

Likelihood fits – Advanced algorithms and tools

Miniworkshop on Statistical Tools

Volker Blobel – Universität Hamburg

1. Introduction
2. Likelihood functions for physics analysis
3. Optimization

1. Introduction

Likelihood fits in physics analysis: two main steps

- Construction of the Likelihood function (\Rightarrow **part 2**), and
- Minimization of the Likelihood function (Optimization) (\Rightarrow **part 3**)

The construction of the Likelihood function requires knowledge of the physics concepts and of details of the statistics of the measurement process.

- **Training + Optimized packages are required!**

Some recent books on Optimization:

*Jorge Nocedal and Stephen J. Wright, **Numerical Optimization**, Springer Series in Operations Research (1999)*

*J. Frederic Bonnans, J. Charles Gilbert, Claude Lemaréchal, Claudia A. Sagastizábal, **Numerical Optimization** – Theoretical and Practical Aspects, Springer (2006)*

*Per Christian Hansen, **Rank-Deficient and Discrete Ill-posed Problems**, Siam (1998)*

Optimization with bound constraints, equality constraints, and inequality constraints is treated.

...the world of HEP ...

Some concepts/notions/strategies, which are popular in HEP, are not mentioned in the text books:

- χ^2 minimization
- “...we are using an iterative method ...” (often extremely slow steepest descent method)
- determination of parameter dependence by grid interpolation
- simulated annealing, for continuous problems
- inequality constraints (bounds) treated by transformation (MINUIT)
- straight-line fits with MINUIT, higher-order polynomial fit without orthogonalization
- “fit uncertainties”, “fit errors”; individual χ^2 contributions to ..., unbiased residuals ...
- “inappropriateness of using $\Delta\chi^2 = 1$ (*for error estimation*) ... used $\Delta\chi^2 = 50$ for 90 % confidence level”
- proofs and error estimates by Monte Carlo simulation, experimental Mathematics
- unsmearing bin-by-bin correction of measured distributions
- “...solution of linear system of equations requires inversion ...”
- “...in 4200 unknowns: computational infeasible; even worse, non-linear fit won’t converge ...”.

2. Likelihood function for optimal physics analysis

The construction of ML functions \mathcal{L} is sometimes done in packages, but more often done individually for a certain analysis. An inappropriate construction may lead to a **bias** in the result.

Negative log-likelihood function $F(\mathbf{x})$, requiring minimization w.r.t. the unknowns \mathbf{x}

$$F(\mathbf{x}) = -\log \mathcal{L}(\mathbf{x}) \quad \left(\text{corresponding to } \frac{1}{2} \chi^2(\mathbf{x})\right)$$

allows combination of least squares and maximum-likelihood expressions.

Recommended principles:

- leave data y_j (e.g. counts) unchanged (no corrections to the data)
 - compare unchanged data to expectation t_j , including all corrections, like normalization uncertainty, background.
-
- Do not destroy ideal statistic, e.g. Poisson statistic for counts.
 - Do not introduce correlations between data points from systematic effects by data corrections.
 - If possible, use parametrization with (*almost*) orthogonal parameters: e.g. σ in $N/(\sqrt{2\pi}\sigma) \exp(-(\mu - y)^2/(2\sigma^2))$ instead of $N \exp(-(\mu - y)^2/(2\sigma^2))$.
 - Use correct statistical expressions for densities in maximum likelihood fits. Example: Normalization uncertainty \Rightarrow .

Normalization uncertainty

The normalization factor is represented by $\alpha = 1 \pm \varepsilon$, i.e. ε is the relative normalization error.

From a paper: “... the following definition for the effective χ^2 is adopted:

$$\chi_{\text{global}}^2 = \sum_n w_n \chi_n^2(x) \qquad \chi_n^2(x) = \left(\frac{\alpha_n - 1}{\varepsilon_n} \right)^2 + \sum_{\ell} \left(\frac{\boxed{\alpha_n y_{n\ell}} - t_{n\ell}(x)}{\sigma_{n\ell}^y} \right)^2$$

(n labels different experiments).

