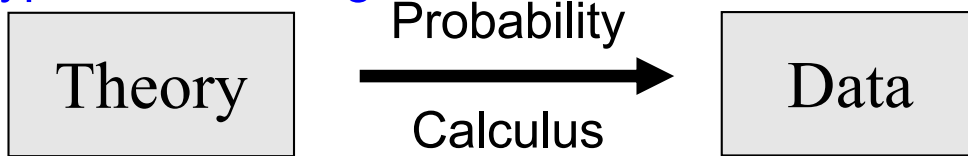


## Parameter estimation: general concepts

Hypothesis testing:



Given predictions, what can one say about data?

Given data, what can one say about parameters or properties as well as about correctness of predictions?

Parameter estimation:



**estimator** = procedure giving a value for a parameter or a property of distribution (pdf) from actual data values

notation: estimator for  $\theta$  is  $\hat{\theta}$  (a **hat** indicates estimator)

**estimate** = observed value of an estimator (often  $\hat{\theta}_{\text{obs}}$ )

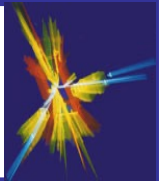
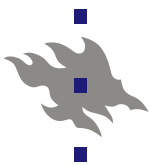
How does one construct an estimator  $\hat{\theta}(\bar{x})$ ?

**Exists no golden rule how to construct an estimator!**

Examples of estimators are arithmetic mean & variance:

$$\hat{\mu}(\{x\}) = \frac{1}{N} \sum_i x_i \quad \hat{V}(\{x\}) = \frac{1}{N-1} \sum_i (x_i - \hat{\mu})^2$$

N.B.  $\hat{\theta}(\bar{x})$  function of random variables & random variable itself, characterized by a pdf  $g(\hat{\theta}; \theta, n)$ , which depends on (true value of)  $\theta$  & has expectation value, variance, etc...



Often start by requiring **consistency**:  $\lim_{n \rightarrow \infty} \hat{\theta} = \theta$  i.e. as sample size increases, estimate converges to true value:

$$\text{for any } \varepsilon > 0, \quad \lim_{n \rightarrow \infty} P(|\hat{\theta} - \theta| > \varepsilon) = 0$$

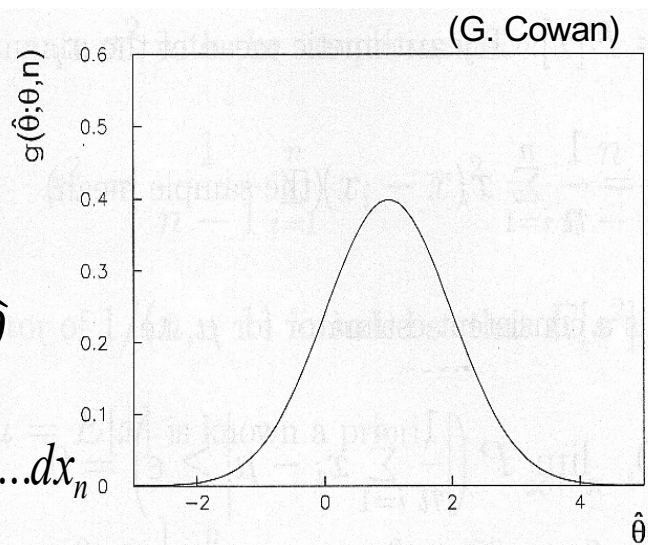
NB! convergence in sense of probability, i.e. no guaranty that any particular  $\hat{\theta}_{\text{obs}}$  will be within given distance of  $\theta$ .

$g(\hat{\theta}; \theta, n)$  is the pdf of  $\hat{\theta}$  for a fixed sample size  $n$ .

Expectation value of  $\hat{\theta}$ :

$$E[\hat{\theta}] = \int \hat{\theta} g(\hat{\theta}; \theta, n) d\hat{\theta}$$

$$\int \dots \int \hat{\theta}(\bar{x}) f(x_1; \theta) \dots f(x_n; \theta) dx_1 \dots dx_n$$



**variance**  $V[\hat{\theta}] = \sigma_{\hat{\theta}}^2$

$\sigma_{\hat{\theta}}$  = "statistical" uncertainty  
 $b$  = "systematic" uncertainty  
(due to construction of  $\hat{\theta}$ )

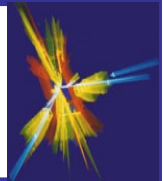
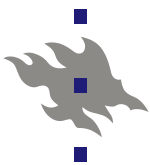
**bias**  $b = E[\hat{\theta}] - \theta$

For most estimators:  $\sigma_{\hat{\theta}} \propto 1/\sqrt{n}$ ,  $b \propto 1/n$

Good estimator: **consistent**, **unbiased** ( $E[\hat{\theta}] = \theta$ ) and **efficient** (i.e. has minimal possible variance = "RCF bound").

**RCF bound:**  $V[\hat{\theta}] \geq \left(1 + \frac{\partial b}{\partial \theta}\right)^2 / E\left[-\frac{\partial^2 \ln L}{\partial \theta^2}\right]$  ( $b$  = bias)

Also "CR (Cramér-Rao) bound" or "information inequality".  
 $L$  = likelihood function (defined on slide 5),  $F$  = Fréchet.



Consider  $n$  measurements of random variable  $x, x_1, \dots, x_n$ .  
Arithmetic mean a natural choice as estimator of  $\mu = E[x]$ :

$$\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (= \text{the sample mean})$$

If  $V[x]$  finite, then  $\bar{x}$  is a consistent estimator for  $\mu$ , i.e.

$$\text{for any } \varepsilon > 0, \quad \lim_{n \rightarrow \infty} P\left(\left|\frac{1}{n} \sum_{i=1}^n x_i - \mu\right| > \varepsilon\right) = 0$$

i.e. the **weak law of large numbers**. Expectation value of  $\bar{x}$ :

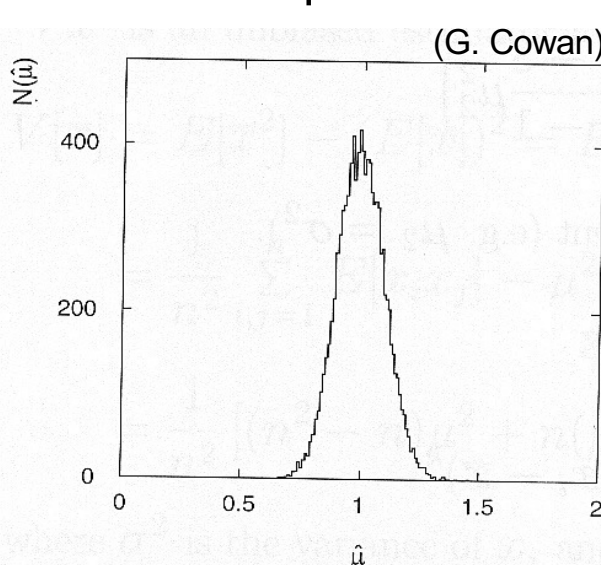
$$E[\bar{x}] = E\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n} \sum_{i=1}^n E[x_i] = \frac{1}{n} \sum_{i=1}^n \mu = \mu$$

→  $\bar{x}$  is an unbiased estimator for  $\mu$ . The variance of  $\bar{x}$  is

$$V[\bar{x}] = E[\bar{x}^2] - (E[\bar{x}])^2 = \frac{1}{n^2} \sum_{i,j=1}^n E[x_i x_j] - \mu^2 = \frac{\sigma^2}{n},$$

$\sigma^2 = \text{variance of } x, E[x_i x_j] = \mu^2 \text{ for } i \neq j \text{ and } E[x_i^2] = \mu^2 + \sigma^2$

Example of estimator for mean: take samples of  $n = 100$  values from a Gaussian MC generator with  $\mu = 1$  &  $\sigma^2 = 1$ . Calculate sample mean & repeat procedure many times.



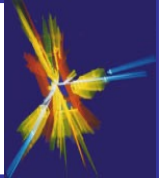
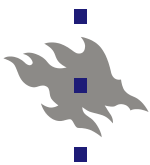
Enter values into a histogram.

$$\bar{\hat{\mu}} = 0.9981 \quad (\hat{\mu} \text{ unbiased})$$

Sample standard deviation of

$$\hat{\mu} \text{ values} = 0.0995 \approx \frac{\sigma}{\sqrt{n}}$$

NB! pdf of  $\hat{\mu} \approx \text{Gaussian}$   
(result of central limit theorem).



Suppose mean  $\mu$  and variance  $V[x] = \sigma^2$  both unknown.  
Then estimate  $\sigma^2$  using sample variance

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{n}{n-1} (\overline{x^2} - \bar{x}^2)$$

Factor  $1/(n-1)$  introduced to have  $E[s^2] = \sigma^2$  (unbiased).  
If mean  $\mu = E[x]$  known then estimate  $\sigma^2$  using statistic  $S^2$

$$S^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 = \overline{x^2} - \mu^2 \quad \text{also } E[S^2] = \sigma^2 \text{ (i.e. unbiased estimator).}$$

Variance of  $s^2(S^2)$  calculated with  $k^{\text{th}}$  central moments  $\mu_k$ .  
 $\mu_k$ 's estimated from corresponding estimator  $m_k$  or  $M_k$ .

$$V[s^2] = \frac{1}{n} \left( \mu_4 - \frac{n-3}{n-1} \mu_2^2 \right) \quad V[S^2] = \frac{1}{n} (\mu_4 - \mu_2^2)$$

$$\hat{\mu}_k = m_k = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^k \quad \hat{\mu}_k = M_k = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^k$$

A natural estimator for standard deviation,  $\sigma$ , then

$$\hat{\sigma} = s = \sqrt{s^2} \quad \left( \text{or in case } \mu \text{ known } \hat{\sigma} = S = \sqrt{S^2} \right)$$

For variance of estimator for standard deviation:

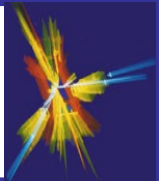
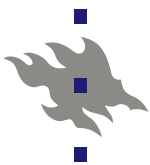
$$V[\hat{s}^2] = (d\sigma^2/d\sigma)^2 V[\hat{\sigma}] = 4\sigma^2 V[\hat{\sigma}] \Rightarrow V[\hat{\sigma}] = V[\hat{s}^2]/4\sigma^2$$

For a Gaussian pdf  $\hat{s}^2(\sigma) = \sigma^2 \Rightarrow d\hat{s}^2/d\sigma = 2\sigma$

$$E[\hat{\mu}_4] = 3\sigma^4 \Rightarrow V[s^2] = 2\sigma^4/(n-1) \Rightarrow \sigma_s = \sigma / \sqrt{2(n-1)}$$

Another quality measure of estimator: **mean square error**

$$MSE = E[(\hat{\theta} - \theta)^2] = V[\hat{\theta}] + b^2 \quad \text{(used as measure in e.g. unfolding methods)}$$



Random variable  $x$  distributed according to pdf  $f(x, \theta)$ .  
Assume functional form of  $f$  known but not parameter  $\theta$ .  
**maximum likelihood method** (suurimman uskottavuuden menetelmä) technique for estimating  $\theta$  from a data sample.

If  $f(x, \theta)$  correct pdf hypothesis, then

$$P(x_i \text{ found in } [x_i, x_i + dx_i] \text{ for all } i) = \prod_{i=1}^n f(x_i, \theta) dx_i$$

If hypothesis (including value of  $\theta$ ) correct (= true)

→ expect higher probability for the data

If hypothesized functional form wrong or  $\theta$  value far away

→ expect lower probability for the data

⇒ higher value of **likelihood function** close to true  $\theta$

$$L(\theta) = \prod_{i=1}^n f(x_i, \theta)$$

NB!  $L(\theta) = f_{\text{sample}}(\bar{x}; \theta)$ , but  $L(\theta)$  regarded only a function of  $\theta$ , measurements  $x_i$ 's constants, "experiment" finished.

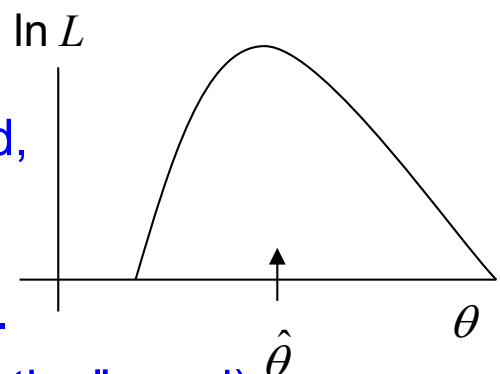
Define ML estimator  $\hat{\theta}$  as value of  $\theta$  that maximizes  $L(\theta)$ .

