

# DATA ANALYSIS

---

Toni Šćulac

Faculty of Science, University of Split, Croatia

CERN School of Computing 2022, Krakow, Poland

# LECTURES OUTLINE

---

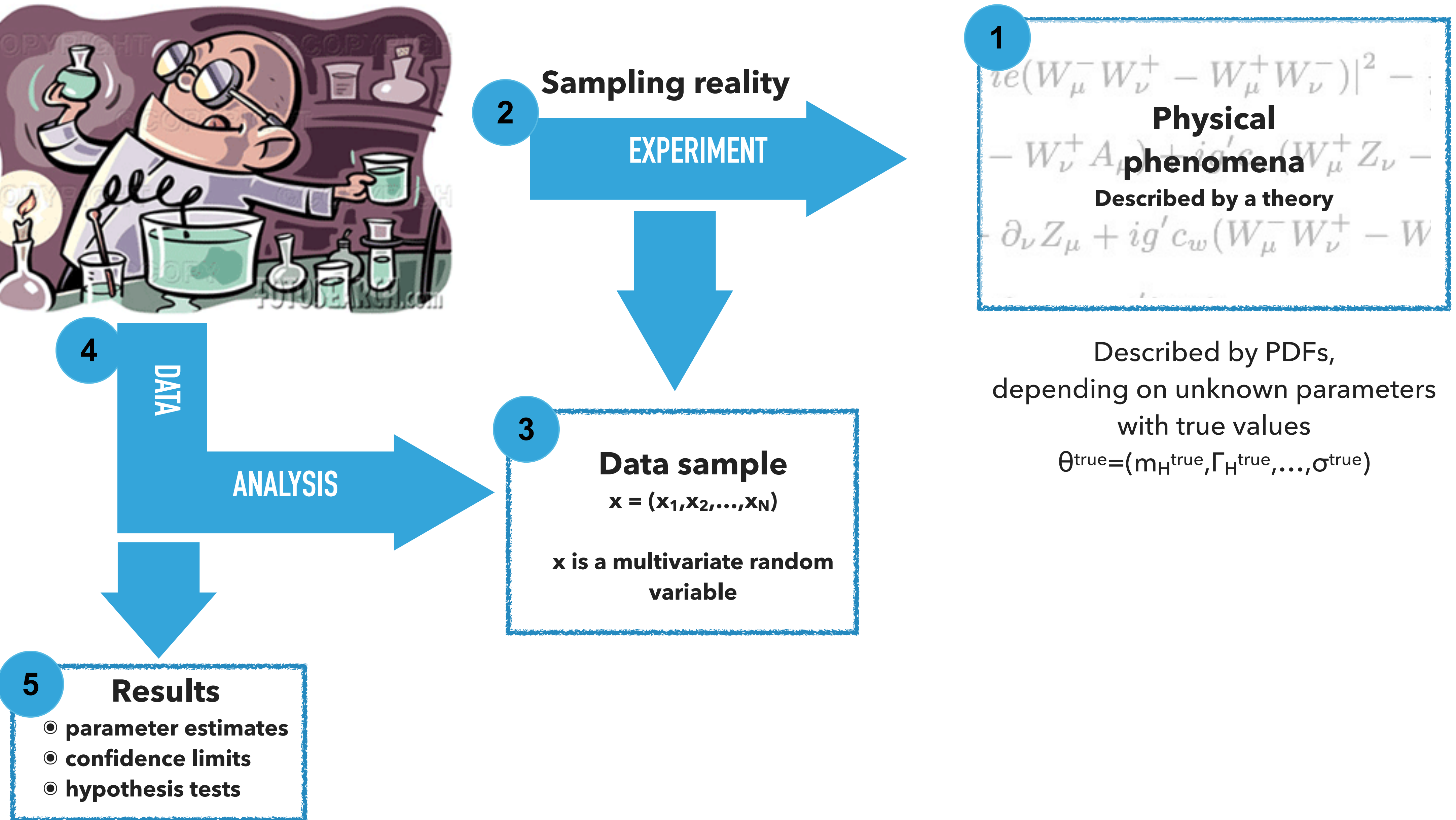
2

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

# PARAMETER ESTIMATION AND CONFIDENCE INTERVALS

# GENERAL PICTURE REMINDER

4



- The parameters of a PDF are constants that characterise its shape:

$$f(x; \theta) = \frac{1}{\theta} e^{-\frac{x}{\theta}}$$

- where  $x$  is measured data, and  $\theta$  are parameters that we are trying to estimate (measure)
- Suppose we have a sample of observed values  $\vec{x} = (x_1, x_1, \dots, x_n)$
- Our goal is to find some function of the data to estimate the parameter(s)
  - we write the **parameter estimator** with a hat  $\hat{\theta}(\vec{x})$
  - we usually call the procedure of estimating parameter(s): **parameter fitting**

- Task: find the average height of all students in a university on the basis of an (honestly selected) sample of  $N$  students
- Some possible ways of getting the result:
  - 1) Add up all the heights and divide by  $N$
  - 2) Add up the first 10 heights and divide by 10. Ignore the rest
  - 3) Add up all the heights and divide by  $N-1$
  - 4) Throw away the data and give the answer as 1.8 m
  - 5) Multiply all the heights and take the  $N$ -th root
  - 6) Choose the most popular height (the mode)
  - 7) Add up the tallest and shortest height and divide by 2
  - 8) Add up the second, fourth, etc. and divide by  $N/2$  for  $N$  even or by  $(N-1)/2$  for  $N$  odd



## ● Consistent

- Estimate converges to the true value as amount of data increases

$$\hat{\theta} \xrightarrow{\text{more data}} \theta^{true}$$

## ● Unbiased

- Bias is the difference between expected value of the estimator and the true value of the parameter

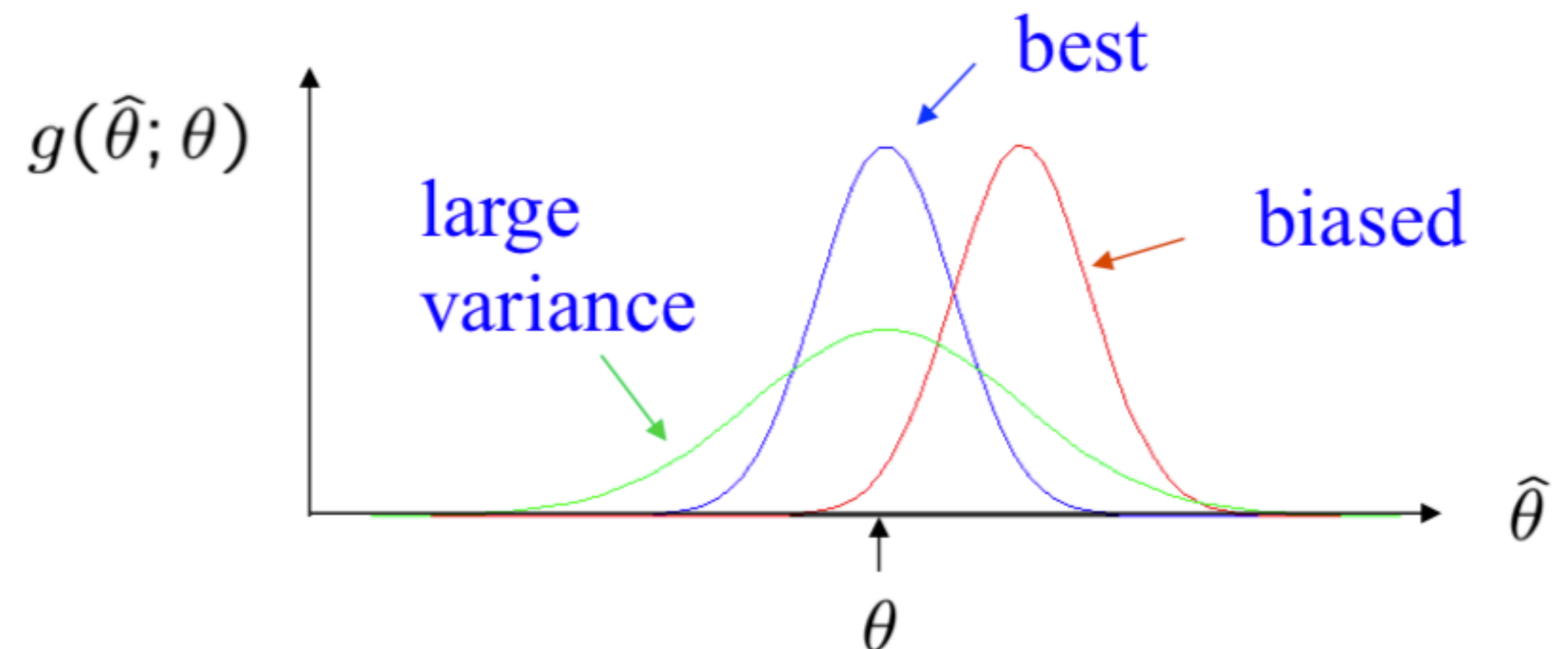
$$b = E(\hat{\theta}) - \theta^{true} = 0$$

## ● Efficient

- Its variance is small

## ● Robust

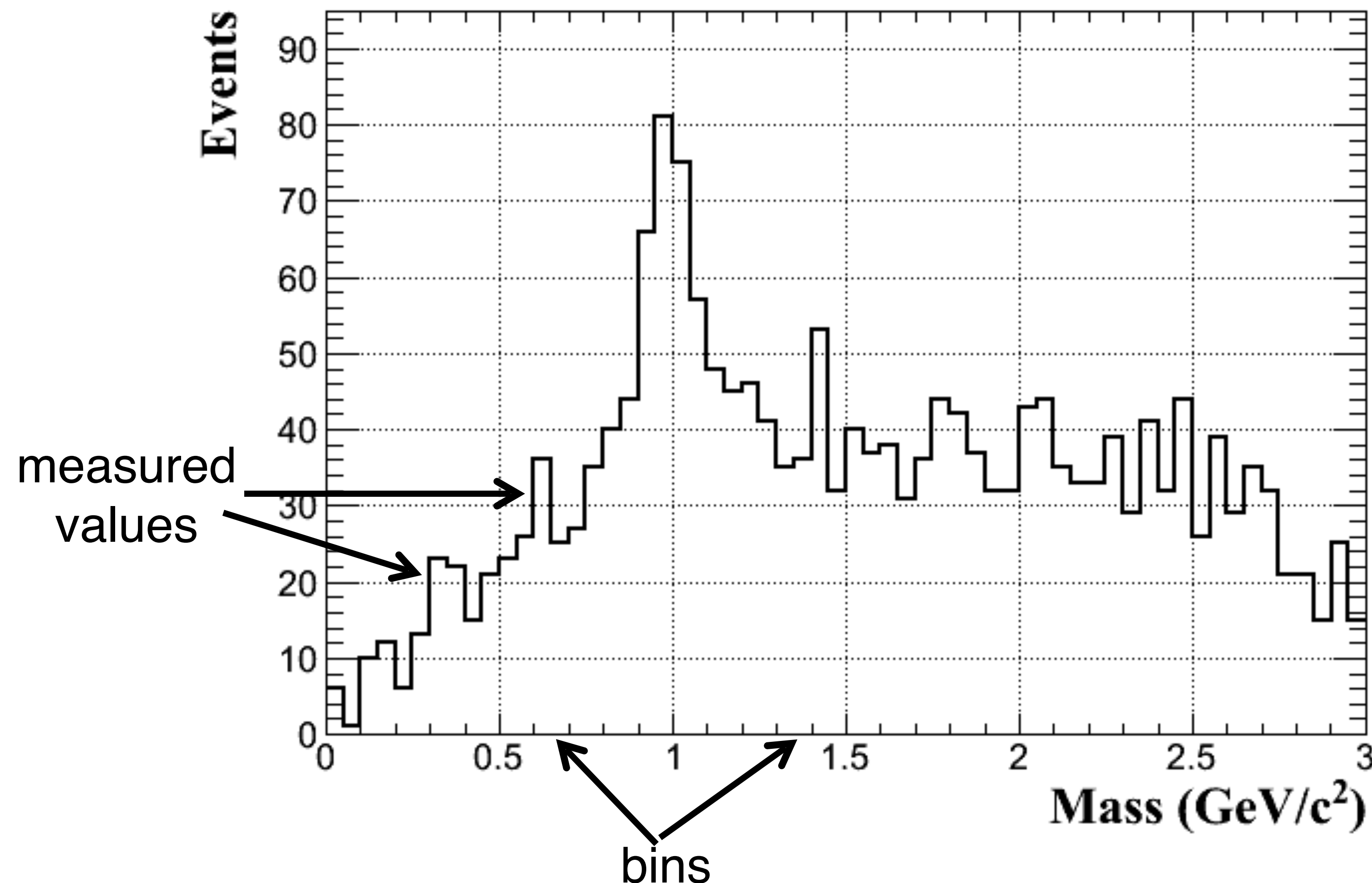
- Insensitive to departures from assumptions in the PDF