The normalization parameter α_n should be applied as a factor to the expectation $t_{n\ell}(x)$, instead of the data $y_{n\ell}$; otherwise a normalization bias is introduced.

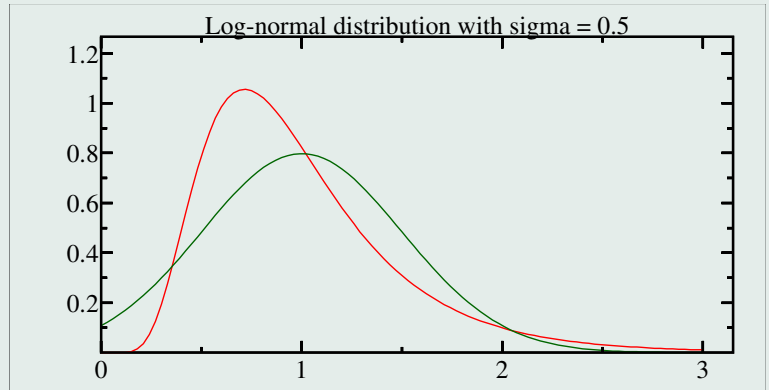
$$\alpha_n y_{n\ell} - t_{n\ell}(x) \quad \Rightarrow \quad y_{n\ell} - \alpha_n t_{n\ell}(x)$$

... and what about weight factor w_n ?

Normalization error = uncertainty in the factor, used to convert *event numbers* to *cross sections*.
Because of its origin – product of many factors, each with some uncertainty – the factor perhaps follows a **log-normal distribution** due the multiplicative central limit theorem.

For a log-normal distribution of a random variable α with $E[\alpha] = 1$ and standard deviation of ε the contribution to the ML-function is

$$\Delta F^{\text{norm}} = \frac{1}{2} \ln \alpha \left(3 + \frac{\ln \alpha}{\ln(1 + \varepsilon^2)} \right) \\ \rightarrow \frac{1}{2} \frac{(\alpha - 1)^2}{\varepsilon^2} \quad \text{for small } \varepsilon$$



The normal and the log-normal distribution, both with mean 1 and standard deviation $\varepsilon = 0.5$.

Log-normal density of a random factor $\alpha = 1 \pm \varepsilon$:

$$p(\alpha) = \exp \left[-\frac{1}{2} \ln \alpha \left(3 + \frac{\ln \alpha}{\ln(1 + \varepsilon^2)} \right) + \frac{\varepsilon^2}{8} - \ln \varepsilon - \ln \sqrt{2\pi} \right]$$

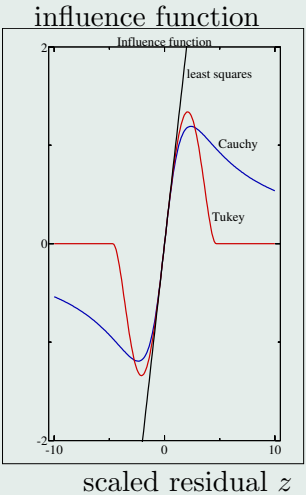
The objective function in least squares is the sum of **squares** of scaled residuals z , with **larger influence** for larger residuals (outliers) \Rightarrow the presence of outliers in the data will deteriorate the fit result.

The **square** is replaced in M-estimates by a dependence with **reduced influence** for larger residuals (non-gaussian maximum likelihood).

Least squares $F(\mathbf{x}) = \frac{1}{2} \sum_i z_i^2$ $z_i = \left(\frac{y_i - t_i}{\sigma_i} \right)$

e.g. Cauchy $\frac{1}{2} z_i^2 \rightarrow \frac{c^2}{2} \ln(1 + (z/c)^2)$

	$\rho(z) = -\ln \text{pdf}(z)$	influence function $\psi(z) = d\rho(z)/dz$	add. weight $\omega(z) = \psi(z)/z$
Least squares	$= \frac{1}{2} z^2$	$= z$	$= 1$
Cauchy($c = 2.3849$)	$= \frac{c^2}{2} \ln(1 + (z/c)^2)$	$= \frac{z}{1 + (z/c)^2}$	$= \frac{1}{1 + (z/c)^2}$
Huber $\begin{cases} \text{if } z \leq c = 1.345 \\ \text{if } z > c = 1.345 \end{cases}$	$= \begin{cases} z^2/2 \\ c(z - c/2) \end{cases}$	$= \begin{cases} z \\ c \cdot \text{sign}(z) \end{cases}$	$= \begin{cases} 1 \\ c/ z \end{cases}$