For  $m$  parameters, usually find solution  $\hat{\theta}_1, \dots, \hat{\theta}_m$  by solving

$$\frac{\partial L}{\partial \theta_i} = 0, \quad i = 1, \dots, m$$

In practice maximize  $\ln L(\theta)$  instead,  
Can then add individual  $\ln P(x_i)$ 's.

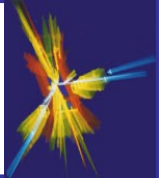
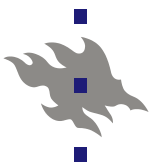
In  $L(\theta)$  might have more than one local maximum → take highest one.



N.B. No binning of data ("all information" used).

N.B. Definition of ML estimators don't guarantee optimality  
→ investigate properties such as bias, variance ...

In most cases, especially for sufficient large data samples,  
ML estimators generally most optimal estimator choice.



Suppose proper decay times of a certain type of unstable states measured for  $n$  decays,  $t_1, \dots, t_n$ . Choose as hypothesis for  $t$  distribution an exponential pdf with mean  $\tau$ .

$$f(t; \tau) = e^{-t/\tau} / \tau$$

Task to estimate value of  $\tau$ . Use **log-likelihood function** instead to find parameter value giving maximum value for function. Equivalent since logarithm a monotonic function ( $\rightarrow$  maximum at same value). In addition, products in  $L$  becomes sums in  $\ln L$  and exponentials becomes factors.

$$\ln L(\tau) = \sum_{i=1}^n \ln f(t_i; \tau) = \sum_{i=1}^n (-\ln \tau - t_i / \tau)$$

$$\text{set } \frac{\partial \ln L}{\partial \tau} = 0 \text{ and solve for } \tau \rightarrow \hat{\tau} = \frac{1}{n} \sum_{i=1}^n t_i$$

How find out whether  $\hat{\tau}$  is an unbiased estimator for  $\tau$ ?

i) Find pdf  $g(\hat{\tau}; \tau, n)$  and compute  $b = E[\hat{\tau}] - \tau$

ii) Compute  $E[\hat{\tau}(t_1, \dots, t_n)] = \int \dots \int \hat{\tau}(\bar{t}) f_{\text{joint}}(\bar{t}; \tau) dt_1 \dots dt_n =$

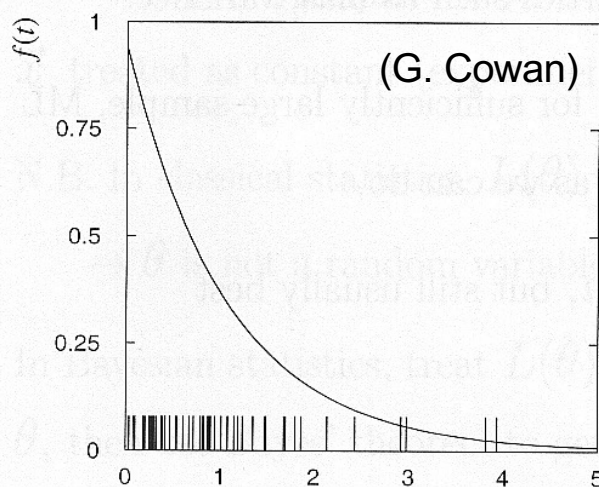
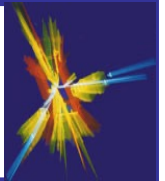
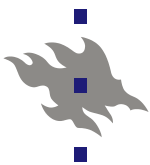
$$\int \dots \int \left( \sum_{i=1}^n t_i \right) \frac{e^{-t_1/\tau}}{\tau} \dots \frac{e^{-t_n/\tau}}{\tau} \frac{dt_1 \dots dt_n}{n} = \frac{1}{n} \sum_{i=1}^n \left( \int \frac{t_i}{\tau} e^{-t_i/\tau} dt_i \prod_{j \neq i} \frac{e^{-t_j/\tau}}{\tau} dt_j \right)$$

$$= \sum_{i=1}^n \tau / n = \tau \rightarrow \hat{\tau} \text{ unbiased estimator for } \tau !!$$

iii) Could make same conclusion without any calculation based on the fact that the sample mean an unbiased estimator for  $E[t]$  and for exponential pdf  $E[t] = \tau$ .

Suppose that one is interested in decay constant  $\lambda = \ln 2 / \tau$  instead of mean lifetime  $\tau$ . ML estimator for  $\lambda$ ?





A sample of 50 observations of proper time,  $t$ , ("ticks" on x-axis), generated using MC assuming exponential distribution with mean  $\tau = 1.0$ . Curve result of a maximum likelihood fit to observations, giving  $\hat{\tau} = 1.062$ .

Given a function  $a(\theta)$  of some parameter  $\theta$ , one has

$$\frac{\partial L}{\partial \theta} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial \theta} = 0 \quad \Rightarrow \quad \frac{\partial L}{\partial a} = 0 \bigg|_{a=a(\theta)} \quad \text{unless} \quad \frac{\partial a}{\partial \theta} = 0$$

So  $a$  maximizing  $L_a(a)$  is  $a(\hat{\theta})$ , where  $\hat{\theta}$  maximizes  $L_{\theta}(\theta)$ .

→ ML estimator of function  $a(\theta)$  is  $\hat{a} = a(\hat{\theta})$

This is called invariance. ML estimators are invariant.

So for decay constant, one gets  $\hat{\lambda} = \frac{\ln 2}{\hat{\tau}} = \ln 2 \cdot n / \sum_{i=1}^n t_i$

Is  $\hat{\lambda}$  an unbiased estimator for  $\lambda$ ?

For  $\hat{\lambda}$  one can show that  $E[\hat{\lambda}] = \frac{n\lambda}{n-1} = \frac{\ln 2}{\tau} \frac{n}{n-1}$

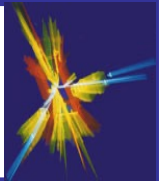
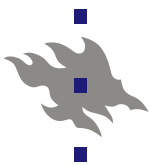
→  $\hat{\lambda}$  has a bias that goes to zero for  $n \rightarrow \infty$ .

Above true for ML estimators ( $b \rightarrow 0$ , when  $n \rightarrow \infty$ ).

Example where ML fails: assume taxis numbered 1 to  $N_{\text{taxi}}$ , ML estimator for  $N_{\text{taxi}}$  from  $m$  taxi number observations?

$$f(n) = \frac{1}{(N_{\text{taxi}} - 1)} \Rightarrow L = \frac{1}{(N_{\text{taxi}} - 1)^m} \Rightarrow \frac{\partial \ln L}{\partial N_{\text{taxi}}} \quad \text{no local maxima} \\ \text{so } \hat{N}_{\text{taxi}} \text{ undefined.}$$

Ansatz:  $\hat{N}_{\text{taxi}} = 2\bar{n} - 1$ , where  $\bar{n} = \text{mean}$ ;  $E[\hat{N}_{\text{taxi}}] = N_{\text{taxi}}$ .



$n$  measurement of a variable  $x$  assumed to be Gaussian distributed with unknown  $\mu$  &  $\sigma$ . Log-likelihood function:

$$\ln L(\mu, \sigma) = \sum_{i=1}^n \ln f(x_i; \mu, \sigma) = \sum_{i=1}^n \left( -\frac{1}{2} (\ln 2\pi + \ln \sigma^2) - \frac{(x_i - \mu)^2}{2\sigma^2} \right)$$

set  $\frac{\partial \ln L}{\partial \mu} = 0$  &  $\frac{\partial \ln L}{\partial \sigma^2} = 0$  and solve for  $\mu$  &  $\sigma^2 \rightarrow$

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

Already known that  $\hat{\mu}$  is an unbiased estimator for  $\mu$ .

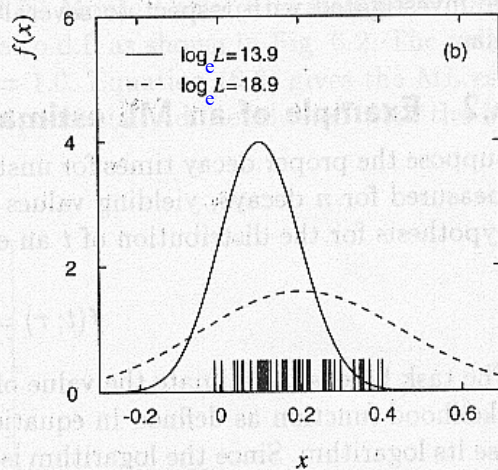
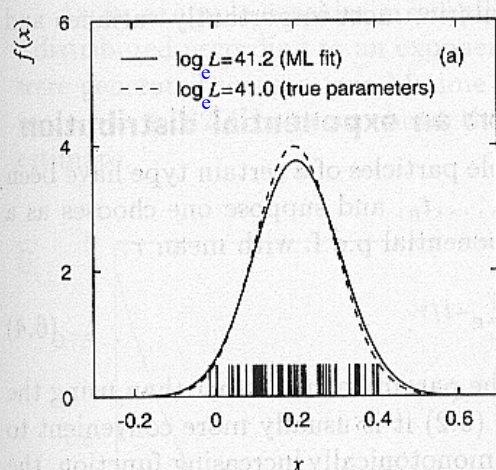
What about  $\hat{\sigma}^2$ ?  $E[\hat{\sigma}^2] = \frac{n-1}{n} \sigma^2$

So ML estimator for  $\sigma^2$  has a bias, but  $b \rightarrow 0$  for  $n \rightarrow \infty$ .

Recall, however, that the sample variance

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu})^2 = \frac{n}{n-1} (\overline{x^2} - \hat{\mu}^2)$$

is an unbiased estimator for the variance of any pdf.



(G. Cowan)

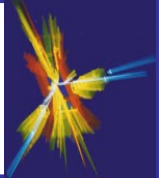
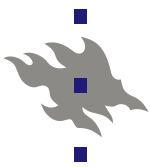
50 observations of gaussian variable  $x$ ;  $\mu_x = 0.2$  &  $\sigma_x = 0.1$ .

(a) pdf of parameters maximizing  $\ln L$  & true parameters.

(b) pdf of parameters far from true ones  $\rightarrow$  low  $\ln L$  values.

What about statistical uncertainty of ML estimates?





## Variance of ML estimators: the analytic method

A direct way of estimating uncertainty of estimate is to

compute variance of estimator, e.g. when  $\hat{\tau} = \frac{1}{n} \sum_{i=1}^n t_i$   
i.e. width of the pdf  $g(\hat{\tau}; \tau, n)$ :

$$\begin{aligned} V[\hat{\tau}] &= E[\hat{\tau}^2] - (E[\hat{\tau}])^2 = \int \dots \int \left( \frac{1}{n} \sum_{i=1}^n t_i \right)^2 \frac{e^{-t_1/\tau}}{\tau} \dots \frac{e^{-t_n/\tau}}{\tau} dt_1 \dots dt_n - (E[\hat{\tau}])^2 \\ &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j \neq i} \left( \int (t_i^2 + t_i t_j) \frac{e^{-t_i/\tau}}{\tau} dt_i \frac{e^{-t_j/\tau}}{\tau} dt_j \prod_{k \neq i, j} \int \frac{e^{-t_k/\tau}}{\tau} dt_k \right) - (E[\hat{\tau}])^2 \\ &= \frac{\tau^2 (2n + n(n-1))}{n^2} - \tau^2 = \frac{\tau^2}{n} \rightarrow V[\hat{\tau}] \quad \begin{array}{l} n \text{ times smaller} \\ \text{than } V[t] \end{array} \end{aligned}$$

(in fact this result was obvious, since here  $\hat{\tau} = \bar{t}$ )

N.B.  $V[\hat{\tau}]$  &  $\sigma_{\hat{\tau}}$  are functions of true (& unknown!)  $\tau$ .

Estimate standard deviation using  $\hat{\sigma}_{\hat{\tau}} = \hat{\tau} / \sqrt{n}$

Estimated standard deviation often quoted as "statistical

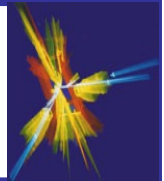
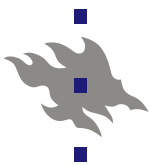
error" of a measurement e.g.  $\hat{\tau} \pm \hat{\sigma}_{\hat{\tau}} = 1.062 \pm 0.150$

should be interpreted as: ML estimate for  $\tau$  is 1.062.

ML estimate for  $\sigma$  of  $g(\hat{\tau}; \tau, n)$  is 0.150.

