

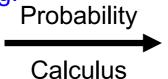
Parameter estimation



Parameter estimation: general concepts

Hypothesis testing:

Theory



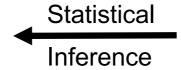
Data

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:

Theory



Data

<u>estimator</u> = procedure giving a value for a parameter or a property of distribution (pdf) from actual data values notation: estimator for θ is $\hat{\theta}$ (a **hat** indicates estimator) <u>estimate</u> = observed value of an estimator (often $\hat{\theta}_{obs}$) How does one construct an estimator $\hat{\theta}(\bar{x})$?

Exists no golden rule how to construct an estimator!

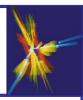
Examples of estimators are arithmetic mean & variance:

$$\hat{\mu}(\{x\}) = \frac{1}{N} \sum_{i} x_{i} \qquad \hat{V}(\{x\}) = \frac{1}{N-1} \sum_{i} (x_{i} - \hat{\mu})^{2}$$

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) θ & has expectation value, variance, etc...



Properties of estimators



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

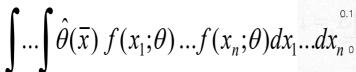
for any
$$\varepsilon > 0$$
, $\lim_{n \to \infty} P(|\hat{\theta} - \theta| > \varepsilon) = 0$

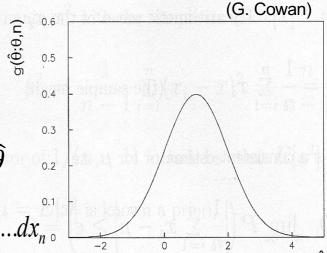
NB! convergence in sense of probability, i.e. no guaranty that any particular $\hat{\theta}_{obs}$ will be within given distance of θ .

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

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

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





variance $V[\hat{ heta}] = \sigma_{\widehat{ heta}}^2$

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

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

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

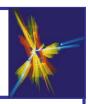
Good estimator: **consistent, unbiased** $(E[\hat{\theta}] = \theta)$ and **efficient** (i.e. has minimal possible variance = "RCF bound").

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

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



Estimator for the mean



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

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

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

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

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

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

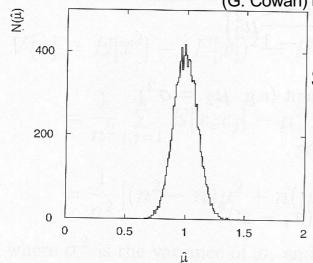
 $ightarrow \bar{x}$ is an unbiased estimator for μ . The variance of \bar{x} is

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

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

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.



(G. Cowan) Enter values into a histogram. $\frac{\bar{\hat{\mu}}}{\hat{\mu}} = 0.9981 \quad (\hat{\mu} \text{ unbiased})$

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)



HELSINGFORS UNIVERSITET Estimators for the variance



Suppose mean μ and variance $V[x] = \sigma^2$ both unknown. Then estimate σ^2 using <u>sample variance</u>

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

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

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

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

$$V[s^{2}] = \frac{1}{n} \left(\mu_{4} - \frac{n-3}{n-1} \mu_{2}^{2} \right) \qquad V[S^{2}] = \frac{1}{n} \left(\mu_{4} - \mu_{2}^{2} \right)$$

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

A natural estimator for standard deviation, σ , then

$$\hat{\sigma} = s = \sqrt{s^2}$$
 (or in case μ known $\hat{\sigma} = S = \sqrt{S^2}$)

For variance of estimator for standard deviation:

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

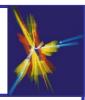
For a Gaussian pdf $\hat{s}^2(\sigma) = \sigma^2 \Longrightarrow d\hat{s}^2/d\sigma = 2\sigma$

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

Another quality measure of estimator: mean square error

$$MSE = E[(\hat{\theta} - \theta)^2] = V[\hat{\theta}] + b^2$$
 (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 θ . **maximum likelihood method** (suurimman uskottavuuden menetelmä) technique for estimating θ 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 θ) correct (= true)

→ expect higher probability for the data

If hypothesized functional form wrong or θ value far away

→ expect lower probability for the data

ln L

 \Rightarrow higher value of **likelihood function** close to true θ

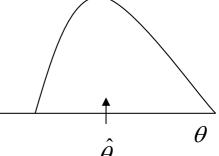
$$L(\theta) = \prod_{i=1}^{n} f(x_i, \theta)$$

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

Define ML estimator $\hat{\theta}$ as value of θ that maximizes $L(\theta)$. For m parameters, usually find solution $\hat{\theta}_1, \dots, \hat{\theta}_m$ by solving

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

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 \rightarrow take highest one.