Nuisance parameters and profile likelihood

Statistical analysis of subset of parameters, e.g. one or two or ... parameters requires minimization w.r.t. to all other parameters (called nuisance parameters) of problem, i.e. eliminating the other parameters, but accounting for the extra uncertainty due to the *fitted* nuisance parameters.

One-parameter analysis: minimize $F(x_1, \mathbf{x}_r)$ for set of fixed values for x_1 w.r.t. all other parameters \mathbf{x}_r (option MINOS in MINUIT):

$$P(x_1) = \min_{\mathbf{x}_r} F(x_1, \mathbf{x}_r)$$

Two-parameter analysis: minimize $F(x_1, x_2, \mathbf{x}_r)$ for fixed pairs (x_1, x_2) w.r.t. all other parameters \mathbf{x}_r (option CONTUR in MINUIT):

$$P(x_1, x_2) = \min_{\mathbf{x}_r} F(x_1, x_2, \mathbf{x}_r)$$

Partitioning: Solution for subset \mathbf{x}_1 by partitioning of \mathbf{x} into $(\mathbf{x}_1, \mathbf{x}_2)$ and elimination of \mathbf{x}_2 \Rightarrow (next part).

3. Optimization

Standard NEWTON method

- Construct an **Objective Function** $F(\mathbf{x})$ from **statistical** considerations, depending on the parameter n -vector \mathbf{x} , to be minimized. $\mathbf{x} \in \mathcal{R}^n$ with start value \mathbf{x}_0 .

- **Step 1:** construct quadratic model $\tilde{F}_k(\mathbf{d}) = F_k + \mathbf{d}^T \nabla F_k + \frac{1}{2} \mathbf{d}^T \mathbf{C}_k \mathbf{d}$

function value: F_k gradient: ∇F_k sec. der. matrix (Hessian): $\mathbf{C}_k \approx \nabla^2 F_k$ at \mathbf{x}_k

- **Step 2:** determine minimum of quadratic model function: step vector \mathbf{d}_k

solve $\mathbf{C}_k \mathbf{d}_k = \mathbf{b}_k \equiv -\nabla F_k$ for \mathbf{d}_k and update: $\mathbf{x}_{k+1} = \mathbf{x}_k + 1 \cdot \mathbf{d}_k$

- repeat steps 1 and 2 (i.e. iterate) for non-linear problems (only single step for $\mathbf{C} = \text{constant}$)
- covariance matrix of parameter n -vector \mathbf{x} is given by inverse matrix: $\mathbf{B} = \mathbf{C}^{-1}$
(i.e. **error propagation**) (equivalent to “ $\Delta\chi^2 = 1$ ”)

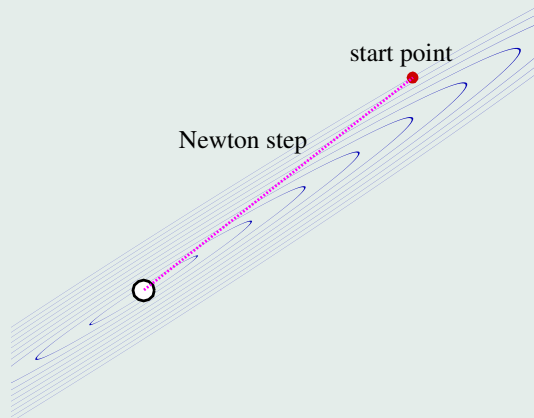
Improved convergence by

line-search $\min \varphi(\alpha) = F(\mathbf{x}_k + \alpha \cdot \mathbf{d}_k)$ w.r.t. factor α

or trust-region $\min \tilde{F}_k(\mathbf{d}) = F_k + \mathbf{d}^T \nabla F_k + \frac{1}{2} \mathbf{d}^T \mathbf{C}_k \mathbf{d}$ w.r.t. \mathbf{d} with constraint $\|\mathbf{d}\| \leq \Delta_k$