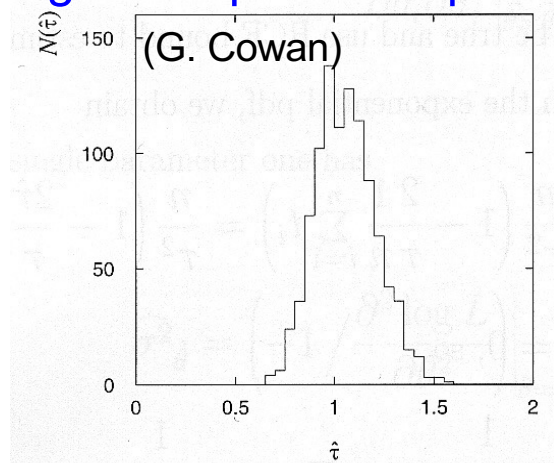
If pdf  $g(\hat{\tau}; \tau, n)$  Gaussian, then  $[\hat{\tau} - \hat{\sigma}_{\hat{\tau}}, \hat{\tau} + \hat{\sigma}_{\hat{\tau}}]$  equivalent to "68.3 % confidence interval" for  $\hat{\tau}$ , generally accepted way to quote uncertainty even when errors non-Gaussian.

NB! very seldom variance explicitly computable as above!!



## Variance of ML estimators: the Monte Carlo method

Cases too difficult to solve analytically (or  $g(\hat{\tau}; \tau, n)$  not known), ML estimator distribution investigated with MC. Simulate large number of **pseudoexperiments**, compute ML estimates each time, resulting distribution  $\approx g(\hat{\tau}; \tau, n)$ . Use experimental  $\mu$  as "true" value & MC to get distribution of sample means, width = unbiased variance estimator. E.g. for exponential pdf  $\hat{\tau} = 1.062$ , used as "true" MC value.



A histogram with ML estimates from 1000 MC experiments with 50 observations each time. MC used  $\tau = 1.062$  as true lifetime. Calculated standard deviation of histogram entries,  $s$ , is 0.151.

Similar to analytical estimate  $\hat{\tau}/\sqrt{n} = 1.062/\sqrt{50} = 0.150$

NB!  $g(\hat{\tau}; \tau, n)$  approx. Gaussian ( $\Leftrightarrow$  central limit theorem)  $\rightarrow$  true in general for ML estimators in large sample limit.

## Variance of ML estimators: the RCF bound

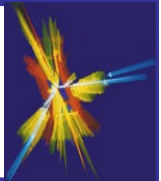
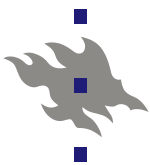
A lower bound on variance of any estimator (not just ML)

$$V[\hat{\theta}] \geq \left(1 + \frac{\partial b}{\partial \theta}\right)^2 / E\left[-\frac{\partial^2 \ln L}{\partial \theta^2}\right] \quad (b = \text{bias})$$

Rao-Cramer-Frechet bound (or "information inequality").

If equality true, then corresponding estimator **efficient**.

ML estimators efficient in large sample limit. So, assume estimator efficient & use RCF bound to estimate  $V[\hat{\theta}]$ .



For the example with exponential pdf, one obtains

$$\frac{\partial^2 \ln L}{\partial \tau^2} = \frac{n}{\tau^2} \left( 1 - \frac{2}{n\tau} \sum_{i=1}^n t_i \right) = \frac{n}{\tau^2} \left( 1 - \frac{2\hat{\tau}}{\tau} \right) \quad \& \quad b=0 \quad \text{so}$$

$$V[\hat{\tau}] \geq \left( E \left[ -\frac{n}{\tau^2} \left( 1 - \frac{2\hat{\tau}}{\tau} \right) \right] \right)^{-1} = \left( -\frac{n}{\tau^2} \left( 1 - \frac{2E[\hat{\tau}]}{\tau} \right) \right)^{-1} = \frac{\tau^2}{n}$$

same variance as obtained from analytical calculation

→ ML  $\hat{\tau}$  an efficient estimator for parameter  $\tau$  for any  $n$ .

For  $\bar{\theta} = (\theta_1, \dots, \theta_m)$  with efficient estimators and zero bias

$$(V^{-1})_{ij} = E \left[ -\frac{\partial^2 \ln L}{\partial \theta_i \partial \theta_j} \right] \Rightarrow (V^{-1})_{ij} = -\frac{\partial^2 \ln L}{\partial \theta_i \partial \theta_j} \Big|_{\bar{\theta} = \hat{\theta}}$$

Impractical to compute RCF bound analytically. In case of sufficiently large data sample, estimate  $V^{-1}$  by evaluating 2<sup>nd</sup> derivate at the ML estimates with the measured data.

Procedure: 1<sup>st</sup> numerically maximize  $\ln L$ , then determine matrix of 2<sup>nd</sup> derivates using finite differences evaluated at ML estimates, finally invert result to find covariance matrix.

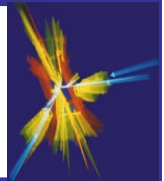
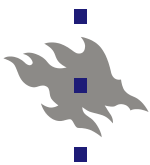
Variance of ML estimators: the graphical method

Extension of RCF bound technique leads to a graphical technique for obtaining the variance of ML estimators.

Expand  $\ln L(\theta)$  around ML estimate  $\hat{\theta}$  of parameter  $\theta$ :

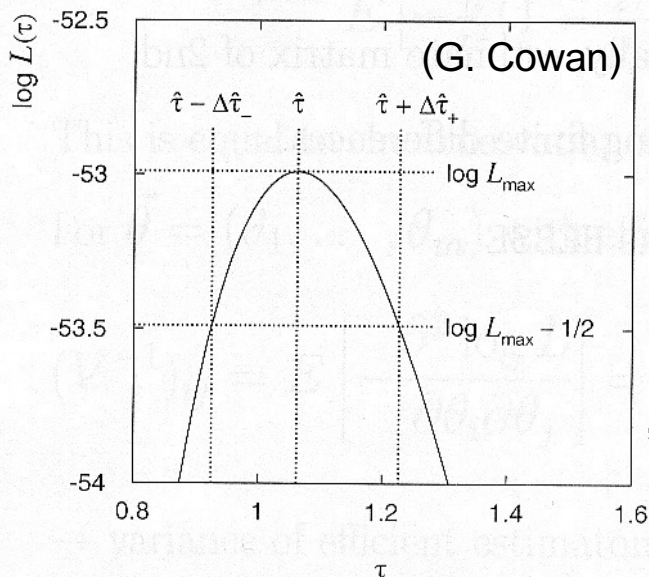
$$\ln L(\theta) = \ln L(\hat{\theta}) + \left[ \frac{\partial \ln L}{\partial \theta} \right]_{\theta=\hat{\theta}} (\theta - \hat{\theta}) + \frac{1}{2!} \left[ \frac{\partial^2 \ln L}{\partial \theta^2} \right]_{\theta=\hat{\theta}} (\theta - \hat{\theta})^2 + \dots$$

$$\text{now } \ln L(\hat{\theta}) = \ln L_{\max} \quad \& \quad \left[ \frac{\partial \ln L}{\partial \theta} \right]_{\theta=\hat{\theta}} = 0 \Rightarrow$$



$$\ln L(\theta) \approx \ln L_{\max} - \frac{(\theta - \hat{\theta})^2}{2\sigma_{\hat{\theta}}^2} \Rightarrow \ln L(\hat{\theta} \pm \sigma_{\hat{\theta}}) \approx \ln L_{\max} - \frac{1}{2}$$

1 (2) standard deviation change of  $\theta$  from its ML estimate leads to a  $\ln L(\theta)$  decrease by 0.5 (2.0) from  $\ln L_{\max}$ .



The exponential distribution example  $\hat{\tau} = 1.062$

$$\Delta \hat{\tau}_- = 0.137, \Delta \hat{\tau}_+ = 0.165$$

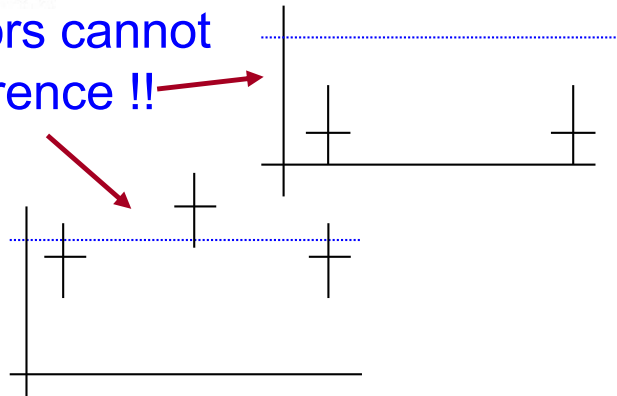
Usually set  $\hat{\tau} = 1.062_{-0.137}^{+0.165}$

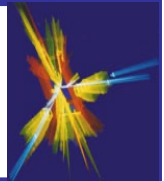
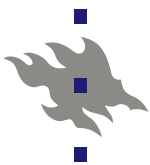
Interval  $[\hat{\tau} - \Delta \hat{\tau}_-, \hat{\tau} + \Delta \hat{\tau}_+]$  interpreted as estimate for "68.3 % confidence interval"

### Summary on ML estimators:

ML estimators cannot tell the difference !!

- consistent.
- invariant.
- biased for small  $n$ .
- not "right", just sensible.
- don't give "most likely value of  $\theta$ " but value of  $\theta$  for which the data is "most likely" (highest likelihood).
- efficient for large  $n$  (saturates the RCF bound).
- often imply the use of numerical methods. (analytical  $\ln L$  maximisation for  $>1$  free variables in practice impractical)
- the  $\ln L_{\max}$  value in itself contains no valuable information  $\Rightarrow$  won't indicate if chosen function for pdf correct or not.





Upto now, normalisation has been fixed; can also leave normalisation free e.g. treat  $n$  as Poisson random variable with mean  $\nu \rightarrow$  result of experiment:  $n$  &  $n$   $x$ -values  $x_1 \dots x_n$

**Extended likelihood function:**  $L(\nu, \bar{\theta}) = \frac{\nu^n}{n!} e^{-\nu} \prod_{i=1}^n f(x_i; \bar{\theta})$

2 separate cases: either  $\nu$  independent or a function of  $\bar{\theta}$   
theory/model gives  $\nu = \nu(\bar{\theta})$ , dropping constant terms  $\rightarrow$

$$\ln L(\bar{\theta}) = n \ln \nu - \nu + \sum_{i=1}^n \ln f(x_i) = \sum_{i=1}^n \ln(\nu(\bar{\theta}) f(x_i; \bar{\theta})) - \nu(\bar{\theta})$$

Now more information used  $\rightarrow$  smaller variances for  $\hat{\bar{\theta}}$

Example: particle scattering, expected number of events

$\nu = \varepsilon \sigma \int \mathcal{L} dt$ , where  $\varepsilon$  detection efficiency,  $\sigma$  scattering cross section ("probability" given by theory) &  $\mathcal{L}$  luminosity (= "particle flux"). Use not only event variables but also  $\sigma$ .

NB! before 'repetition of experiment' = same number of events. Here, means same  $\nu$  (e.g. same integrated  $\mathcal{L}$ ).

Suppose  $\nu$  &  $\bar{\theta}$  are (functionally) independent:

$$\frac{\partial \ln L}{\partial \nu} = 0 \rightarrow \hat{\nu} = n; \quad \frac{\partial \ln L}{\partial \theta_j} = 0 \rightarrow \text{usual ML } \hat{\bar{\theta}}$$

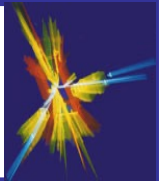
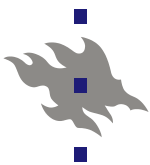
i.e. an additional fluctuation source. Useful sometimes e.g. when  $f(x; \bar{\theta})$  superposition of known components

$f(x; \bar{\theta}) = \sum_{j=1}^m \theta_j f_j(x)$ . Can use usual ML with constrain but then not all  $\theta_j$  independent &

$\theta_m = 1 - \sum_{j=1}^{m-1} \theta_j$  different  $\theta_j$  treated differently.

With extended ML avoid all of that.





