

PHYSTAT2011 – Workshop on Unfolding

Thursday Jan 20th 2011

Volker Blobel – Universität Hamburg

Unfolding methods for particle physics

1. Unfolding – direct and inverse process
 2. Naive unfolding
 3. Convolution/deconvolution
 4. Orthogonalization
 5. Truncation and correlations
 6. Regularization
 7. DCT and projection methods
 8. Iterative unfolding
- + Appendix

*Unfolding is a complex **statistical** and mathematical problem.*

1. Unfolding – direct and inverse processes

The process of the transition between the **true distribution** $f(t)$ and the **measured distribution** $g(s)$ for linear inverse problems is described by the Fredholm integral equation of the first kind:

$$\int_{\Omega} K(s, t) f(t) dt + b(s) = g(s)$$

($b(s)$ = background contribution). Two types of processes are based on the integral equation:

| | |
|------------------------------------|---|
| direct process (MC) | true/MC dist. $f(t) \implies g(s)$ measured dist. |
| inverse process (unfolding) | measured dist. $g(s) \implies f(t)$ true dist. |

Discretization: the integral equation becomes an (usually ill-posed) linear system of equations:

$$\mathbf{Ax} = \mathbf{y} \qquad y_i = \int_{s_{i-1}}^{s_i} g(s) ds \quad i = 1, 2, \dots, m$$

(assuming a case without background contribution) with the representation

| | |
|---|--|
| true distribution $f(t) \Rightarrow \mathbf{x}$ | n -vector of unknowns |
| measured distribution $g(s) \Rightarrow \mathbf{y}$ | m -vector of measured data |
| Kernel $K(s, t) \Rightarrow \mathbf{A}$ | rectangular m -by- n response matrix . |

The variables s , t and vectors \mathbf{x} , \mathbf{y} can be multi-dimensional. Elements of the response matrix \mathbf{A} are (positive) probabilities, and include efficiency.

Discretization

- The response matrix \mathbf{A} has to be determined in particle physics by MC ;
- several different methods can be applied in the discretization: simple binning, quadrature methods (weighted sums), B-splines (allows re-weighting of MC) . . . ;
- unfolding is independent of the assumed $f(t)^{MC}$, if the Fredholm equation above is correct,
- but for complex physical measurements this assumption may not be true, and the response function is influenced by the assumed $f(t)^{MC}$. In this case the Fredholm integral equation should be rewritten in the form,

$$\int_{\Omega} K(s, t; f) f(t) dt + b(s) = g(s) ,$$

which represents a more difficult *nonlinear* inverse problem; $f(t)^{MC}$ should be close to the expected result.

Problems in Particle Physics differ from problems in other fields:

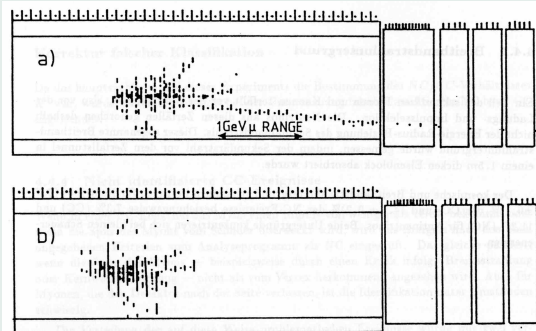
- input errors are well-known (Poisson data, . . . covariance matrix) \mathbf{V}_y ;
- covariance matrix of result is required, no bias, small correlations;
- dimension parameters are small compared to other fields (with e.g. 10^6 parameters)

Unfolding is more general than “data correction”

Neutral current neutrino interactions (narrow band beam): CHARM collaboration

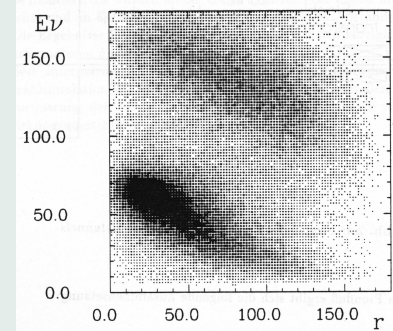
Measurement of $E_{\text{had}}, \vartheta_{\text{had}}, r_{\text{beam}}$

CC: $\nu N \rightarrow [\text{hadrons}] \mu$
 $E_{\nu, \text{in}} = E_{\text{had}} + E_{\mu}$



NC: $\nu N \rightarrow [\text{hadrons}] \nu$
 $E_{\nu, \text{in}} = E_{\text{had}} + \text{???}$

beam flux $\phi(E_{\nu}, r_{\text{beam}})$



“Indirect” determination of $d\sigma/dy$ and $d\sigma/dx$:

direct problem (by MC)

$d\sigma/dx; d\sigma/dy \Rightarrow$ predicted dist. in $E_{\text{had}}, \vartheta_{\text{had}}, r_{\text{beam}}$

inverse problem (by unfolding)

$d\sigma/dx; d\sigma/dy \Leftarrow$ measured dist. in $E_{\text{had}}, \vartheta_{\text{had}}, r_{\text{beam}}$

Unfolding is more general than “data correction” to correct migration effects.

2. Naive unfolding

The data errors are represented by a m -vector \mathbf{e} , and the actually measured distribution \mathbf{y} is given by

$$\text{measured distribution } \mathbf{y} = \mathbf{y}_{\text{exact}} + \mathbf{e} = \mathbf{A} \mathbf{x}_{\text{exact}} + \mathbf{e} \quad \mathbf{e} = \text{data errors}$$

$$\text{Unfolding: } \mathbf{y}; \mathbf{V}_y \xrightarrow{\mathbf{A}} \mathbf{x}; \mathbf{V}_x \quad \min_{\mathbf{x}} \{ \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2 \}$$

The n -by- m pseudoinverse \mathbf{A}^+ is a generalization of the inverse matrix (also called Moore-Penrose *generalized inverse*); it satisfies the relation $\mathbf{A}^+ \mathbf{A} = \mathbf{I}$, and allows the least squares solution by

$$\begin{aligned} \hat{\mathbf{x}} &= \mathbf{A}^+ \mathbf{y} & \mathbf{A}^+ &= (\mathbf{A}^T \mathbf{V}_y^{-1} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{V}_y^{-1} \\ \mathbf{V}_x &= \mathbf{A}^+ \mathbf{V}_y \mathbf{A}^{+T} = (\mathbf{A}^T \mathbf{V}_y^{-1} \mathbf{A})^{-1} \end{aligned}$$

The response matrix \mathbf{A} and the pseudoinverse \mathbf{A}^+ do not depend on any assumption about \mathbf{x} .

$$\begin{aligned} \text{estimate } \hat{\mathbf{x}} &= \mathbf{A}^+ \mathbf{y} = \mathbf{A}^+ \mathbf{y}_{\text{exact}} + \mathbf{A}^+ \mathbf{e} = \mathbf{A}^+ \mathbf{A} \mathbf{x}_{\text{exact}} + \mathbf{A}^+ \mathbf{e} \\ &= \mathbf{x}_{\text{exact}} + \underbrace{(\mathbf{A}^+ \mathbf{A} - \mathbf{I}) \mathbf{x}_{\text{exact}}}_{\text{systematic error}} + \underbrace{\mathbf{A}^+ \mathbf{e}}_{\text{statistical error}} \end{aligned}$$

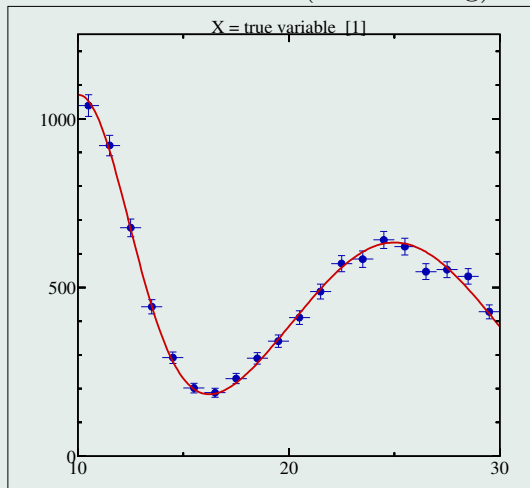
The systematic error (**bias**) $(\mathbf{A}^+ \mathbf{A} - \mathbf{I}) \mathbf{x}_{\text{exact}}$ depends in $\mathbf{x}_{\text{exact}}$!

$$\text{Least squares } \mathbf{A}^+ \mathbf{A} = (\mathbf{A}^T \mathbf{V}_y^{-1} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{V}_y^{-1} \mathbf{A} = \mathbf{I} \quad \text{no bias!} \quad \text{Good or bad?}$$

Naive result with narrow bins

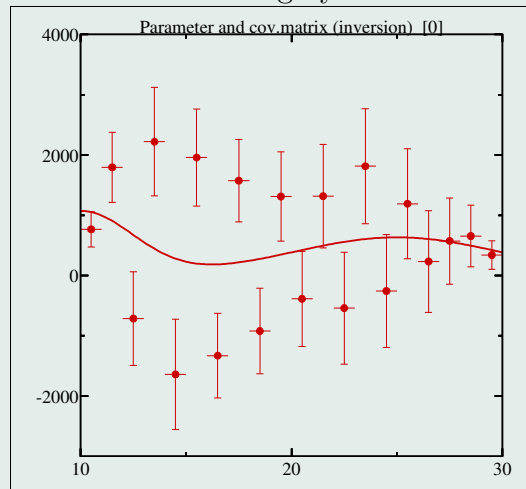
Example of unfolding problem with narrow bins

Perfect resolution (no smearing)



Histogram
for sample with 10 000 entries.

