



Diagonalizing the Unfolding Problem

- figure out what can be learned and what not -

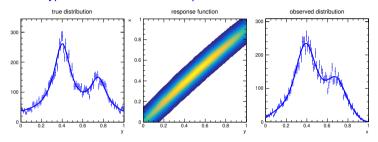
Michael Schmelling / MPI for Nuclear Physics

Disclaimer

This presentation is about ongoing work, i.e. math and story line are still evolving. Only the simplest case of a 1-dim linear unfolding problem is discussed. The talk focuses on results. Mathematical details are given in the backup slides.

Introduction

sketch of a typical 1-dim measurement problem



- given: noisy data that estimate the observable distibution
- known: detector response, i.e. efficiency, biases and smearing
- wanted: an estimate of the true distribution

in the following: linear detector response and counting statistics ->

The linear model

Fredholm integral equation of 1st kind

$$g(x) = \int dy \, R(x,y) f(y)$$

- f(y): true density unknown
 - ightharpoonup integration over the full range allowed for y
 - lacktriangleright nature generates values y according to f(y)
- \blacksquare R(x, y): response function of the detector known
 - ightharpoonup density in x for fixed y
- g(x): parent distribution of the actual observations approximately known
 - lacktriangleright an experiment records IID measurements $x_i, i=1,\ldots,N$ drawn from g(x)
 - \blacktriangleright the sample size N is a Poisson distributed random variable

The full story?

* real life is more complicated than Fredholm's equation

$$egin{aligned} g(x) &= \int dy \ R_0(x) \ &+ \int dy \ R_1(x,y) f(y) \ &+ \int dy \ dy' \ R_2(x,y,y') f(y) f(y') \ &+ \int dy \ dy' dy'' \ R_3(x,y,y',y'') f(y) f(y') f(y'') \ &+ \dots \end{aligned}$$

back to looking only at the 2nd term →

often used: histogram approximation

integral equation → matrix equation

$$a_k = \sum_{l=1}^{n_b} R_{kl} b_l \quad ext{or} \quad a = R \ b$$

with

$$a_k = \int_{\Delta x_k} \!\!\! dx \, g(x) \;,\; b_l = \int_{\Delta y_l} \!\!\! dy \, f(y) \quad ext{and} \quad R_{kl} = rac{1}{\Delta y_l} \int_{\Delta x_k} \!\!\! dx \int_{\Delta y_l} \!\!\! dy \, R(x,y)$$

- a: histogram of the observed distribution
- b: histogram of the true distribution
- R: response matrix
- numerically attractive because it requires to handle only finite dimensional matrices
- model dependence from choice of binning



Ansatz: expansion into a complete set of basis functions

 \diamond starting point: symmetric positive definite operator T(y,y')

$$T(y,y') = \int dx \ dx' \ R(x,y) \ R(x',y') \ w(x,x')$$

- appropriate functional form to address the unfolding problem
- lacksquare any symmetric positive definite weight function $w(x,x^\prime)$ is allowed
- lacksquare simplest choice: $w(x,x')=\delta(x-x')$
- optimum choice for IID data with Poisson distributed sample size:

$$w(x,x') = rac{\delta(x-x')}{g(x)}$$

- ightharpoonup in practice use a Maximum Likelihood estimate $\hat{q}(x)$ instead of q(x)
- ightharpoonup asymptotic properties are obtained by setting $\hat{q}(x) = q(x)$

lacktriangle properties of the eigenfunctions $b_k(y)$ of T(y,y')

$$\int dy' \ T(y,y') b_k(y') = \lambda_k \ b_k(y)$$

- lacksquare ordered positive eigenvalues $\lambda_k > \lambda_{k+1} > 0 \ \forall \ k$
- $lacksquare w(x,x') \propto 1/N \Rightarrow \lambda_k \propto 1/N$
- orthonormality