Same with extended ML,

$$\ln L(\nu, \bar{\theta}) = \sum_{i=1}^n \ln \left( \sum_{j=1}^m \nu \theta_j f_j(x_i) \right) - \nu$$

define  $\mu_j = \nu \theta_j$  as expected number of events of type  $j$ ,

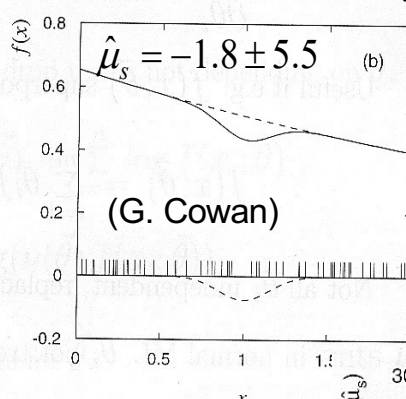
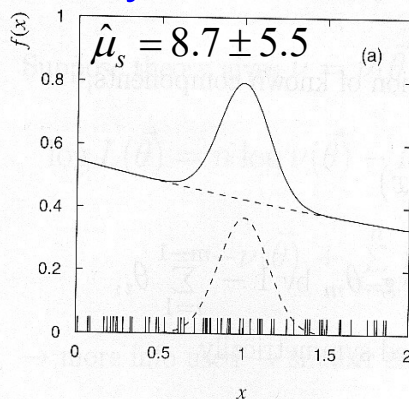
$$\ln L(\bar{\mu}) = \sum_{i=1}^n \ln \left( \sum_{j=1}^m \mu_j f_j(x_i) \right) - \sum_{j=1}^m \mu_j$$

Now all parameters treated symmetrically. Often  $\mu_j$ 's more closely related to wanted final result e.g. production cross section for type  $j$  events. NB! ML fitted  $\mu_j$ 's can be  $< 0$  & in case  $\mu_j < 0$  unphysical, must decide how to treat that case. Example: 2 types of events, signal( $s$ ) & background( $b$ ).

$$f(x) = \frac{\mu_s}{\mu_s + \mu_b} f_s(x) + \frac{\mu_b}{\mu_s + \mu_b} f_b(x)$$

Assume  $f_s(x)$  &  $f_b(x)$  known, estimate  $\mu_s$  &  $\mu_b$

mostly 'works' well ... but sometimes  $\hat{\mu}_s < 0$

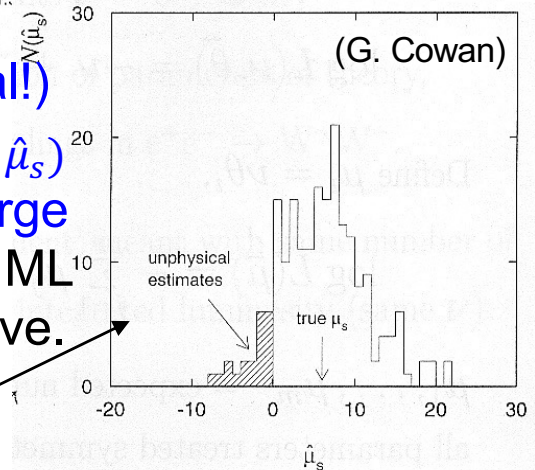


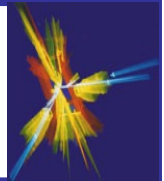
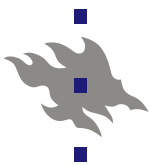
extended ML fit to 2 MC samples generated with  $\mu_s = 6$  &  $\mu_b = 60$ .  $f_s(x)$  Gaussian &  $f_b(x)$  exponential.

Can report negative  $\hat{\mu}_s$  (unphysical!) or take as estimator  $\hat{\mu}_s^{\text{phys}} = \max(0, \hat{\mu}_s)$  (biased!). Can be a problem for large samples. Example: 200 extended ML fits on similar MC samples as above.

$$\hat{\mu}_s = 6.1 \pm 0.4 \quad (\text{using all } \hat{\mu}_s)$$

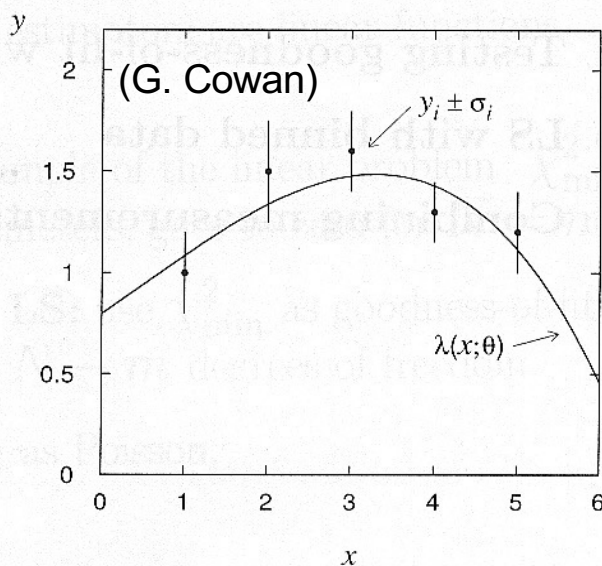
$$\hat{\mu}_s = 6.4 \pm 0.4 \quad (\text{using only } \hat{\mu}_s^{\text{phys}})$$





## The method of least squares (pienimmän neliösumman menetelmä)

$N$  independent Gaussian random variables  $y_i$ ,  $i = 1, \dots, N$  with expectation values  $E[y_i] = \lambda_i = \lambda(x_i; \bar{\theta})$  that depend on unknown parameters  $\bar{\theta}$ . Each  $y_i$  related to a known  $x_i$  & known  $V[y_i] = \sigma_i^2$  (e.g. estimated measurement uncertainty). E.g. temperature measurement  $T$  at positions  $x_i$ .



Least squares problem in a nut shell:  $N$  values  $y_1, \dots, y_N$  measured with uncertainties  $\sigma_1, \dots, \sigma_N$  at  $x$ -values  $x_1, \dots, x_N$  (known without uncertainties). Each value  $\lambda_i$  of  $y_i$  given by function  $\lambda(x_i, \theta)$ . Goal is to minimize  $\chi^2$  sum by adjusting parameters  $\theta$  (i.e. to find most optimal curve through points).

Joint pdf for  $N$  independent Gaussians  $y_i$  is product of  $N$  Gaussians:

$$g(\bar{y}; \bar{\lambda}, \bar{\sigma}) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{(y_i - \lambda_i)^2}{2\sigma_i^2}\right)$$

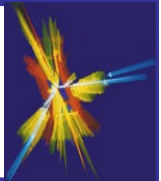
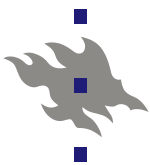
Then log-likelihood function (drop terms independent of  $\bar{\theta}$ ):

$$\ln L(\bar{\theta}) = -\frac{1}{2} \sum_{i=1}^N (y_i - \lambda(x_i, \bar{\theta}))^2 / \sigma_i^2$$

However maximizing  $\ln L(\bar{\theta})$  equivalent to minimizing

$$\chi^2(\bar{\theta}) = \sum_{i=1}^N (y_i - \lambda(x_i, \bar{\theta}))^2 / \sigma_i^2$$

Basis of **method of least squares**: minimize quadratic sum of difference between measured & hypothesized values, weighted by inverse of measurement variance.



So ML justifies somehow method of least-squares (LS).  
What "proves" validity of ML? Nothing, only an assumption  
LS estimators have particularly desirable properties when

$\lambda(x; \bar{\theta})$  linear function of  $\theta$ :  $\lambda(x; \bar{\theta}) = \sum_{j=1}^m a_j(x) \theta_j$ ,  
where  $a_j(x)$  are any linearly independent functions of  $x$  (i.e. one term can't be given as linear combination of the others)  $\Rightarrow$

- LS estimators have zero bias & are efficient (i.e. have minimal variance) for any  $N$  ("Gauss-Markov" theorem).
  - LS estimators & their variances can be found analytically although one may still prefer to estimate them numerically.
- Variance (also **MC method** valid):

- Calculate **analytically** covariance matrix  $V_{ij} = \text{cov}[\hat{\theta}_i, \hat{\theta}_j]$
- Alternatively one can estimate elements **numerically**:

$$(V^{-1})_{ij} = \frac{1}{2} \left[ \frac{\partial^2 \chi^2}{\partial \theta_i \partial \theta_j} \right]_{\bar{\theta} = \hat{\theta}} \rightarrow \text{coincides with RCF bound if } \lambda \text{ linear function of } \theta \text{ \& } V[y_i] \text{ known and Gaussian, then } \ln L = -\chi^2/2.$$

- Since  $\lambda(x; \bar{\theta})$  linear in parameters  $\bar{\theta}$ ,  $\chi^2$  quadratic in

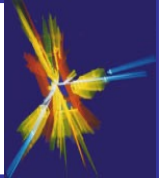
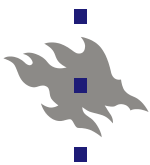
$$\chi^2(\bar{\theta}) = \chi^2(\hat{\bar{\theta}}) + \frac{1}{2} \sum_{i,j=1}^m \left[ \frac{\partial^2 \chi^2}{\partial \theta_i \partial \theta_j} \right]_{\bar{\theta} = \hat{\bar{\theta}}} (\theta_i - \hat{\theta}_i)(\theta_j - \hat{\theta}_j)$$

(an expansion of the  $\chi^2$  function around its minima)

$\rightarrow$  1 standard deviation contour in parameter space given by curve whose tangents are  $\hat{\theta}_i \pm \hat{\sigma}_{\hat{\theta}_i}$  that corresponds to

$$\chi^2(\bar{\theta}) = \chi^2(\hat{\bar{\theta}}) + 1 = \chi_{\min}^2 + 1 \quad \text{graphical method}$$

Even when  $\lambda(x; \bar{\theta})$  not linear in  $\bar{\theta}$  & formula not really valid, region  $\chi^2(\bar{\theta}) \leq \chi_{\min}^2 + 1$  can still be interpreted as a "confidence region" with a given probability of containing true  $\bar{\theta}$ .



A usual application of LS method is a polynomial fit.  
Independent data  $y \pm \Delta y$  measured at different  $x$ -values:

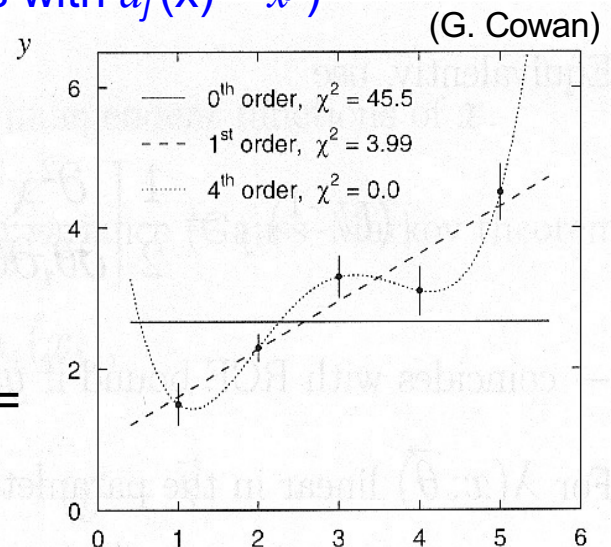
$$\lambda(x; \theta_0, \dots, \theta_m) = \sum_{j=0}^m \theta_j x^j \quad (m+1 \text{ free parameters})$$

(special case of linear LS fits with  $a_j(x) = x^j$ )

LS polynomial fit to 5 data points (with  $\chi^2_{\min}$  indicated)

- 0th order (1 parameter)
- 1st order (2 parameters)
- 4th order (5 parameters)

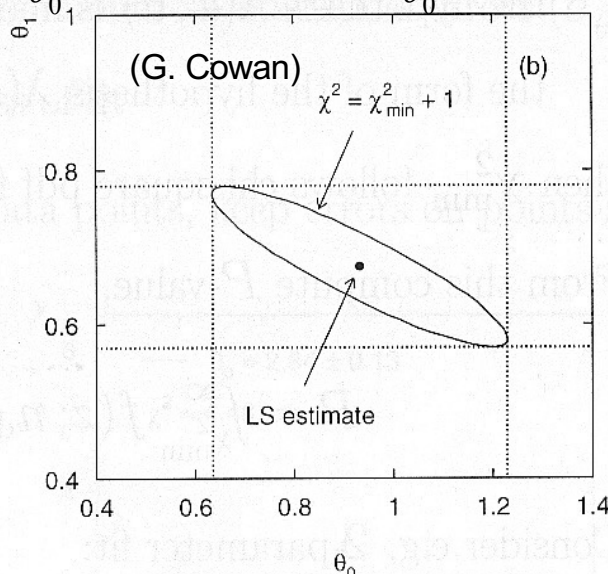
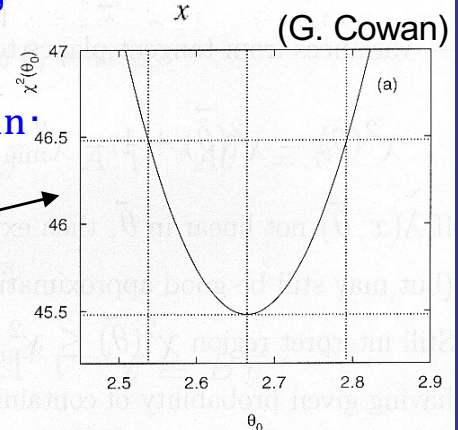
NB! last case meaningless since number of parameters = number of data points



Uncertainties & covariances estimated using any of the 4 methods, all related to  $\chi^2$  change when parameters changed from those giving  $\chi^2_{\min}$ .