Naive unfolding by inversion



Huge fluctuations, due to large negative correlations: neighbour bin -95% (second $+85\%$).

True curve $f(x)$ is shown in red.

Orthogonalization!

3. Convolution/deconvolution with Gaussian response

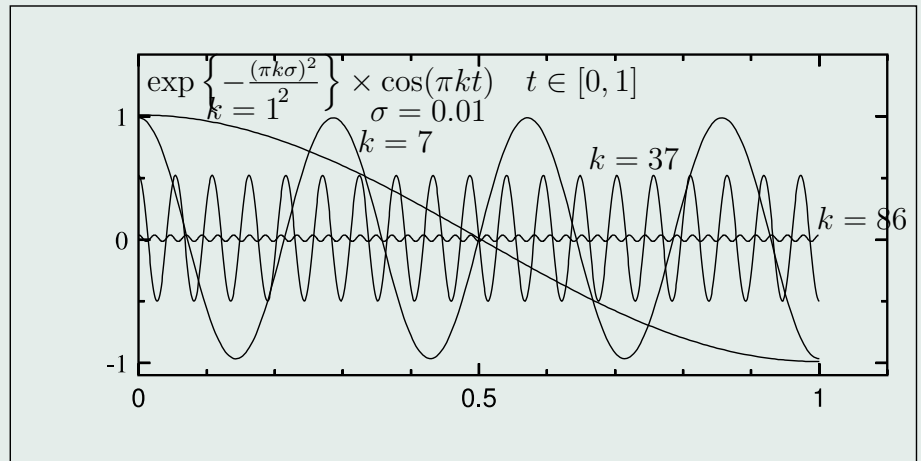
An even function $f(t)$ with period 1 can be approximated by a sum with n terms

$$f(t) \approx a_0 + \sum_{k=1}^{n-1} a_k \cos(\pi kt) \quad g(s) \approx \alpha_0 + \sum_{k=1}^{n-1} \alpha_k \cos(\pi ks)$$

Convolution of functions $\cos(\pi kt)$ by a kernel function $K(s, t) \equiv K(s - t)$ given by a Gaussian resolution function (standard deviation σ):

$$\int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(s-t)^2}{2\sigma^2}\right) \times \cos(\pi kt) dt = \exp\left(-\frac{(\pi k\sigma)^2}{2}\right) \times \cos(\pi ks),$$

The amplitude is attenuated by an **exponential factor**, which will become $\ll 1$ for larger values of k .



...and deconvolution

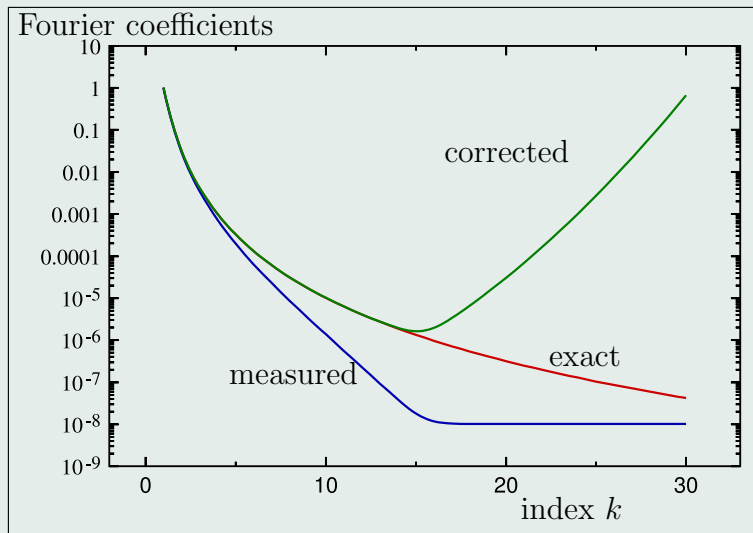
Deconvolution: expand convoluted (“measured”) function $g(s)$ to obtain coefficients α_k :

$$\begin{aligned} \text{Correct coefficients back: } \hat{a}_k &= \exp\left(\frac{\pi^2 k^2 \sigma^2}{2}\right) \times \alpha_k && \text{for small } k \\ &= \exp\left(\frac{\pi^2 k^2 \sigma^2}{2}\right) \times \sqrt{\alpha_k^2 + \epsilon^2} \xrightarrow{k \rightarrow \infty} \epsilon \times \exp\left(\frac{\pi^2 k^2 \sigma^2}{2}\right) \end{aligned}$$

Fourier coefficients α_k are below the exact ones a_k and reach for $k = 15$ a level of about 10^{-8} due to **round-of-errors** ϵ . The deconvoluted Fourier coefficients, labelled *corrected*, are correct only up to $k = 15$.

Deconvolution with $k \gg 15$ will have result dominated by noise.

For bin width $w = \sigma$ the factor is > 100 . A factor below 10 is reached for $w > 1.5\sigma$.



4. Orthogonalization

Singular value decomposition, applied to rectangular m -by- n matrix \mathbf{A} :

$(m \geq n)$

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \sum_{i=1}^n \sigma_i \mathbf{u}_i \mathbf{v}_i^T \quad \mathbf{U}^T \mathbf{U} = \mathbf{V}^T \mathbf{V} = \mathbf{V}\mathbf{V}^T = \mathbf{I}$$

$$\begin{pmatrix} & \text{\scriptsize (n)} \\ \text{\scriptsize (m)} & \mathbf{A} \end{pmatrix} = \begin{pmatrix} & \text{\scriptsize (n)} \\ \text{\scriptsize (m)} & \mathbf{U} \end{pmatrix} \cdot \begin{pmatrix} \text{\scriptsize (n)} \\ \mathbf{\Sigma} \end{pmatrix} \begin{pmatrix} \text{\scriptsize (n)} \\ \mathbf{V}^T \end{pmatrix}$$

Matrix $\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_n)$ with ordered singular values $\sigma_1 \geq \sigma_2 \geq \dots \sigma_n \geq 0$, assuming pre-whitening of eq. $\mathbf{A}\mathbf{x} \simeq \mathbf{y}$, i.e. $\mathbf{V}_y = \mathbf{I}$.

Diagonalization of symmetric LS matrix (or Hessian of log-Likelihood functions):

$$\mathbf{C} = \mathbf{A}^T \mathbf{A} = (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^T \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \mathbf{V}\mathbf{\Sigma}^2 \mathbf{V}^T = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$$

Eigenvalues λ_j of symmetric matrix \mathbf{C} are equal to squared singular values σ_j , and eigenvectors are equal to the singular vectors of matrix \mathbf{V} .

Least squares solution

SVD ... “a new way to see into the heart of a matrix” (Gilbert Strang)

$$\text{Matrix product using SVD: } \mathbf{Ax} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \mathbf{x} = \sum_{j=1}^n \sigma_j (\mathbf{v}_j^T \mathbf{x}) \mathbf{u}_j = \mathbf{y}$$

Fourier coefficients $\mathbf{c} = \mathbf{U}^T \mathbf{y}$ with $\mathbf{V}_c = \mathbf{I}$ represent measurement \mathbf{y} .