$$\int dy \ b_k(y) \ b_l(y) = \delta_{kl}$$

completeness

$$\sum_{k=0}^{\infty} b_k(y) \ b_k(y') = \delta(y-y')$$

 \diamond the unfolding problem in the basis $b_k(y)$

$$f(y) = \sum_{k=0}^\infty eta_k b_k(y) \quad ext{and} \quad g(x) = \sum_{k=0}^\infty eta_k \ a_k(x) \quad ext{with} \quad a_k(x) = \int dy \ R(x,y) \ b_k(y)$$

note: from now on specialize to

$$w(x,x')=rac{\delta(x-x')}{\hat{g}(x)}$$

which implies

$$\int dx \, rac{a_k(x) \, a_l(x)}{\hat{g}(x)} = \lambda_k \delta_{kl}$$
 and it follows $\int dx \, g(x) rac{a_k(x)}{\hat{g}(x)} = \lambda_k eta_k$

employ this to estimate $\beta_k \rightarrow$

❖ use the Law of Large Numbers and that expectation values are integrals

$$\lambda_k eta_k = \int dx \ g(x) rac{a_k(x)}{\hat{g}(x)} = N \int dx \ rac{g(x)}{N} rac{a_k(x)}{\hat{g}(x)} = N \left\langle rac{a_k(x)}{\hat{g}(x)}
ight
angle$$

 \blacksquare estimates from a sample of N measurements

$$N\left\langle rac{a_k(x)}{\hat{g}(x)}
ight
angle pprox \sum_{i=1}^N rac{a_k(x_i)}{\hat{g}(x_i)} = \lambda_k \hat{eta}_k$$

lacksquare covariance matrix of \hat{eta}_k in the limit $\hat{g}(x) o g(x)$

$$C_{kl}(\hat{eta}) = rac{\delta_{kl}}{\lambda_k}$$

- the expansion coefficients contribute independent information ("diagonalization")
- ▶ the variance grows with higher orders

Information content of the expansion coefficients

- lacksquare counting experiments: information = number of events = significance² = $(n/\sqrt{n})^2$
- \blacksquare significance² defines the equivalent number of events in β_k :

$$N_k = rac{eta_k^2}{C_{kk}(eta)} = eta_k^2 \lambda_k$$

 \blacksquare sum rule in the limit $\hat{g}(x) = g(x)$:

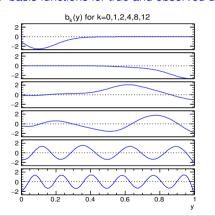
$$\sum_{k=0}^{\infty} N_k = N$$

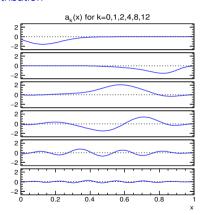
- \blacktriangleright the observed N events are distributed into independent measurements of the β_k
- ▶ the individual shares are proportional to β_k^2 and λ_k
- only the low-order coefficient get significant shares

numerical studies

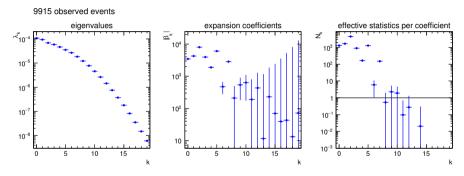
Properties of the introductory example

* basis functions for true and observed distribution



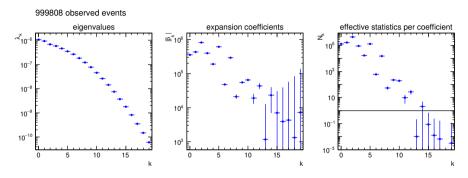


coefficients for a low statistics measurement



→ the experimental information is contained in 10 coefficients

coefficients for a high statistics measurement



→ the experimental information is contained in 14 coefficients

Interim summary

- a recipe for an optimal basis $b_k(y)$ to parametrize f(y) has been given
 - the expansion coefficients are independent
 - they are estimated by simple sums of weights over the data
- all information of the data about the true distribution is contained in a small number of leading order coefficients
 - for smooth functions the higher order terms are small $\beta_k \ll 1$
 - higher order terms are suppressed by $\lambda_k \ll 1$
 - \blacktriangleright the variance of the higher order terms grows with $1/\lambda_k$
- unfolding has been factorized into two problems
 - identifying the available information
 - \blacktriangleright constructing an estimate $\hat{f}(y)$ of the true distribution f(y)

