## Applied Statistics

#### The Chi-Square Distribution, Fit & Test

The Chi-Square fit is also (originally) known as Method of Least Squares, though this method does not include uncertainties on the data points involved.





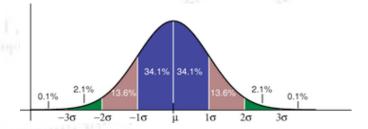






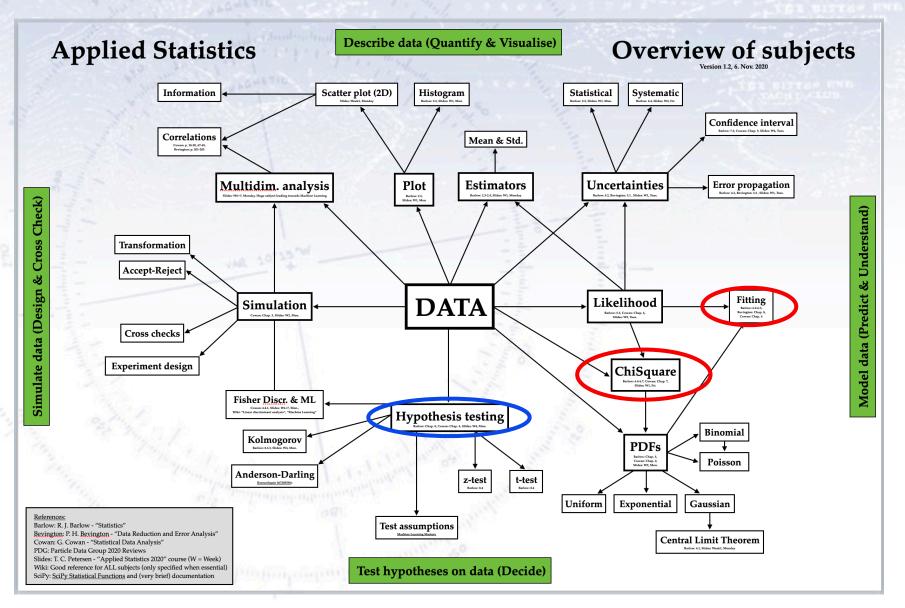


Troels C. Petersen (NBI)



"Statistics is merely a quantisation of common sense"

## ChiSquare



# The discovery of Ceres

Dwarf planet and the largest astroid (r=487km)

Theta Ophiuchi



# The discovery of Ceres

Dwarf planet and the largest astroid (r=487km)



On the 1st of January 1801 Giuseppe Piazzi discovered "new light" and could follow this comet/planet until 11th of February. He published the positions, but due to Ceres being behind the sun, it would be out of sight until the following winter. Following the calculations of a 24 year old mathematician/physicist, it was recovered on the 31st of December 1801 by von Zach and H. Olbers.

The young man's name was Carl Friedrich Gauss, and the method he used/invented for this was...

# The discovery of Ceres

Dwarf planet and the largest astroid (r=487km)



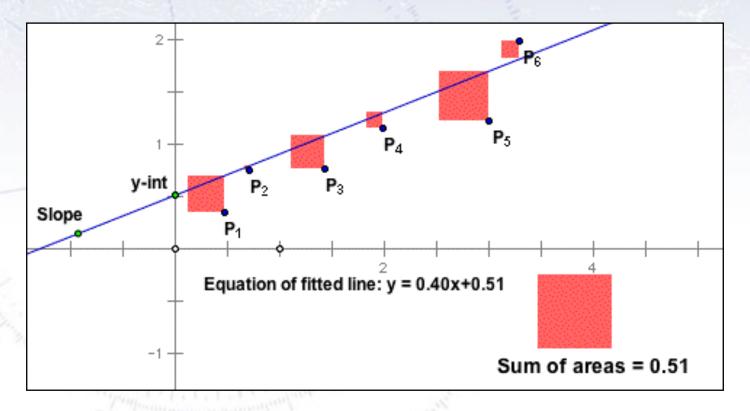
On the 1st of January 1801 Giuseppe Piazzi discovered "new light" and could follow this comet/planet until 11th of February. He published the positions, but due to Ceres being behind the sun, it would be out of sight until the following winter. Following the calculations of a 24 year old mathematician/physicist, it was recovered on the 31st of December 1801 by von Zach and H. Olbers.

The young man's name was Carl Friedrich Gauss, and the method he used/invented for this was...

...method of least squares!

South

The problem at hand is determining the curve that best fitted data:

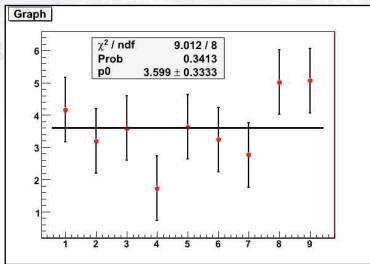


The "best fit" is found by minimising the sum of the squares...

Originally, uncertainties were not included (not "invented" yet!)

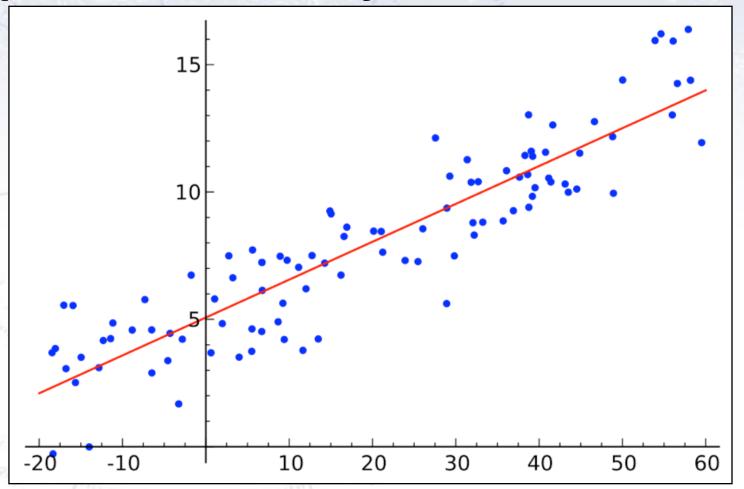
The method of least squares is a standard approach to the approximate solution of **overdetermined systems**, i.e. sets of equations in which there are **more equations than unknowns**.

"Least squares" means that the overall solution minimises the sum of the squares of the errors made in solving every single equation.



The most important application is in **data fitting**. The best fit in the least-squares sense minimises the **sum of squared residuals**, a residual being the difference between an observed value and the fitted value provided by a model.

The problem at hand is determining the curve that best fitted data:



Originally, uncertainties were not included (not "invented" yet!)

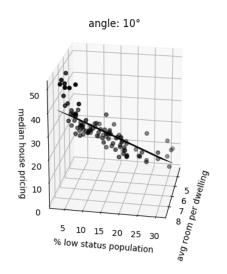
## MultiDim Linear Regression

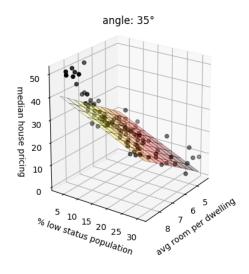
The advantage of the linear regression (without uncertainties) is that it can be done "easily" in multiple dimensions.