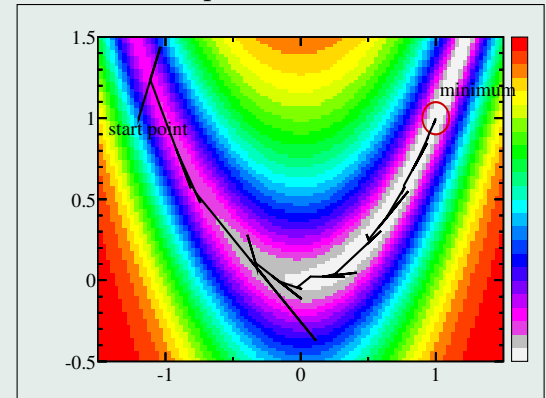
“... extremely important, and might supersede line-searches, sooner or later.”

... large correlation



Quadratic function is minimized by a single NEWTON step; steepest-descent would be inefficient due to high correlation.

... non-quadratic behaviour



Non-quadratic function: minimization requires to follow very narrow curved valley.

- NEWTON method based on second-derivative matrix has quadratic convergence rate.
- Simple methods like steepest-descent with only linear convergence rate are slow; convergence may not occur as the iteration stagnates (can be misinterpreted as indication for convergence).

Derivative calculation

Gauss-Newton methods, to calculate gradient and Hessian approximation of the log ML-function $F(\mathbf{x})$ from *first* derivatives of the parametrization function f :
(only 6 statements to accumulate \mathbf{b} and \mathbf{C})

- Least squares: fit of function f to data \mathbf{y} (Gauss-Newton, f'' terms ignored)

$$\mathbf{b}_j = \sum_i w_i \left(\frac{\partial f_i}{\partial x_j} \right) (y_i - f_i) \quad \mathbf{C}_{jk} \approx \sum_i w_i \left(\frac{\partial f_i}{\partial x_j} \right) \left(\frac{\partial f_i}{\partial x_k} \right)$$

- ML-estimation (Poisson): fit of (bin-integrated) function f to counts y_i (f'' terms ignored)

$$\mathbf{b}_j = \sum_i \left(\frac{\partial f_i}{\partial x_j} \right) \frac{y_i - f_i}{f_i} \quad \mathbf{C}_{jk} \equiv (\nabla^2 F_{jk}) \approx \sum_i \left(\frac{y_i}{f_i^2} \right) \left(\frac{\partial f_i}{\partial x_j} \right) \left(\frac{\partial f_i}{\partial x_k} \right)$$

numerical differentiation: by finite-differences in a number of repeated function evaluations;

first der.: $2n$ second der.: $n(n-1)/2 \propto n^2$ function evaluations

automatic differentiation: takes as input a computer code able to calculate a function value –
produces as output another computer code able to calculate the derivatives;

symbolic differentiation: algebraic specification of the function is manipulated to produce new
algebraic expressions for the derivatives.

Solution by partitioning

Interest in subset \mathbf{x}_1 of total parameter vector $(\mathbf{x}_1, \mathbf{x}_2)$

Partitioning the matrix equation $\mathbf{C}\mathbf{x} = \mathbf{b}$ (\mathbf{C} symmetric, \mathbf{C}_{11} and \mathbf{C}_{22} are square matrices):

$$\left(\begin{array}{c|c} \mathbf{C}_{11} & \mathbf{C}_{12} \\ \hline \mathbf{C}_{12}^T & \mathbf{C}_{22} \end{array} \right) \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{pmatrix}$$

with $\mathbf{C}_{12} \neq 0$. Solve matrix equation for \mathbf{x}_1 , e.g. by inversion:

$$\begin{pmatrix} \mathbf{C}_{11} - \mathbf{C}_{12}\mathbf{C}_{22}^{-1}\mathbf{C}_{12}^T \end{pmatrix} \begin{pmatrix} \mathbf{x}_1 \end{pmatrix} = \begin{pmatrix} \mathbf{b}_1 - \mathbf{C}_{12}[\mathbf{C}_{22}^{-1} \mathbf{b}_2] \end{pmatrix}$$

covariance matrix $\mathbf{V}_{11}(\mathbf{x}_1) = \begin{pmatrix} \mathbf{C}_{11} - \mathbf{C}_{12}\mathbf{C}_{22}^{-1}\mathbf{C}_{12}^T \end{pmatrix}^{-1}$