# EXAMPLE IN HEP - HISTOGRAM FITTING

8

- In counting experiments we usually represent data in histograms
- In the following example we will study a particle mass histogram



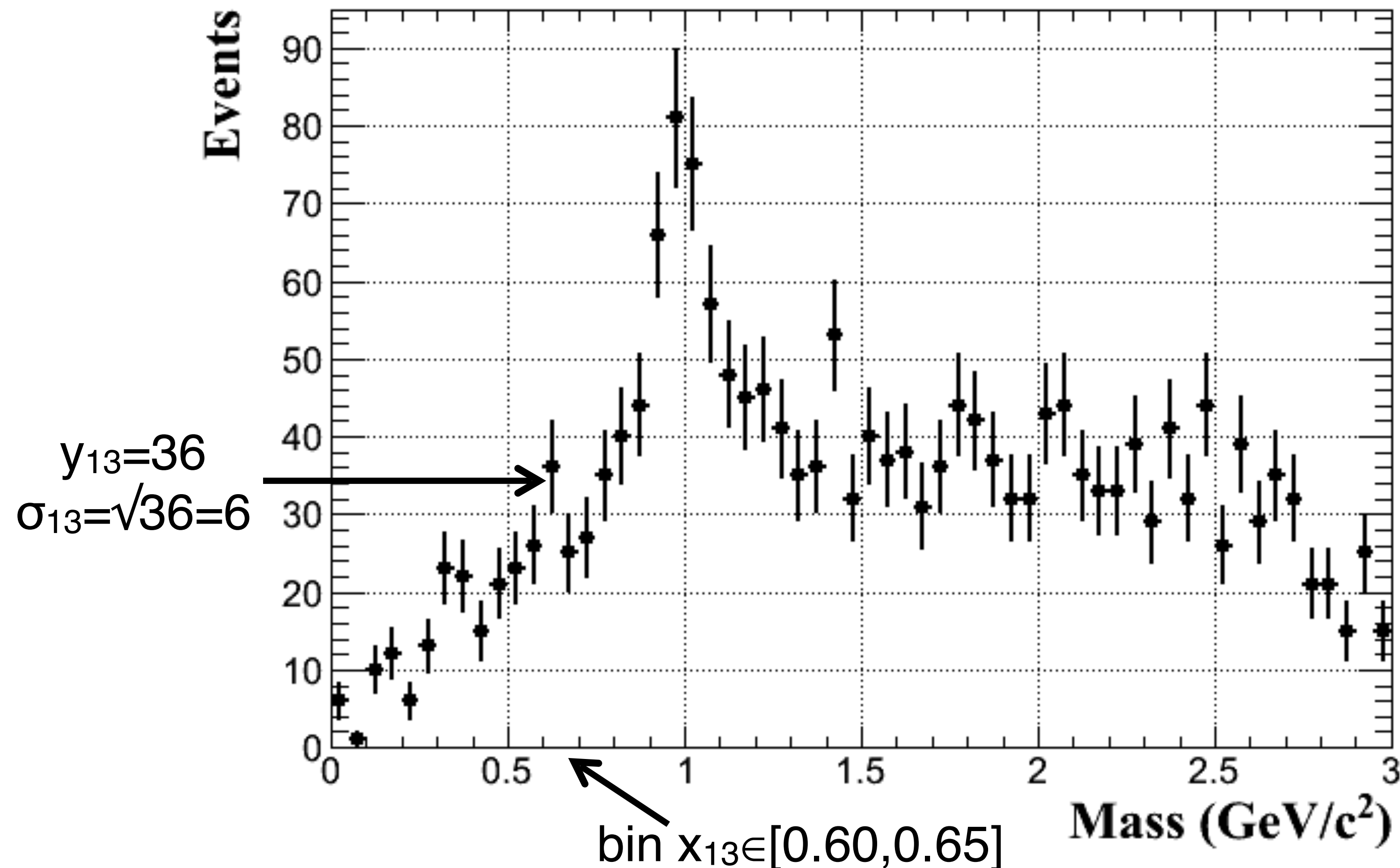
Root:  
`histo->Draw();`



# EXAMPLE IN HEP - HISTOGRAM FITTING

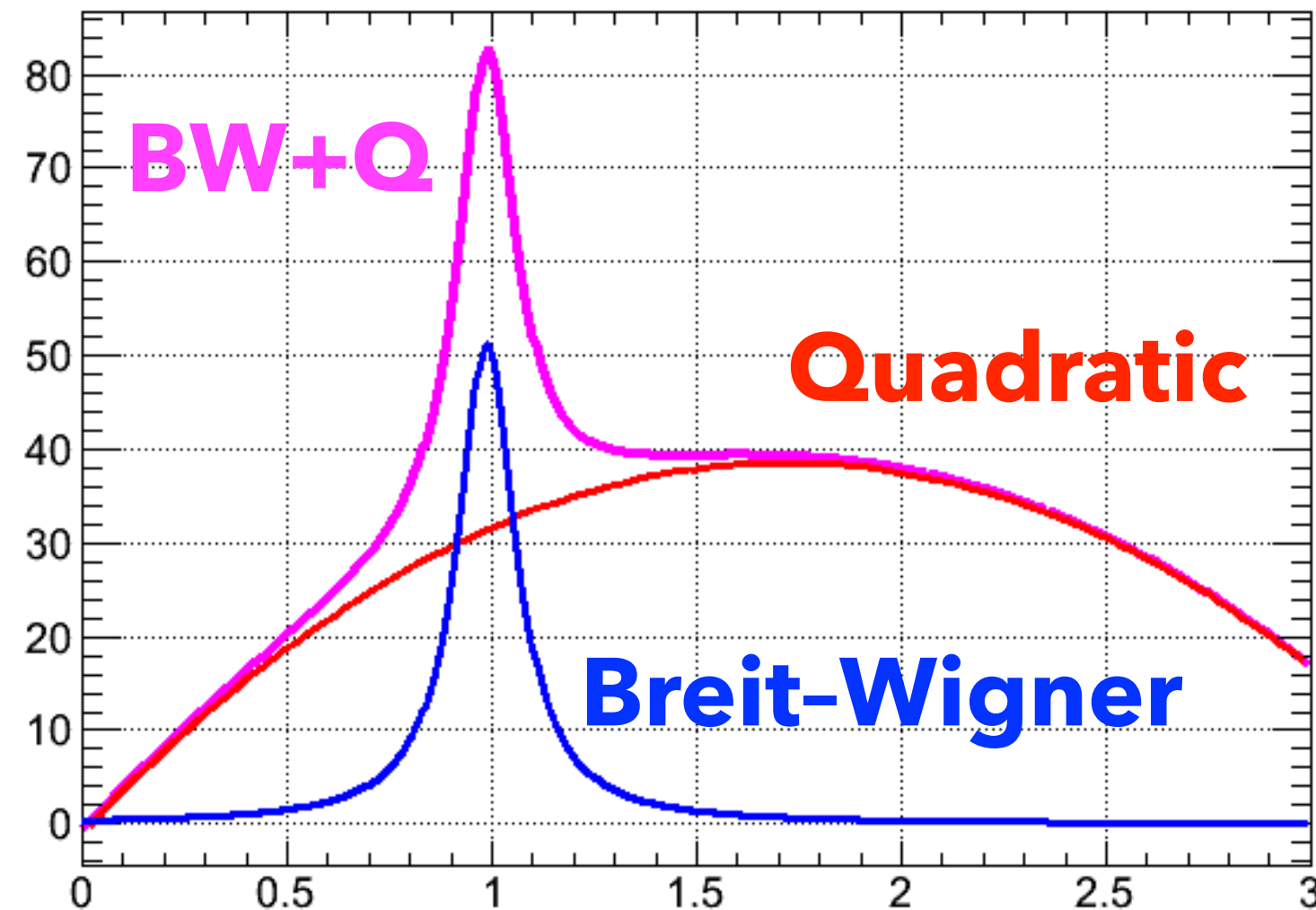
9

- Measured values have statistical uncertainties so it is better to represent them with points and error bars
  - each bin has a Poisson uncertainty



Root:  
`histo->Draw("ep");`

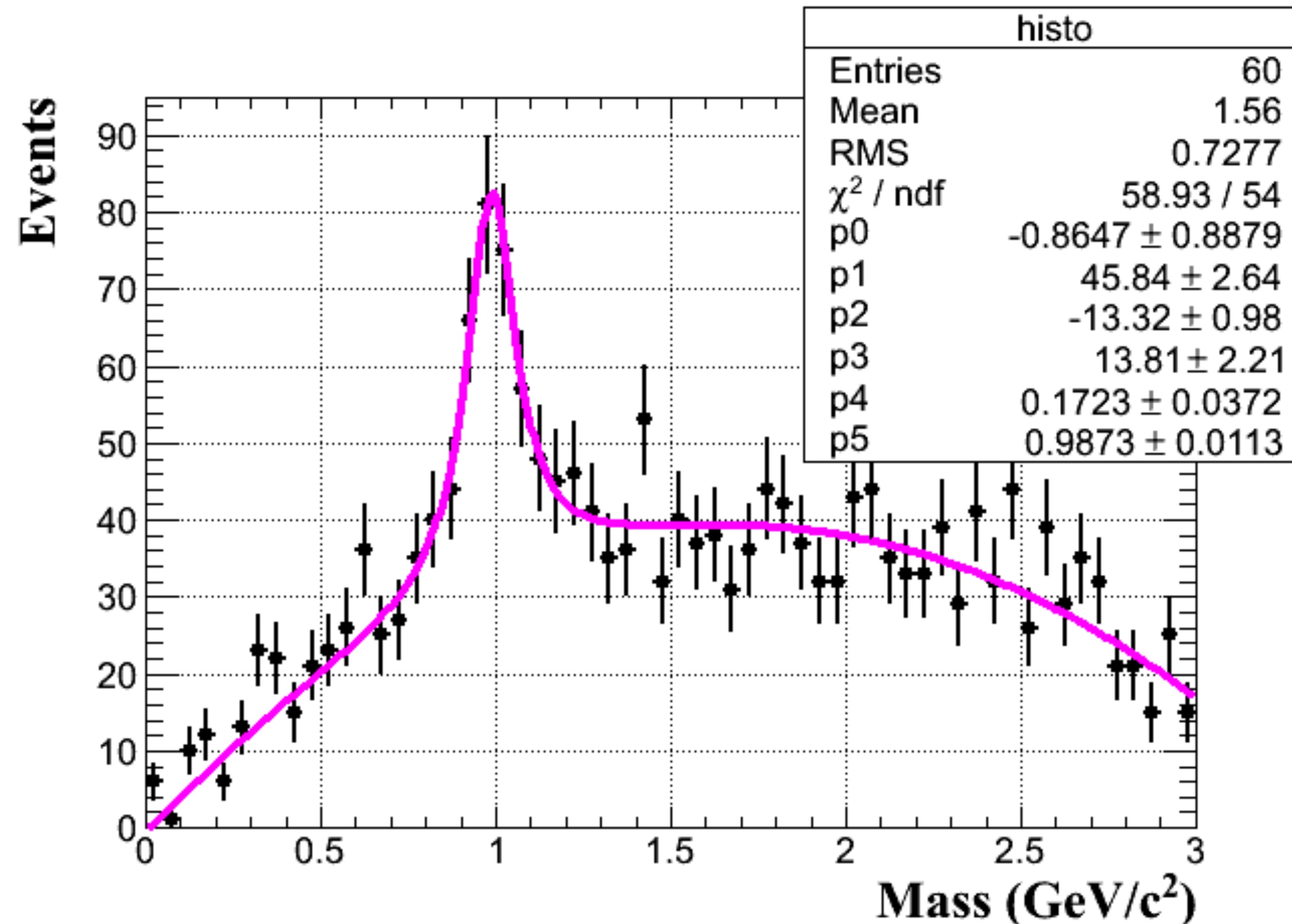
- Therefore we have
  - a set of precisely known values  $\mathbf{x} = (x_1, \dots, x_N)$  - **histograms bins**
  - At each  $x_i$ 
    - a measured value  $y_i$  - **number of events in a given bin**
    - a corresponding **error on measured value**  $\sigma_i$
- We are missing a theoretical PDF  $f(x_i; \theta^{true})$  with true parameters  $\theta^{true}$  so we can calculate **parameter estimator**  $\hat{\theta}$



$$BW(x; D, \Gamma, M) \approx \frac{D\Gamma}{(x^2 - M^2)^2 + 0.25\Gamma^2}$$

$$Q(x; A, B, C) = A + Bx + Cx^2$$

$$f(x_i, \theta^{true}) = f(x_i; D, \Gamma, M, A, B, C) = BW(x_i; D, \Gamma, M) + Q(x_i; A, B, C)$$