0th order:  $\hat{\theta}_0 = 2.66 \pm 0.13$ ,  $\chi^2_{\min} = 45.5$

$\sigma_{\hat{\theta}_0}$  from  $\chi^2(\hat{\theta}_0 \pm \sigma_{\hat{\theta}_0}) \leq \chi^2_{\min} + 1$



1st order:

$$\hat{\theta}_0 = 0.93 \pm 0.30,$$

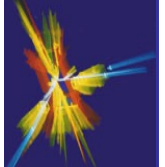
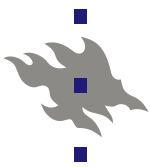
$$\hat{\theta}_1 = 0.68 \pm 0.10,$$

$$\chi^2_{\min} = 3.99,$$

$$\text{cov}[\hat{\theta}_0, \hat{\theta}_1] = -0.028, \rho =$$

$$\text{cov}[\hat{\theta}_0, \hat{\theta}_1] / \sigma_{\hat{\theta}_0} \sigma_{\hat{\theta}_1} = -0.93$$





Let's examine our 1st order polynomial fit in more detail.  
How would one determine the LS estimators for the 0th & 1st order terms in this case. The expression to minimize is

$$\chi^2(\bar{\theta}) = \sum_{i=1}^N (y_i - \theta_0 - \theta_1 x_i)^2 / \sigma_i^2$$

To get the LS estimators one has to look for a local minima:  
let's further make the simplification that all  $\sigma_i$ 's are equal ( $= \sigma$ ).

$$\frac{\partial \chi^2}{\partial \theta_0} = 0 \rightarrow \sum_{i=1}^N (y_i - \hat{\theta}_0 - \hat{\theta}_1 x_i) = \bar{y} - \hat{\theta}_0 - \hat{\theta}_1 \bar{x} = 0$$

$$\frac{\partial \chi^2}{\partial \theta_1} = 0 \rightarrow \sum_{i=1}^N x_i (y_i - \hat{\theta}_0 - \hat{\theta}_1 x_i) = \bar{x}\bar{y} - \hat{\theta}_0 \bar{x} - \hat{\theta}_1 \bar{x}^2 = 0$$

Combining the two equations, one obtains for the LS estimators:

$$\hat{\theta}_1 = (\bar{x}\bar{y} - \bar{x} \bar{y}) / (\bar{x}^2 - \bar{x}^2) = \text{cov}[x, y] / V[x]$$

$$\hat{\theta}_0 = (\bar{x}^2 \bar{y} - \bar{x} \bar{x}\bar{y}) / (\bar{x}^2 - \bar{x}^2) = \bar{y} - \hat{\theta}_1 \bar{x}$$

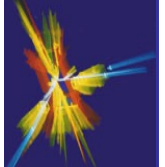
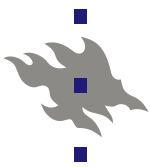
These are standard equations to get the slope  $\theta_1$  & the intercept with the  $y$ -axis  $\theta_0$  for a straight line fit. With different  $\sigma_i$ 's for each point, the formula should only be modified such that each point is given a weight  $1/\sigma_i^2$  & normalisation = the total weight  $\Sigma 1/\sigma_i^2$ .

Obtained  $\chi_{\min}^2$  value used to estimate probability of hypothesis, if true, would give observed data.  $(y_i - \lambda(x_i; \theta)) / \sigma_i$  – a measure of agreement btwn observed data & hypothesis.

$\chi_{\min}^2$  obtained from a LS fit follows the standard chi-square distribution with the degrees-of-freedom =  $N - m$  if:

- $y_i, i = 1, \dots, N$  Gaussian random variables with known covariance matrix  $V_{ij}$  or independent  $y_i$ 's with known  $\sigma_i$ 's.
- hypothesis  $\lambda(x; \bar{\theta})$  is linear in parameters  $\theta_i, i = 1, \dots, m$ .
- functional form of the hypothesis  $\lambda(x; \bar{\theta})$  is correct.





If all previous satisfied, one can calculate the  $P$ -value:

$$P = \int_{\chi_{\min}^2}^{\infty} f(z; n_d) dz$$

Example: consider our 1st order polynomial fit ( $m = 2$ )

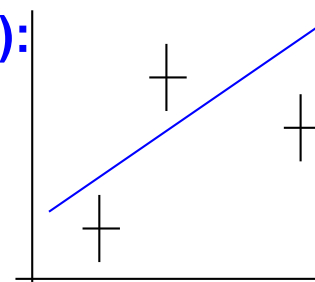
$$\chi_{\min}^2 = 3.99, N - m = 3 \rightarrow P = 0.263$$

i.e. if it is true that  $\lambda(x)$  would be a straight line and if the experiment would be repeated many times, then in 26.3 % of the cases, one would obtain a worse (i.e. higher)  $\chi_{\min}^2$ .

NB!  $E[\chi_{\min}^2] = n_d$  (number degrees-of-freedom)  $\rightarrow$  each data point should contribute  $\approx 1$  to the  $\chi^2$

$\chi_{\min}^2/n_d \gg 1$  (or a very small  $P$ -value):

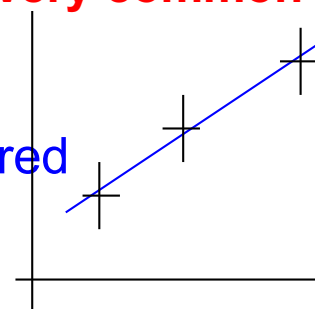
- hypothesis ("= function") wrong
- or measurement wrong/bad
- or uncertainties underestimated
- or extremely bad luck (unlikely!)



Very common pit fall!

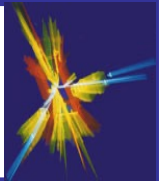
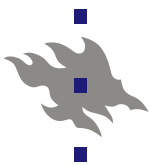
$\chi_{\min}^2/n_d \ll 1$  (or a large  $P$ -value):

- the uncertainties overestimated
- or correlations of uncertainties ignored
- or extremely good luck (unlikely!)



Note distinction btwn small statistical uncertainty on LS estimator & a good LS fit (i.e. small  $\chi_{\min}^2$ )  
Statistical uncertainty estimated from change of  $\chi^2$  near its minimum is independent of the absolute value of the  $\chi_{\min}^2$ .

Variance (statistical uncertainty) of estimator tells us: If experiment repeated many times, how wide is distribution of estimates (doesn't tell whether hypothesis correct or not).



Very common use of LS method is to combine a number of measurements of the same quantity. Then one has:

$y_i$  = result of measurement  $i$ ,  $i = 1, \dots, N$ ;  
 $\sigma_i^2 = V[y_i]$ , assumed to be known;  
 $\lambda$  = true value (takes role of  $\theta$ , no  $x$ -dependence).

For independent  $y_i$ 's, minimize:  $\chi^2(\lambda) = \sum_{i=1}^N \frac{(y_i - \lambda)^2}{\sigma_i^2}$

$$\frac{\partial \chi^2}{\partial \lambda} = 0 \rightarrow \hat{\lambda} = \sum_{i=1}^N \frac{y_i}{\sigma_i^2} / \sum_{i=1}^N \frac{1}{\sigma_i^2}$$

$$V[\hat{\lambda}] = \left( \frac{1}{2} \left[ \frac{\partial^2 \chi^2}{\partial^2 \lambda} \right]_{\lambda=\hat{\lambda}} \right)^{-1} = 1 / \sum_{i=1}^N \frac{1}{\sigma_i^2}$$

Well-known formula for **weighted average**. Variance of average < variances of individual measurements. More precise measurements (i.e. having smaller variances) have larger weight. Generalized to non-independent measurements  $y_i$  i.e.  $\text{cov}[y_i, y_j] = V_{ij}$ . Then have to minimize:

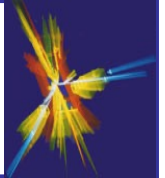
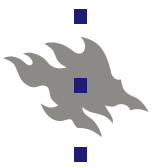
$$\chi^2(\lambda) = \sum_{i,j=1}^N (y_i - \lambda)(V^{-1})_{ij}(y_j - \lambda)$$

$$\frac{\partial \chi^2}{\partial \lambda} = 0 \rightarrow \hat{\lambda} = \sum_{i=1}^N w_i y_i, \quad w_i = \sum_{j=1}^N (V^{-1})_{ij} / \sum_{k,l=1}^N (V^{-1})_{kl}$$

$$V[\hat{\lambda}] = \left( \frac{1}{2} \left[ \frac{\partial^2 \chi^2}{\partial^2 \lambda} \right]_{\lambda=\hat{\lambda}} \right)^{-1} = \sum_{i,j=1}^N w_i V_{ij} w_j = \mathbf{w}^T \mathbf{V} \mathbf{w}$$

$$\sum_{i=1}^N w_i = 1 \rightarrow E[\hat{\lambda}] = \sum_{i=1}^N w_i E[y_i] = \lambda \sum_{i=1}^N w_i = \lambda, \text{ i.e. unbiased}$$

Assumption: individual  $y_i$ 's unbiased. The weights in LS prescription give RCF bound for variance ("Gauss-Markov").



Averaging 2 correlated measurements: measurements  $y_1$  &  $y_2$  with covariance matrix  $V$ :

$$V = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix} \rightarrow V^{-1} = \frac{1}{1-\rho^2} \begin{pmatrix} 1/\sigma_1^2 & -\rho/\sigma_1\sigma_2 \\ -\rho/\sigma_1\sigma_2 & 1/\sigma_2^2 \end{pmatrix}$$

then  $\hat{\lambda} = wy_1 + (1-w)y_2$ ,  $w = \frac{\sigma_2^2 - \rho\sigma_1\sigma_2}{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2}$  and

$$V[\hat{\lambda}] = \frac{(1-\rho^2)\sigma_1^2\sigma_2^2}{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2} \equiv \sigma^2 \Leftrightarrow \frac{1}{\sigma^2} = \frac{1}{1-\rho^2} \left( \frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2} - \frac{2\rho}{\sigma_1\sigma_2} \right)$$

Increase of inverse variance due to the 2<sup>nd</sup> measurement:

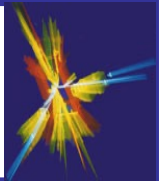
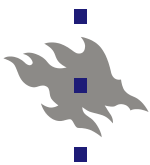
$$\frac{1}{\sigma^2} - \frac{1}{\sigma_1^2} = \frac{1}{1-\rho^2} \left( \frac{\rho}{\sigma_1} - \frac{1}{\sigma_2} \right)^2 \geq 0$$

$\rho \leq 1 \Rightarrow$  2<sup>nd</sup> measurement only beneficial for average (i.e. new combined variance  $\leq$ ). No variance change when  $\rho = \sigma_1/\sigma_2$  (incl.  $\rho = 1$  &  $\sigma_1 = \sigma_2$ , i.e. same measurement twice).

If  $\rho > \sigma_2/\sigma_1$  then  $w < 0$  & weighted average not btwn  $y_1$  &  $y_2$ , due to a large positive correlation btwn  $y_1$  &  $y_2$ . Can happen in the case of common normalisation uncertainty.

Usable for calibrating common variable see e.g. G.Cowan: *Statistical Data Analysis*, page 109, where temperature ( $T$ ) estimate improved using measurements at same  $T$  with two rulers having different thermal expansion coefficients.

Overlapping samples can be used to estimate covariance matrix if not known. Either use MC generated samples or use real data by dividing data sample into a large number of subsamples & determining estimators  $y_1, \dots, y_N$  for each subsample & from these correlations coefficient matrix.