Least squares solution with SVD:

$$\begin{aligned} \mathbf{U}^T \cdot | \quad \mathbf{Ax} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \hat{\mathbf{x}} = \mathbf{y} \\ \mathbf{\Sigma}\mathbf{V}^T \hat{\mathbf{x}} = \mathbf{U}^T \mathbf{y} = \mathbf{c} \quad \text{Fourier coefficients } c_j = \mathbf{u}_j^T \mathbf{y} \pm 1 \\ \hat{\mathbf{x}} = \mathbf{V}\mathbf{\Sigma}^{-1} \mathbf{c} \\ \mathbf{V}_x = \mathbf{V}\mathbf{\Sigma}^{-1} \mathbf{\Sigma}^{-1} \mathbf{V}^T \quad \text{error propagation} \end{aligned}$$

$$\hat{\mathbf{x}} = \sum_{j=1}^n \left(\frac{1}{\sigma_j} \right) c_j \mathbf{v}_j \quad \mathbf{V}_x = \sum_{j=1}^n \left(\frac{1}{\sigma_j^2} \right) \mathbf{v}_j \mathbf{v}_j^T$$

Problem with zero or very small singular values!

Vanishing singular values and truncation

Assumption: p non-zero singular values of total n values with contribution $d_j = c_j/\sigma_j$.

$$\hat{\mathbf{x}} = \underbrace{\sum_{j=1}^p d_j \mathbf{v}_j}_{\mathbf{x}_{\text{range}} \in \mathbb{R}^p} + \underbrace{\sum_{j=p+1}^n \tilde{d}_j \mathbf{v}_j}_{\mathbf{x}_{\text{null}} \in \mathbb{R}^{n-p}} \quad \mathbf{A}\hat{\mathbf{x}} = \sum_{j=1}^p \sigma_j d_j \mathbf{v}_j + \underbrace{\sum_{j=p+1}^n \sigma_j \tilde{d}_j \mathbf{v}_j}_{=0} = \mathbf{y}$$

$(n - p)$ contributions \tilde{d}_j are arbitrary: $\mathbf{x}_{\text{null}} \in \mathbb{R}^{n-p}$ without influence on measured distribution \mathbf{y}

Alternatives for unfolding solution:

- \mathbf{x}_{null} = plausible contribution; in iterative methods the initial contributions $\mathbf{x}_{\text{null}} \in \mathbb{R}^{n-p}$ remain unchanged;
- minimum norm solution: $\mathbf{x}_{\text{null}} = 0$ with $n > p$;

$$\hat{\mathbf{x}} = \mathbf{x}_{\text{range}} + \mathbf{x}_{\text{null}} \quad \|\mathbf{x}_{\text{range}} + \mathbf{x}_{\text{null}}\|^2 = \|\mathbf{x}_{\text{range}}\|^2 + \|\mathbf{x}_{\text{null}}\|^2$$

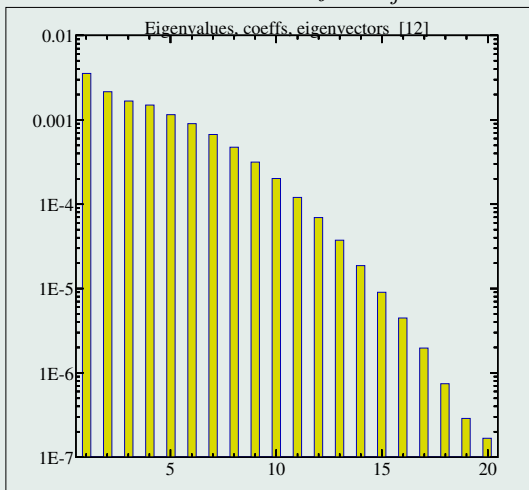
- reduction of dimension: $n' = p$ with full-rank covariance matrix.

What is the “best” approach?

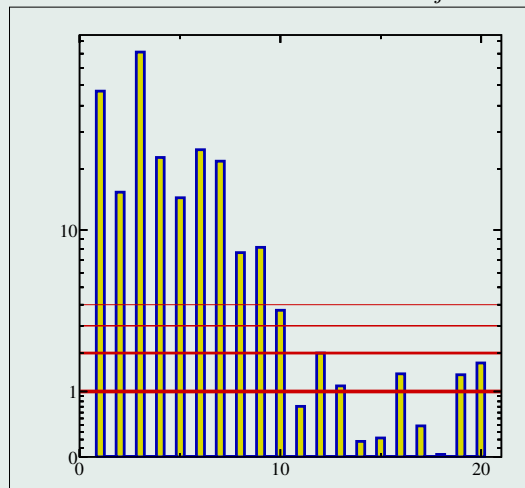
Eigenvalues and Fourier coefficients

Example of unfolding problem with narrow bins

Eigenvalues $\lambda_j = \sigma_j^2$



Normalized coefficients c_j



Eigenvalues decrease by 4 orders of magnitude.

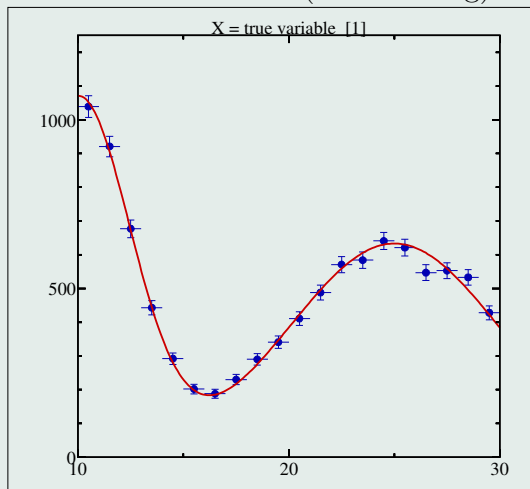
Only 10 of the 20 coefficients are significant.

Red lines are for 1, 2, 3 and 4 standard deviations.
Statistical errors are 1 for all coefficients.

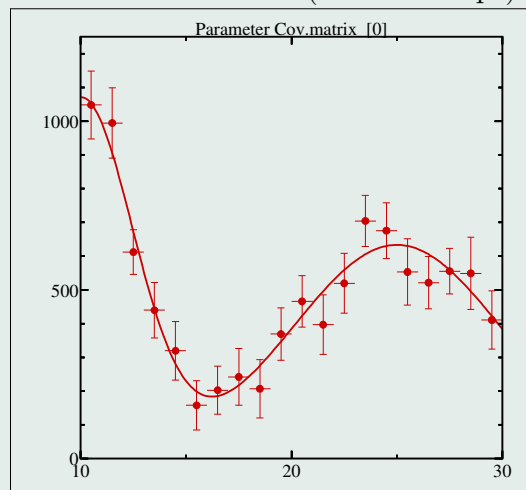
Example with truncation

Example of unfolding problem with narrow bins

Perfect resolution (no smearing)



Truncation method (15 terms kept)



Histogram
for sample with 10 000 entries.

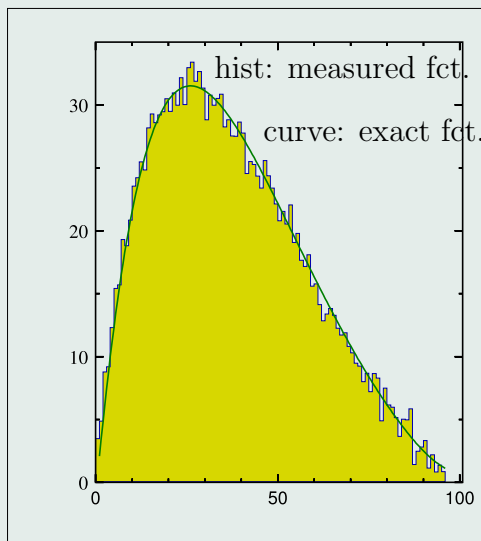
Reduced fluctuations

True curve $f(x)$ is shown in red.

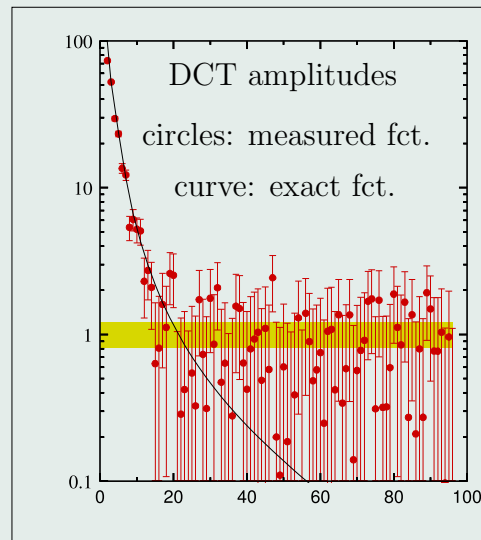
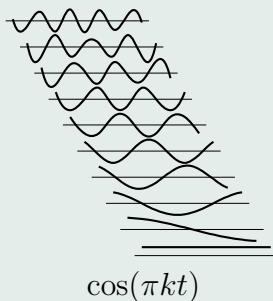
5. Truncation and positive correlations

“Can unfolding dist have $\sigma < \sqrt{n}$?”