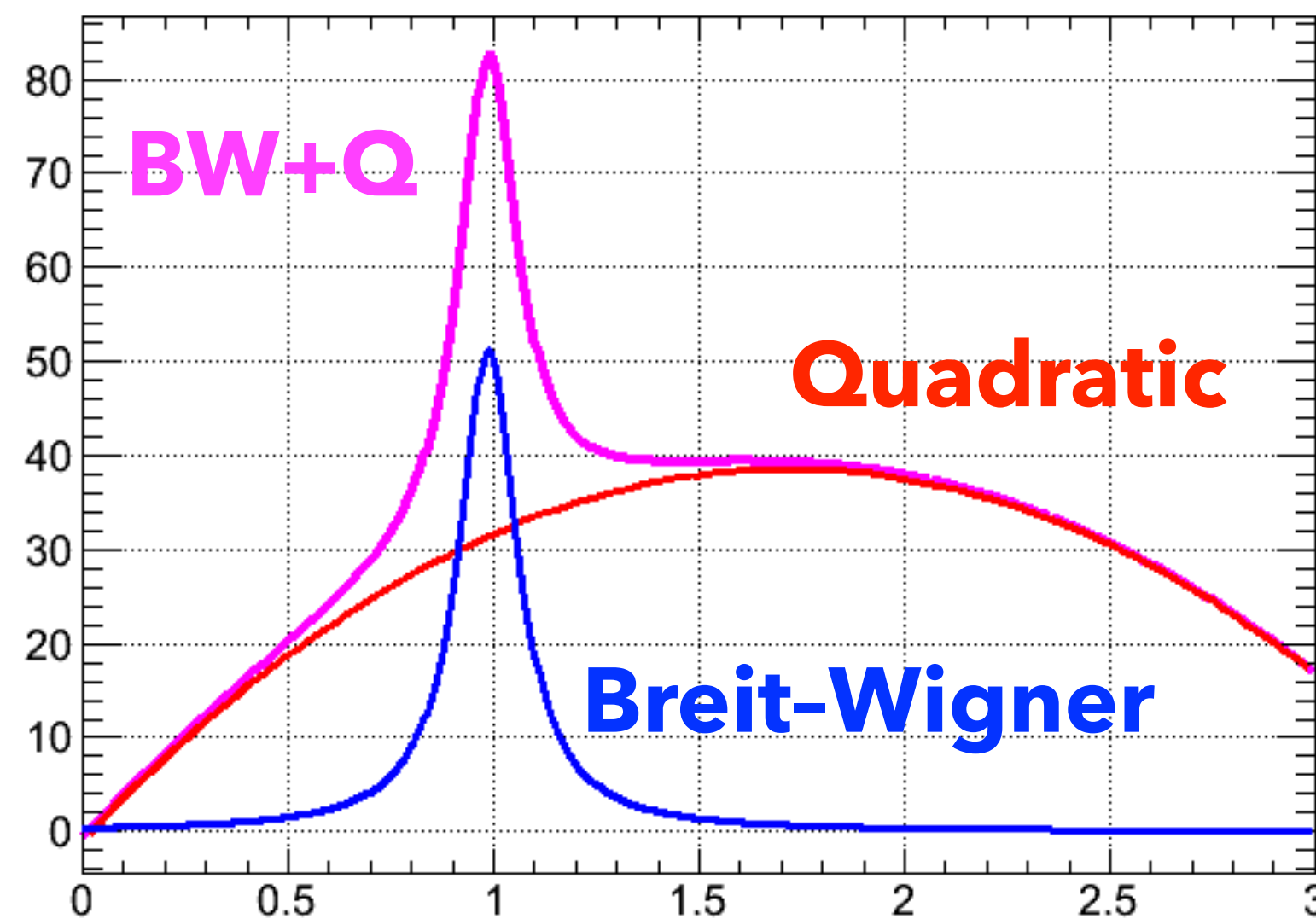
# EXAMPLE IN HEP - HISTOGRAM FITTING

12

1

**Physical phenomena**

Described by a theory



2

**Sampling reality**

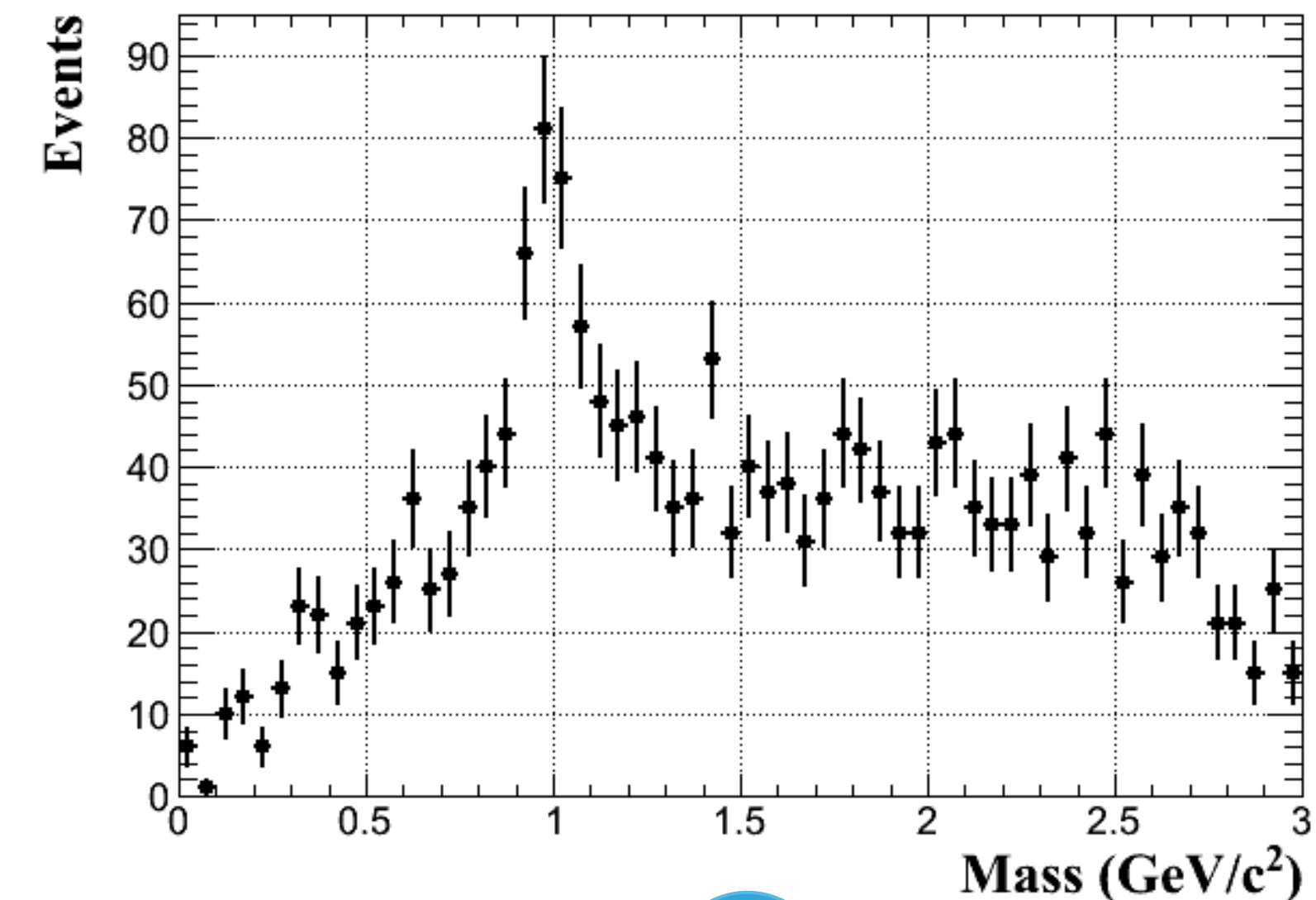
EXPERIMENT

3

**Data sample**

$\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$

$\mathbf{x}$  is a multivariate random variable



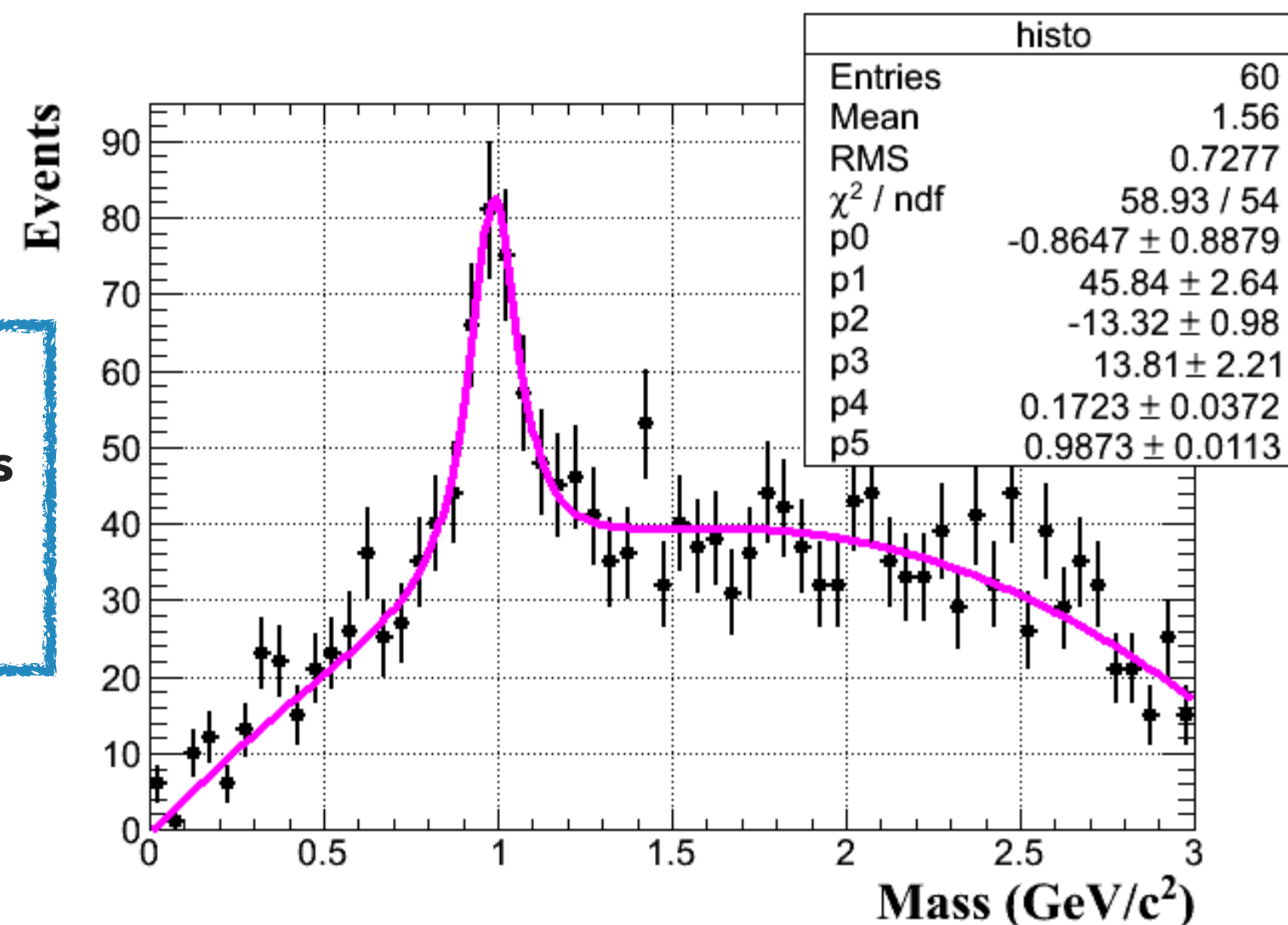
4

DATA ANALYSIS

5

**Results**

- parameter estimates
- confidence limits
- hypothesis tests





- 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  $T = T(x_1, x_2, \dots, x_N)$ 
  - For example, the sample mean  $\bar{x} = \frac{1}{N} \sum x_i$  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 example of statistics (plural of statistic!): sample median, variance, standard deviation, t-statistic, chi-square statistic, kurtosis, skewness, ...