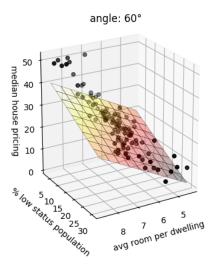
This has been used in several fields of science, in particular economy, but also medicin.

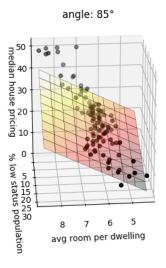
The model typically used is given below, and the analysis called "ANOVA".

$$f(x) = ax_1 + bx_2 + c$$



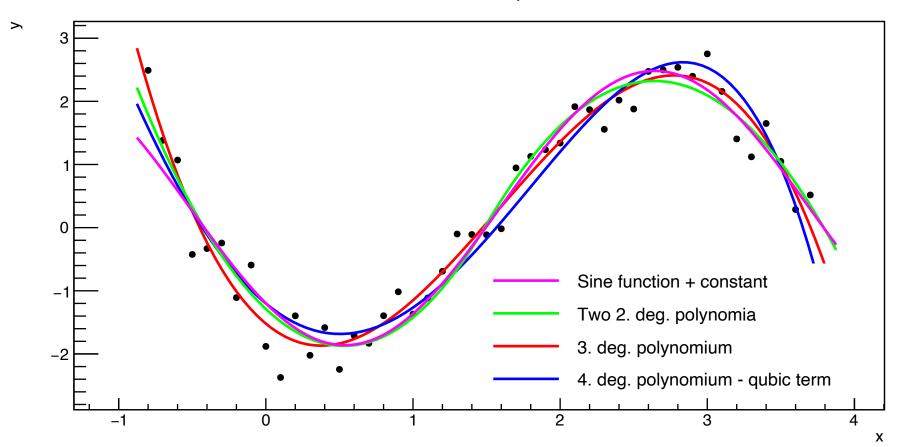






Look at the figure below, and determine which curve fits best...

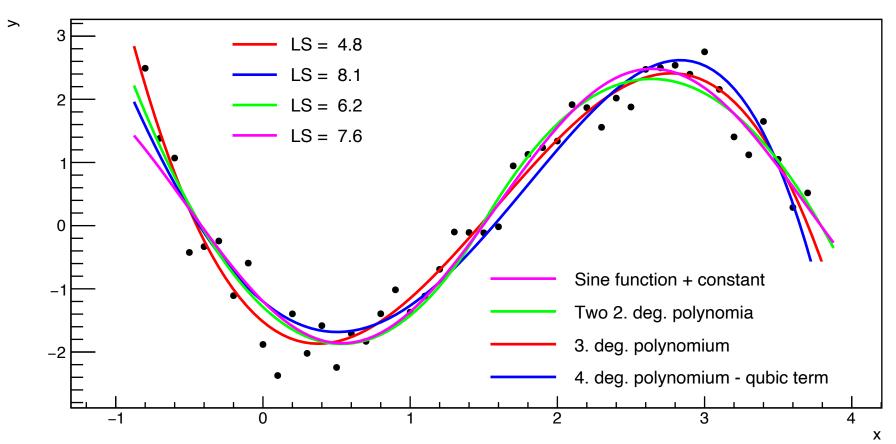
Illustration of Least Squares' Method



Well, what do you define as "best"?

Look at the figure below, and determine which curve fits best...

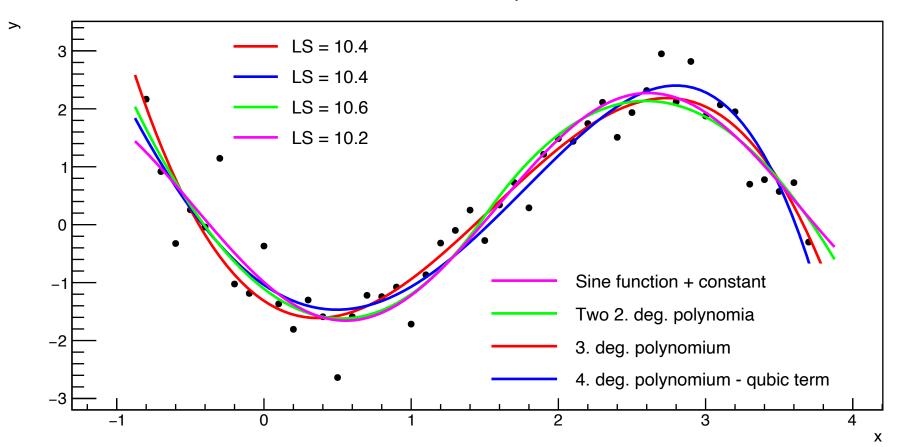
Illustration of Least Squares' Method



Well, what do you define as "best"? And how good is it?!?

Look at the figure below, and determine which curve fits best...

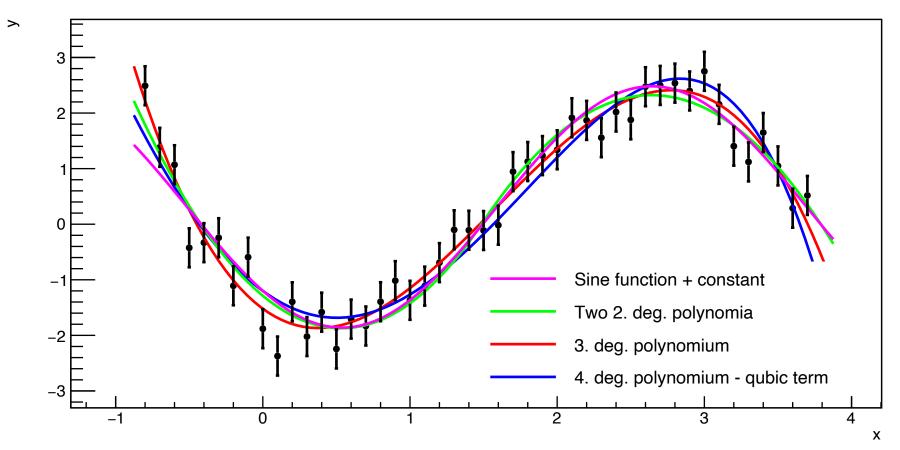
Illustration of Least Squares' Method



Well, what do you define as "best"? And how good is it?!?

Look at the figure below, and determine which curve fits best...

Illustration of ChiSquare Method



Well, what do you define as "best"?

#### Defining the Chi-Square

Problem Statement: Given N data points  $(x,y,\sigma_y)$ , adjust the parameter(s)  $\theta$  of a model, such that it fits data best.