Assume distribution \mathbf{y} without migration: discrete cosine transformation $\mathbf{c} = \mathbf{U}_{\text{DCT}}^T \mathbf{y}$



Histogram with unit cov. matrix



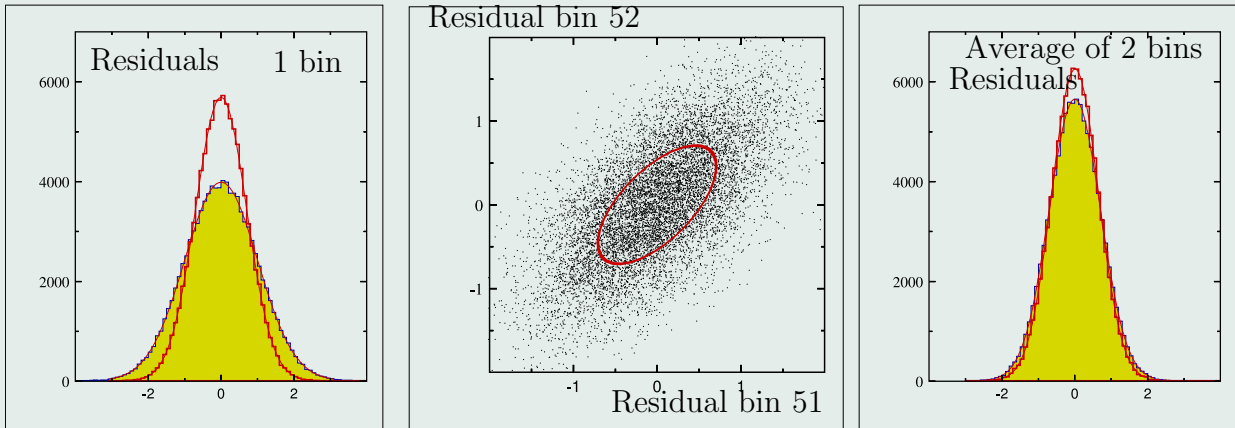
DCT amplitudes with unit cov. matrix

All high-frequency coefficients with index > 20 compatible with zero: truncation of second half of contributions (coeffs. and elements of cov. matrix) and back-transformation $\hat{\mathbf{y}} = \mathbf{U}_{\text{DCT}} \mathbf{c}$ reduces noise.

Note: $\Delta\chi^2 = c_j^2$ for single coefficient c_j !

Increase of accuracy ...

- Histograms of differences to exact bin content: single bin, bin 51 vs. bin 52, 2-bin average;
- broader distribution from unmodified histogram, narrower from histogram with truncation after half the coefficients;
- truncation results in higher local accuracy (single bins), no bias, but positive correlations;
- perfect agreement of width with calculated errors and correlations.



Truncation of coefficients of high frequency contributions

- no bias introduced; cov.matrix calculation is accurate (but rank defect);
- higher precision for single bins only, due to introduction of positive correlations;
- neither higher total precision by truncation nor distortion of data or errors.

6. Regularization

Norm regularization: $\min_x \{ \|\mathbf{Ax} - \mathbf{y}\|^2 + \tau \|\mathbf{x}\|^2 \} \implies \mathbf{x} = \mathbf{A}^\# \mathbf{y} = \left[(\mathbf{A}^\top \mathbf{A} + \tau \mathbf{I})^{-1} \mathbf{A}^\top \right] \mathbf{y}$

Solved for fixed τ by inversion, or (better) ...

$$\mathbf{A}^\# \mathbf{A} \neq \mathbf{I}$$

using SVD $\bar{\mathbf{x}} = \mathbf{V} \underbrace{\left[(\boldsymbol{\Sigma}^2 + \tau \mathbf{I})^{-1} \boldsymbol{\Sigma}^2 \right]}_{\text{filter factor } \varphi} \underbrace{\boldsymbol{\Sigma}^{-1} (\mathbf{U}^\top \mathbf{y})}_{\text{coeff. } \mathbf{c}}$

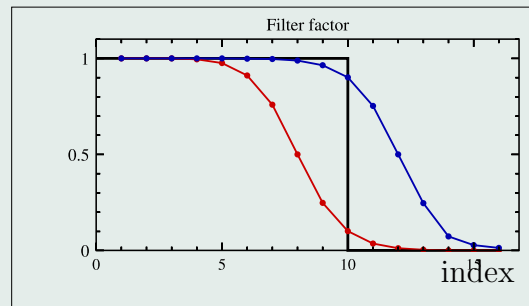
Effect of regularization: introduction of filter factor φ_j :

$$\hat{\mathbf{x}} = \sum_{j=1}^n \frac{1}{\sigma_j} \mathbf{c}_j \mathbf{v}_j \implies \hat{\mathbf{x}} = \sum_{j=1}^n \frac{1}{\sigma_j} \varphi_j \mathbf{c}_j \mathbf{v}_j$$

$$\mathbf{V}_x = \sum_{j=1}^n \frac{1}{\sigma_j^2} \mathbf{v}_j \mathbf{v}_j^\top \implies \mathbf{V}_x = \sum_{j=1}^n \frac{1}{\sigma_j^2} \varphi_j^2 \mathbf{v}_j \mathbf{v}_j^\top$$

$$\text{filterfactor } \varphi_j = \left(\frac{\sigma_j^2}{\sigma_j^2 + \tau} \right) = \begin{cases} 1 & \text{for } \sigma_j^2 \gg \tau \\ 1/2 & \text{for } \sigma_j^2 = \tau \\ 0 & \text{for } \sigma_j^2 \ll \tau \end{cases}$$

Different dependence possible: $1 / (1 + (\tau / \sigma_j^2)^\alpha)$



Regularization with differential operator

$$\mathbf{L} \text{ regularization: } \min_{\mathbf{x}} \{ \|\mathbf{A}\mathbf{x} - \mathbf{y}\|^2 + \tau \|\mathbf{L}\mathbf{x}\|^2 \} \implies \mathbf{A}^\# = \left[(\mathbf{A}^\top \mathbf{A} + \tau \mathbf{L}^\top \mathbf{L})^{-1} \mathbf{A}^\top \right]$$

Solution by more complicated mathematical operation (“generalized singular value decomposition”), but almost identical unfolding formalism with filter factors φ_j .

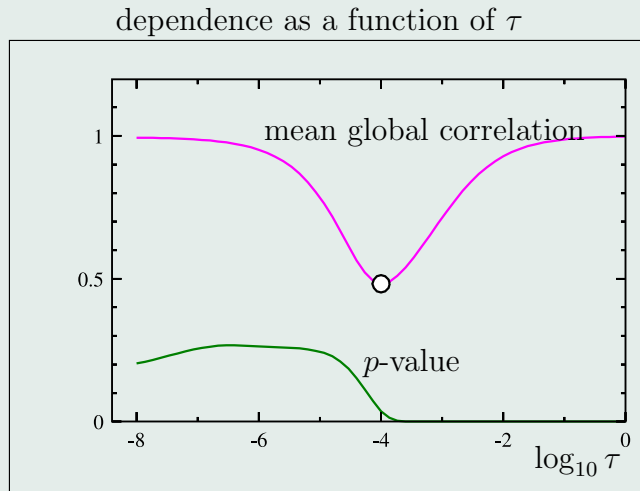
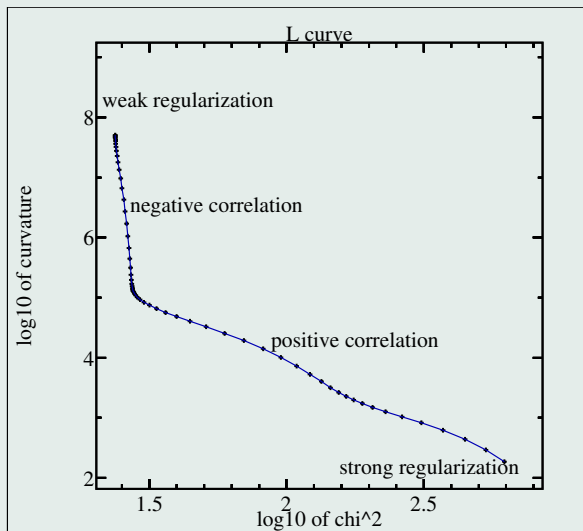
$$\text{most popular: sec.der. } \mathbf{L}_2^r = \begin{pmatrix} 1 & -1 & 0 & 0 & \cdots & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & \cdots & 0 & 0 & 0 \\ 0 & -1 & 2 & -1 & \cdots & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & -1 & 2 & -1 \\ 0 & 0 & 0 & 0 & \cdots & 0 & -1 & 1 \end{pmatrix} = \mathbf{U}_{\text{DCT}} \mathbf{\Lambda} \mathbf{U}_{\text{DCT}}^\top$$

My opinion about regularization:

- A result without regularization with many data points represents noise – nothing else;
- removing coefficient c_j (by filtering) means $\Delta\chi^2 = c_j^2$;
- regularization will not introduce unwanted bias, and allows to reduce or suppress insignificant contributions (noise), that would destroy the unfolding result.

