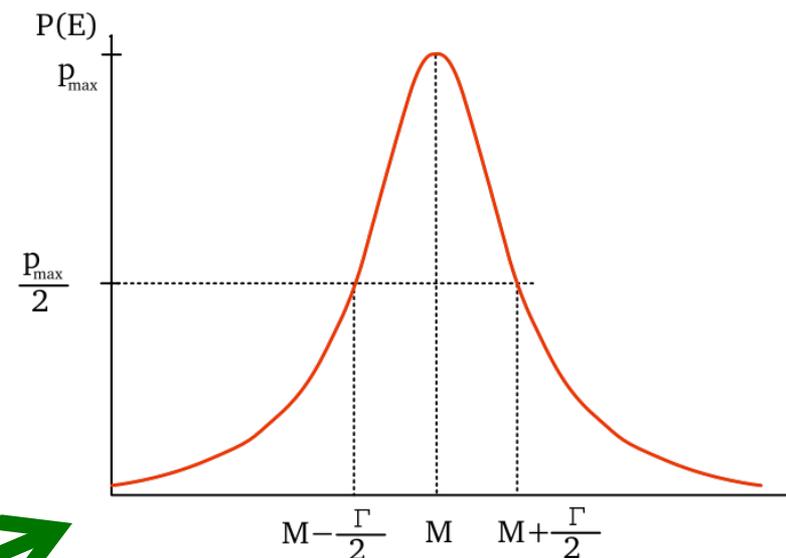
- Used in Bayesian statistics

## ▶ Gamma distribution

- Probability model for waiting time

## ▶ Cauchy or Lorentz or Breit-Wigner distribution

- A solution to the differential equation describing a **resonance**
- Energy distribution of a resonance

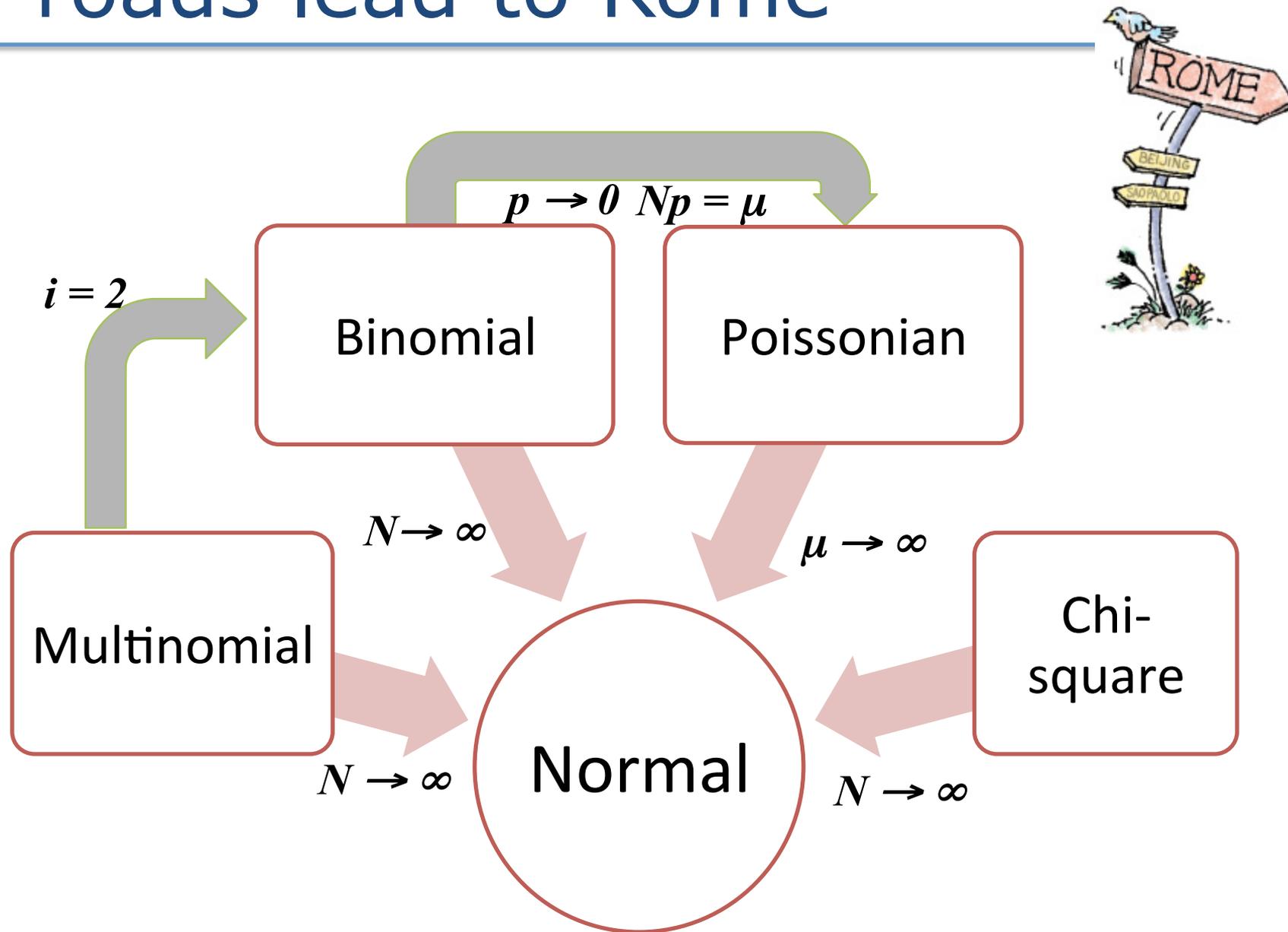


$$P(E) \sim \frac{1}{(E^2 - M^2)^2 + M^2 \Gamma^2}$$

## ▶ Log-Normal distribution

- Used when including systematic errors in the analysis
- If  $x$  is Log-Normally distributed, then  $\log(x)$  is Normally distributed

# All roads lead to Rome



# References

- ▶ F. James, *Statistical Methods in Experimental Physics*, World Scientific 2006
- ▶ R. J. Barlow, *Statistics – A guide to the Use of Statistical Methods in Physical Sciences*, Wiley 1999
- ▶ G. Cowan, *Statistical Data Analysis*, Oxford Univ. Press, 1998
- ▶ D. S. Sivia, *Data Analysis – A Bayesian Tutorial*, Oxford University Press, 2008
- ▶ L. Lyons, *Statistics for nuclear and particle physicists*, Cambridge University Press 1992
- ▶ PDG, *The Review of Particle Physics*, J. Beringer et al., Phys. Rev. D86, 010001 (2012), <http://pdg.lbl.gov/>
  - Chapter 35: *Probability*
  - Chapter 36: *Statistics*
  - Chapter 37: *Monte Carlo Techniques*
  - And references therein
- ▶ S. Baker and R. D. Cousins, *Clarification of the use of chi-square and likelihood functions in fits to histograms*, Nucl.Instrum.Meth.221:437-442,1984.
- ▶ ROOT Users Guide, <http://root.cern.ch/drupal/content/users-guide>
- ▶ Luca Lista, *Statistical methods for data analysis*, <http://people.na.infn.it/~lista/Statistics/>
- ▶ M. Liendl, *Experiment Simulation*, CERN School of Computing 2006
- ▶ M. Liendl, A. Heikkinen, *Experiment Simulation*, CERN School of Computing 2008
- ▶ I. Puljak, A. Heikkinen, *Data analysis*, CERN School of Computing 2010