The best way to do this, given uncertainties  $\sigma_i$  on  $y_i$  is by minimising:

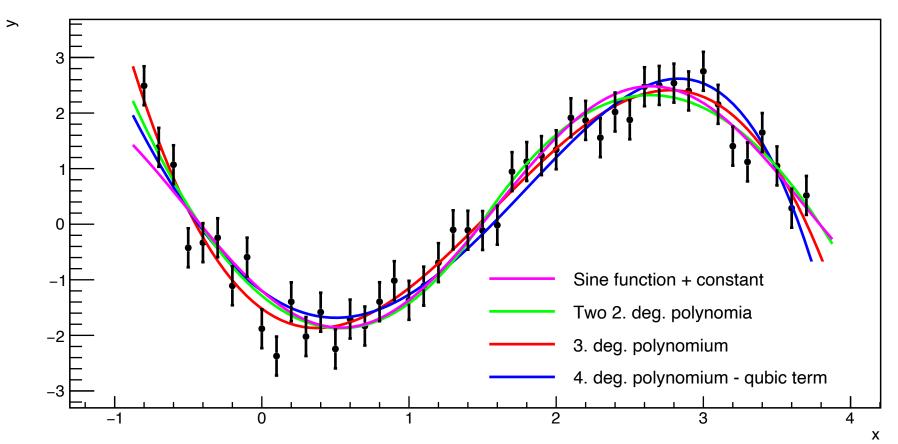
$$\chi^2(\theta) = \sum_{i}^{N} \frac{(y_i - f(x_i, \theta))^2}{\sigma_i^2}$$

#### The power of this method is hard to overstate!

Not only does it provide a simple, elegant and unique way of fitting data, but more importantly it provides a **goodness-of-fit measure**.

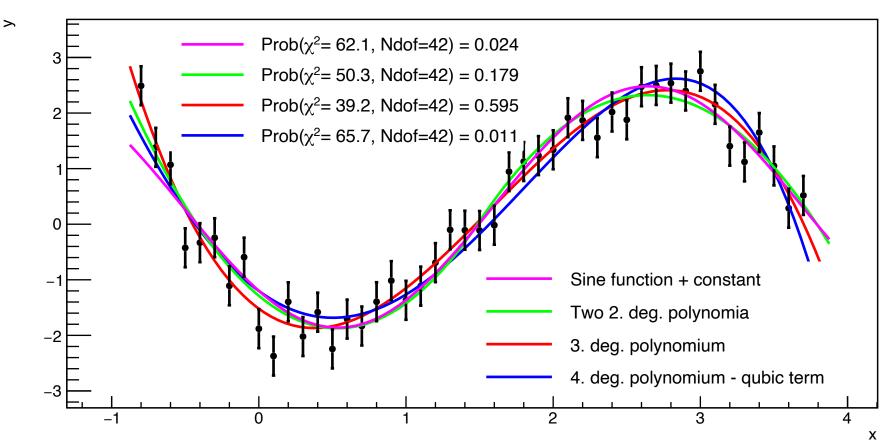
#### This is the Chi-Square test!

Look at the figure below, and determine which curve fits best...



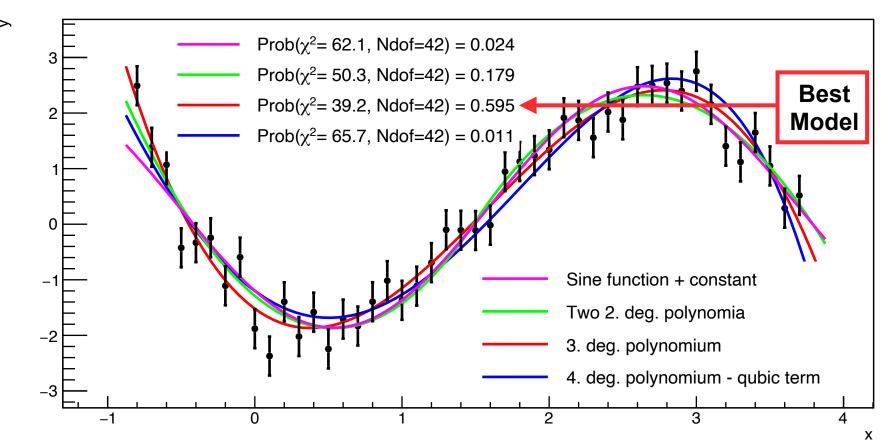
Well, what do you define as "best"?

Look at the figure below, and determine which curve fits best...



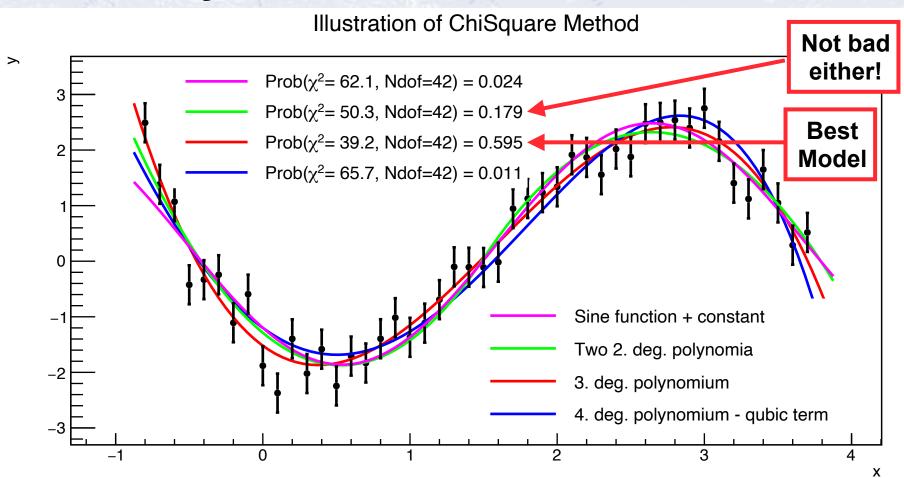
Well, what do you define as "best"? The Chi2 quantifies this!

Look at the figure below, and determine which curve fits best...



Well, what do you define as "best"? The Chi2 quantifies this!

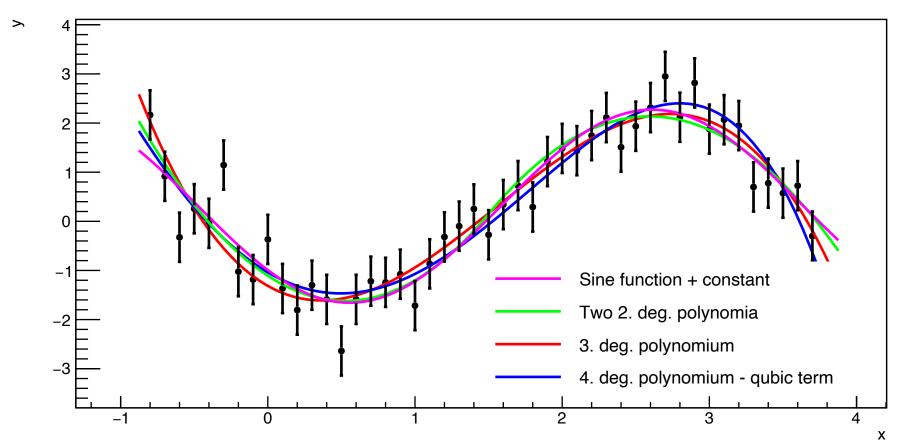
Look at the figure below, and determine which curve fits best...



