

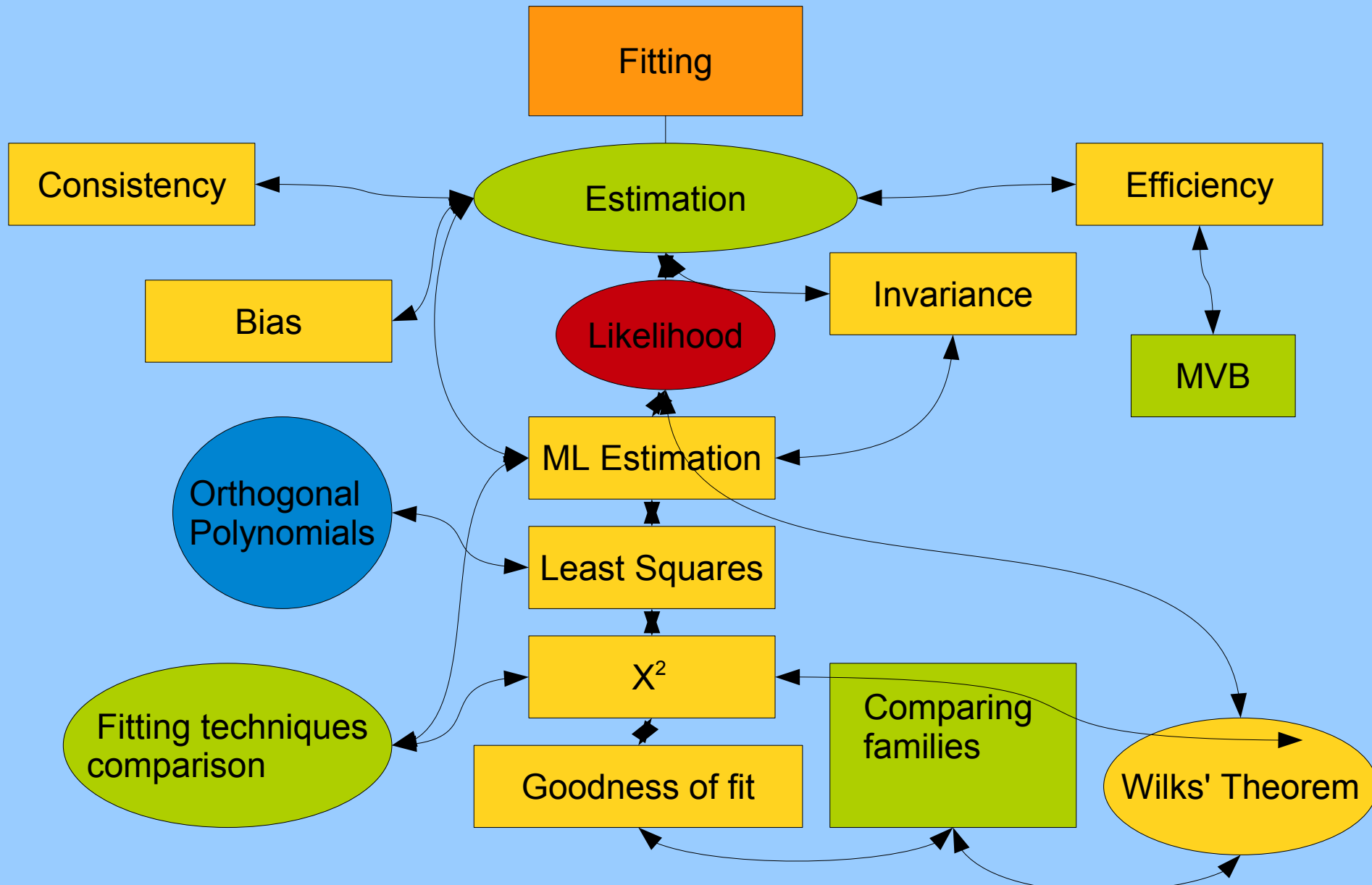
# Statistics (2)

## Fitting

*Roger Barlow*  
*Manchester University*

IDPASC school  
Sesimbra  
14<sup>th</sup> December 2010

# Summary



# Fitting and Estimation

Data sample  $\{x_1, x_2, x_3, \dots\}$  confronts theory – pdf  $P(x; a)$

( $a$  may be multidimensional)

Estimator  $\hat{a}(x_1, x_2, x_3, \dots)$  is a process returning a value for  $a$ .