Determination of regularization parameter

- Regularization corresponds technically to “weak” a-priori information about \mathbf{x} (like “measurement” with standard deviation $\approx \sqrt{1/\tau}$ for norm regularization);
- if $\tau > \sigma_j^2$ of significant Fourier coefficients $c_j \Rightarrow$ **no bias**;
- evaluate several quantities as function of τ for wide range: $\tau \in (\tau_L, \tau_R)$;
- effective number of degrees of freedom: $\sum \varphi_j$



Presentation of regularized result I

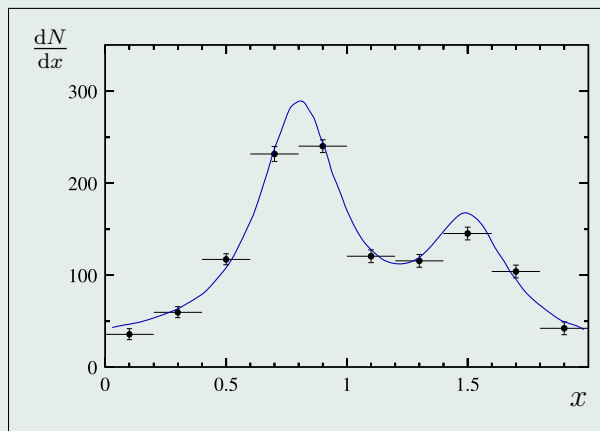
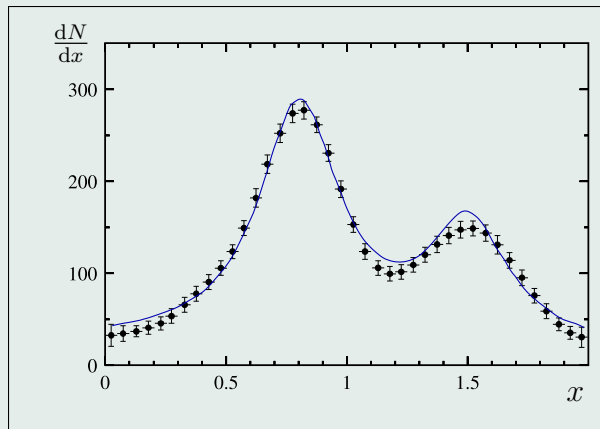
... with MC input

(from A. Hoecker and V. Kartvelishvili: NIM A 372)

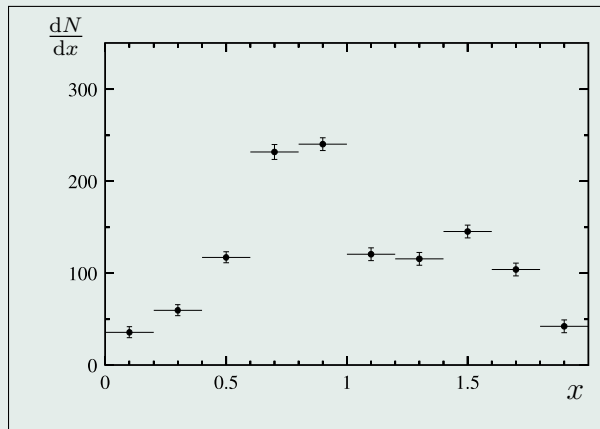
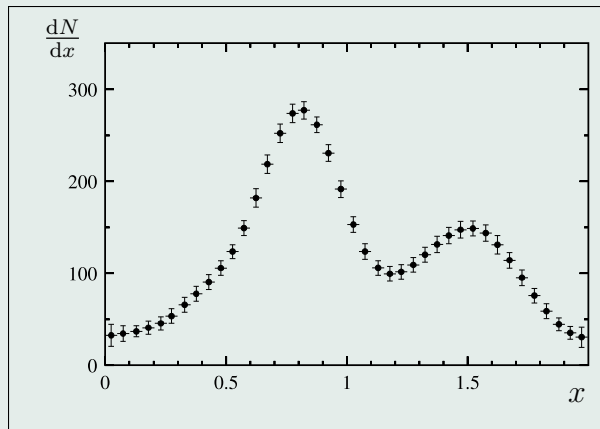
- Result with 40 data points, almost like a “band” representing result
- constructed from 10 significant parameters, with 40-by-40 covariance matrix, singular with rank 10; use “effective” weight matrix for fits

$$[\mathbf{V}_x]^{-1} = \sum_{j=1}^n \sigma_j^2 \mathbf{v}_j \mathbf{v}_j^T$$

- large positive correlations, therefore few sign-changes of the residuals to MC input distribution.
- (Same) result with 10 data points, each point represents a bin average of the result
- 10-by-10 covariance matrix non-singular, inverse is weight matrix.
- small and negligible correlations.



- Plots (without the true MC distribution): both data sets represent the same information!
- What is the better data presentation in Particle Physics?
- Try to calculate “limits” from the data!
- Note: not obtained by different methods, but (only) different presentation of identical information.



Example: two unknowns, three measured values

The measured vector \mathbf{y} and the response matrix \mathbf{A} are given by

$$\mathbf{y} = \underbrace{\mathbf{A} \begin{pmatrix} 1 \\ 1 \end{pmatrix}}_{\text{exact}} + \underbrace{\begin{pmatrix} 0.01 \\ -0.03 \\ 0.02 \end{pmatrix}}_{\text{error}} = \underbrace{\begin{pmatrix} 0.27 \\ 0.25 \\ 3.33 \end{pmatrix}}_{\text{measured}} \quad \mathbf{A} = \begin{pmatrix} 0.16 & 0.10 \\ 0.17 & 0.11 \\ 2.02 & 1.29 \end{pmatrix}$$

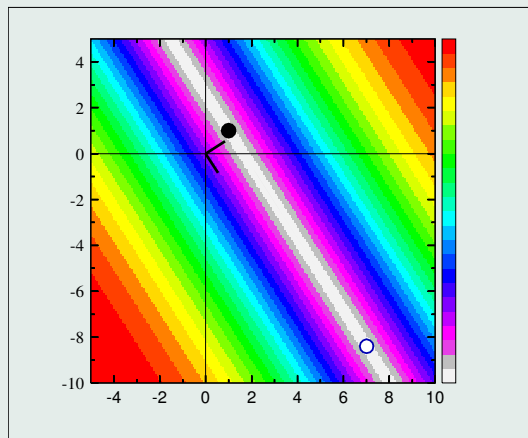
If analysed by the singular value decomposition the singular values are

$$\sigma_1 = 2.4127 \quad \sigma_2 = 0.0022$$

The least-squares solution $\hat{\mathbf{x}}_{\text{LS}}$, assuming a standard deviations of 0.02 for the three elements of \mathbf{y}

$$\hat{\mathbf{x}}_{\text{exact}} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad \hat{\mathbf{x}}_{\text{LS}} = \begin{pmatrix} 7.01 \pm 4.90 \\ -8.40 \pm 7.67 \end{pmatrix}$$

$$\text{correlation coeff. } \rho = -0.999998$$



Example from Per Christian Hansen: [Discrete Inverse Problems – Insight and Algorithms](#)

Example with regularization

- Regularization term $\tau \|\mathbf{x}\|^2$ is equivalent to assumption of measurement information

$$(\mathbf{x})_j = 0 \pm \sqrt{\frac{1}{\tau}}$$

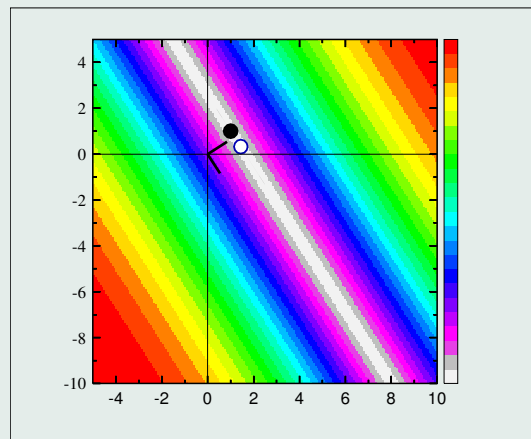
- Example: assume solution to have norm $\|\mathbf{x}\|$ of the order of 1.
- Regularization parameter value $\tau = 1/4$ corresponds to

$$(\mathbf{x})_j = 0 \pm 2 \quad j = 1, 2$$

Regularized solution of example:

$$\hat{\mathbf{x}}_{exact} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad \hat{\mathbf{x}}_{reg LS} = \begin{pmatrix} 1.44 \pm 1.05 \\ 0.33 \pm 1.65 \end{pmatrix}$$

correlation coeff. $\rho = -0.99996$



Calculation using normal fit program `aplcon`.