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

N.B. Definition of ML estimators don't guarantee optimality

 \rightarrow investigate properties such as bias, variance ...

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



HELSINGFORS UNIVERSITET Example of ML estimator UNIVERSITY OF HELSINKI



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

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

Task to estimate value of τ . 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.

$$\begin{split} & \ln L(\tau) = \sum_{i=1}^n \ln f(t_i;\tau) = \sum_{i=1}^n \left(-\ln \tau - t_i/\tau\right) \\ & \text{set} \quad \frac{\partial \ln L}{\partial \tau} = 0 \quad \text{and solve for } \tau \quad \rightarrow \quad \hat{\tau} = \frac{1}{n} \sum_{i=1}^n t_i \end{split}$$

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

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

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

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

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

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 $\lambda = \ln 2$ / τ instead of mean lifetime τ . ML estimator for λ ?

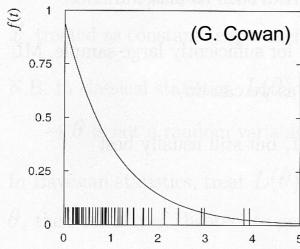


HELSINGFORS UNIVERSIT Example of ML estimators (cont.) UNIVERSITY OF HELSINKI



S 1 Addid May Approximated R

HELSINGIN YLIOPISTO



(G. Cowan) 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 θ , one has

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

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

 \rightarrow 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 λ ?

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

 $\rightarrow \hat{\lambda}$ 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 $N_{\rm taxi}$, ML estimator for $N_{\rm taxi}$ from m taxi number observations?

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

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



HELSINGFORS UNIVERSITE Example of ML estimators (cont.)





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

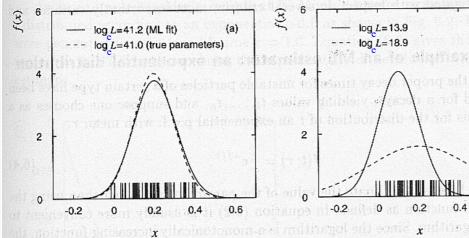
Already known that $\hat{\mu}$ is an unbiased estimator for μ .

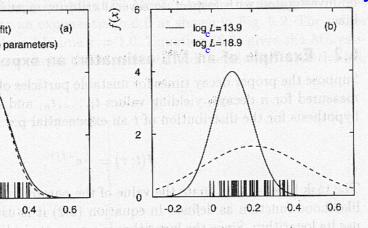
What about
$$\widehat{\sigma^2}$$
? $E[\widehat{\sigma^2}] = \frac{n-1}{n} \sigma^2$

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



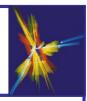


(G. Cowan)

- 50 observations of gaussian variable x; μ_x = 0.2 & σ_x = 0.1
- (a) pdf of parameters maximizing $\ln L$ & true parameters.
- (b) pdf of parameters far from true ones \rightarrow low In 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 $g(\hat{\tau}; \tau, n)$:

$$V[\hat{\tau}] = E[\hat{\tau}^2] - (E[\hat{\tau}])^2 = \int ... \int \left(\frac{1}{n} \sum_{i=1}^n t_i\right)^2 \frac{e^{-t_1/\tau}}{\tau} ... \frac{e^{-t_n/\tau}}{\tau} dt_1 ... dt_n - (E[\hat{\tau}])^2$$

$$= \frac{1}{n^2} \sum_{i=1}^{n} \sum_{j \neq i} \left(\int (t_i^2 + t_i t_j) \frac{e^{-t_i/\tau}}{\tau} dt_i \frac{e^{-t_j/\tau}}{\tau} dt_j \prod_{k \neq i,j} \int \frac{e^{-t_k/\tau}}{\tau} dt_k \right) - (E[\hat{\tau}])^2$$

$$= \frac{\tau^2(2n + n(n-1))}{n^2} - \tau^2 = \frac{\tau^2}{n} \rightarrow V[\hat{\tau}] \qquad n \text{ times smaller}$$
than $V[t]$

(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!) τ .

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. $\hat{\tau} \pm \hat{\sigma}_{\hat{\tau}} = 1.062 \pm 0.150$

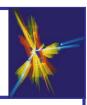
should be interpreted as: ML estimate for τ is 1.062.