A 'good' estimator is

- Consistent
- Unbiased
- Invariant
- Efficient

*These qualities depend on the estimator and on the particular pdf*

Explanations follow. Introduce (again) the Likelihood

$$L(x_1, x_2, x_3, \dots; a) = P(x_1; a) P(x_2; a) P(x_3; a) \dots$$

## Introduce the Expectation value

$$\langle f \rangle = \iiint \dots f(x_1, x_2, x_3, \dots) L(x_1, x_2, x_3, \dots, a) dx_1 dx_2 dx_3, \dots$$

Integrating over the space of results but not over  $a$ .

It is the average you would get from a large number of samples. Analogous to Quantum Mechanics.

Consistency requires:  $\lim_{N \rightarrow \infty} \langle \hat{a} \rangle = a$

i.e. given more and more data, the estimator will tend to the right answer

This is normally quite easy to establish

Require  $\langle \hat{a} \rangle = a$  (even for finite sample sizes)

If a bias is known, it can be corrected for

Standard example: estimate mean and variance of pdf from data sample

$$\hat{\mu} = \frac{1}{N} \sum x_i \qquad \hat{V} = \frac{1}{N} \sum (x_i - \hat{\mu})^2$$

This tends to underestimate V. Correct by factor  $N/(N-1)$

Desirable to have a procedure which is transparent to the form of  $a$ , i.e. need not worry about the difference between  $\hat{a}^2$  and  $\widehat{a^2}$

This is incompatible with unbiasedness.  
The well known formula (previous slide) is unbiased for  $V$  but biased for  $\sigma$

Minimise  $\langle (\hat{a} - a)^2 \rangle$

The spread of results of your estimator about the true value

Remarkable fact: there is a limit on this (Minimum Variance Bound, or Cramer-Rao bound)

$$V(\hat{a}) \geq \frac{-1}{\left\langle \frac{d^2 \ln L}{da^2} \right\rangle}$$

## Repeated Gaussian measurements

**Bias** 
$$\hat{\mu} = \frac{1}{N} \sum x_i$$

$$\iiint dx_1 dx_2 dx_3 \left( \frac{(x_1 - \mu)}{N} + \dots \right) \frac{e^{-(x_1 - \mu)^2 / 2\sigma^2}}{\sigma \sqrt{2\pi}} \dots = 0$$

**Variance** 
$$\iiint dx_1 dx_2 dx_3 \left( \frac{(x_1 - \mu)^2}{N^2} + \dots \right) \frac{e^{-(x_1 - \mu)^2 / 2\sigma^2}}{\sigma \sqrt{2\pi}} \dots = \frac{\sigma^2}{N}$$

**MVB** 
$$\ln L = \sum \frac{-(x_i - \mu)^2}{2\sigma^2} - N \ln(\sigma \sqrt{2\pi}); \quad \frac{d^2 \ln L}{d\mu^2} = \frac{-N}{\sigma^2}$$



# More examples

Centre of a top hat function:  $\frac{1}{2}(\max + \min)$

$$\sigma^2 = \frac{W}{2(N+1)(N+2)}$$

More efficient than the mean.

Several Gaussian measurements with different  $\sigma$ : weight each measurement by  $(1/\sigma)^2$ . - normalised

But don't weight Poisson measurements by their value.