first steps ->

Truncated expansion with only the well measured coefficients

$$f(y) = \sum_{k=0}^{\infty} \beta_k b_k(y) \quad \Rightarrow \quad \hat{f}(y) = \sum_{k=0}^{m} \hat{\beta}_k b_k(y)$$

relation between $\hat{f}(y)$ and f(y) if $\hat{\beta}_k = \beta_k$

$$\hat{f}(y) = \int dy' \, P(y,y') f(y') \quad ext{with} \quad P(y,y') = \sum_{k=0}^m b_k(y) b_k(y')$$

- ightharpoonup P(y,y') projects f(y) to the subspace spanned by the first m functions
- completeness of the basis $b_k(y)$ implies $P(y,y') = \delta(y-y')$ for $m \to \infty$
- different points y and y' are correlated with covariance

$$C(y,y') = \sum_{k=0}^m rac{1}{\lambda_k} b_k(y) b_k(y')$$

Concluding discussion

- every unfolding result $\hat{f}(y)$ is characterized by a posterior response P(y, y')
- $\hat{f}(y)$ will be biased even if the coefficients β_k are unbiased
- In It of the figure of the structures α in the β_k can cause unphysical structures α in $\hat{f}(y)$
 - \triangleright estimating f(y) by the truncated expansion may give unsatisfactory results
- \blacksquare the basis $b_k(y)$ is detector specific
 - comparison or combination of experiments not obvious
 - \triangleright check of theoretical predictions can be done for any basis $b_{k}(u)$
- more sophisticated regularisation schemes try to inject missing information
 - adjust expansion coefficients within uncertainties
 - invent higher order terms according to some prior knowledge
 - modify posterior response and correlation function
 - but the fundamental limitations of the truncated expansion always apply

Backup

Orthonormality of the basis $b_k(y)$

Application of the eigenvalue equation shows

$$I_{kl} = \int dy \; b_l(y) \; b_k(y) = rac{1}{\lambda_k} \int dy \; dy' \; b_l(y) \; T(y,y') \; b_k(y') = rac{\lambda_l}{\lambda_k} \int dy' \; b_l(y) \; b_k(y')$$

and thus

$$I_{kl} = rac{\lambda_l}{\lambda_k} I_{kl}$$

and it follows

- $\lambda_{k} \neq \lambda_{l}$: $I_{kl} = \delta_{kl}$
- $\lambda_k = \lambda_l$: orthogonal linear combinations of $b_k(y)$ and $b_l(y)$ exist

and with appropriate normalization:

$$\int dy \ b_l(y) \ b_k(y) = \delta_{kl}$$

Completeness of the basis $b_k(y)$

Given an arbitrary function f(y) and it's expansion into $b_k(y)$

$$f(y) = \sum_{k=0}^{\infty} \beta_k b_k(y)$$
,

the coefficients β_k are obtained, using the orthonormality of the $b_k(y)$, by

$$\int dy \; b_k(y) f(y) = \int dy \; b_k(y) \sum_{l=0}^{\infty} eta_l b_l(y) = \sum_{l=0}^{\infty} eta_l \int dy \; b_k(y) \; b_l(y) = \sum_{l=0}^{\infty} eta_l \; \delta_{kl} = eta_k \; .$$