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

If pdf $g(\hat{\tau}; \tau, n)$ Gaussian, then $[\hat{\tau} - \widehat{\sigma}_{\hat{\tau}}, \hat{\tau} + \widehat{\sigma}_{\hat{\tau}}]$ 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!!



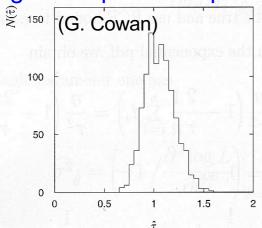


Variance of ML estimators: the Monte Carlo method

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 μ 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 $\hat{\tau}/\sqrt{n} = 1.062/\sqrt{50} = 0.150$

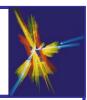
NB! $g(\hat{\tau}; \tau, n)$ approx. Gaussian (\Leftrightarrow central limit theorem) \rightarrow true in general for ML estimators in large sample limit. Variance of ML estimators: **the RCF bound**

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

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

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

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

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

same variance as obtained from analytical calculation \rightarrow ML $\hat{\tau}$ an efficient estimator for parameter τ for any n.

For $\theta = (\theta_1, ..., \theta_m)$ with efficient estimators and zero bias

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

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

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

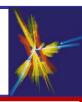
Variance of ML estimators: the graphical method

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 θ :

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

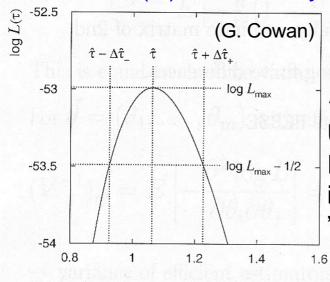
$$\operatorname{now} \ln L(\hat{\theta}) = \ln L_{\max} \quad \& \quad \left[\frac{\partial \ln L}{\partial \theta} \right]_{\theta = \hat{\theta}} = 0 \quad \Rightarrow \quad$$





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

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



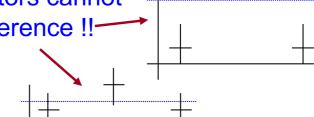
(G. Cowan) The exponential distribution example $\hat{ au}=1.062$

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

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:

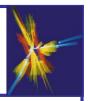
ML estimators cannot tell the difference!!



- consistent.
- invariant.
- biased for small n.
- not "right", just sensible.
- don't give "most likely value of θ " but value of θ 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 lnL maximisation for >1 free variables in practice impractical)
- the lnL_{max} value in itself contains no valuable information \Rightarrow won't indicate if choosen function for pdf correct or not.



Extended ML



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:
$$L(v, \overline{\theta}) = \frac{v^n}{n!} e^{-v} \prod_{i=1}^n f(x_i; \overline{\theta})$$

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

$$\ln L(\overline{\theta}) = n \ln v - v + \sum_{i=1}^{n} \ln f(x_i) = \sum_{i=1}^{n} \ln \left(v(\overline{\theta}) f(x_i; \overline{\theta}) \right) - v(\overline{\theta})$$

Now more information used \rightarrow smaller variances for $\bar{\theta}$ Example: particle scattering, expected number of events $v = \varepsilon \, \sigma \int \mathcal{L} \, dt$, where ε detection efficiency, σ scattering cross section ("probability" given by theory) & \mathcal{L} luminosity ("flux"). Use not only event describing variables but also σ . N.B.before `repetition of experiment' = same number of events. Here, meaning same v (e.g. same integrated \mathcal{L}) Suppose v & $\bar{\theta}$ are (functionally) independent:

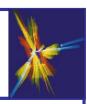
$$\frac{\partial \ln L}{\partial v} = 0 \rightarrow \hat{v} = n; \quad \frac{\partial \ln L}{\partial \theta_j} = 0 \rightarrow \text{ usual ML } \hat{\bar{\theta}}$$

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

With extended ML avoid all of that.



Extended ML (cont.)