7. DCT and projection methods

Discrete cosine transformation DCT = orthonormal transformation by matrix U_{DCT} with $U_{\text{DCT}} U_{\text{DCT}}^T = I$,

$$\text{transformations: } \mathbf{X} = U_{\text{DCT}}^T \mathbf{x} \quad \mathbf{x} = U_{\text{DCT}} \mathbf{X}$$

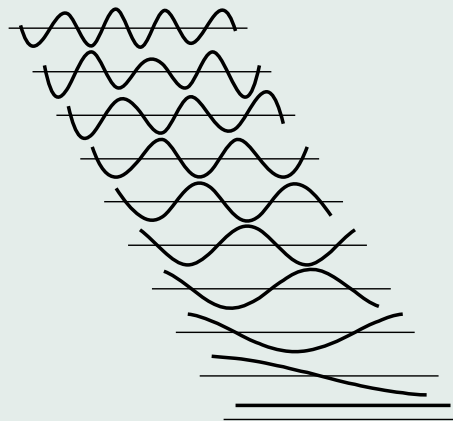
$$U_{jk} = \begin{cases} \sqrt{1/n} & k = 0 \\ \sqrt{2/n} \cos[\pi k(j + 1/2)/n] & k = 1, 2, \dots, n-1 \end{cases}$$

Relation to derivative regularization matrix $L = L_2^r$:

$$L = U_{\text{DCT}} \Lambda U_{\text{DCT}}^T \quad L^T L = U_{\text{DCT}} \Lambda^2 U_{\text{DCT}}^T$$

with eigenvalues: $\lambda_k = 4 \sin^2\left(\frac{k\pi}{2n}\right) \quad k = 0, 1, \dots$

- DCT is purely real, and concentrates “energy” into lower order coefficients better than discrete Fourier transformation;
- separability – perform DCT in any of the directions first and then apply to second direction: coefficients will not change: 2-D DCT is 2×1 -D DCT;
- used in modern coding standards like JPEG, MPEG.



$$u_k(t) = \cos(\pi k t)$$

$$k = 0, 1, \dots, n-1$$

$$t \in (0, 1)$$

Projection methods

So far: approximation of solution $\mathbf{x}_{\text{exact}} \in \mathbb{R}^n$ in low-dimensional subspace of low-frequency components:

- by truncation, keeping the first p Fourier coefficients of SVD method (TSVD), or
- by regularization using e.g. second-derivative matrix \mathbf{L}_2^r with regularization parameter τ .

Alternative: use of fixed projection n -by- p matrix $\mathbf{W}_{(p)} = \{\mathbf{v}_1, \mathbf{v}_2 \dots \mathbf{v}_p\}$ with $\bar{\mathbf{x}} \in \mathbb{R}^p$:

$$\mathbf{x} = \mathbf{W}_{(p)} \bar{\mathbf{x}} \qquad \bar{\mathbf{A}} = \mathbf{A} \mathbf{W}_{(p)}$$

$$\min_{\mathbf{x}} Q(\mathbf{x}) \quad \text{with} \quad Q(\mathbf{x}) = \|\mathbf{A} \mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{A} \mathbf{W}_{(p)} \bar{\mathbf{x}} - \mathbf{y}\|^2 = \|\bar{\mathbf{A}} \bar{\mathbf{x}} - \mathbf{y}\|^2$$

- Basis vectors $\mathbf{v}_j =$ e.g. **DCT eigenvectors**, or SVD singular vectors;
- solutions $\bar{\mathbf{x}}$ and $\mathbf{x} = \mathbf{W}_{(p)} \bar{\mathbf{x}}$ dominated by low-frequency components;
- truncation or filtering with filter functions φ_j possible;
- fixed one- or two-dimensional transformation, to be used with different MC data sets to study MC differences.

8. Iterative unfolding

Why **iterative methods**? What is the advantage compared to a direct solution?

- Unfolding with extremely large dimensions (e.g. 10^6 parameters in picture deblurring) requires iterative methods;
- dimension parameters in Particle Physics problems are small, no cpu-time problems for direct analytic methods, like SVD;
- certain iterative methods for inverse problems are popular in Particle Physics: they have semi-convergence, with **implicit regularization**.

Gauß had used iterative methods at least since the year 1823:

You will in future hardly eliminate directly, at least not when you have more than two unknowns. The indirect procedure can be done while one is half asleep, or is thinking about other things. [Carl Friedrich Gauß, Werke IX, p.278]

Landweber iteration

Formula for one Landweber⁺ iteration step $k = 0, 1, 2, \dots$:

$$\mathbf{x}^{[k+1]} := \mathbf{x}^{[k]} + \omega \mathbf{A}^T (\mathbf{y} - \mathbf{A} \mathbf{x}^{[k]}) \quad 0 < \omega < 2/\sigma_1^2$$

Semi-convergence corresponds to regularization ($j = 1, 2, \dots, n$): iteration number k

$$\mathbf{x}^{[k]} = \sum_{j=1}^n \frac{\varphi_j^{[k]}}{\sigma_j} (\mathbf{u}_j^T \mathbf{y}) \mathbf{v}_j \quad \varphi_j^{[k]} \approx 1 - (1 - \omega \sigma_j^2)^k \approx \begin{cases} 1 & \text{for large } \sigma_j^2 \\ k (\omega \sigma_j^2) & \text{for } \sigma_j^2 \ll 1/\omega \\ \rightarrow 0 & \text{for } \sigma_j^2 \rightarrow 0 \end{cases}$$

Iteration number k plays the role of a regularization parameter:

- convergence fast for components with a large singular values σ_j ;
- very slow for components with small singular value: components unchanged after few iterations;
- after a large number of iterations the unique (oscillating) solution of linear system;
- no prescription for covariance matrix calculation.

⁺ Other names associated with the algorithm are Richardson, Fridman, Picard and Cimmino.

Iterative methods in particle physics

In general there is the attempt, by an iterative tuning process, to use the “correct” distribution in the MC simulation, i.e. that distribution that should be extracted from the measured data $\mathbf{y} \simeq \mathbf{A}\mathbf{x}$.

Iterative improvement of a matrix $\mathbf{M}_x^{[k]}$, which depends on the solution \mathbf{x} and which should allow to extract the solution

$$\mathbf{x}^{[k+1]} = \mathbf{M}_x^{[k]} \mathbf{y} \implies \mathbf{x}^{[k+1]} = \left(\mathbf{M}_x^{[k]} \mathbf{A} \right) \mathbf{x} \quad \text{but} \quad \left(\mathbf{M}_x^{[k]} \mathbf{A} \right) \neq \mathbf{I}$$

“equation” valid only for certain \mathbf{x}

bin-by-bin correction factor: $\mathbf{M}_x = \text{diagonal}$

other iterative methods: $\mathbf{M}_x = \text{matrix with non-negative elements}$

Advantage: popular and accepted by collaborations; simple mathematics: no “complicated” operations like SVD;

Disadvantage: no direct error propagation with matrix \mathbf{M}_x possible; unknown regularization strength; correlations unknown and ignored in bin-by-bin correction method; questionable statistically and mathematically; applicable only to “correct” for migration effects, no general unfolding.

Summary

- Folding is a direct process – it is robust and simple, but a folded model prediction does not show the sensitivity of the measurement (and is insensitive to contribution \mathbf{x}_{null}).
- Unfolding is the inverse process – it is a discrete ill-posed problem, mathematically complex, with a response matrix with large condition number, but it allows to study the sensitivity of the measurement:
 - test of black-box algorithms on a few selected examples are not sufficient;
 - it is essential to understand the statistical and mathematical properties of the algorithm, and to understand the detector;
 - orthogonalization and regularization are generally accepted concepts in other fields;
 - there are still several open questions to apply successfully unfolding for particle physics experiment.