# Maximum Likelihood

Estimate  $a$  by choosing the value which maximises  $L(x_1, x_2, x_3, \dots, a)$ . Or, for convenience,  $\ln L = \sum \ln P(x_i, a)$

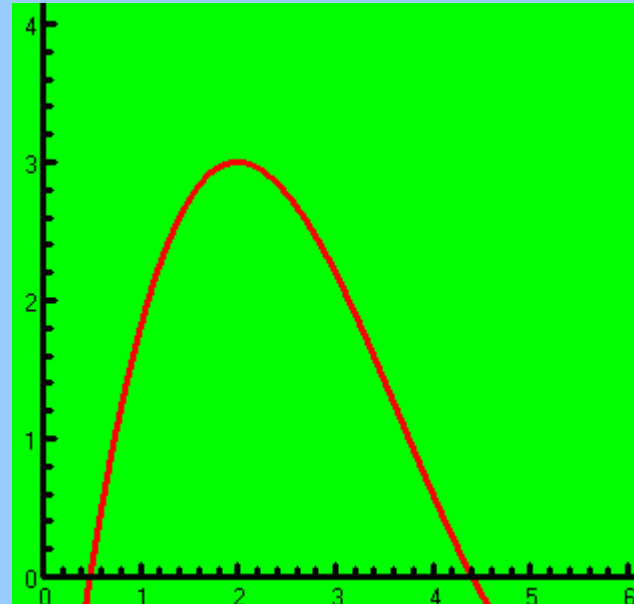
Consistency	Yes
Bias-free	No
Invariance	Yes
Efficiency	Yes, in large N limit

This is a technique, but not the only one.

Use by algebra in simple cases or numerically in tougher ones

Adjust  $a$  to maximise  
 $\ln L$

If you have a form  
for  $(d \ln L / da)$  that  
helps a lot.



Use MINUIT or ROOT or....., especially if  $a$  is  
multidimensional

Maximising: requires  $\sum d \ln P(x_i, a) / da = 0$

This leads to fractions with no nice solution  
 – unless  $P$  is exponential.

Given set of  $x_i$ , measured  $y_i$ , predictions  $f(x_i)$

subject to Gaussian smearing – Max  
 likelihood mean minimising

$$\chi^2 = \frac{\sum (y_i - f(x_i; a))^2}{\sigma_i^2}$$

Classic example: straight line fit  $f(x) = mx + c$

$$m = \frac{\overline{xy} - \bar{x}\bar{y}}{\overline{x^2} - \bar{x}^2}; \quad c = \bar{y} - m\bar{x}$$

# The Normal Equations

If  $f$  is linear function of  $a_1, a_2, a_3 \dots a_M$

$$-f_i = f(x_i) = \sum_j a_j g_j(x_i)$$

Maximum Likelihood = Minimum  $\chi^2$

$$\sum_i 2(y_i - \sum_j a_j g_j(x_i)) g_k(x_i) = 0$$

$$\sum_i y_i g_k(x_i) = \sum_j a_j \sum_i g_j(x_i) g_k(x_i)$$

Solve for the coefficients  $a_j$  by inverting  
matrix

# Orthogonal Polynomials

Good trick: construct the  $g(x)$  functions so that the matrix is diagonal

If fitting polynomial up to 5<sup>th</sup> power (say), can use  $1, x, x^2, x^3, x^4, x^5$  or  $1, x, 2x^2-1, 4x^3-3x, 8x^4-8x^2+1, 16x^5-20x^3+5x$ , or whatever

Choose  $g_0 = 1$

Choose  $g_1 = x - (\sum x)/N$  so that makes  $\sum g_0 g_1 = 0$

And so on iteratively  $g_r(x) = x^r + \sum_{rs} c_{rs} g_s(x)$

$$c_{rs} = -\sum_i x_i^r g_s(x_i) / \sum_i g_s^2(x_i)$$

These polynomials are orthogonal over a specific dataset

# Fitting histograms

Raw data  $\{x_1, x_2, x_3, \dots, x_N\}$

Often sorted into bins  $\{n_1, n_2, n_3, \dots, n_m\}$

Number of entries in bin is Poisson

$$\chi^2 = \sum \frac{(n_i - f(x_i; a))^2}{\sigma_i^2} \rightarrow \sum \frac{(n_i - f(x_i; a))^2}{f(x_i; a)} \rightarrow \sum \frac{(n_i - f(x_i; a))^2}{n_i}$$

Last form sometimes used as a definition for  $\chi^2$ , though really only an approximation

Fit function to histogram by minimising  $\chi^2$ .