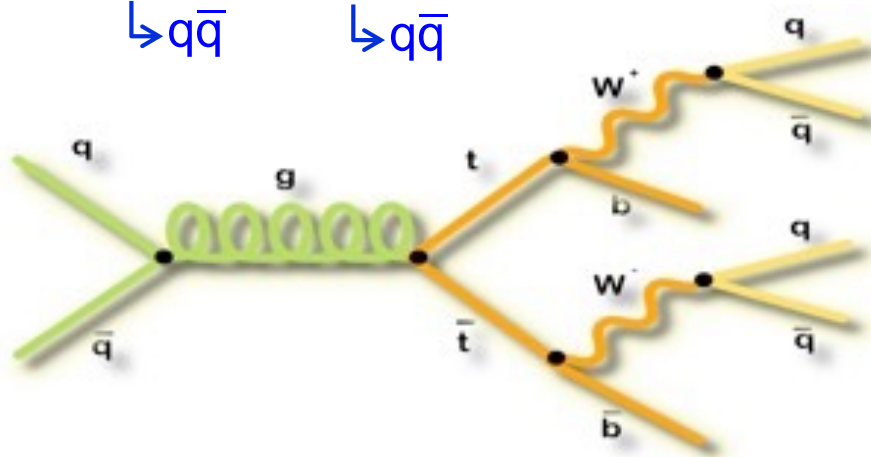
LS fit with constraints (penalty functions/nuisance terms) to improve data quality: Sometimes inputs ( $x_i$ 's) suffer from significant uncertainties or there are some scalings (nuisance terms) involved that directly effects our result. Then a possible solution can be to include additional terms in the  $\chi^2$  sum and look for the global  $\chi^2_{\min}$  that minimizes everything including the variation of inputs (with respect to their uncertainties) or the scalings.

Example: top mass reconstruction at CDF experiment

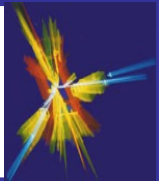
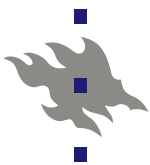
Proton-antiproton collisions create top-antitop pair:

$p\bar{p} \rightarrow t (\rightarrow W^+ b) \bar{t} (\rightarrow W^- \bar{b}) X$

$\hookrightarrow q\bar{q} \quad \hookrightarrow q\bar{q}$



experimental challenge: quarks ( $q$ ,  $b$ ) hadronize and make jets whose energy (& momentum) measured with poor resolution (10-20 %)  $\Rightarrow$  if raw energy (& momentum) measurement used the top mass poorly reconstructed. Make use of additional constraints: two of the jets should make a  $W$  resonance ( $\mu_W = 80.42 \pm 0.03$  GeV for the  $W$  mass) with a width  $\Gamma_W = 2.1$  GeV + combined with the third jet should make a top quark with same  $\mu_{\text{top}}$  (= top mass) as the other top (= "triplet" of jets) within a width  $\Gamma_{\text{top}} = 1.6$  GeV.



i) The masses are first determined by using the formulae given:

$$m_i = \sqrt{E_i^2 - P_{x,i}^2 - P_{y,i}^2 - P_{z,i}^2}$$

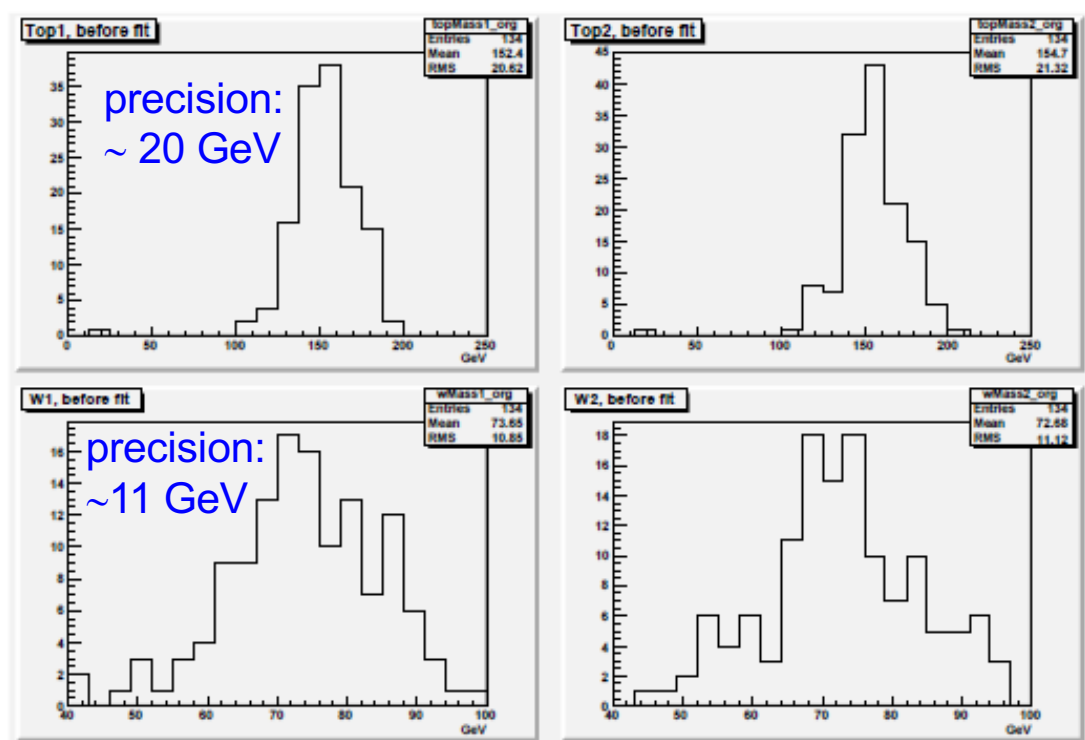
The energies of the  $j$ :th ( $j = 1, 2$ ) top quark and the W-boson are (from the channel given):

$E_{bj(qj)}$   $\equiv$  energy of b-jet (q-jet)  $j$   
 $P_{x(y,z),bj(qj)} \equiv x(y,z)$  momentum  
 component of b-jet (q-jet)  $j$

$$E_{t,j} = E_{bj} + E_{qj1} + E_{qj2}$$

(same for  $P_x$ ,  
 $P_y$  and  $P_z$ )

$$E_{w,j} = E_{qj1} + E_{qj2}$$



ii) Now we test the hypothesis given

$$\chi^2 = \frac{(m_{top1} - m_{top2})^2}{2\Gamma_{top}^2} + \frac{(m_{W1} - m_{Wnom})^2}{\Gamma_W^2} + \frac{(m_{W2} - m_{Wnom})^2}{\Gamma_W^2},$$

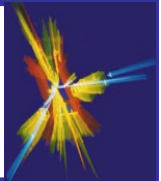
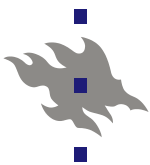
$$\Gamma_{top} = 1.6 \text{ GeV}, \Gamma_W = 2.1 \text{ GeV}, m_{Wnom} = 80.42 \text{ GeV}$$

$\chi^2$  large ( $> 100$ ). Cure: allow energy measurements to vary

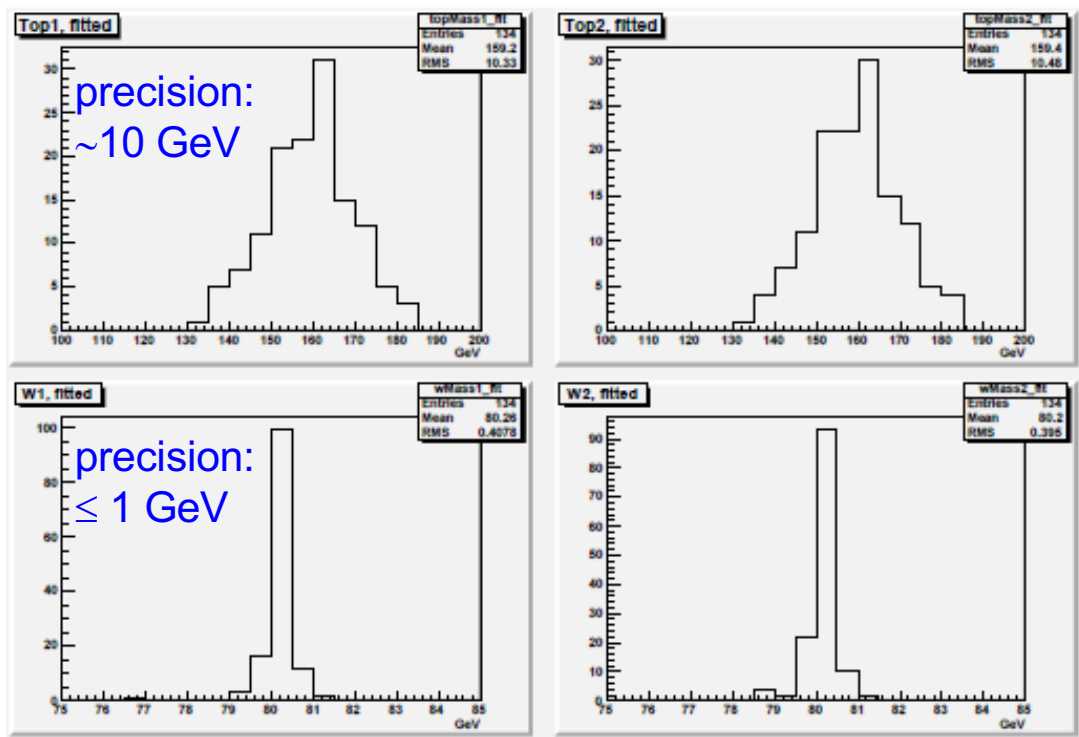
$$\chi_{new}^2 = \chi^2 + \sum_{i=b1,b2,q11,q12,q21,q22} \frac{(E_i - c_i \cdot E_i)^2}{\sigma_{E_i}^2}$$

Let each  $c_i$  vary in order to minimize whole  $\chi_{new}^2$

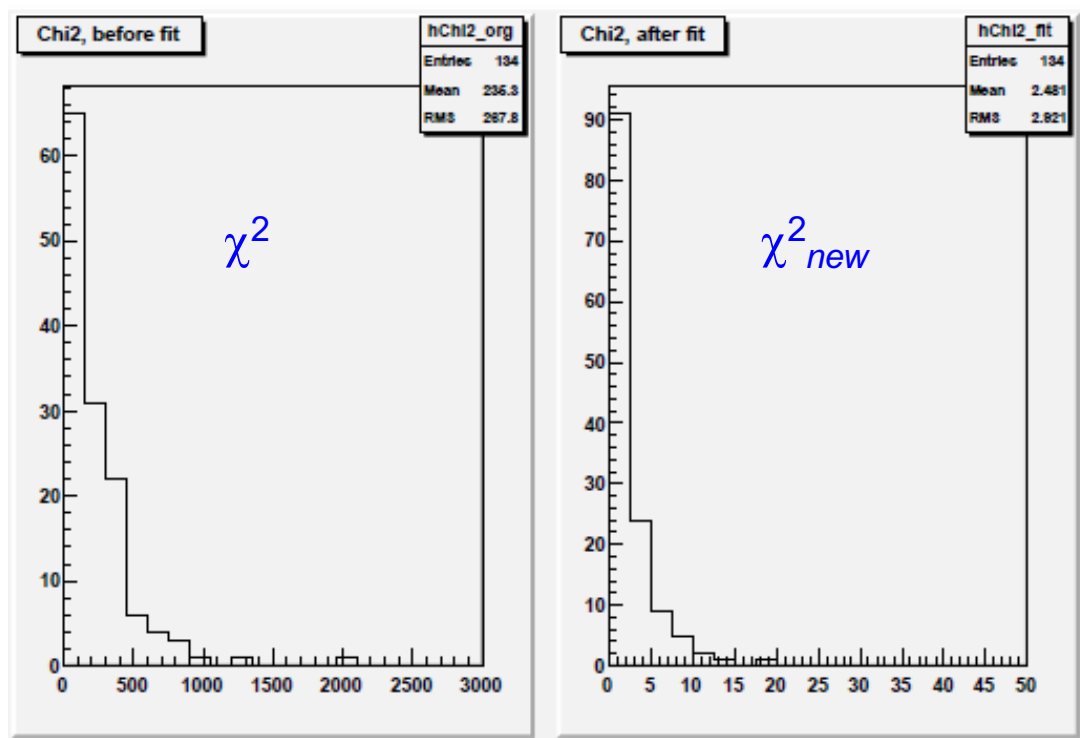




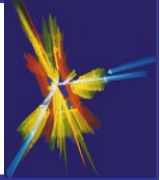
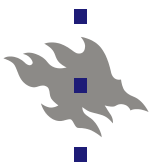
## Result:



global  $\chi^2$  values significantly improved (by a factor  $\sim 100$ !)