Well, what do you define as "best"? The Chi2 quantifies this!

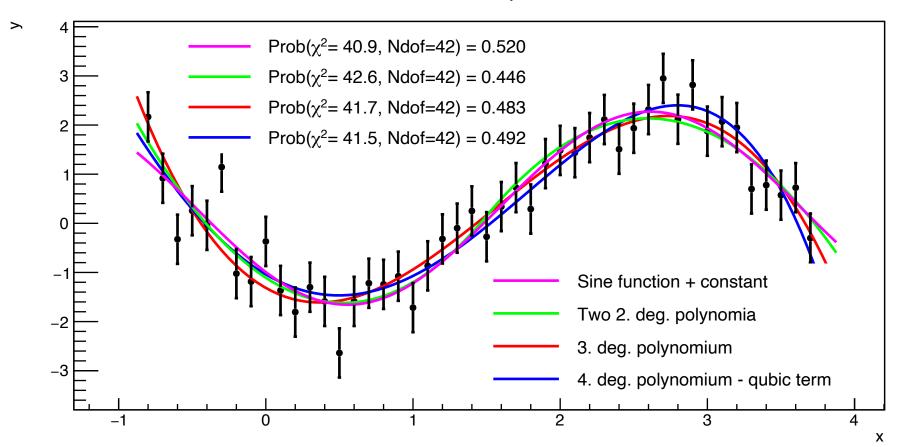
Look at the figure below, and determine which curve fits best...

Illustration of ChiSquare Method



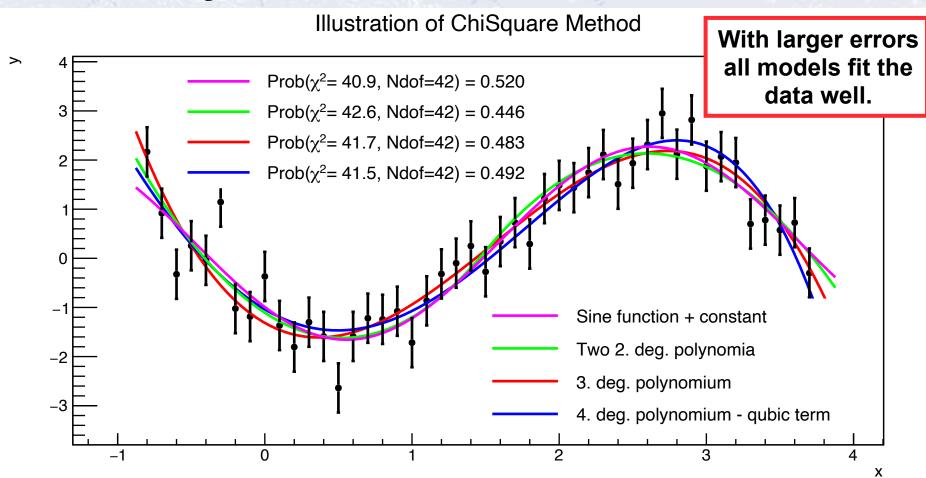
What about now with larger errors?

Look at the figure below, and determine which curve fits best...



What about now with larger errors?

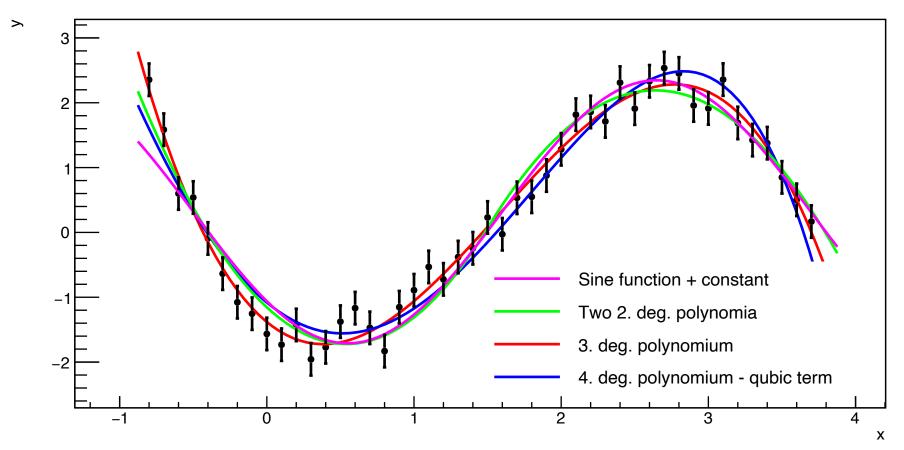
Look at the figure below, and determine which curve fits best...



What about now with larger errors?

Look at the figure below, and determine which curve fits best...

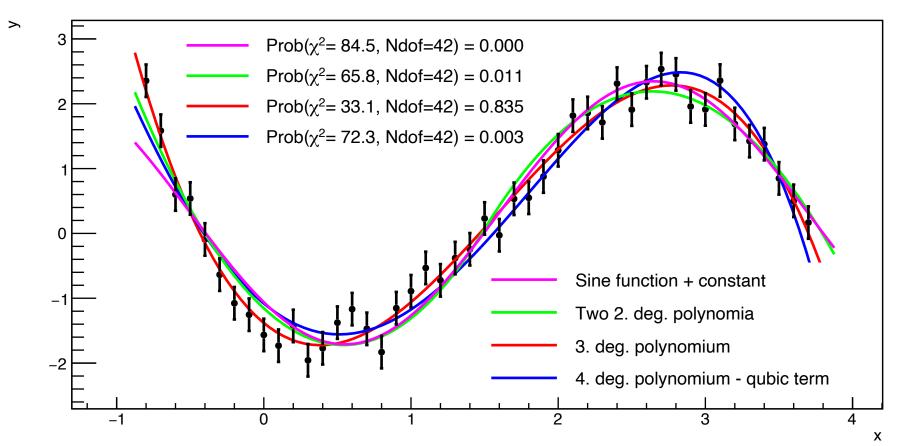
Illustration of ChiSquare Method



What does **smaller** errors do?

Look at the figure below, and determine which curve fits best...

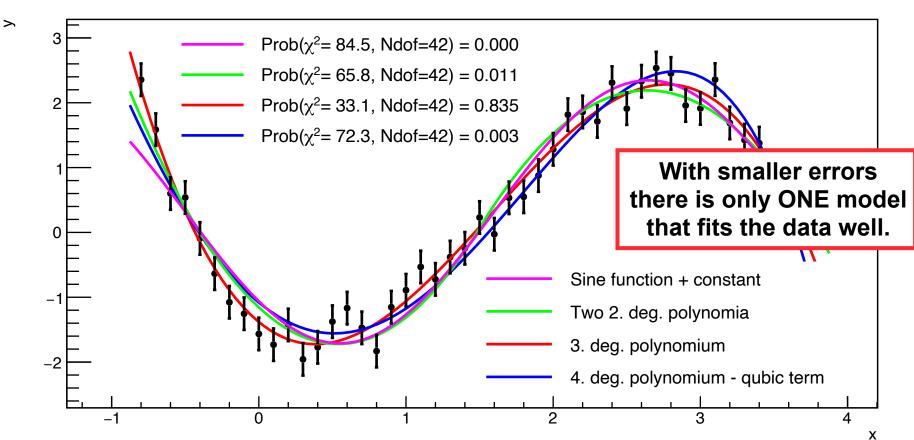
Illustration of ChiSquare Method



What does **smaller** errors do?