## THE MAXIMUM LIKELIHOOD METHOD

- Gives consistent and asymptotically unbiased estimators
- Widely used in practice

## THE LEAST SQUARES (CHI-SQUARE) METHOD

- Gives consistent estimator
- Linear Chi-Square estimator is unbiased
- Frequently used in histogram fitting

## THE METHOD OF MOMENTS

- Gives consistent and asymptotically unbiased estimators
- Not as efficient as the Maximum Likelihood method



- 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** is a product of all single events probabilities:
  - $P(x; \theta) = P(x_1; \theta)P(x_2; \theta) \cdots P(x_N; \theta) = \prod P(x_i; \theta)$
- When the variable **x is replaced by the observed** data  $x^{\text{OBS}}$ , then P is no longer a PDF
- It is usual to denote it by L and called  $L(x^{\text{OBS}}; \theta)$  **the likelihood function**
  - Which is now a function of  $\theta$  only  $L(\theta) = P(x^{\text{OBS}}; \theta)$
- Often in the literature, it's convenient to keep X as a variable and continue to use notation  $L(X; \theta)$

- The probability that all  $N$  independent events will happen is given by the likelihood function  $L(x; \theta) = \prod f(x_i; \theta)$
- The principle of maximum likelihood (ML) says: **The maximum likelihood estimator  $\hat{\theta}$  is the value of  $\theta$  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**  
$$\ln L(x; \theta) = \sum \ln f(x_i; \theta)$$
- For  $p$  parameters we get a set of  $p$  **likelihood equations**: 
$$\frac{\partial \ln L(x; \theta)}{\partial \theta_j} = 0$$
- It is often more convenient the **minimise  $-\ln L$  or  $-2\ln L$**

- Consider the lifetime pdf  $f(t; \tau) = \frac{1}{\tau} e^{(-\frac{t}{\tau})}$
- Suppose we have **measured data**  $t(t_1, \dots, t_N)$
- Our **likelihood function** is defined as  $L(\tau) = \prod f(t_i; \tau)$
- The value of  $\tau$  for which  $L(\tau)$  is maximum also gives the maximum value of its **log-likelihood function**  $\ln L(\tau) = \sum \ln f(t_i; \tau) = \sum (\ln \frac{1}{\tau} - \frac{t_i}{\tau})$
- Solving one likelihood equation  $\frac{\partial \ln L(\tau)}{\partial \tau} = 0$  gives  $\hat{\tau} = \frac{1}{N} \sum t_i$
- Try generating 100 Monte Carlo toys for  $\tau = 1$  and estimating  $\hat{\tau}$  using the ML method

- ML estimator is **consistent**
- ML estimate is approximately **unbiased** and **efficient** for large samples
  - Usually biased for small samples
- ML estimate is **invariant**
  - A transformation of parameter won't change the answer
  - Keep in mind that invariance comes at the cost of a bias!
- Extra care to be taken when the best value of parameters are near imposed limits
- **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 the **unbinned maximum likelihood**
- ML has many advantages, but a few drawbacks too

- In Bayesian statistics, both  $\theta$  and  $x$  are random variables
- We want to know the conditional PDF for  $\theta$  given the data  $x$ :

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

- where  $\pi(\theta)$  is the prior probability density for  $\theta$ , reflecting the stage of knowledge of  $\theta$  before measuring the data  $x$ 
  - If we choose “prior ignorance”  $\pi(\theta) = \text{const}$ , then  $\hat{\theta}_{\text{Bayes}} = \hat{\theta}_{\text{ML}}$
  - No golden rule on how to define  $\pi(\theta)$
- In Bayesian statistics all our knowledge about  $\theta$  is in  $p(\theta | x)$ 
  - It is often a very complicated multidimensional function that is hard to report
  - Summarised using an estimator  $\hat{\theta}_{\text{Bayes}}$  which is often defined as the mode of  $p(\theta | x)$



- Likelihood function ( $L$ ) is constructed by replacing the variable  $x$  by the observed data in a product of single events probabilities
- Maximising (minimising) the  $\ln L$  ( $-2 \ln L$ ) function gives the parameter estimate  $\hat{\theta}_{ML}$
- $\hat{\theta}_{ML}$  does not mean that the estimate is the “most likely” value of  $\theta$ , it is the value that makes your data most likely
- ML estimate is unbiased and efficient for large samples, be careful if you want to use it for small samples
- ML can be used to fit binned data
- ML can be extended to deal with the case where the number of expected events is not a fixed number but a random number



- Suppose you have a set of precisely known (without error) values  $x(x_1, \dots, x_N)$  with a corresponding set of measured values  $y(y_1, \dots, y_N)$  with corresponding uncertainties  $\sigma(\sigma_1, \dots, \sigma_N)$ 
  - For example  $x_i$  histogram mass bins with  $y_i$  events with Poissonian uncertainty  $\sigma_i$
- Suppose you also know a function  $f(x; \theta)$  which predicts the value of  $y_i$  for any  $x_i$ . It depends on an unknown parameter  $\theta$ , which you are trying to determine.
  - In our example function  $f(x; \theta)$  would be theoretical prediction for number of events at a given mass
- To find best estimate of  $\theta$  we minimise the suitably weighted sum of squared differences between measured and predicted values, the so called “**least squares**” or “**chi-square**”:

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

● Estimator is found by finding the value which minimises  $\chi^2 : \frac{\partial \chi^2}{\partial \theta} = 0$

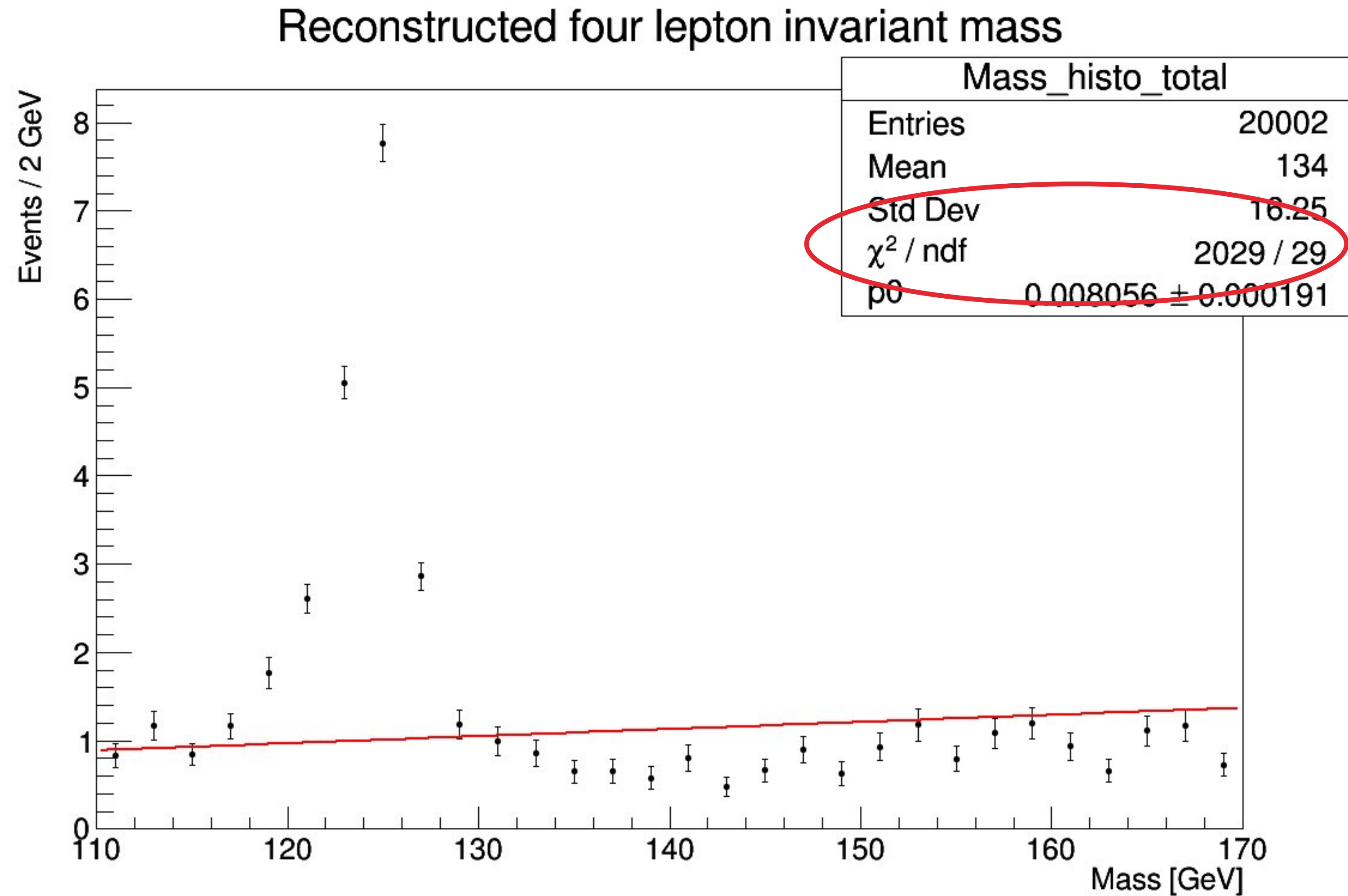
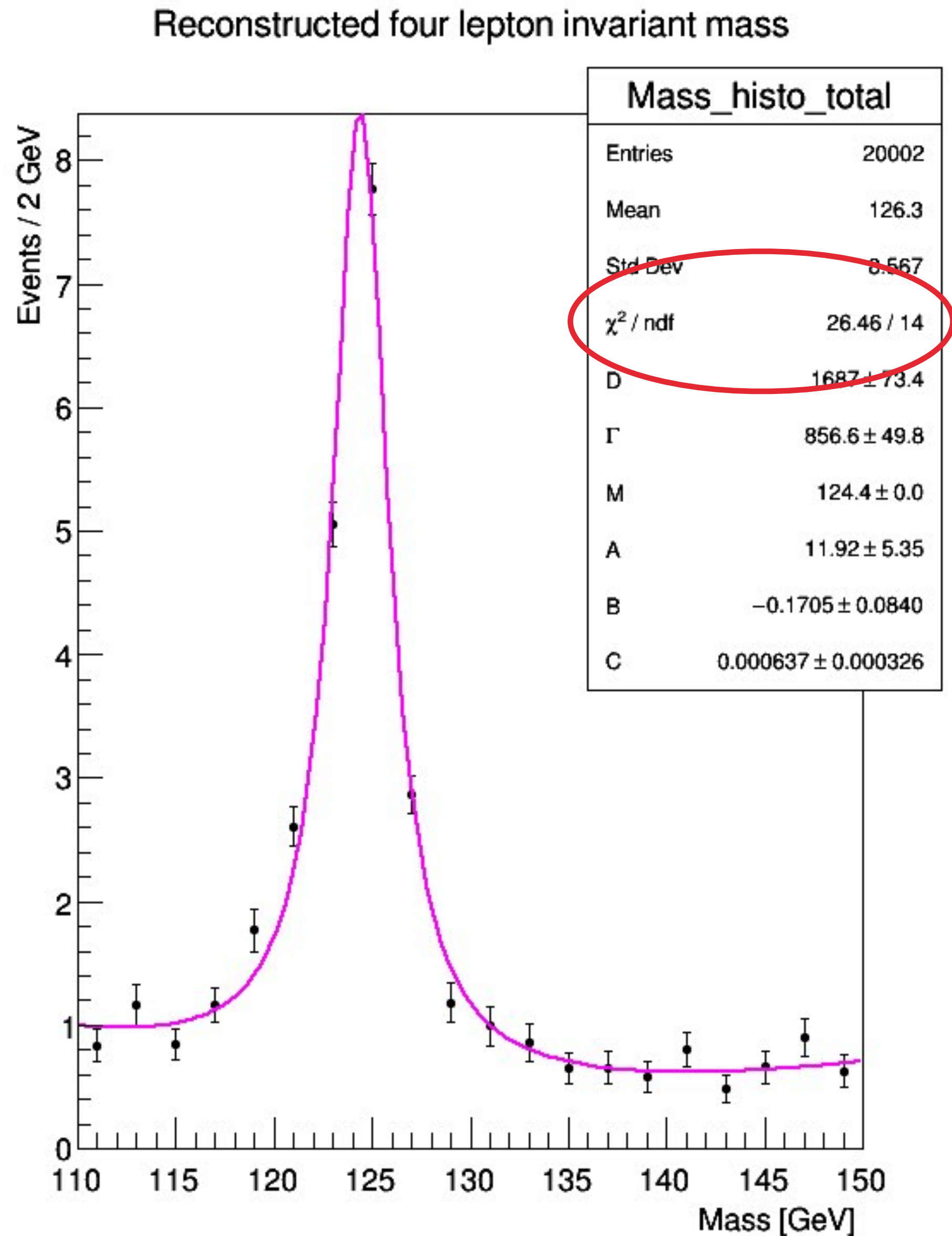
● The quantity  $\chi^2 = \sum_{i=1}^N \frac{(y_i^{data} - y_i^{ideal})^2}{(expected\ error)^2}$  gives information about the fit quality