Note: “local solution” $\mathbf{x}_2^* \equiv [\mathbf{C}_{22}^{-1} \mathbf{b}_2]$

Applications:

- Track detector alignment (Millepede): simultaneous fit of alignment parameters (\mathbf{x}_1) and an unlimited number of parameters of tracks (\mathbf{x}_2 , many).
- Barlow fraction fit of histograms, using MC histograms with limited statistics.

MINUIT: “used by people in many fields, many applications; lifetime long compared with computers, accelerators” [**MINUIT** Upgrade Project 2002]

From the **MINUIT** User’s Guide (2004):

“...suited to handle difficult problems, including those which may require guidance in order to find the correct solution.”

“...not intended for the repeated solution of identical parametrized problems ...”

“...no limit on the number of parameters ...however the ‘technological’ limitations of MINUIT can be seen around a maximum of 15 free parameters ...”

number of parameters	small number of cases	large number of cases
small number	MINUIT	—
derivatives available	[MINUIT]	L-BFGS , N
constraints required	Newton + Lagrange	
large number	L-BFGS	—
ill-conditioned	Regularization	

Additional methods needed for special problems:
large number of parameters, large number of cases,
constraints, ill-conditioned problems.

Symmetric matrices

Space $\propto n^2$ and cpu-time $\propto n^3$ for standard (in-place) solution of a matrix equation of dimension n :

dimension n	space		cpu-time	
10	0.8	kB	10^{-5}	sec
1 000	8	MB	10	sec
100 000	80	GB	120	days

Method acceptable for n up to 1000 (or even 10 000 – few hours cpu-time), if matrix condition is sufficient, and robust algorithm is used.

Sparse matrix: large matrices are often *sparse* – can be stored in $\ll n^2$ locations, but inverse of sparse matrix is full matrix!

Generalized residual minimization: (related to method of conjugate gradients);

system of linear equations $\mathbf{C} \mathbf{x} = \mathbf{y}$ equivalent to $\min \|\mathbf{C} \mathbf{x} - \mathbf{y}\|_2$

\mathbf{C} may be indefinite, very large and sparse. Needs only product of sparse matrix \mathbf{C} \times vector, e.g. MINRES*. **Example:** $n = 50\,558$ with cpu-time of < 3 minutes instead of 15 days.

Band matrix, arrow matrix, skyline matrix: cpu-time only $\propto n^1$ by Cholesky method with decomposition $\mathbf{C} = \mathbf{L} \mathbf{D} \mathbf{L}^T$, forward $\mathbf{L} \mathbf{z} = \mathbf{y}$ and backward solution $\mathbf{L}^T \mathbf{x} = \mathbf{D}^{-1} \mathbf{z}$. Even elements of inverse matrix within band can be calculated fast.

Specialized parametrization of problem: \Rightarrow e.g. band matrix msc track fit by broken lines (large nr of parameters) **10 \times faster** than Kalman.

*) C. C. Paige and M. A. Saunders (1975), Solution of sparse indefinite systems of linear equations, SIAM J. Numer. Anal. 12(4), pp. 617-629.