Appendix

- Types of unfolding problems
- Open questions
- Is the unfolding result allowed to depend on the MC input dependence?
- Comparison of codes
- Unfolding program *RUN*
- Literature
- Correlations
- Averaging correlated data
- Smoothing by truncation of DCT amplitudes
- Low-pass regularization
- Example with low-pass filter

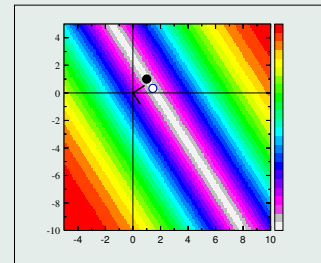
Types of unfolding problems

Parametrized unfolding

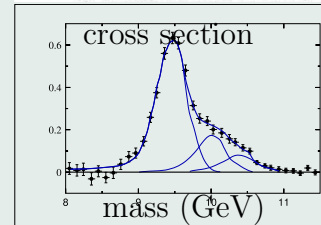
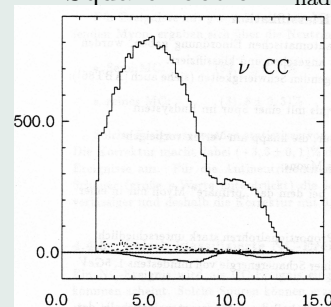
parameterized expression with few parameters no regularization necessary – using normal fit program.

Parameter-free unfolding bins or classes, no parametrization

- Classification problems, small number of bins: no reduction of number of bins possible;
- Correction of distributions with small migration effects;
- Unfolding with transformation, evtl. > 1 measured distribution;
- Structureless distributions (small number of Fourier coefficients):
 - Regularization with possibility to reduce number of bins;
 - Steeply falling: try transformation to remove steepness, e.g. $\sqrt{E_{\text{had}}}$, evtl. use parametrization
- Distributions with narrow structures, peaks: try to use parametrization e.g. Gaussian peaks)



Square root of E_{had}



Open questions

- discretization: continuous B-splines of higher order (necessary for re-weighting of MC – RUN) instead of simple discontinuous bins
- input of n -tuples instead of vectors/matrices
- strategies for steeply falling distributions and for class data (e.g. 2-, 3- ... jet events)
- alternative low-pass filter strategies
- data presentation with small or large number of data points
- projection methods (DCT) in one- and two-dimensions
- case of limited acceptance regions (e.g. low p_T – using constraints (Lagrange multiplier – RUN)) in one- and two-dimensions
- using general fit program (like aplcon) for unfolding? (flexible: e.g. fit may include uncertainty of \mathbf{A} or background fit)
- uncertainties of response matrix elements
- algorithm for automatic variable bin definition (RUN)
- statistical comparison of information content for different number of degrees of freedom
- use of detailed estimates of the unfolding result (multiplicative, additive)

Is the unfolding result allowed to depend on the MC input dependence?

My opinion: **NO**.

Do ‘blind analysis’, whenever possible. Avoid any bias w.r.t. an expected result.

From a paper on the CFM: *‘The correction of the detector acceptance using Monte Carlo modelling requires the cross section model used in the simulation to be sufficiently close to the data, such that migration between the bins are well reproduced. . . . In practice this is achieved using an iterative MC event reweighting procedure which converges after one iteration for the measurement region.’*

In an iterative method: *‘. . . it gives the best results (in terms of its ability to reproduce the true distribution) if one makes a realistic guess about the distribution that the true values follow’*

What happens in the case of a completely insensitive detector?

Regularization methods will not be able to get any result!

In the iterative method: *‘One finds then that the final probabilities are equal to the initial ones’*

Comparison of codes

| Method | RUN | GURU | Tunfold | Iterative |
|-------------------|-----|------|---------|-----------------|
| Input: matrix | | ✓ | ✓ | ✓ |
| Input: n -tuple | ✓ | | | |
| Orthogonalization | ✓ | ✓ | | |
| Input errors | ✓ | ✓ | ✓ | |
| Least squares | | ✓ | ✓ | ? |
| MaxLik (Poisson) | ✓ | | | |
| Regularization | ✓ | ✓ | ✓ | implicit |
| iterative | | | | ✓ |
| automatic binning | ✓ | | | |
| Cov.mat. by prop. | ✓ | ✓ | ✓ | |
| MC re-weighting | ✓ | | | |

GURU (Fortran, SVD, by Andreas Hoecker and Vato Kartvelishvili) and TUNFOLD (by Stefan Schmitt) in RooUnfold (ROOT Unfolding Framework, by Tim Adye et al.)

RUN (Fortran), converted to C++ by Natalie Milke (Uni Dortmund)

Unfolding program *RUN*

Development started 1979/1980 for neutral current neutrino experiment CHARM:

- First problem was reconstruction of cross section $d\sigma/dy$ from the measured values E_{had} and radius $r_{\text{interaction}}$ (only these 2 quantities were measurable).
- Input are n -tuples, and 1-dim. or 2-dim. or 3-dim. measured histograms.
- Instead of Least Squares the ML method with Poisson statistic was used (sometimes only few entries/bin for > 1 -dim. histograms), with diagonalization of Hessian.
- For the intermediate result cubic B -splines were used to avoid discontinuities.
- A special option allows to check the consistency of MC simulation.
- Used in other experiments:
Neutrino physics, 2-photon-physics at e^+e^- colliders, in astrophysics, and still used in 2010 (LHC, D0 at FNAL).
- Conversion to C++ by Natalie Milke (Uni Dortmund).

Literature

Per Christian Hansen.

Rank-deficient and Discrete Ill-Posed Problems: Numerical Aspects of Linear Inversion,
Volume 160 of *SIAM monographs on mathematical modeling and computation*. Society for Industrial and Applied Mathematics, Philadelphia, 1997.

Discrete Inverse Problems – Insight and Algorithms,

Volume 160 of *Fundamentals of Algorithms*. Society for Industrial and Applied Mathematics, Philadelphia, 2010.

Jari Kaipio and Erkki Somersalo.

Statistical and Computational Inverse Problems,

Volume 160 of *Applied Mathematical Science*. Springer, 2004.

Statistical inverse problems: Discretization, model reduction and inverse crimes,

Journal of Computational and Applied Mathematics, 198:493 – 504, 2007.

Curtis R. Vogel.

Computational Methods for Inverse Problems,

Volume 160 of *SIAM Frontiers in Applied mathematics*. Society for Industrial and Applied Mathematics, Philadelphia, 2002.

Andreas Rieder.

Keine Probleme mit Inversen Problemen.

Vieweg, 2003.

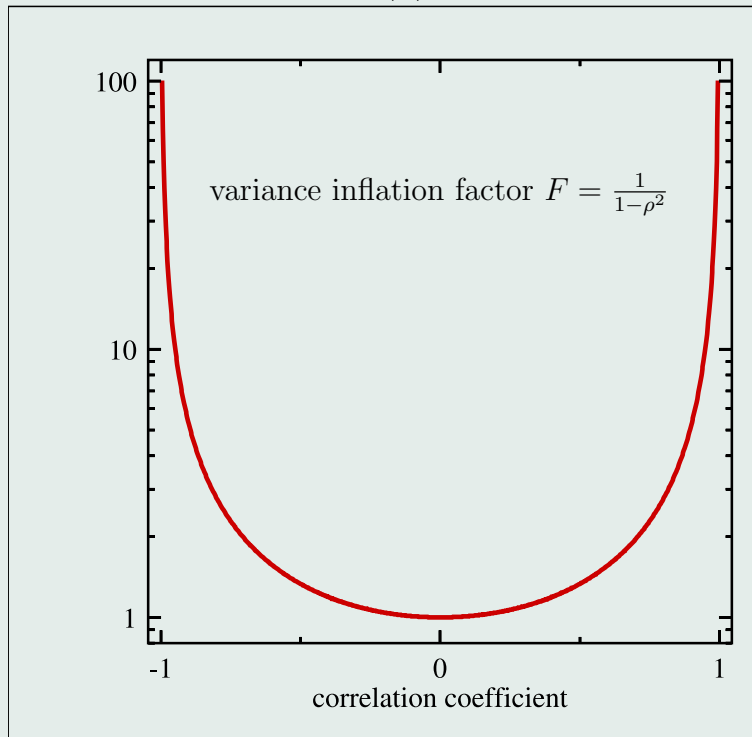
Ake Björk.

Numerical Methods for Least Squares Problems.

Society for Industrial and Applied Mathematics, Philadelphia, 1996.