| small $\chi^2$       | large $\chi^2$        |
|----------------------|-----------------------|
| good fit             | bad fit (bad model)   |
| overestimated errors | underestimated errors |

● Since  $\langle \chi^2 \rangle = N$ , easy way to estimate the fit quality is to check if  $\frac{\chi^2}{N.D.O.F} \approx 1$ , N.D.O.F is calculated as (N - free parameters)

# CHI-SQUARE FIT TEST - EXAMPLE

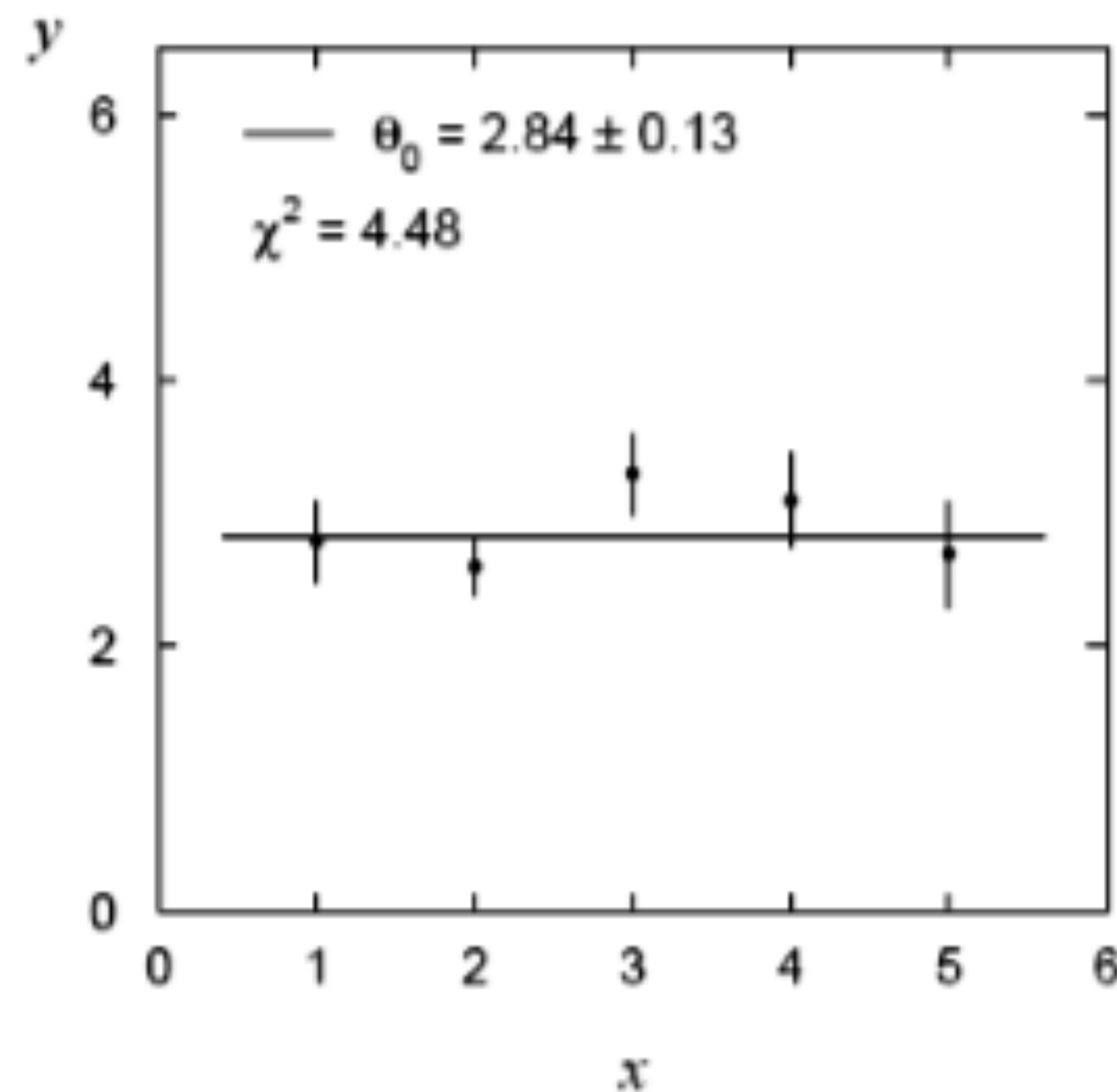
23



- LS has particularly desirable properties if  $f(x; \theta)$  is a linear function of  $\theta$ :

$$f(x; \theta) = \sum_{j=1}^m a_j(x) \theta_j, \text{ where } a_j(x) \text{ are linearly independent functions of } x$$

- estimators and their variances can be found analytically
- the estimators have zero bias and minimum variance

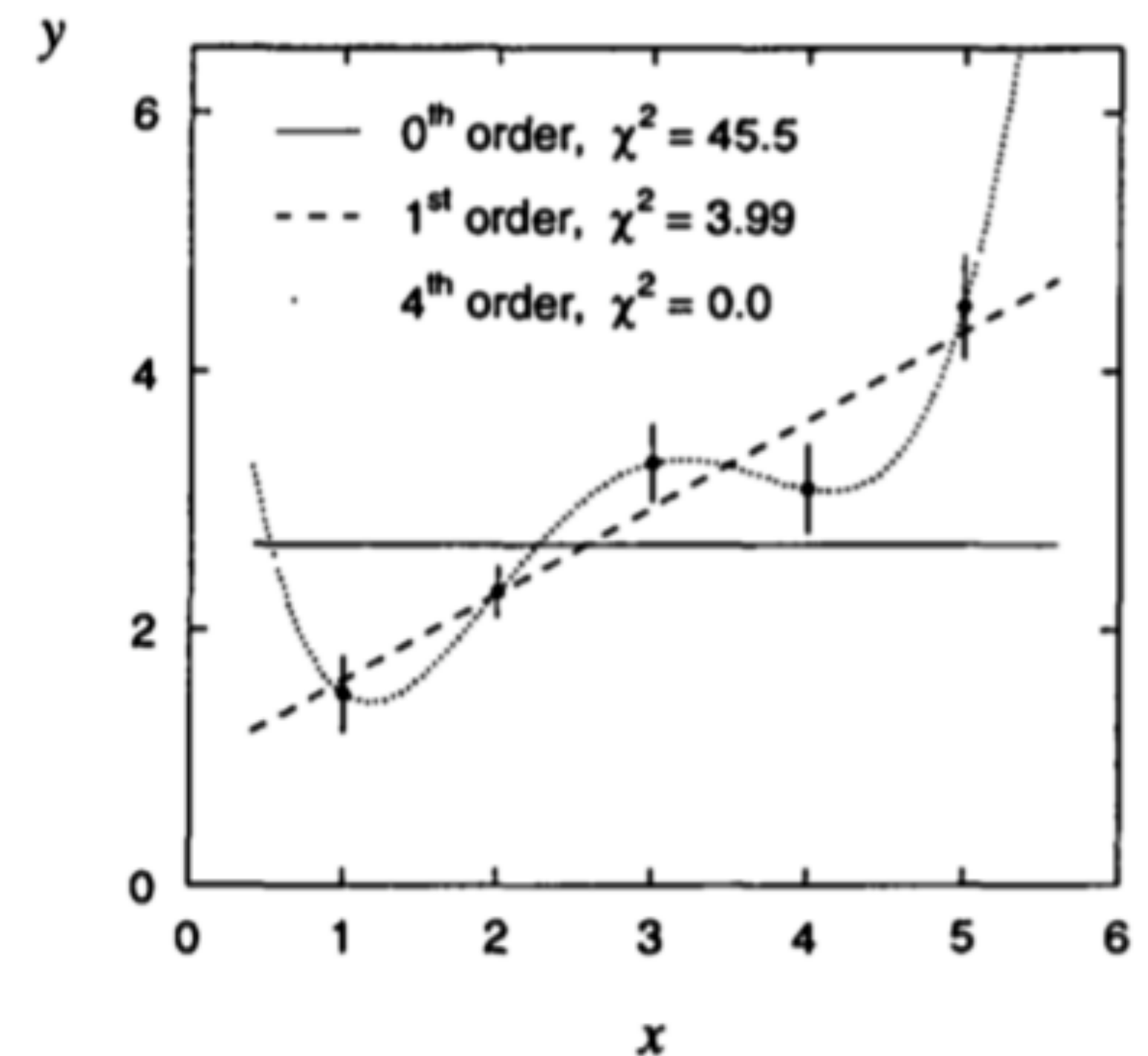




- Assume we measure 5 values of a quantity  $y$ , measured with errors  $\sigma_y$  at different values of  $x$

- For the fit function we try polynomial of order  $m$ : 
$$f(x; \theta) = \sum_{j=0}^m x^j \theta_j$$

- 0-th order: the weighted average
- 1-st order: a very good description
- 4-th order: equal number of parameters as points
- For Gaussian distributed  $y$  LS = ML!