Look at the figure below, and determine which curve fits best...

Illustration of ChiSquare Method



What does **smaller** errors do?

#### Defining the Chi-Square

Problem Statement: Given N data points  $(x,y,\sigma_y)$ , adjust the parameter(s)  $\theta$  of a model, such that it fits data best.

The best way to do this, given uncertainties  $\sigma_i$  on  $y_i$  is by minimising:

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

#### The power of this method is hard to overstate!

Not only does it provide a simple, elegant and unique way of fitting data, but more importantly it provides a **goodness-of-fit measure**.

#### This is the Chi-Square test!

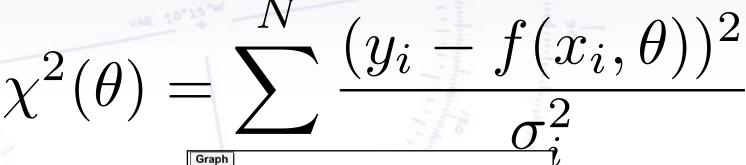
#### Defining the Chi-Square

Problem St Note that when doing a weighted mean, arameter(s) one should check if the measurements agree with each other!

The best w

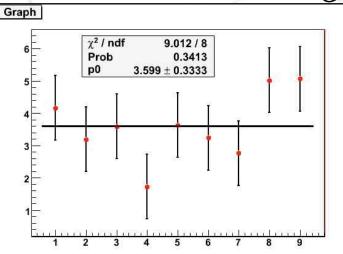
This can be done with a ChiSquare test.

imising:



#### The power of

Not only does it pro data, but more imp



#### to overstate!

que way of fitting of-fit measure.

test!

#### Weighted mean & ChiSquare

The weighted mean is actually an **analytical ChiSquare minimisation to a constant**. The result is the same, and one can then calculate Prob( $\chi^2$ , Ndof).

#### Example:

Data (from pendulum experiment) could be four length measurement (in mm):

 $d: [17.8 \pm 0.5, 18.1 \pm 0.3, 17.7 \pm 0.5, 17.7 \pm 0.2]$ 

The output from the above data is (many digits for *checks only*):

Mean = 17.8098 mm

Error on mean = 0.15057 mm

ChiSquare = 1.28574

Ndof = 3

Probability = 0.7325213

NOTE: This seems a very nice (and precise) result, and it may very well be. BUT, it might also be, that we all four estimated it from the same photo or similarly, which could be biased by an angled view. Then we would be fooling ourselves. We will discuss such "systematic uncertainties" more!

#### Weighted mean & ChiSquare

The weighted mean is actually an **analytical ChiSquare minimisation to a constant**. The result is the same, and one can then calculate Prob( $\chi^2$ , Ndof).

#### Example:

Data (from pendulum experiment) could be four length measurement (in mm):

$$d: [17.8 \pm 0.5, 18.1 \pm 0.3, 17.7 \pm 0.5, 17.7 \pm 0.2]$$

The output from the above data is (many digits for *checks only*):

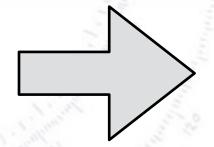
Mean = 17.8098 mm

Error on mean = 0.15057 mm

ChiSquare = 1.28574

Ndof = 3

Probability = 0.7325213

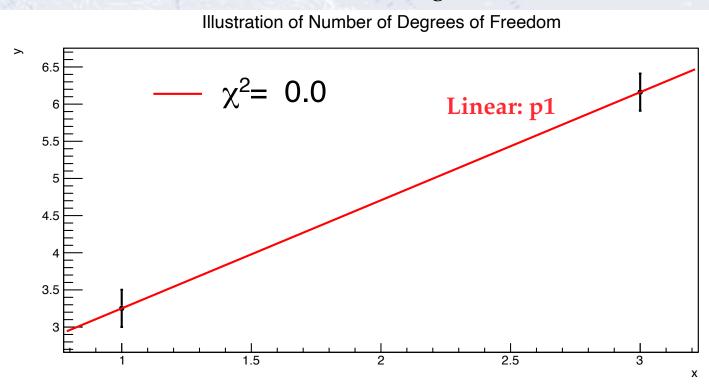


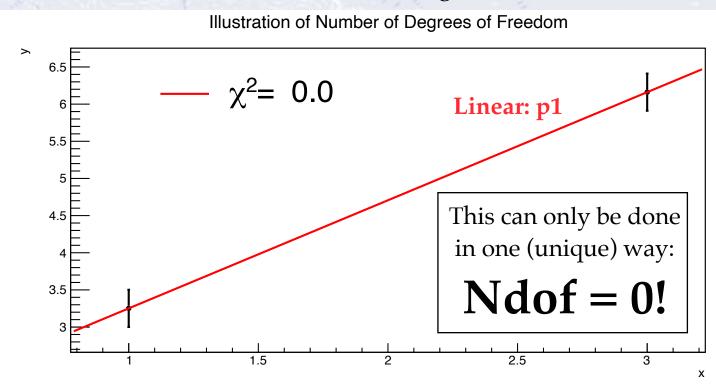
 $d = (17.81 \pm 0.15) \text{ mm}$  $p(\chi^2=1.3, N_{dof}=3) = 0.73$ 

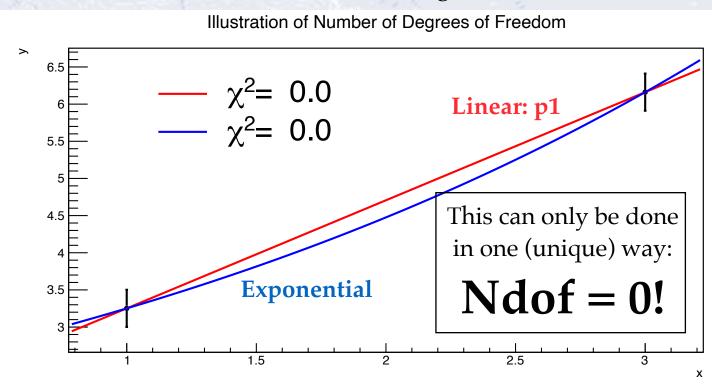
NOTE: This seems a very nice (and precise) result, and it may very well be. BUT, it might also be, that we all four estimated it from the same photo or similarly, which could be biased by an angled view. Then we would be fooling ourselves. We will discuss such "systematic uncertainties" more!