Same with extended ML,

$$\ln L(\nu, \overline{\theta}) = \sum_{i=1}^{n} \ln \left(\sum_{j=1}^{m} \nu \theta_{j} f_{j}(x_{i}) \right) - \nu$$

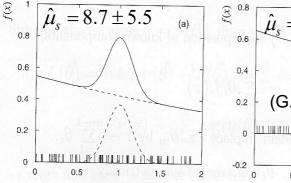
define $\mu_i = \nu \theta_i$ as expected number of events of type j,

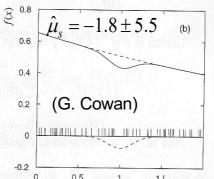
$$\ln L(\overline{\mu}) = \sum_{i=1}^{n} \ln \left(\sum_{j=1}^{m} \mu_j f_j(x_i) \right) - \sum_{j=1}^{m} \mu_j$$

Now all parameters treated symmetrically. Often μ_j 's more closely related to wanted final result e.g. production cross section for type j events. NB! ML fitted μ_j 's can be < 0 & in case μ_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 μ_s & μ_b

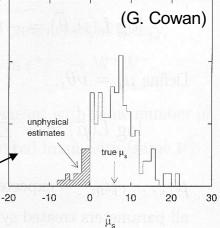
mostly 'works' well ... but sometimes $\hat{\mu}_s < 0$





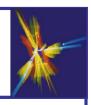
extended ML fit to 2 MC samples generated with μ_s = 6 & μ_b = 60. $f_s(x)$ Gaussian & $f_b(x)$ exponential.

Can report negative $\hat{\mu}_s$ (unphysical!) or take as estimator $\hat{\mu}_s^{\text{phys}} = \max(0, \hat{\mu}_s)^{20}$ (biased!). Can be a problem for large samples. Example: 200 extended ML¹⁰ fits on similar MC samples as above. $\hat{\mu}_s = 6.1 \pm 0.4$ (using all $\hat{\mu}_s$)



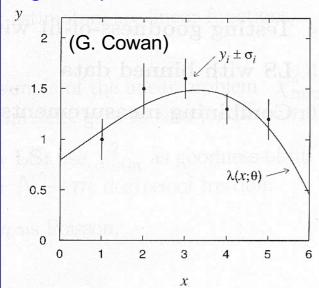


Method of least squares



The method of least squares (pienimmän neliösumman menetelmä)

N independent Gaussian random variables y_i , i = 1,..., 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 $y_1, ..., y_N$ measured with uncertainties $\sigma_1, ..., \sigma_N$ at x-values $x_1, ..., x_N$ (known without uncertainties). Each value λ_i of y_i given by function λ (x_i , θ). Goal is to minimize χ^2 sum by adjusting parameters θ (i.e. to find most optimal curve through points).

Joint pdf for N independent Gaussians y_i is product of N Gaussians: $N = 1 \quad (y_i = y_i^2)^2$

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

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

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

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

$$\chi^{2}(\overline{\theta}) = \sum_{i=1}^{N} (y_{i} - \lambda(x_{i}, \overline{\theta}))^{2} / \sigma_{i}^{2}$$

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



HELSINGFORS UNIVERSIFIE ar LS estimators & their variances



UNIVERSITY OF HELSINKI

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

$$\lambda(x; \overline{\theta})$$
 linear function of $\theta: \lambda(x; \overline{\theta}) = \sum_{j=1}^{m} a_j(x)\theta_j$,

where $a_i(x)$ are any linearly independent functions of x (i.e. one term can't be given as linear combination of the others) =

- 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:

$$(V^{-1})_{ij} = \frac{1}{2} \left[\frac{\partial^2 \chi^2}{\partial \theta_i \partial \theta_j} \right]_{\bar{\theta} = \bar{\hat{\theta}}}$$

 $(V^{-1})_{ij} = rac{1}{2} \left[rac{\partial^2 \chi^2}{\partial \theta_i \partial \theta_j}
ight]_{\bar{\theta} = \bar{\theta}}
ight]$ ightarrow coincides with KUF bound if λ linear function of $\theta \& V[y_i]$ known and Gaussian, then $L = -\chi^2/2$.