Correlations

Correlation coefficients $|\rho| > 0.5$ are dangerous:

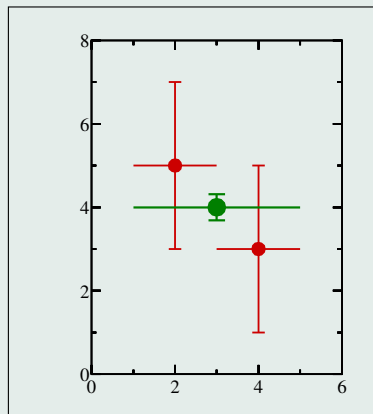
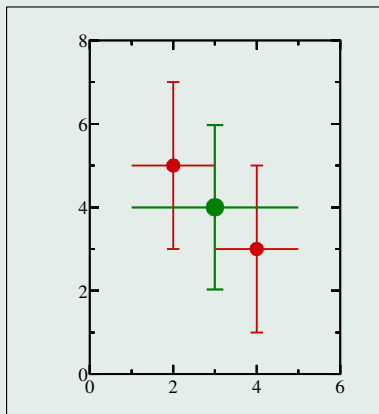


Averaging correlated data

The Figures show two adjacent data points d_1 and d_2 with a large positive (left) and negative (right) correlation coefficient, assuming the same standard deviations: $\sigma = \sigma_1 = \sigma_2$.

$$\mathbf{V} = \begin{pmatrix} \sigma^2 & \rho_{12}\sigma^2 \\ \rho_{12}\sigma^2 & \sigma^2 \end{pmatrix} \quad \text{with} \quad \rho_{12} = \pm 0.95$$

Average: the average value $\bar{d} = \frac{1}{2}(d_1 + d_2)$ has a variance of $\mathbf{V}_{\bar{d}} = \frac{1}{2}(1 + \rho_{12})\sigma^2$ (see middle point in figures).

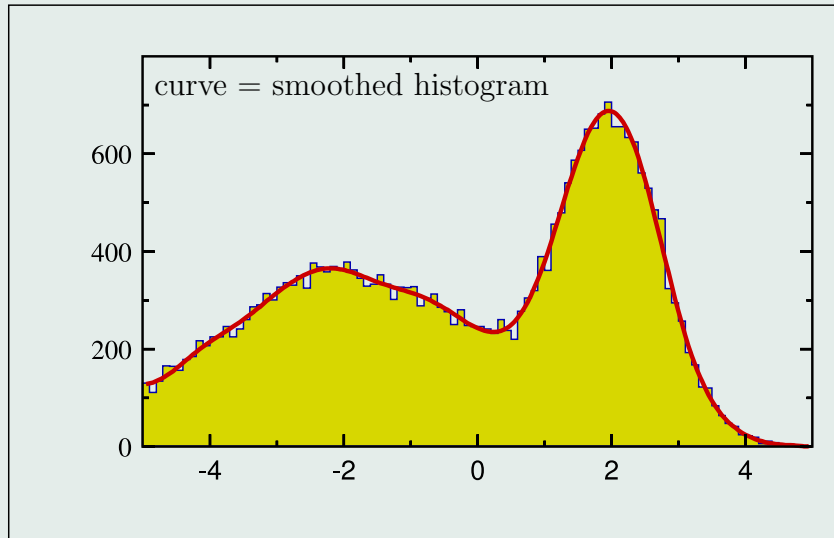


For highly correlated data the properties for averaging and χ^2 -comparison with predicted values are not intuitive, but have to be performed with an explicit calculation, based on the inverse covariance matrix $\mathbf{V}_x^{-1} = \mathbf{W}_x$ (weight matrix).

Smoothing by truncation of DCT amplitudes

Example for smoothing of histogram by DCT and truncation:

- Transformation ($\sqrt{n_j}$) of bin content (Poisson) to stable variance;
- discrete cosine transformation and truncation of insignificant coefficients;
- back transformations.



Low-pass regularization

Without regularization there are large bin-to-bin fluctuations due to negative correlations between neighbour bins. These fluctuations can be suppressed in a low-pass filter by averaging 3-to-1 bins:

$$\bar{x}_j = \frac{1}{4}x_{j-1} + \frac{1}{2}x_j + \frac{1}{4}x_{j+1} \quad \text{or general} \quad \bar{x}_j = a_j x_{j-1} + (1 - 2a_j) x_j + a_j x_{j+1}$$

The factor a_j can be chosen to minimize¹⁾ the variance of \bar{x}_j , using the known matrix \mathbf{V}_x .

Pro: Fluctuations are really suppressed and the true dependence is clearer visible.

No bias, if number of bins large and no strong structure.

Con: In regions of larger second-derivatives a bias is introduced, because the above filter assumes an almost linear dependence over 3-point regions.

First and last bins disappear.

The general averaging algorithm for this ‘local’-regularization method:

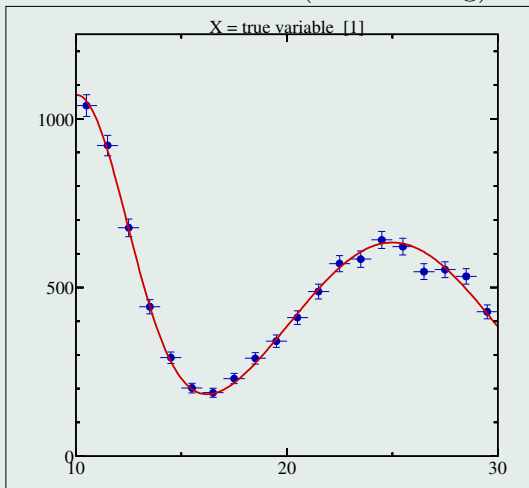
$$\begin{aligned} \bar{\mathbf{x}} &= \mathbf{T}\mathbf{x} \\ \mathbf{V}_{\bar{\mathbf{x}}} &= \mathbf{T}\mathbf{V}_x\mathbf{T}^T \end{aligned} \quad \text{with} \quad \mathbf{T} = \frac{1}{4} \begin{pmatrix} 1 & 2 & 1 & 0 & \cdots & 0 & 0 & 0 \\ 0 & 1 & 2 & 1 & \cdots & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 1 & 2 & 1 \end{pmatrix}$$

1) O. Helene et al., NIM A 523 (2004) 186; NIM A 580 (2007) 1466 - 1473

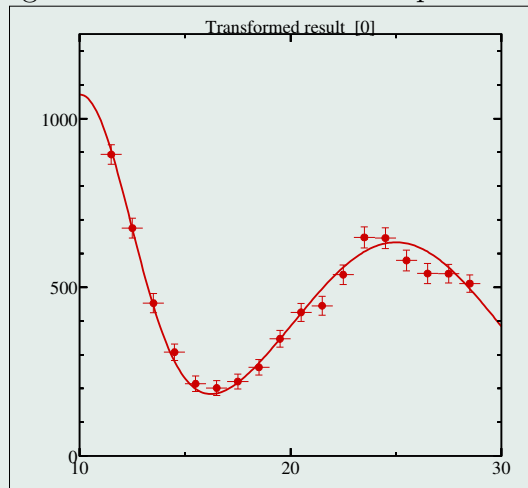
Example with low-pass filter

Example of unfolding problem with narrow bins

Perfect resolution (no smearing)



regularization $\tau = 10^{-5}$ + low-pass filter



Histogram
for sample with 10 000 entries.

Reduced correlations: neighbour bin -10%
and second neighbour -30% .

True curve $f(x)$ is shown in red.

Contents

| | | | |
|---|-----------|---|-----------|
| 1. Unfolding – direct and inverse processes | 2 | 7. DCT and projection methods | 24 |
| Discretization | 3 | Projection methods | 25 |
| Unfolding is more general than “data correction” | 4 | 8. Iterative unfolding | 26 |
| 2. Naive unfolding | 5 | Landweber iteration | 27 |
| Naive result with narrow bins | 6 | Iterative methods in particle physics | 28 |
| 3. Convolution/deconvolution | 8 | Summary | 29 |
| . . . and deconvolution | 9 | Appendix | 30 |
| 4. Orthogonalization | 10 | Types of unfolding problems | 31 |
| Least squares solution | 11 | Open questions | 32 |
| Vanishing singular values and truncation | 12 | Is the unfolding result allowed to depend | |
| Eigenvalues and Fourier coefficients | 13 | on the MC input dependence? | 33 |
| Example with truncation | 14 | Comparison of codes | 34 |
| 5. Truncation and positive correlations | 15 | Unfolding program <i>RUN</i> | 35 |
| Increase of accuracy | 16 | Literature | 36 |
| 6. Regularization | 17 | Correlations | 37 |
| Regularization with differential operator | 18 | Averaging correlated data | 38 |
| Determination of regularization parameter | 19 | Smoothing by truncation of DCT amplitudes | 39 |
| Presentation of regularized result I | 20 | Low-pass regularization | 40 |
| Presentation of regularized result II | 21 | Example with low-pass filter | 41 |
| Example: two unknowns, three measured values | 22 | Table of contents | 42 |
| Example with regularization | 23 | | |