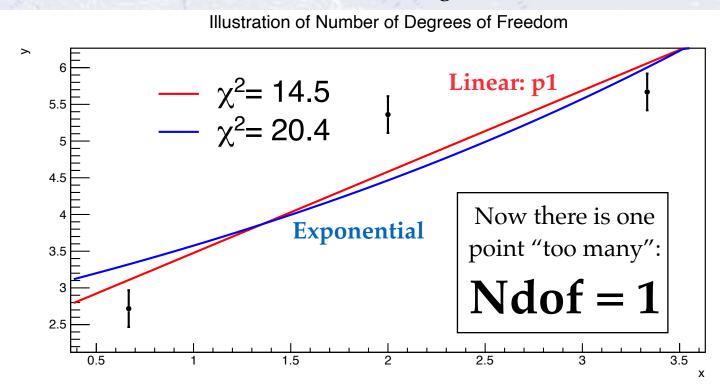
## Why the ChiSquare is great

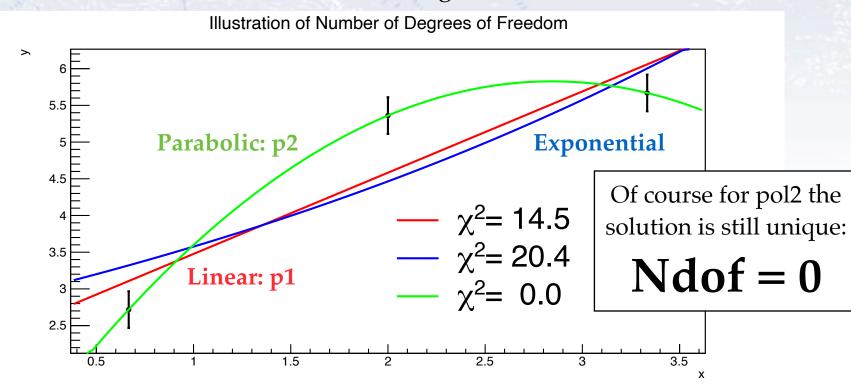
...but not its magic





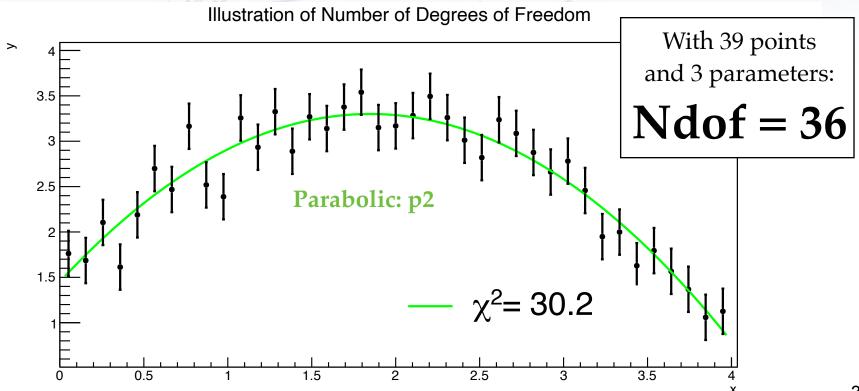






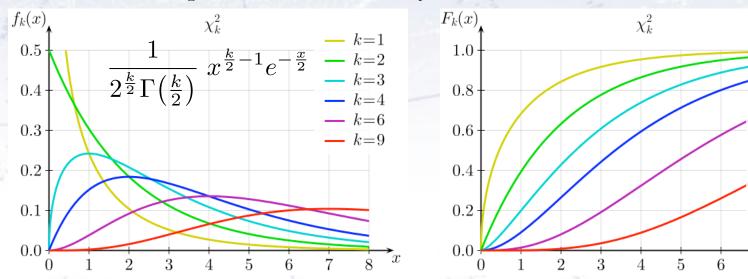
The number of degrees-of-freedom, Ndof, can be calculated as the number of points in the fit minus the number of parameters in the fit function:

$$N_{\text{dof}} = N_{\text{data points}} - N_{\text{fit variables}}$$



#### The Chi-Square distribution and test

The **Chi-Square distribution** for  $N_{dof}$  degrees of freedom is the distribution of the sum of the squares of  $N_{dof}$  normally distributed random variables.



The **Chi-Square test** consists of comparing the Chi-Square value obtained from a fit with the PDF of expected Chi-Square values. This allows the calculation of the *probability* of observing something with the same Chi-Square value or higher...

Rule of thumb: Chi-Square should roughly match N<sub>dof</sub>

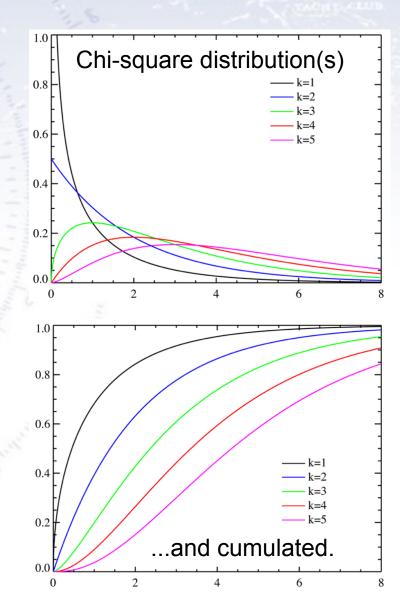
Given a **Chi-square value** and a **number of degrees of freedom** (Ndof), one can obtain a **"goodness-of-fit"**.

It is known, what Chi-square values to expect given the Ndof. One can therefore compare to this (Chi-square) distribution, and see...

what is the probability of getting this Chi-square value or something worse, assuming this is the correct fit function!

#### Example:

A fit gave the Chi-square 7.1 with 5 dof. The chance of getting this Chi-square or worse is... (reading the pink bottom curve (Ndof = k = 5) at 7.1)...



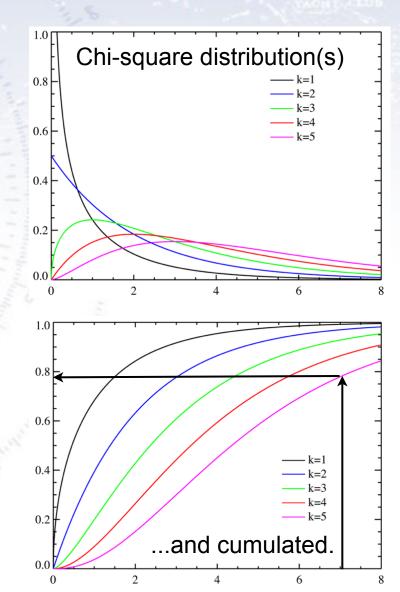
Given a **Chi-square value** and a **number of degrees of freedom** (Ndof), one can obtain a **"goodness-of-fit"**.

It is known, what Chi-square values to expect given the Ndof. One can therefore compare to this (Chi-square) distribution, and see...

what is the probability of getting this Chi-square value or something worse, assuming this is the correct fit function!

#### Example:

A fit gave the Chi-square 7.1 with 5 dof. The chance of getting this Chi-square or worse is... (reading the pink bottom curve (Ndof = k = 5) at 7.1)... 1 - 0.78 = 22%



In the table below, one can get a quick estimate for low  $N_{\text{dof}}$ .