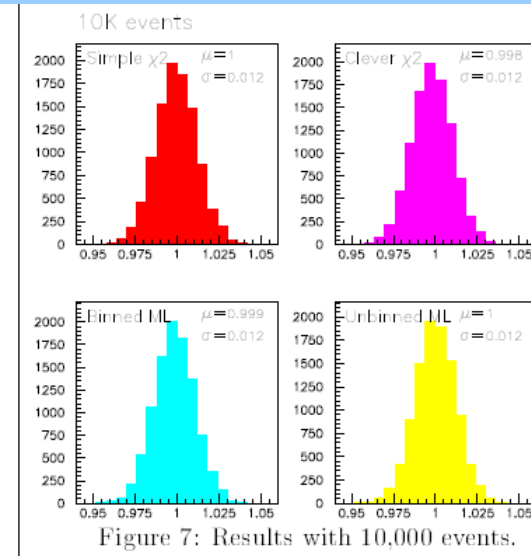
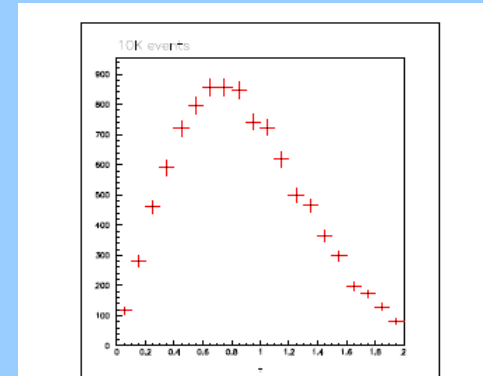
# 4 Techniques

- 1) Minimise naïve  $\chi^2$ . Computationally easy as problem linear
- 2) Minimise full  $\chi^2$ . Slower as problem nonlinear due to terms in the denominator
- 3) Binned Maximum Likelihood. Write the Poisson probability for each bin  $e^{-f_i} f_i^{n_i}/n_i!$  and maximise the sum of logs
- 4) Full maximum likelihood without binning



Fit  $f(x) = \frac{1}{2a} x e^{-ax^2}$

Try (many times)  
 with 10,000 events  
 All methods give same  
 results



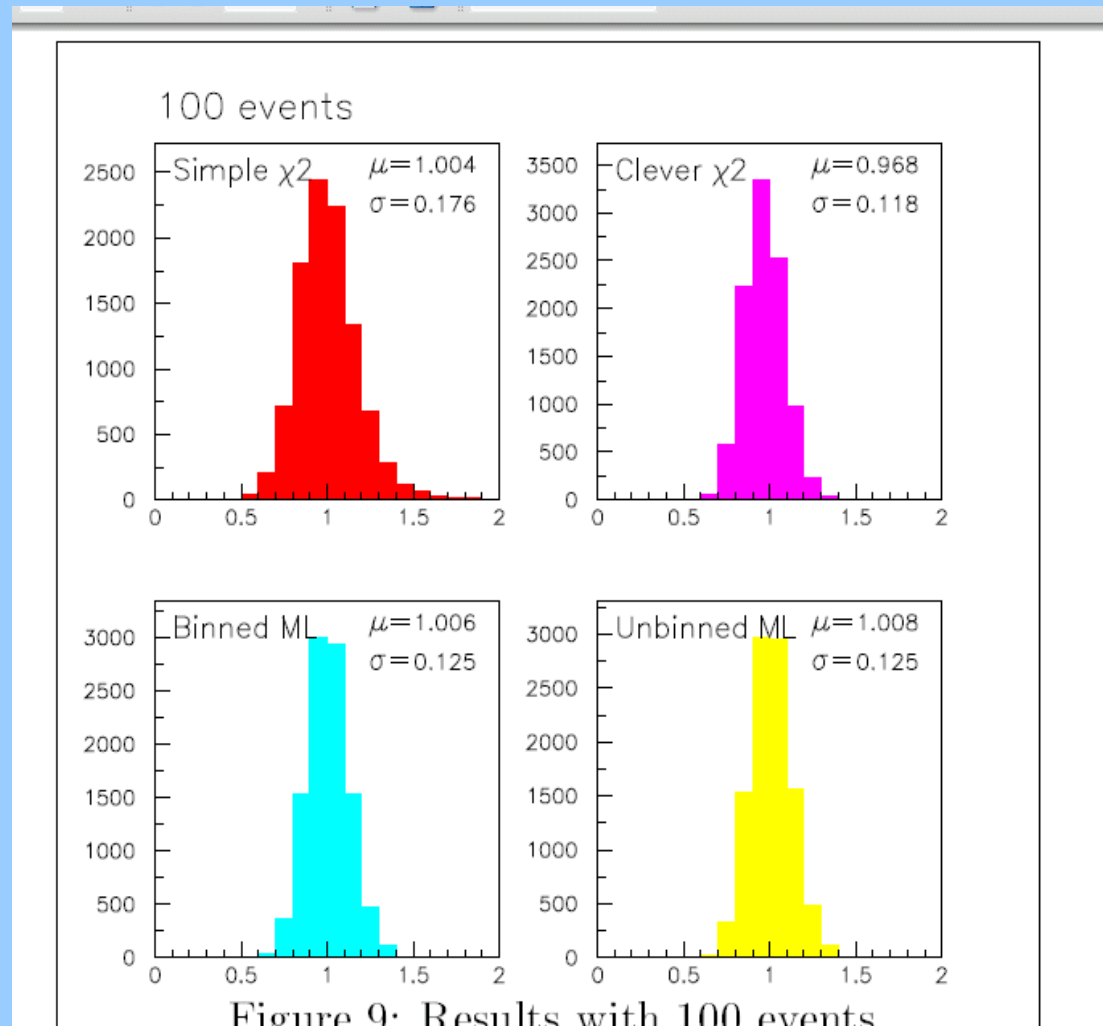
# Histogram fitting (contd)