Invented in the mid 1950's (physicist W.C. Davidon – unpublished)

Start with simple (diagonal) matrix \mathbf{B}_0 (\mathbf{C}_0^{-1}), calculate gradient compute step \mathbf{d} by product

$$\mathbf{d} = -\mathbf{B}_k \times \nabla F_k$$
$$\text{line-search} \quad \min_{\alpha} \varphi(\alpha) \equiv F(\mathbf{x}_k + \alpha \mathbf{d}) \quad \Rightarrow \quad \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha \mathbf{d},$$

calculate differences $\mathbf{s}_k = \mathbf{x}_{k+1} - \mathbf{x}_k$ $\mathbf{y}_k = \nabla F_{k+1} - \nabla F_k$, and update matrix \mathbf{B}_k such that $\mathbf{B}_{k+1} \mathbf{y}_k = \mathbf{s}_k$ (BFGS updating formula, 1970)

$$\mathbf{B}_{k+1} = (\mathbf{I} - \rho_k \mathbf{s}_k \mathbf{y}_k^T) \mathbf{B}_k (\mathbf{I} - \rho_k \mathbf{y}_k \mathbf{s}_k^T) + \rho_k \mathbf{s}_k \mathbf{s}_k^T \quad \text{with } \rho_k = \frac{1}{\mathbf{y}_k^T \mathbf{s}_k}$$

\Rightarrow correct inverse Hessian \mathbf{B} reached after n (linear independent) steps for quadratic problem - but **may diverge** for small $\mathbf{y}_k^T \mathbf{s}_k$!

Requires modern line-search algorithm with sufficient decrease in function value and slope (strong Wolfe conditions) (J.J. Moré and D.J. Thuente (1994)) to have self-correcting properties. (1976)

\Rightarrow a minimization algorithm, that needs only the gradient (no second derivatives).

Limited memory BFGS (L-BFGS)

BFGS algorithm requires matrix with $\propto n^2$ locations:

What to do, if the number of parameters is 100 000 or larger?

Note: the inverse Hessian \mathbf{B}_k is used only in the product $\mathbf{d} = -\mathbf{B}_k \times \nabla F_k$

BFGS: \mathbf{B}_k determined by \mathbf{B}_0 and k pairs $\{\mathbf{s}_0, \mathbf{y}_0\} \dots \{\mathbf{s}_{k-1}, \mathbf{y}_{k-1}\}$.

In J. Nocedal's "Limited memory" algorithm (1980) the vector \mathbf{d} is calculated from

L-BFGS: diagonal \mathbf{B}_0 and the last m pairs $\{\mathbf{s}_{k-m}, \mathbf{y}_{k-m}\} \dots \{\mathbf{s}_{k-1}, \mathbf{y}_{k-1}\}$.

with fixed $m = 7 \dots 29$, i.e. space $\propto n^1$, no matrix storage required! Algorithm often more efficient than original BFGS algorithm, and requires e.g. only 60 Mbyte for 100 000 parameters.

"... algorithm is excellent, it is at present the best choice, often the only possible one, for large-scale problems ..."

Comparison LVMINI/MINUIT

Minimization package LVMINI, using L-BFGS, developed for up to several 100 000 parameters, needs gradient ∇F ; error calculation included.

	Straight-line least squares fit			Rosenbrock fct	
	LVMINI	MINUIT	MINUIT-MINOS	LVMINI	MINUIT
no gradient	-	131	392	-	286
gradient	12 + 4	70	159	71	144
+ diagonal	10 + 4				

2-parameter fits;

in some straight-line fits MINUIT parabolic error and correlation wrong.

	LVMINI	MINUIT	MINUIT-MINOS
no gradient	-	170	638
gradient	14 + 8	118	289
+ diagonal	11 + 8		

4-parameter fits:

Gaussian histogram peak + background

Initial values		Number of function evaluation		
		LVMINI	MINUIT	MINUIT-MINOS
good	no gradient		2074	22 417
	gradient	558 + 40	–	–
	+ diagonal	29 + 40		
bad	no gradient		[2865]	[17 199]
	gradient	601 + 40	–	–
	+ diagonal	104 + 40		

20-parameter fits:

five peaks, described by Students distribution, plus linear background.

Note: errors calculated by LVMINI are as accurate as MINOS errors in this test. Number in [] = minimum not found.