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

$$\chi^{2}(\overline{\theta}) = \chi^{2}(\overline{\hat{\theta}}) + \frac{1}{2} \sum_{i,j=1}^{m} \left[\frac{\partial^{2} \chi^{2}}{\partial \theta_{i} \partial \theta_{j}} \right]_{\overline{\theta} = \overline{\hat{\theta}}} (\theta_{i} - \hat{\theta}_{i}) (\theta_{j} - \hat{\theta}_{j})$$

(an expansion of the χ^2 function around its minima) → 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

$$\chi^2(\overline{\theta}) = \chi^2(\overline{\hat{\theta}}) + 1 = \chi^2_{\min} + 1$$

graphical method

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



LS fits of polynomials



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

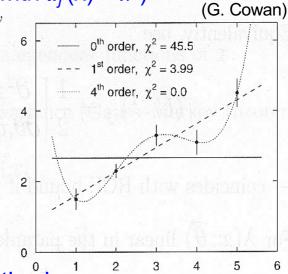
$$\lambda(x; \theta_0, ..., \theta_m) = \sum_{j=0}^{m} \theta_j x^j \quad (m+1 \text{ free parameters})$$

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

LS polynomial fit to 5 data points (with χ^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



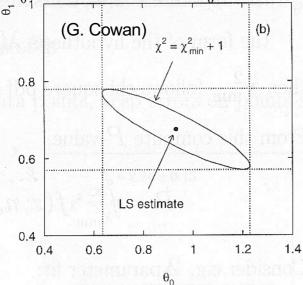
Uncertainties & covariances

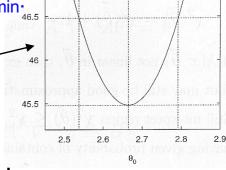
estimated using any of the 4 methods,

all related to χ^2 change when parameters changed from those giving χ^2_{\min}

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

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





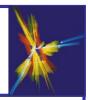
(G. Cowan)

1st order:

$$\begin{split} \hat{\theta}_0 &= 0.93 \pm 0.30, \\ \hat{\theta}_1 &= 0.68 \pm 0.10, \\ \chi^2_{\min} &= 3.99, \\ \cos[\hat{\theta}_0, \hat{\theta}_1] &= -0.028, \rho = \\ \cos[\hat{\theta}_0, \hat{\theta}_1] / \sigma_{\hat{\theta}_0} \sigma_{\hat{\theta}_1} = -0.93 \end{split}$$



Deriving LS estimators



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

$$\chi^{2}(\overline{\theta}) = \sum_{i=1}^{N} (y_{i} - \theta_{0} - \theta_{1}x_{i})^{2} / \sigma_{i}^{2}$$

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

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

$$\frac{\partial \chi^2}{\partial \theta_1} = 0 \to \sum_{i=1}^N x_i \left(y_i - \hat{\theta}_0 - \hat{\theta}_1 x_i \right) = \overline{xy} - \hat{\theta}_0 \overline{x} - \hat{\theta}_1 \overline{x^2} = 0$$

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

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

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

These are standard equations to get the slope θ_1 & the intercept with the *y*-axis θ_0 for a straight line fit. With different σ_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 χ^2_{\min} 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.

 χ^2_{\min} obtained from a LS fit follows the standard chisquare distribution with the degrees-of-freedom = N-m if:

- y_i , i = 1, ..., N Gaussian random variables with known covariance matrix V_{ij} or independent y_i 's with known σ_i 's.
- hypothesis $\lambda(x; \bar{\theta})$ is linear in parameters θ_i , i = 1, ..., 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)

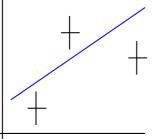
$$\chi_{\min}^2 = 3.99, N - m = 3 \rightarrow P = 0.263$$

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) χ^2_{\min} .

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

 $\chi^2_{\min}/n_{\rm d}$ »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!

 χ^2_{\min}/n_d «1 (or a large *P*-value):

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

ed

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

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).



HELSINGFORS UNIVERS Combining measurements with LS UNIVERSITY OF HELSINKI



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, ..., N ; $\sigma_i^2 = V[y_i]$, assumed to be known;

 λ = true value (takes role of θ , no x-dependence).

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

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

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

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

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. $cov[y_i,y_i]=V_{ii}$. Then have to minimize:

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

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

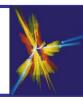
$$V[\hat{\lambda}] = \left(\frac{1}{2} \left[\frac{\partial \chi^2}{\partial^2 \lambda}\right]_{\lambda = \hat{\lambda}}\right)^{-1} = \sum_{i,j=1}^{N} w_i V_{ij} w_j = \mathbf{w}^T \mathbf{V} \mathbf{w}$$

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

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



HELSINGFORS UN COMBINING measurements with LS (cont.) UNIVERSITY OF HELSINKI



Averaging 2 correlated measurements: measurements y_1 & y_2 with covariance matrix \mathbf{V} :

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

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 2nd measurement:

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

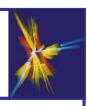
 $ho \le 1 \Longrightarrow 2^{\text{nd}}$ measurement only beneficial for average (i.e new combined variance \le). No variance change when $\rho = \sigma_1/\sigma_2$ (incl. $\rho = 1 \& \sigma_1 = \sigma_2$, i.e. same measurement twice).