With small sample  
(100 events)

Simple  $\chi^2$  goes  
bad due to bins  
with zeros

Full  $\chi^2$  not good as  
Poisson is not  
Gaussian

Two ML methods  
OK



Each term is clearly of order 1.

$$\chi^2 = \sum \left( \frac{y_i - f(x_i; a)}{\sigma_i} \right)^2$$

Full treatment by integrating multi-d gaussian gives  $\chi^2$  distribution  $P(\chi^2, N)$

$$\int_{\chi^2}^{\infty} P(\chi'^2; N) d\chi'^2$$

Mean indeed  $N$ . Shapes vary

Is a p value.  
Often called “ $\chi^2$  probability”

If the fit is bad,  $\chi^2$  is large

Large  $\chi^2 \gg N$ , low p value means:

- The theory is wrong
- The data are wrong
- The errors are wrong
- You are unlucky

If you histogram the p values from many cases (e.g. kinematic fits) the distribution should be flat.

This is obvious if you think about it in the right way

Small  $\chi^2 \ll N$ , p value  $\sim 1$  means:

- The errors are wrong
- You are lucky

Exact  $\chi^2 = N$  means the errors have been calculated from this test, and it says nothing about goodness of fit

## Nice extra feature

If one (or more) of the parameters in the function have been fitted to the data, this improves the  $\chi^2$  by an amount which corresponds to 1 less data point

Hence 'degrees of freedom'  $N_D = N - N_P$

No test available, sorry

Take a 'Toy Monte Carlo' which simulates your data many times, fit and find the likelihood.

Use this distribution to obtain a p value for your likelihood

This is not in most books as it is computationally inefficient. But who cares these days?

Often stated that  $\Delta \ln L = -2 \chi^2$

This strictly relates to changes in likelihood caused by an extra term in model. Valid for relative comparisons within families

E.g. Fit data to straight line.  $\chi^2$  sort of OK

Fit using parabola.  $\chi^2$  improves. If this improvement is  $\gg 1$  the parabola is doing a better job. If only  $\sim 1$  there is no reason to use it



Wilks' theorem lets you compare the merit of adding a further term to your parametrisation: yardstick for whether improved likelihood is significant. Does not report absolute merit as  $\chi^2$  does

Caution! Not to be used for adding bumps at arbitrary positions in the data.

# Summary