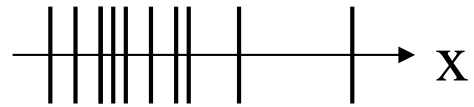


Note: here problem a bit simplified, in reality need also to scale  $P$ 's



Data of very large samples often given as  $N$ -binned histograms:

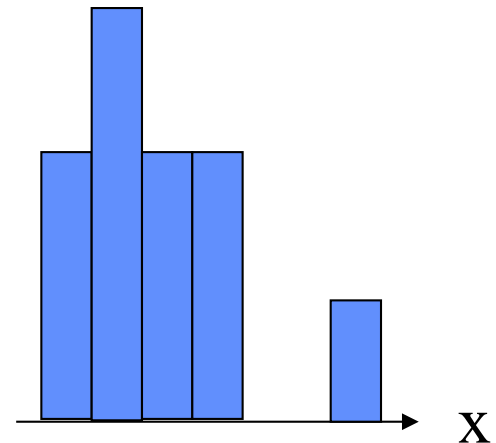
$$\bar{n} = (n_1, \dots, n_N), n_{\text{tot}} = \sum_{i=1}^N n_i$$



The hypothesis for each bin can be expressed as a Poissonian:

$$\bar{v} = (v_1, \dots, v_N), v_{\text{tot}} = \sum_{i=1}^N v_i,$$

$$\text{where } v_i(\bar{\theta}) = v_{\text{tot}} \int_{x_i^{\min}}^{x_i^{\max}} f(x; \bar{\theta})$$



Three different techniques for fitting histograms:

- Binned ML (either leaving the normalisation free or not)

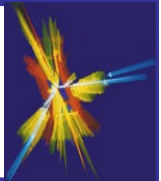
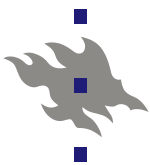
$$\ln L(v_{\text{tot}}, \bar{\theta}) = \sum_{i=1}^N (n_i \ln v_i(v_{\text{tot}}, \bar{\theta}) - v_i(v_{\text{tot}}, \bar{\theta}))$$

$$\ln L(\bar{\theta}) = \sum_{i=1}^N (n_i \ln v_i(\bar{\theta}))$$

**Suggested technique!**

- Proper LS  $\chi^2(\bar{\theta}) = \sum_{i=1}^N (n_i - v_i(\bar{\theta}))^2 / v_i(\bar{\theta}),$
- Modified LS  $\chi^2(\bar{\theta}) = \sum_{i=1}^N (n_i - v_i(\bar{\theta}))^2 / n_i,$

Binned ML gives same result as unbinned ML (unless bin size > feature size) & especially recommended for small statistics. Both LS methods are biased & less efficient if content per bin small (due to asymmetry of Poissonian). LS methods give  $\chi^2$  value directly but  $\chi^2$  value can also be extracted in case of binned ML (see following pages).



Data of large samples often as  $N$ -binned histograms:

$$\bar{n} = (n_1, \dots, n_N), n_{\text{tot}} = \sum_{i=1}^N n_i, \text{ hypothesis given in similarly:}$$

$$\bar{v} = (v_1, \dots, v_N), v_{\text{tot}} = \sum_{i=1}^N v_i, \text{ where } v_i(\bar{\theta}) = n_{\text{tot}} \int_{x_i^{\min}}^{x_i^{\max}} f(x; \bar{\theta})$$

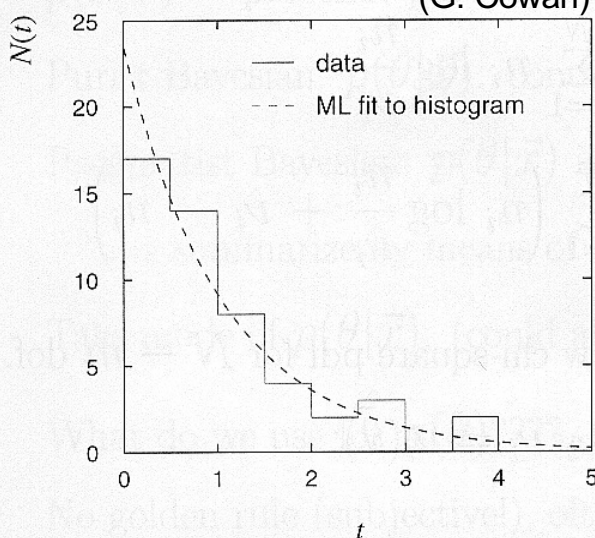
Can regard joint sample pdf as multinomial ( $n_{\text{tot}}$  constant!)

$$f_{\text{joint}}(\bar{n}; \bar{v}) = \frac{n_{\text{tot}}!}{n_1! \dots n_N!} \left( \frac{v_1}{n_{\text{tot}}} \right)^{n_1} \dots \left( \frac{v_N}{n_{\text{tot}}} \right)^{n_N} \rightarrow \ln L(\bar{\theta}) = \sum_{i=1}^N n_i \ln v_i(\bar{\theta})$$

in limit of zero bin width  $\rightarrow$  usual unbinned ML.

NB! 0 or few entries in a bin not a problem for binned ML.

(G. Cowan)



Example with exponential distribution (bin width  $\Delta t = 0.5$ ).

$$\hat{\tau} = 1.07 \pm 0.17$$

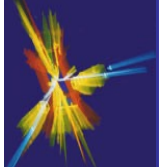
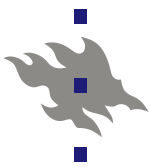
( $1.06 \pm 0.15$  for unbinned ML fit to same sample)

usual result of binning:  
"loose" some information

If total number of entries  $n_{\text{tot}}$  treated as Poissonian variable, one gets for "binned" extended log-likelihood:

$$f_{\text{joint}}(\bar{n}; \bar{v}) = \prod_{i=1}^N \frac{v_i^{n_i}}{n_i!} e^{-v_i} \rightarrow \ln L(v_{\text{tot}}, \bar{\theta}) = \sum_{i=1}^N (n_i \ln v_i(v_{\text{tot}}, \bar{\theta}) - v_i(v_{\text{tot}}, \bar{\theta}))$$

(here pdf a product of a Poisson & a multinomial distribution, same considerations regarding  $n_{\text{tot}}$ 's  $\bar{\theta}$  (non)dependence also valid here as with usual extended ML.) NB!  $v_i(\bar{\theta}) = v_{\text{tot}} \int_{x_i^{\min}}^{x_i^{\max}} f(x; \bar{\theta})$



When using binned ML fits,  $\chi^2$  value estimated using ratio

$$\lambda = \frac{L(\bar{n}, \hat{\bar{v}})}{L(\bar{n}, \bar{n})} = \frac{f_{\text{joint}}(\bar{n}, \hat{\bar{v}})}{f_{\text{joint}}(\bar{n}, \bar{n})} = \prod_{i=1}^N \left( \frac{\hat{v}_i}{n_i} \right)^{n_i} (\cdot e^{n_{\text{tot}} - \hat{v}_{\text{tot}}} \text{ if } n_{\text{tot}} \text{ free})$$

then for e.g. multinomial only ( $M$ , i.e.  $n_{\text{tot}}$  fixed) or multinomial + a Poisson ( $P$ , i.e.  $n_{\text{tot}}$  free), one can use:

$$\chi_M^2 = -2 \ln \lambda_M = 2 \sum_{i=1}^N n_i \ln \frac{n_i}{\hat{v}_i} \quad (N - m - 1 \text{ dof})$$

$$\chi_P^2 = -2 \ln \lambda_P = 2 \sum_{i=1}^N \left( n_i \ln \frac{n_i}{\hat{v}_i} + \hat{v}_i - n_i \right) \quad (N - m \text{ dof})$$

Both variables follow chi-square distribution for number of degrees-of-freedom (dof) indicated in the parenthesis.

Here aim of ML fit is to estimate mean value for each bin:

$$\hat{v}_i = n_{\text{tot}} \int_{x_i^{\min}}^{x_i^{\max}} f(x; \hat{\theta}) dx \quad (\text{fit } m \text{ parameters})$$

Since  $L(\bar{n}, \bar{n})$  doesn't depend on any unknown parameters, then parameters maximizing  $L(\bar{n}, \hat{\bar{v}})$  maximizes also  $\lambda(\hat{\bar{v}})$ .

NB! zero entries in a bin  $i$  is not a problem, put  $(n_i)^{n_i} = 1$ .

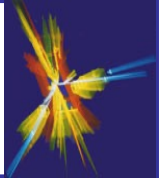
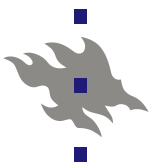
Alternative, use Pearson's  $\chi^2$  test & replace  $v_i$  by  $\hat{v}_i = v_i(\hat{\theta})$

$$\chi^2 = \sum_{i=1}^N (n_i - \hat{v}_i)^2 / \hat{v}_i \quad (\text{Poisson: dof} = N - m)$$

$$\text{NB! } \hat{p}_i = \hat{v}_i / \hat{v}_{\text{tot}}$$

$$\chi^2 = \sum_{i=1}^N (n_i - \hat{p}_i n_{\text{tot}})^2 / \hat{p}_i n_{\text{tot}} \quad (\text{multinomial: dof} = N - m - 1)$$

as usual follows chi-square distribution only if all  $\hat{v}_i \geq 5$ .



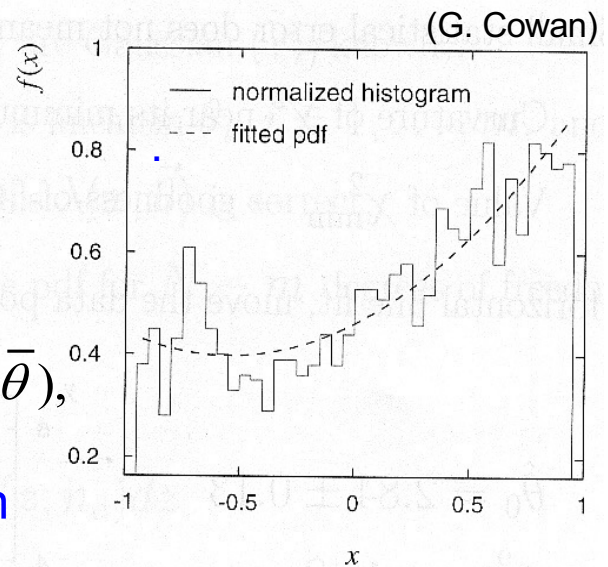
Consider data put in a histogram with  $N$  bins,  $n$  entries & a hypothesized pdf  $f(x; \bar{\theta})$

Define:

$y_i$  = number of entries in bin  $i$

$$\nu_i(\bar{\theta}) = n \int_{x_i^{\min}}^{x_i^{\max}} f(x; \bar{\theta}) dx = n p_i(\bar{\theta}),$$

where  $x_i^{\min}$  &  $x_i^{\max}$  are the bin limits and  $p_i$  the probability to have an entry in bin  $i$ .



**LS method for binned data:** parameters  $\bar{\theta}$  minimizing

$$\chi^2(\bar{\theta}) = \sum_{i=1}^N (y_i - \nu_i(\bar{\theta}))^2 / \sigma_i^2,$$

where  $\sigma_i^2 = V[y_i]$  not known a priori (make assumption).

Most common assumption: treat the  $y_i$ 's as Poisson random variables & in place of true variance put either

$$\sigma_i^2 = \nu_i(\bar{\theta}) \quad (\text{LS method})$$

$$\sigma_i^2 = y_i \quad (\text{modified LS method})$$

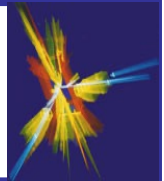
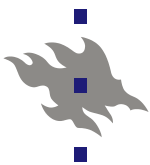
MLS may be easier to deal with computationally, but  $\chi^2_{\min}$  no longer follows standard chi-square distribution (& a priori undefined if bins without any entries exists).