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.

Overlaping 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 & determing estimators $y_1, ..., y_N$ for each subsample & from these correlations coefficient matrix.



HELSINGFORS L'S'FIT With penalty functions (nuisance terms) UNIVERSITY OF HELSINKI



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 χ^2 sum and look for the global χ^2_{min} that minimazes everything including the variation of inputs (with respect to their uncertainties) or the scalings.

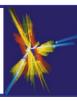
Example: top mass reconstruction at CDF experiment Proton-antiproton collisions create top-antitop pair:

$$\begin{array}{c} p\overline{p} \rightarrow t \ (\rightarrow W^+b)\overline{t} \ (\rightarrow W^-\overline{b}) \ X \\ \downarrow \geqslant q\overline{q} \qquad \qquad \downarrow \geqslant q\overline{q} \\ \hline q \qquad \qquad \qquad \downarrow \qquad \qquad q \\ \hline q \qquad \qquad \qquad \downarrow \qquad \qquad q \\ \hline \end{array}$$

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 (μ_W = 80.42 \pm 0.03 GeV for the W mass) with a width Γ_W = 2.1 GeV + combined with the third jet should make a top quark with same μ_{top} (= top mass) as the other top (= "triplet" of jets) within a width Γ_{top} = 1.6 GeV.



HELSINGFORS LISVIFIT With penalty functions (nuisance terms)



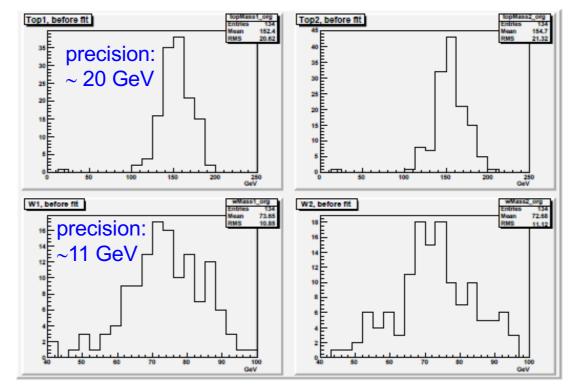
i) The masses are first determined by using the formulae given:

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

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})



ii) Now we test the hypothesis given

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

 Γ_{top} = 1.6 GeV, Γ_W = 2.1 GeV, m_{Wnom} = 80.42 GeV

 χ^2 large (> 100). Cure: allow energy measurements to vary

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

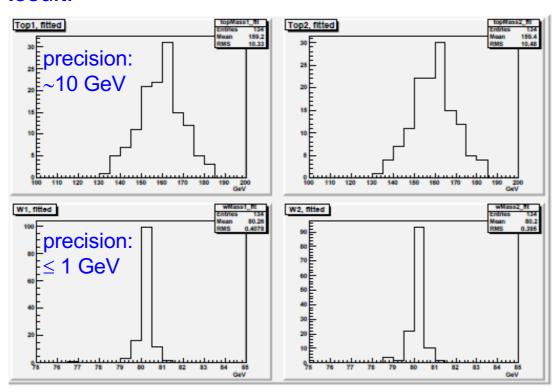
Let each c_i vary in order to minimaze whole χ^2_{new}



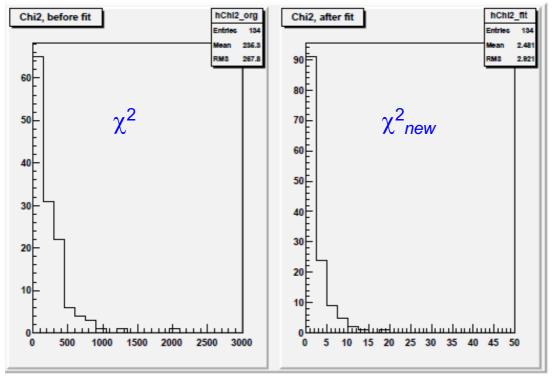
HELSINGFORS LISVIFIC With penalty functions (nuisance terms) UNIVERSITY OF HELSINKI



Result:



global χ^2 values significantly improved (by a factor ~100!)



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