- If  $y_i$  are Poissonian distributed variance is equal to the mean value so there are two choices

- Pearson's Chi-Square is  $\chi^2(\theta) = \sum_{i=1}^N \frac{(y_i - \lambda_i(\theta))^2}{\lambda_i(\theta)}$

- now  $\sigma_i$  depends on parameters  $\theta$  that complicates the minimisation procedure

- Neyman's or modified Chi-Square is  $\chi^2(\theta) = \sum_{i=1}^N \frac{(y_i - \lambda_i(\theta))^2}{y_i}$

- minimisation simpler but errors may be poorly estimated
  - problem for  $y_i = 0$

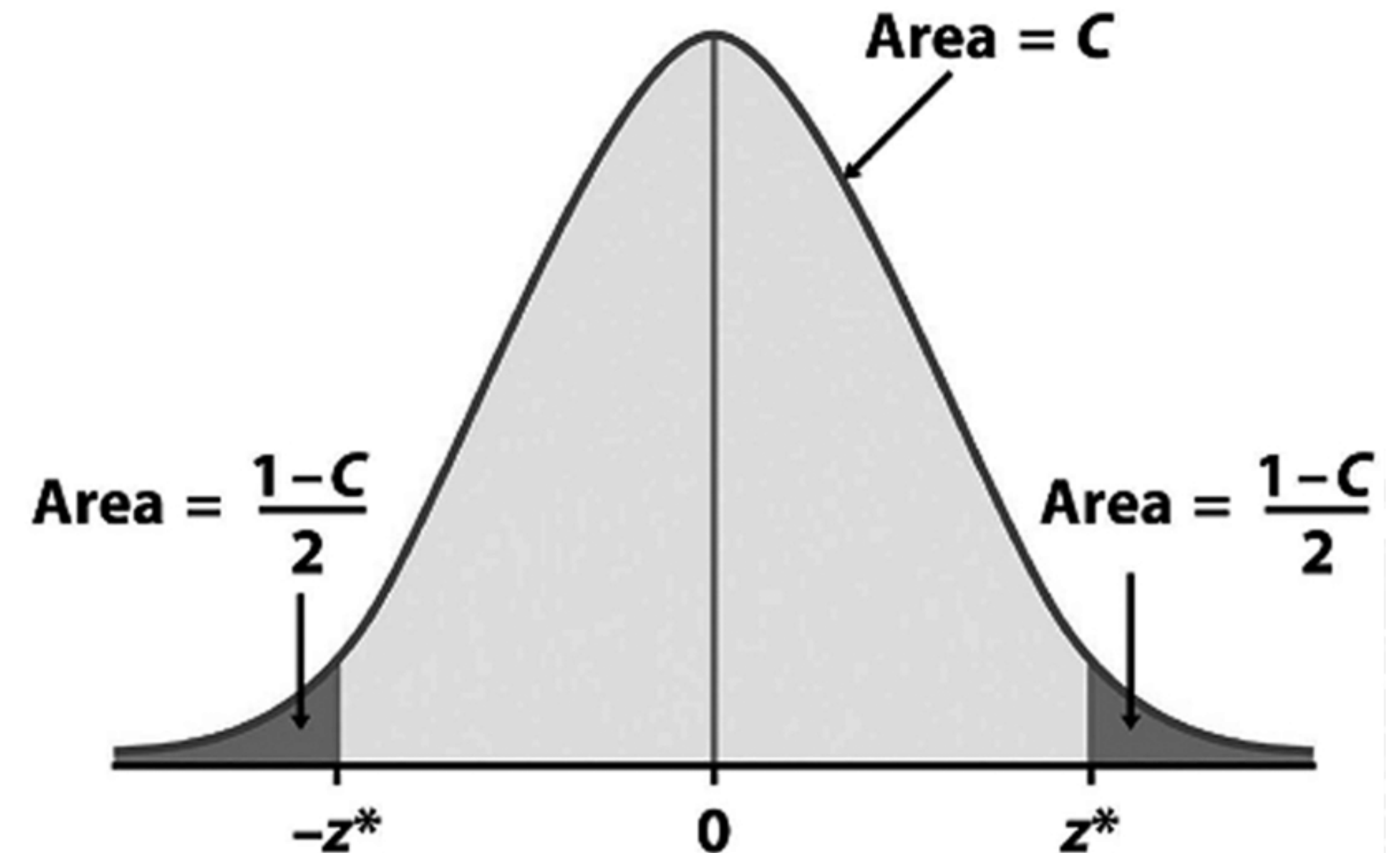


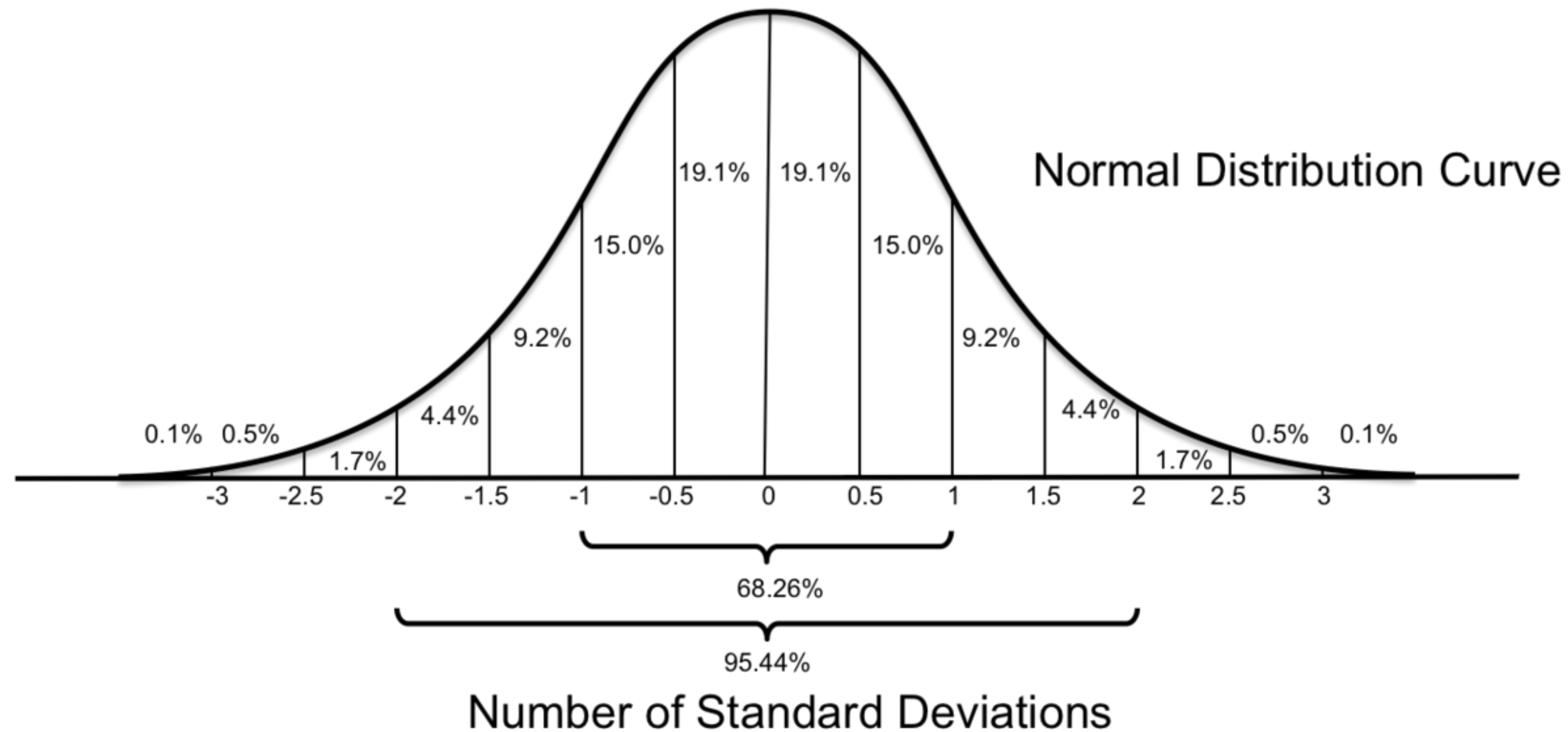
- 
- In addition to a “point estimate” of a parameter we should report an interval reflecting its statistical uncertainty.
  - Desirable properties of such an interval:
    - communicate objectively the result of the experiment
    - have a given probability of containing the true parameter
    - provide information needed to draw conclusions about the parameter
    - communicate incorporated prior beliefs and relevant assumptions
  - Often use  $\pm$  the estimated standard deviation ( $\sigma$ ) of the estimator
  - In some cases, however, this is not adequate:
    - estimate near a physical boundary
    - if the PDF is not Gaussian

- Let some measured quantity be distributed according to some PDF  $f(x; \theta)$ , we can determine the probability that  $x$  lies within some interval, with some confidence  $C$ :

$$P(x_- < x < x_+) = \int_{x_-}^{x_+} f(x; \theta) dx = C$$

- We say that  $x$  lies in the interval  $[x_-, x_+]$  with confidence  $C$





● If  $f(x; \theta)$  is a Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ :

●  $x_{\pm} = \mu \pm 1 \cdot \sigma$      $C = 68 \%$

●  $x_{\pm} = \mu \pm 2 \cdot \sigma$      $C = 95.4 \%$

●  $x_{\pm} = \mu \pm 1.64 \cdot \sigma$      $C = 90 \%$

●  $x_{\pm} = \mu \pm 1.96 \cdot \sigma$      $C = 95 \%$

$$P(x_- < x < x_+) = \int_{x_-}^{x_+} f(x; \theta) dx = C$$

● There are 3 conventional ways to choose an interval around the centre:

1) **Symmetric interval:**  $x_-$  and  $x_+$  equidistant from the mean

2) **Shortest interval:** minimizes  $(x_+ - x_-)$

3) **Central interval:**  $\int_{-\infty}^{x_-} f(x; \theta) dx = \int_{x_+}^{+\infty} f(x; \theta) dx = \frac{1 - C}{2}$

● For the Gaussian, and any symmetric distributions, 3 definitions are equivalent

- So far we have considered only two-tailed intervals, but sometimes one-tailed limits are also useful

- for example in the case of measuring a parameter near a physical boundary

- **Upper limit:**  $x$  lies below  $x_+$  at confidence level  $C$ : 
$$\int_{-\infty}^{x_+} f(x; \theta) dx = C$$

- **Lower limit:**  $x$  lies above  $x_-$  at confidence level  $C$ : 
$$\int_{x_-}^{+\infty} f(x; \theta) dx = C$$

- In a measurement two things involved:
  - True physical parameters:  $\theta^{true}$
  - Measurement of the physical parameter (parameter estimation):  $\hat{\theta}$
- Given the measurement  $\hat{\theta} \pm \sigma_{\theta}$  what can we say about  $\theta^{true}$  ?
- Can we say that  $\theta^{true}$  lies within  $\hat{\theta} \pm \sigma_{\theta}$  with 68% probability?
  - **NO!!!**
  - $\theta^{true}$  is **not a random variable**! It lies in the measured interval or it does not!
- We can say that if we repeat the experiment many times with the same sample size, construct the interval according to the same prescription each time, in 68% of the experiments  $\hat{\theta} \pm \sigma_{\theta}$  interval will cover  $\theta^{true}$ .



- There are two ways to obtain confidence intervals for the parameter estimated by the Maximum Likelihood method

- **Analytical way:**

- If we assume the **Gaussian approximation** we can estimate the confidence interval by matrix inversion:

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

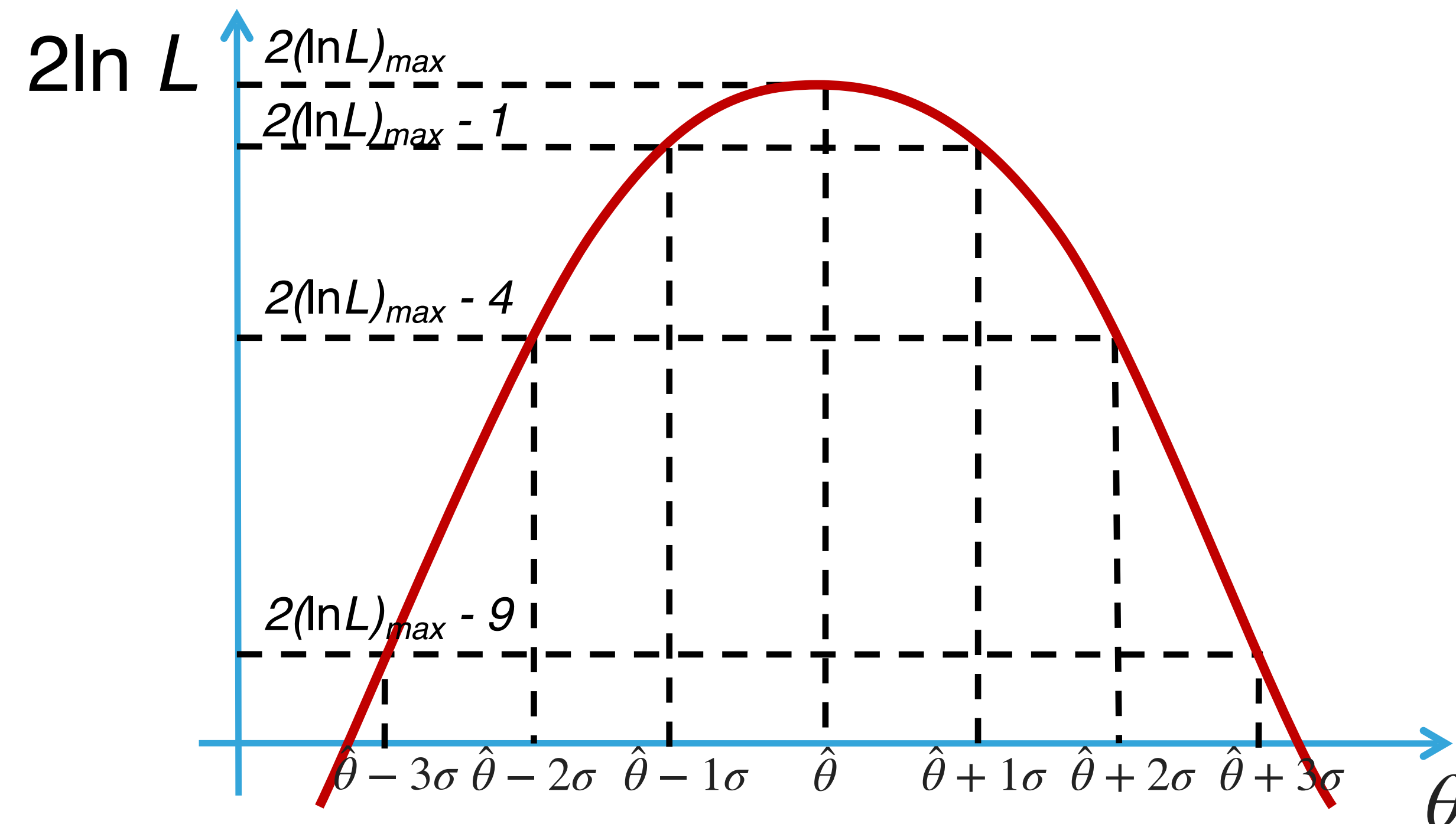
- If the likelihood function is non-Gaussian and in the limit of small number of events this approximation will give symmetrical interval while that might not be the case
  - Possible to solve by hand only for very simple PDF cases, otherwise numerical solution needed
    - Matrix inversion done with HESSE/MINUIT algorithm in ROOT

- **From the Log-Likelihood curve**

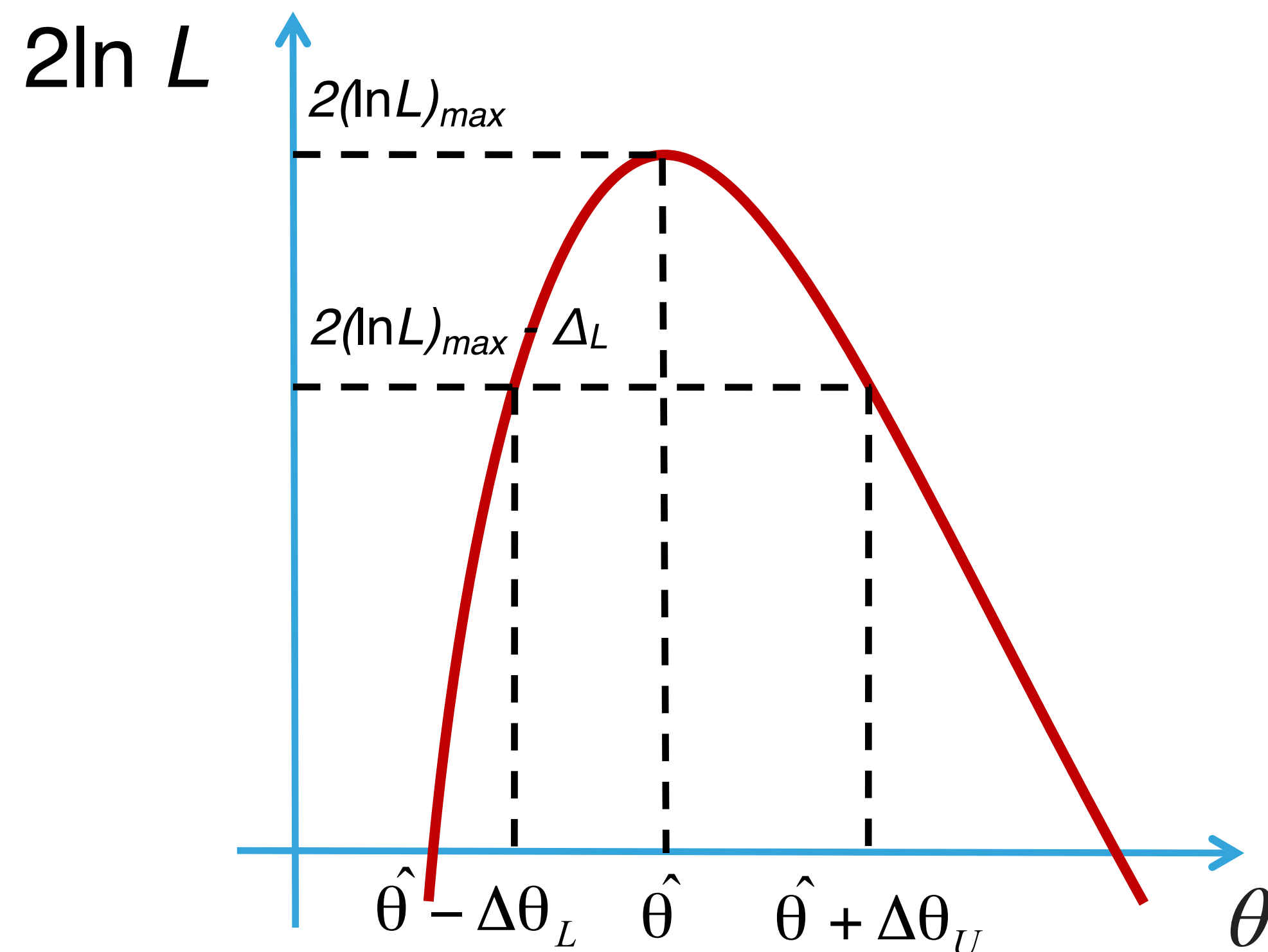
- Extract  $\sigma_{\hat{\theta}}$  from log-likelihood scan using:

$$\ln L(\hat{\theta} \pm N \cdot \sigma_{\hat{\theta}}) = \ln L_{\max} - \frac{N^2}{2}$$

- This is the same as looking for  $2\ln L_{\max} - N^2$



- The Log-Likelihood function can be asymmetric
  - for smaller samples, very non-Gaussian PDFs, non-linear problems,...
- The confidence interval is still extracted from the Log-Likelihood curve using the same prescription
  - This leads to asymmetrical confidence interval that should be used when quoting the final result

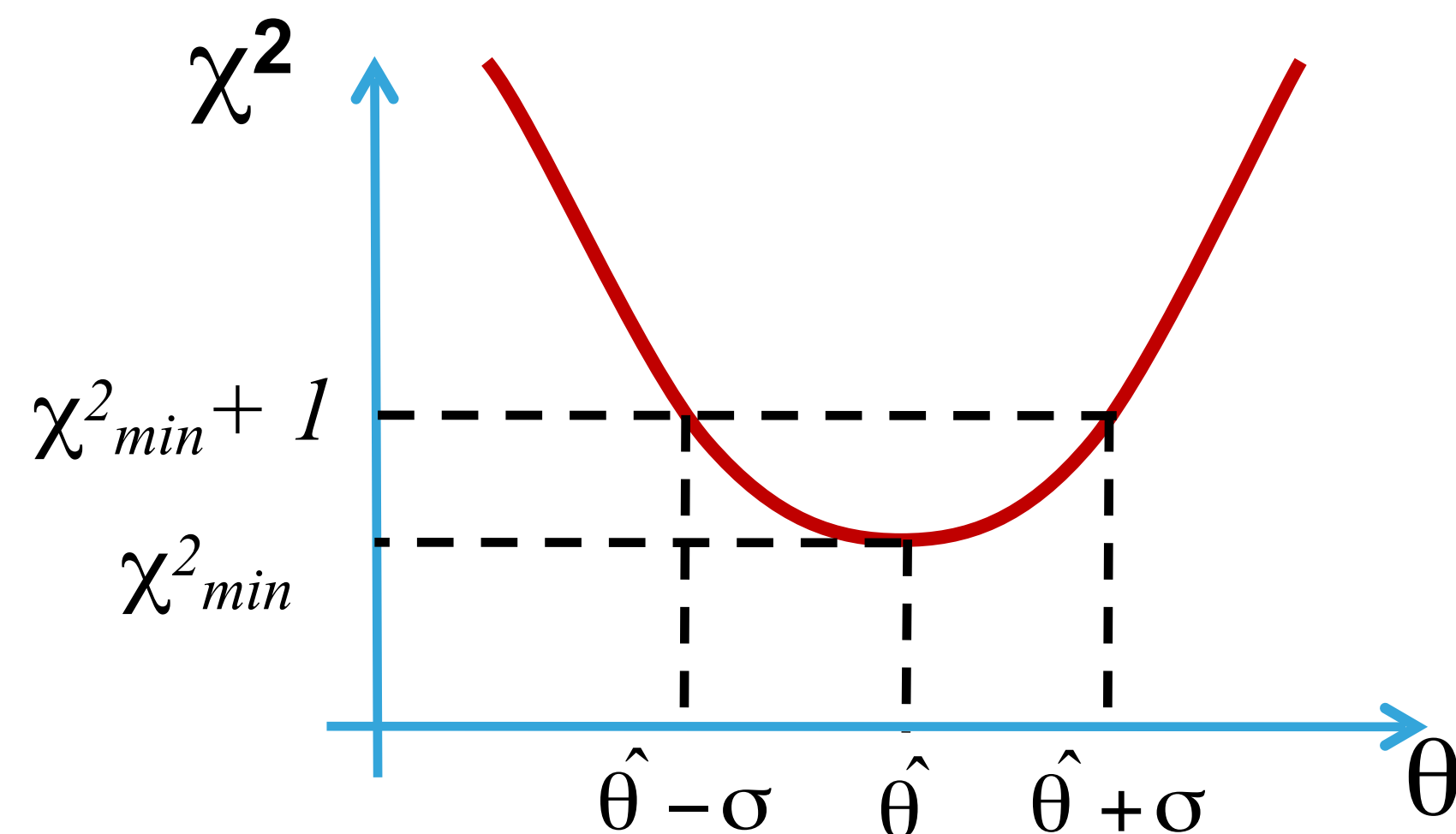


| CL    | $\Delta_L$ |
|-------|------------|
| 68.27 | 1          |
| 95.45 | 4          |
| 99.73 | 9          |

- The confidence intervals for the Least Squares (Chi-Square) method are obtained in the identical way as for the Maximum likelihood method
- **Analytical way of matrix inversion:**
  - Solving analytically (or numerically):

$$\text{cov}^{-1}(\theta_i, \theta_j) = \frac{1}{2} \frac{\partial^2 \chi^2}{\partial \theta_i \partial \theta_j} \bigg|_{\theta=\hat{\theta}}$$