Often when LS method is used for binned data,

normalization  $\nu_{\text{tot}}$   
for total number  
of entries left free

$$\nu_i(\bar{\theta}, \nu) = \nu_{\text{tot}} \int_{x_i^{\min}}^{x_i^{\max}} f(x; \bar{\theta}) dx = \nu_{\text{tot}} p_i(\bar{\theta}),$$





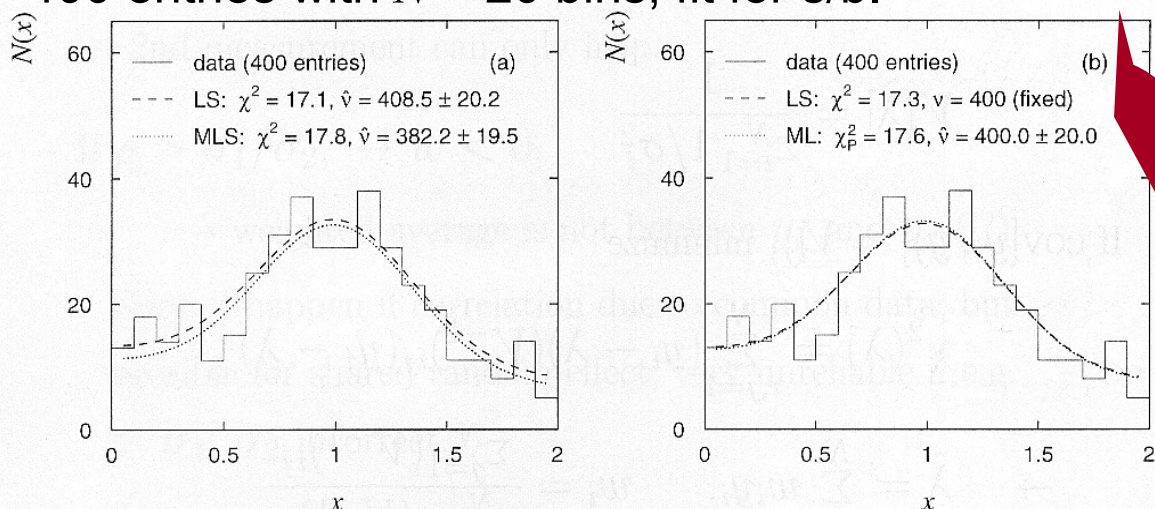
Can also use LS with binned data but need to be aware that  $\hat{v}$  a non-consistent (= bad) estimator for  $n$ . Can show

$$\hat{v}_{\text{LS}} = n + \frac{\chi_{\min}^2}{2} \quad \hat{v}_{\text{MLS}} = n - \chi_{\min}^2$$

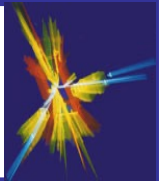
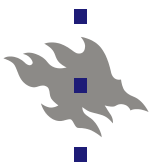
since expect contribution to  $\chi^2$  of  $O(1)$  per bin, the relative fitted number of entries typically  $N/2n$  too high (LS) or  $N/n$  too low (MLS). Taking rule of thumb of  $\geq 5$  entries per bin, can get normalization errors upto 10-20%. Therefore if LS method used, take normalisation from number of entries !!

Recall that ML method for binned data doesn't have such problems, i.e.  $\hat{v}_{\text{ML}} = n$ . Also variances of ML estimators converge faster with  $n$  to RCF bound (= smallest possible variance) than LS or MLS, see e.g. W.T. Eadie *et al.*, Statistical methods in Experimental Physics, North-Holland, 1971  
→ recommended to use ML estimators for binned data !!

Example: gaussian signal( $s$ ) + exponential background( $b$ ).  
MC data in  $[0,2]$  with  $s/b = 1$ ,  $\mu_s = 1$ ,  $\sigma_s = 0.35$  &  $\xi_b = 4$ .  
4 different fits: LS & MLS( $v$  free), LS( $v$  fixed) & ML( $v$  free).  
 $n = 400$  entries with  $N = 20$  bins, fit for  $s/b$ . (G. Cowan)



Although standard deviation same in all cases, 2 later fits are to be preferred, since number of entries are correct.



**Systematic effects**, a general naming for all effects that alters measured quantities / result somehow, including background, efficiency, calibration, bias, extrapolation, resolution, time variations, geometric effects, dead time ... Measured quantities & results corrected for such effects using a procedure (resulting in a correction factor). Uncertainty (or i.e. upper limit of the knowledge) in the procedure or correction factor is called **systematic uncertainty**. (aim should be to correct for systematic effects, correction done at certain precision = "systematical uncertainty").

E.g.: charge particle transverse momentum measurement

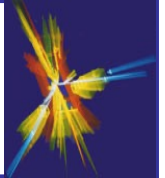
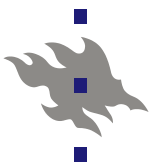
$$p_T [\text{GeV}] = 0.3B\rho [\text{Tm}]$$

- 0.3 a hard number – speed of light  $\times 10^{-9}$ , known exactly.
- radius of curvature  $\rho$  obtained for each particle & subject to uncertainties. If devices aligned, on average balanced.
- magnetic field  $B$  measured with some finite accuracy.

For each particle,  $B$  applied either always too high or low  
 $\Rightarrow$  similar systematic effect always irrespective of statistics.

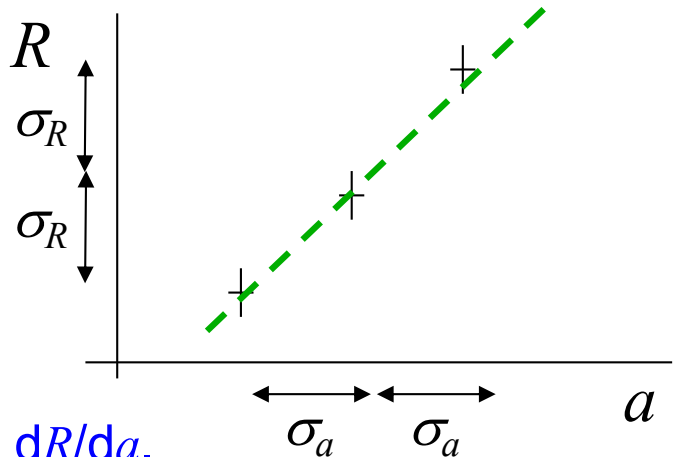
Final result usually given with a statistical & a systematic uncertainty. Total measurement uncertainty often obtained by adding the two in quadrature & hence assuming systematic uncertainties to be Gaussian (mostly not correct !!). Best to quote the two uncertainties apart & give result as:

$A_{FB} = -0.102 \pm 0.012 \pm 0.023$  or  $\sigma = 45 \pm 4 \pm 1$  mb,  
where first uncertainty is statistical & second systematic.  
This eases combination & comparison with other results plus shows relative importance of different contributions.



When direct uncertainty propagation can't be applied must numerically estimate variation of result,  $R$ , on parameter  $a \Rightarrow \sigma_a = dR/da$ . Typical in experimental physics due to e.g. theoretical parameters & experimental calibrations.

If  $a$  known only with some precision  $\sigma_a$ , compute  $R$  for values  $a = a_0 \pm (\text{a few})\sigma_a$ , where  $a_0$  reference value of  $a$  (variation sufficiently large to make a clear measurable change on  $R$ ).



Assume linear & determine  $dR/da$ .  
How treat uncertainties in  $dR/da$ ?

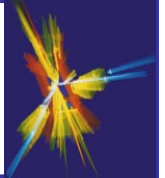
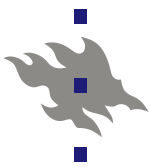
- ignore it since typically  $\sigma_{dR/da} \ll dR/da$  !!

Alternative (for significant uncertainties): Include it into likelihood. Regard  $a$  as another parameter in analysis ML or include term  $(a - a_0)^2 / \sigma_a^2$  as a contribution to  $\chi^2$ . If  $a$  allowed to vary result changes & gives smaller uncertainty. Works if data really depends on  $a$  & fitted  $a$  makes sense.

Searching for “hidden” systematic uncertainties: after estimating all known systematic uncertainties, must assure absence of any further “hidden” systematic uncertainty.

Good practice to repeat the analysis in different forms e.g.

- vary the range of the data used to extract the result
- vary applied cuts to ensure sample purity & data quality
- include & exclude data subsets(taken in different conditions)
- use histograms of different bin size for extracting result
- determine quantities by counting & parameter estimation



Compatibility of two results: Often result,  $R$ , can be extracted from same data using two different methods to investigate for a systematic uncertainty. Assume estimates  $R_1$  &  $R_2$  with statistical errors  $\sigma_1$  &  $\sigma_2$ . How can one decide whether  $R_1$  &  $R_2$  compatible or not compatible? Examine

$$\Delta = R_1 - R_2 \quad \sigma_\Delta^2 = \sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2$$

If  $|\Delta| \leq \sim 1-2\sigma_\Delta$ , assume estimates to agree satisfactorily. Often not at all straight forward to estimate  $\rho$ , sometimes impossible. However can place limits on  $\sigma_\Delta$  using RCF bound. Idea: estimate weighted average of two estimates.

$$R = wR_1 + (1-w)R_2 \Rightarrow w = \frac{\sigma_2^2 - \rho\sigma_1\sigma_2}{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2}$$

$$\Rightarrow \sigma_R^2 = \frac{(1-\rho^2)\sigma_1^2\sigma_2^2}{\sigma_1^2 + \sigma_2^2 - 2\rho\sigma_1\sigma_2}$$

$$\text{RCF bound: } \sigma_R^2 \geq \frac{1}{N \int (d \ln P(x; R) / dR)^2 P(x; R) dx} = \sigma_{R, \min}^2$$

where  $P(x; R)$  is probability density function for distribution.

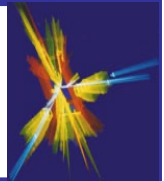
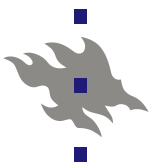
$$\sigma_R \geq \sigma_{R, \min} \Rightarrow \rho = \frac{\sigma_{R, \min}^2 \pm \sqrt{(\sigma_1^2 - \sigma_{R, \min}^2)(\sigma_2^2 - \sigma_{R, \min}^2)}}{\sigma_1\sigma_2}$$

$$\sigma_\Delta^{\max} = \sqrt{(\sigma_1^2 - \sigma_{R, \min}^2)} + \sqrt{(\sigma_2^2 - \sigma_{R, \min}^2)} \quad \sigma_\Delta^{\min} = \left| \sqrt{(\sigma_1^2 - \sigma_{R, \min}^2)} - \sqrt{(\sigma_2^2 - \sigma_{R, \min}^2)} \right|$$

RCF bound can be calculated quite easily in most cases. In practice  $\sigma_\Delta$  range fairly narrow since mostly (at least) 1 of analysis uses an efficient estimator (e.g.  $\sigma_1 = \sigma_{R, \min}$ ).

Hence, sufficient to use

$$\sigma_\Delta = \sqrt{|\sigma_1^2 - \sigma_2^2|}$$



Compatibility example: LEP experiments measured angular distribution of  $e^+e^- \rightarrow \mu^+\mu^-$  events using scattering angle  $\theta$  of produced  $\mu^-$ . Form of probability distribution assumed to be:

$$P(\cos \theta) = 3(1 + \cos^2 \theta)/8 + A \cos \theta,$$

where  $A$  asymmetry between "forward" ( $\cos \theta > 0$ ,  $\theta$  = polar angle between  $\mu^-$  &  $e^-$ ) & "backward" events ( $\cos \theta < 0$ ). With  $\sim 33.5\text{k}$   $e^+e^- \rightarrow \mu^+\mu^-$  events,  $A$  measured both by counting # of forward & backward events & by fit to  $\cos \theta$  distribution.

$$A = 0.0123 \pm 0.0055 \text{ (counting)} \quad A = 0.0084 \pm 0.0051 \text{ (fitting)}$$

The difference is 0.0039. Is this significant? RCF bound:

$$\frac{1}{\sigma_A^2} = N \int_{-1}^{+1} P(x) \left( \frac{d \ln P}{dA} \right)^2 dx = \int_{-1}^{+1} \frac{8Nx^2}{3 + 3x^2 + 8Ax} dx \approx \frac{8N}{3} \int_{-1}^{+1} \frac{x^2}{1 + x^2} dx = \frac{8N}{3} \frac{4 - \pi}{2}$$

$N = 33.5\text{k}$  events  $\Rightarrow \sigma_{A,\min} = 0.0051$ . Hence  $\sigma_\Delta = 0.0021$  & discrepancy between the two results is less than two standard deviations. NB! The fitting method saturates the RCF bound.

### Steps of a traditional experimental physics analysis

1. Devise analysis & get physics result.
2. Do the statistical uncertainty analysis.
3. Alter analysis/input/method to get systematic uncertainties.
4. Repeat step 3 until analysis/input/method exhausted.
5. Combine systematic uncertainties into a total one
6. Search for "hidden" systematics. ☹ identify one
7. Write note/paper & ask for feedback from colleagues.