Degrees of freedom (df)	$\chi^2$ value <sup>[16]</sup>										
1	0.004	0.02	0.06	0.15	0.46	1.07	1.64	2.71	3.84	6.64	10.83
2	0.10	0.21	0.45	0.71	1.39	2.41	3.22	4.60	5.99	9.21	13.82
3	0.35	0.58	1.01	1.42	2.37	3.66	4.64	6.25	7.82	11.34	16.27
4	0.71	1.06	1.65	2.20	3.36	4.88	5.99	7.78	9.49	13.28	18.47
5	1.14	1.61	2.34	3.00	4.35	6.06	7.29	9.24	11.07	15.09	20.52
6	1.63	2.20	3.07	3.83	5.35	7.23	8.56	10.64	12.59	16.81	22.46
7	2.17	2.83	3.82	4.67	6.35	8.38	9.80	12.02	14.07	18.48	24.32
8	2.73	3.49	4.59	5.53	7.34	9.52	11.03	13.36	15.51	20.09	26.12
9	3.32	4.17	5.38	6.39	8.34	10.66	12.24	14.68	16.92	21.67	27.88
10	3.94	4.86	6.18	7.27	9.34	11.78	13.44	15.99	18.31	23.21	29.59
P value (Probability)	0.95	0.90	0.80	0.70	0.50	0.30	0.20	0.10	0.05	0.01	0.001
	Non-significant							Significant			

In the table below, one can get a quick estimate for low  $N_{\text{dof}}$ .

Degrees of freedom	(df)	χ <sup>2</sup> value <sup>[16]</sup>										
1	0.004	0.02	0.06	0.15	0.46	1.07	1.64	2.71	3.84	6.64	10.83	
2	0.10	0.21	0.45	0.71	1.39	2.41	3.22	4.60	5.99	9.21	13.82	
3	0.35	0.50	_			3 66	1 61	6.25	7.82	11.34	16.27	
4 0	<u>Python:</u> chi2 prob = stats.chi2.sf(chi2 value, N <sub>DOF</sub> )							9.49	13.28	18.47		
5	sf (survival function) = 1 - CDF								11.07	15.09	20.52	
6	1.63	2.20	3.07	3.83	5.35	7.23	8.56	10.64	12.59	16.81	22.46	
7	2.17	2.83	3.82	4.67	6.35	8.38	9.80	12.02	14.07	18.48	24.32	
8	2.73	3.49	4.59	5.53	7.34	9.52	11.03	13.36	15.51	20.09	26.12	
9	3.32	4.17	5.38	6.39	8.34	10.66	12.24	14.68	16.92	21.67	27.88	
10	3.94	4.86	6.18	7.27	9.34	11.78	13.44	15.99	18.31	23.21	29.59	
P value (Probabil	ity) 0.95	0.90	0.80	0.70	0.50	0.30	0.20	0.10	0.05	0.01	0.001	
		Non-significant							Sig	Significant		

#### Chi-Square probability interpretation

The Chi-Square probability can **roughly** be interpreted as follows:

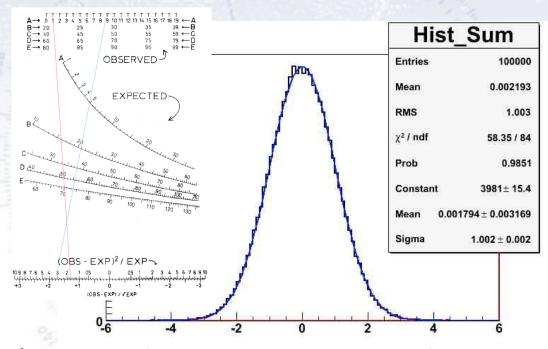
- If  $\chi^2$  / Ndof  $\approx$  1 or more precisely if 0.01 < p( $\chi^2$ ,Ndof) < 0.99, then all is good.
- If  $\chi^2$  / Ndof » 1 or more precisely if  $p(\chi^2, Ndof) < 0.01$ , then your fit is bad, and your hypothesis is probably not correct.
- If  $\chi^2$  / Ndof « 1 or more precisely if 0.99 < p( $\chi^2$ ,Ndof), then your fit is TOO good and you probably overestimated the errors.

If the statistics behind the plot is VERY high (great than 10<sup>6</sup>), then you might have a hard time finding a model, which truly describes all the features in the plot (as now tiny effects become visible), and one hardly ever gets a good Chi-Square probability. However, in this case, one should not worry too much, unless very high precision is wanted.

Anyway, the Chi-Square still allows you to compare several models, and determine which one is the better.

## Chi-Square for binned data

If the data is binned (i.e. put into a histogram), then Pearson's Chi-square applies:



The formula (based on Poisson statistics) is:

$$\chi^2 = \sum_{i \in \text{bin}} \frac{(O_i - E_i)^2}{E_i}$$



## Chi-Square for binned data

While Pearson's Chi-square test is quite useful, it has some limitations, especially when some bins have low statistics.

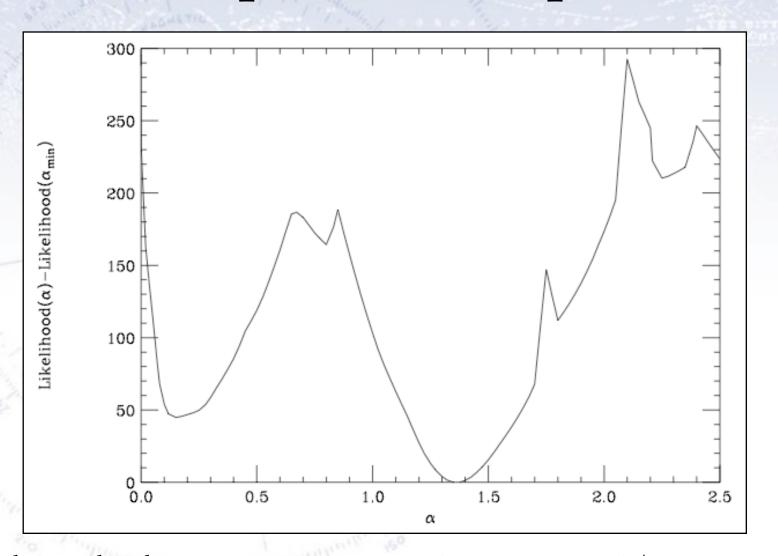
The expected cell count  $(E_i)$  should not be too low. Some require 5 or more, and others require 10 or more. A common rule is 5 or more in 80% of bins, but no cells with zero expected count. When this assumption is not met, Yates's Correction can be applied.

One alternative is to divide by O<sub>i</sub> when O<sub>i</sub> is not 0.

Another alternative is the likelihood fit, which does not suffer under low statistics.