Ill-posed least squares: $\mathbf{Ax} \cong \mathbf{y}$ (with weakly defined degrees of freedom)

Idea of regularization: control the norm of the residuals and the norm of the solution vector \mathbf{x} , simultaneously, by adding regularizing term with regularization parameter $\tau > 0$:

Thikhonov-Phillips: $F_\tau(\mathbf{x}) = \|\mathbf{Ax} - \mathbf{y}\|^2 + \tau \|\mathbf{Lx}\|^2$ with $\mathbf{L} = \mathbf{I}$ or second-derivative term

Matrix equation: $(\mathbf{C} + \tau \mathbf{L}^T \mathbf{L}) \mathbf{x} \equiv (\mathbf{A}^T \mathbf{V}_y^{-1} \mathbf{A} + \tau \mathbf{L}^T \mathbf{L}) \mathbf{x} = \mathbf{A}^T \mathbf{V}_y^{-1} \mathbf{y} \equiv \mathbf{b}$ to be solved

Solution with orthogonalization: eigenvalues λ_j and eigenvectors \mathbf{u}_j

Fourier coefficients $c_j = \frac{1}{\sqrt{\lambda_j}} (\mathbf{b}^T \mathbf{u}_j)$ independent, with error = 1

Solution $\mathbf{x}_\tau = \sum_{j=1}^n f_j \frac{1}{\sqrt{\lambda_j}} c_j \mathbf{u}_j$ with filter factor $f_j = \left(\frac{\lambda_j}{\lambda_j + \tau} \right)$

Regularization parameter τ appears only in filterfactor: scan over a large τ range with selection of τ , that minimizes the **global correlation** between the parameters.

Summary

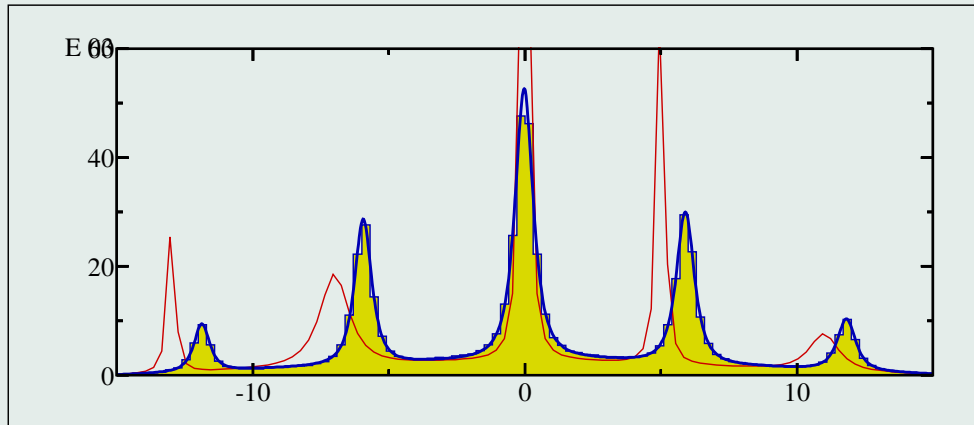
Additional methods needed for different problems: for example

- large number of parameters up to 10^5 ;
- large number of cases;
- solution of ill-posed problems (e.g. unfolding).

New methods developed during the last 40 years: for example

- BFGS with modern line-search (Wolfe conditions);
- limited memory BFGS for large-scale optimization;
- regularization methods with stable solutions

Initial values		Number of function evaluation		
		LVMINI	MINUIT	MINUIT-MINOS
good	no gradient		2074	22 417
	gradient	558 + 40	—	—
	+ diagonal	29 + 40		
bad	no gradient		[2865]	[17 199]
	gradient	601 + 40	—	—
	+ diagonal	104 + 40		



Five peaks, described by Student's distribution, plus linear background. Initial parameter values (bad) corresponds to red line; blue line is fit result. Number in [] = minimum not found.

Contents

1. Introduction	2
... the world of HEP ...	3
2. Likelihood function for optimal physics analysis	4
Normalization uncertainty	5
The log-normal distribution	6
Outlier	7
Nuisance parameters and profile likelihood	8
3. Optimization	9
Minimization	10
Derivative calculation	11
Solution by partitioning	12
Minimization problems	13
Symmetric matrices	14
Quasi-NEWTON algorithm	15
Limited memory BFGS (L-BFGS)	16
Comparison LVMINI/MINUIT	17
Regularization methods	18
Summary	19
Five-peak fit	20
Table of contents	21