- **From the Chi-Square curve**

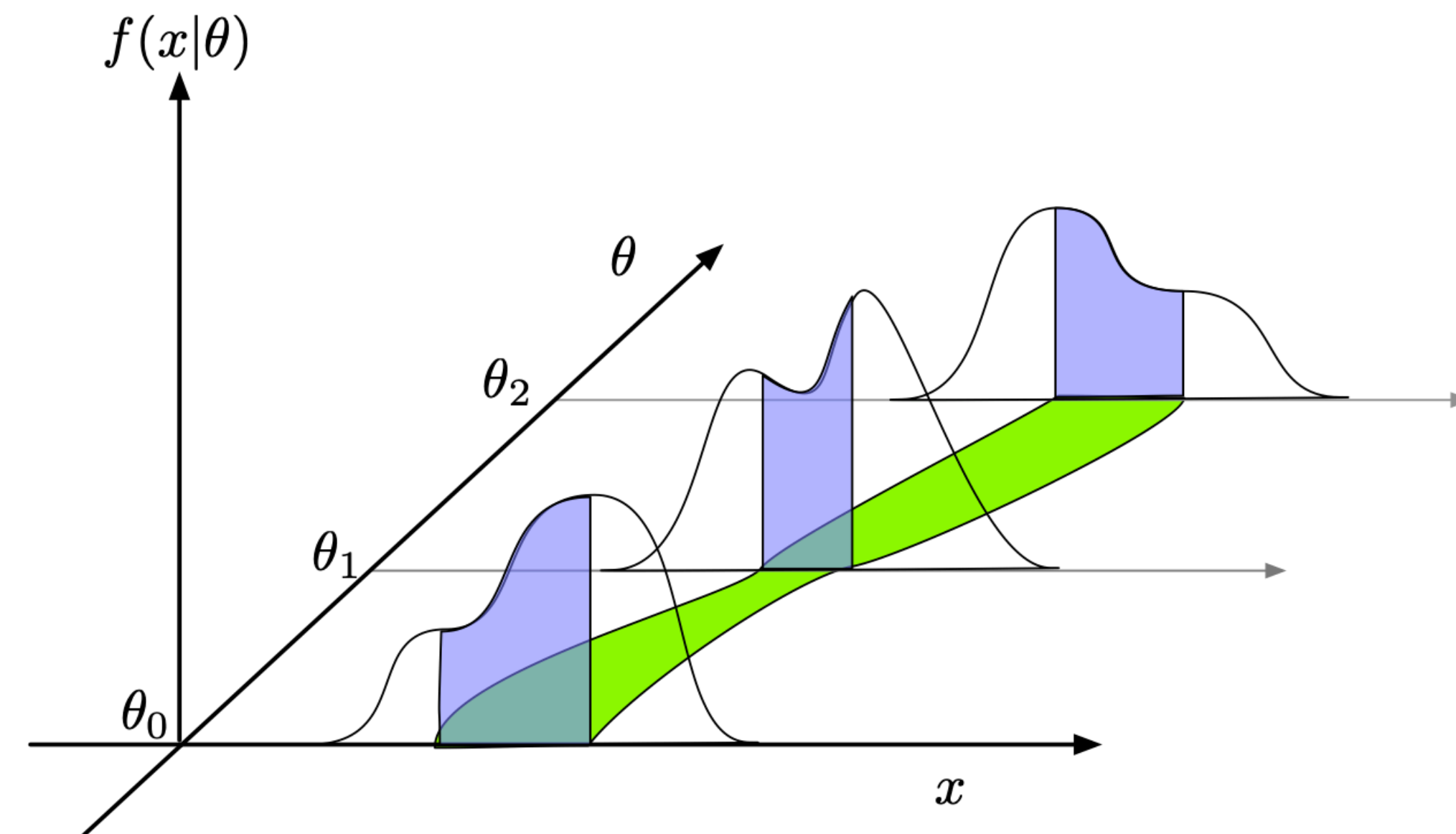


| CL    | $\Delta_L$ |
|-------|------------|
| 68.27 | 1          |
| 95.45 | 4          |
| 99.73 | 9          |

- Using frequentist approach Neyman defines confidence interval of the unknown parameter  $\theta$ :

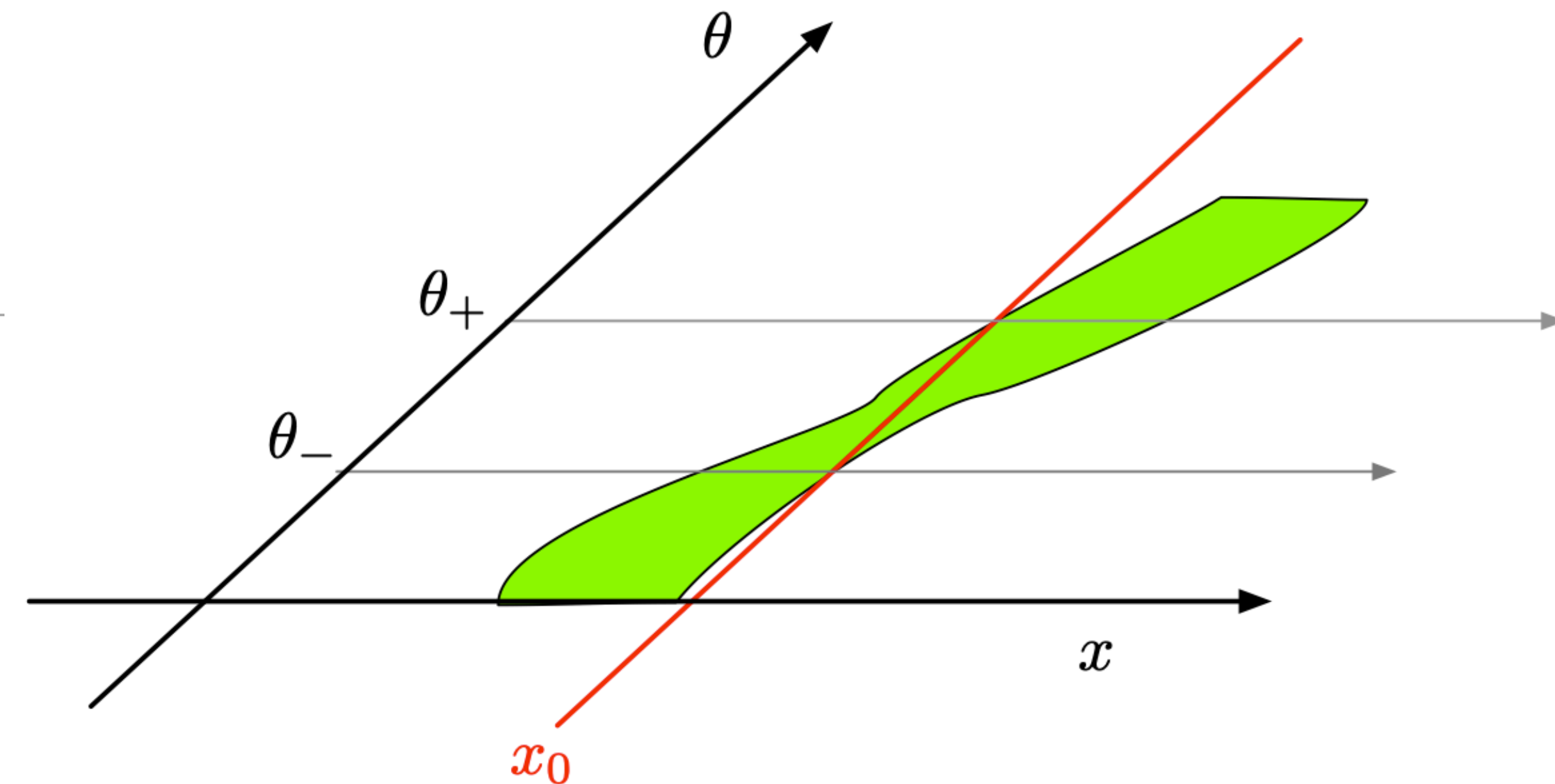
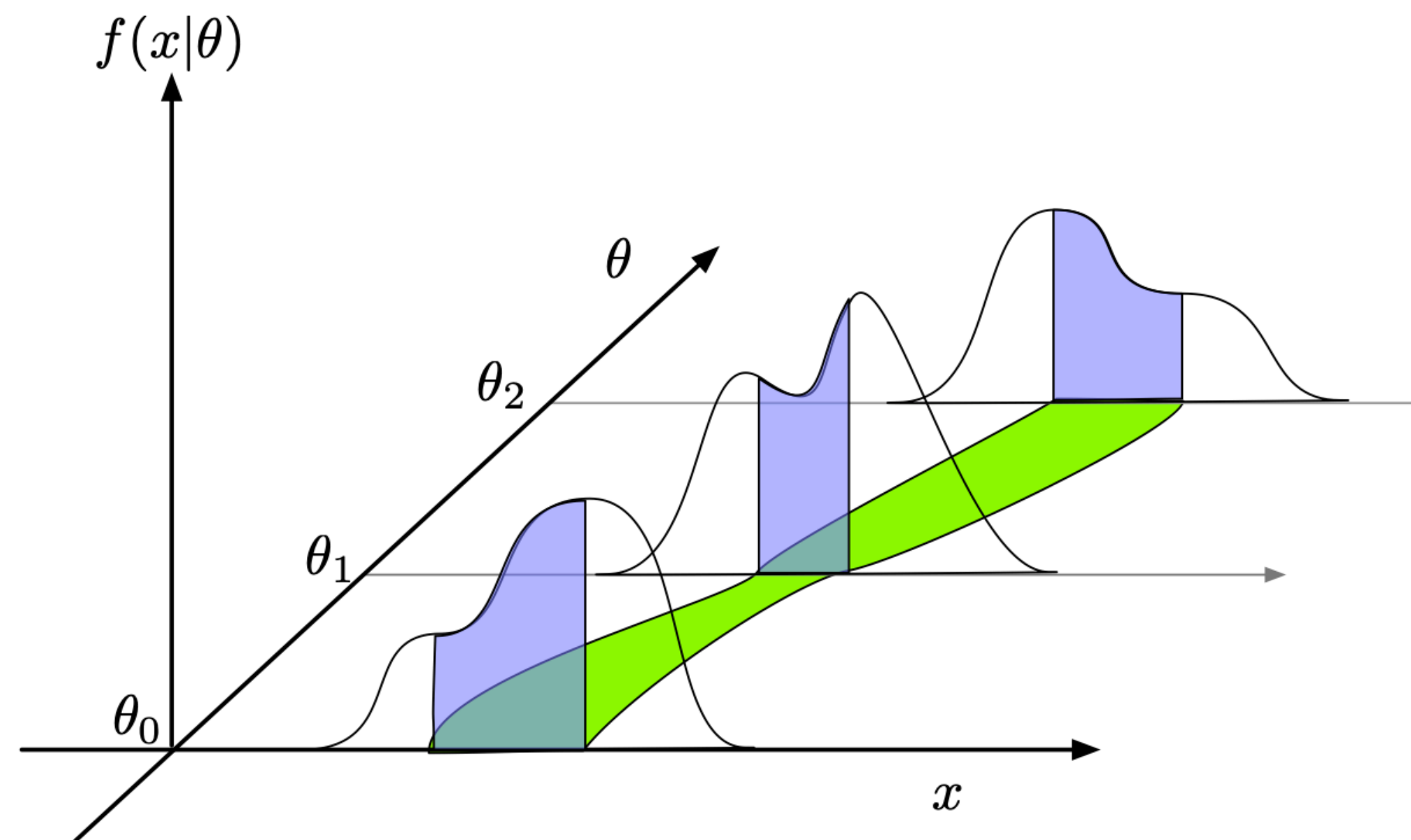
$$P(x_1 < x < x_2; \theta) = \int_{x_1}^{x_2} f(x; \theta) dx = CL$$

- $x$  is the measurement and  $CL$  is predefined confidence level
- Union of  $[x_1, x_2]$  segments for all values of the parameter  $\theta$  is known as the **confidence belt**
- All of these steps are performed **before measuring the data**





- Now we perform the measurement to obtain  $x_0$
- the points  $\theta$  where the belt intersects  $x_0$  are part of the **confidence interval**  $[\theta_-, \theta_+]$  for this measurement
- For every point  $\theta$ , if it were true, the data would fall in its acceptance region with probability CL, so the interval  $[\theta_-, \theta_+]$  covers the true value with probability CL



- Still a frequentist approach!

- In Bayesian statistics, all knowledge about parameter  $\theta$  is contained in the posteriori PDF  $p(\theta | x)$ :

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

- which gives the degree of belief for  $\theta$  to have values in certain region given we observe the data  $x$ 
  - $\pi(\theta)$  is the prior PDF for  $\theta$ , reflecting experimenter's subjective degree of belief about  $\theta$  before the measurement
  - $L(x | \theta)$  is the Likelihood function, i.e. the PDF for the data given a certain value of  $\theta$
  - The dominator simply normalises the posteriori PDF to unity

- We can now use Bayesian statistics to express our degree of belief about  $\theta$  before the measurement:

$$\pi(\theta) = \begin{cases} 0, & m < 0 \\ \text{constant}, & m \geq 0 \end{cases}$$

- assuming a Gaussian PDF we can calculate

$$p(\theta | x) = \frac{e^{-\frac{(x - \theta)^2}{2\sigma^2}}}{\int_0^{\infty} e^{-\frac{(x - \theta')^2}{2\sigma^2}} d\theta'}$$