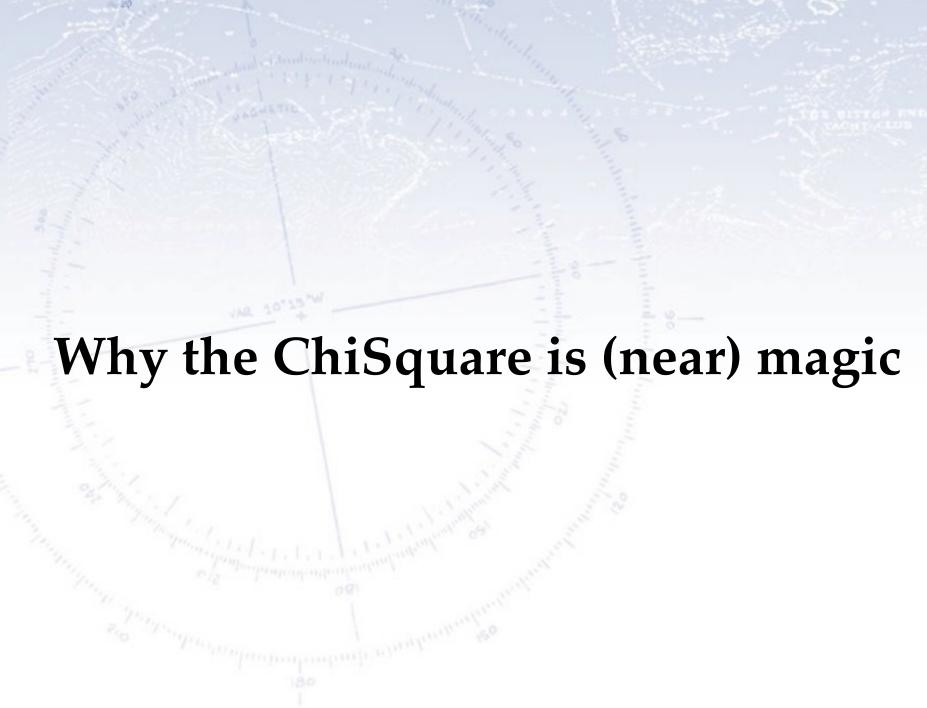
$$\chi^2 = \sum_{i \in \text{bin}} \frac{(O_i - E_i)^2}{E_i}$$

Yet, another alternative is the G-test, which is more robust at low statistics. However, I've never  $G=2\sum_{i\in D}O_i \;\ln(O_i/E_i)$  seen it in use.

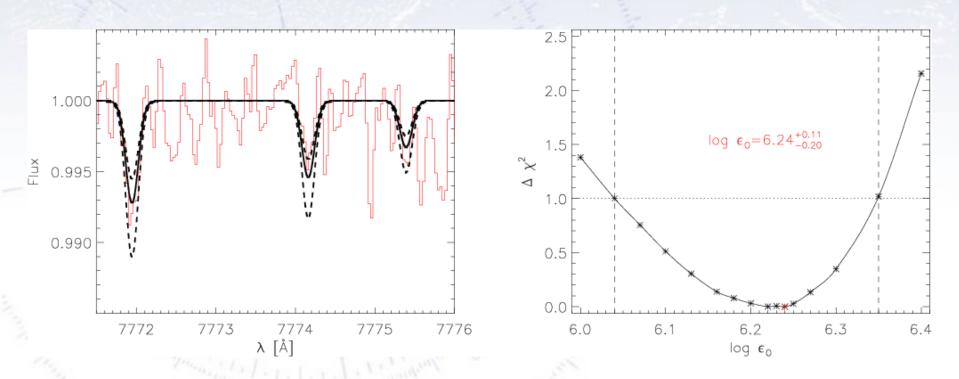


The fact that there are several minima makes fitting difficult/uncertain!

\*Always give good starting values!!!

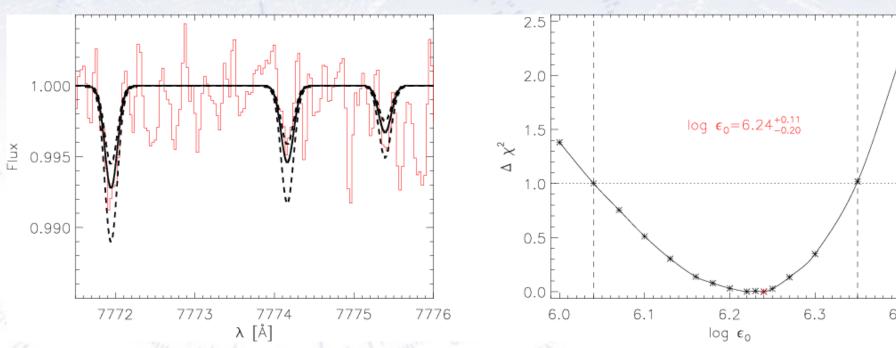


Uncertainties need not always be symmetric (though that is usually better!)



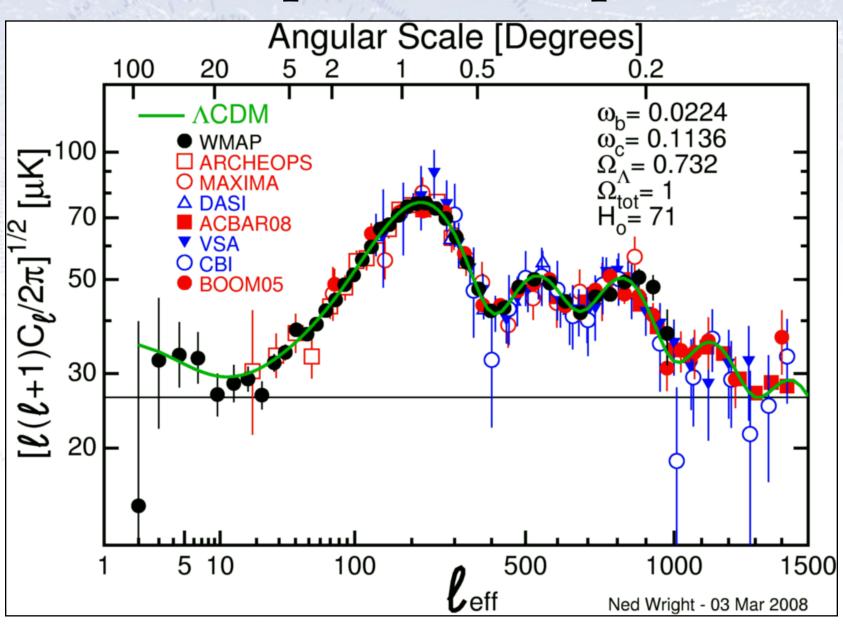
The uncertainty on a parameter is found where the Chi2 has increased by 1 from the minimum.

Uncertainties need not always be symmetric (though that is usually better!)



#### Please commit to memory!

The uncertainty on a parameter is found where the Chi2 has increased by 1 from the minimum.



# Notes on the ChiSquare method

"It was formerly the custom, and is still so in works on the theory of observations, to derive the method of least squares from certain theoretical considerations, the assumed normality of the errors of the observations being one such.

It is however, more than doubtful whether the conditions for the theoretical validity of the method are realised in statistical practice, and the student would do well to regard the method as recommended chiefly by its comparative simplicity and by the fact that it has **stood the test of experience**".

[G.U. Yule and M.G. Kendall 1958]