Inserting the result for β_k then yields

$$f(y) = \sum_{k=0}^{\infty} b_k(y) \int dy' \ b_k(y') f(y') = \int dy' f(y') \sum_{k=0}^{\infty} b_k(y) \ b_k(y')$$

which implies

$$\sum_{k=0}^{\infty} b_k(y) \ b_k(y') = \delta(y-y')$$

Orthogonality of the transformed basis $a_k(x)$

According to the definition of $a_k(x)$ one has

$$\int dx \; dx' \; a_k(x) \, a_l(x') w(x,x') = \int dy \; dy' \; dx \; dx' \, R(x,y) b_k(y) R(x',y') b_l(y') w(x,x')
onumber \ = \int dy \; dy' \; b_k(y) \, T(y,y') b_l(y') = \lambda_l \int dy \; b_k(y) \; b_l(y) = \lambda_l \delta_{kl}$$

which for the special case

$$w(x,x')=rac{\delta(x-x')}{\hat{g}(x)} \quad ext{simplifies to} \quad \int dx \, rac{a_k(x) \, a_l(x)}{\hat{g}(x)} = \lambda_k \delta_{kl}$$

and which with one finds

$$\int dx \ g(x) \ \frac{a_k(x)}{\hat{g}(x)} = \int dx \ \sum_{l=0}^\infty \beta_l \ a_l(x) \frac{a_k(x)}{\hat{g}(x)} = \sum_{l=0}^\infty \beta_l \int dx \ \frac{a_l(x) a_k(x)}{\hat{g}(x)} = \lambda_k \beta_k$$

Decomposition of the response matrix

Using the completeness relation and the definition of $a_k(x)$ one finds:

$$egin{aligned} R(x,y) &= \int dy' \, R(x,y') \, \delta(y-y') \ &= \int dy' \, R(x,y) \, \sum_{k=0}^\infty b_k(y) b_k(y') \ &= \sum_{k=0}^\infty b_k(y) \int dy' \, R(x,y') b_k(y') \ &= \sum_{k=0}^\infty b_k(y) a_k(x) = \sum_{k=0}^\infty a_k(x) b_k(y) \end{aligned}$$

Covariance matrix of $\hat{\beta}_k$

With the weight functions $u_k(x)$

$$u_k(x) = rac{1}{\lambda_k} \, rac{a_k(x)}{\hat{g}(x)} \quad ext{and thus} \quad \hat{eta}_k = \sum_{i=1}^N u_k(x_i)$$

one finds, for Poisson distributed samples sizes N with $\langle N(N-1)\rangle = \langle N\rangle^2$,

$$egin{aligned} C_{kl}(\hat{eta}) &= \left<\hat{eta}_k\hat{eta}_l
ight> - \left<\hat{eta}_k
ight> \left<\hat{eta}_l
ight> \ &= \left< N(N-1)
ight> \left< u_k
ight> \left< u_l
ight> + \left< N
ight> \left< u_k u_l
ight> - \left< N
ight>^2 \left< u_k
ight> \left< u_l
ight> \ &= \left< N
ight> \left< u_k u_l
ight> &= \left< N
ight> \int dx rac{g(x)}{\langle N
ight>} rac{a_k(x) a_l(x)}{\hat{g}^2(x) \lambda_k \lambda_l} &= \int dx rac{a_k(x) a_l(x)}{\hat{g}(x) \lambda_k \lambda_l} &= rac{\delta_{kl}}{\lambda_k} \ . \end{aligned}$$

The last two steps follow in the limit $\hat{g}(x) = g(x)$ and from the orthogonality of the $a_k(x)$.

Sum rule

Given the number of equivalent events in β_k

$$N_k = rac{eta_k^2}{C_{kk}(eta)} = eta_k^2 \lambda_k$$

In the limit $\hat{q}(x) = q(x)$ with

$$\lambda_k eta_k = \int dx \ g(x) rac{a_k(x)}{\hat{g}(x)} = \int dx \ a_k(x)$$

one finds

$$\sum_{k=0}^\infty N_k = \sum_{k=0}^\infty eta_k(eta_k\lambda_k) = \sum_{k=0}^\infty eta_k \int dx \; a_k(x) = \int dx \; \sum_{k=0}^\infty eta_k a_k(x) = \int dx \; g(x) = N \; ,$$

i.e. effectively the total data set of N events is split into disjoint measurements of the β_k .

Covariance matrix of unfolded distribution

The covariance between two points u and u' is

$$C(y,y') = \left\langle \hat{f}(y)\hat{f}(y')\right\rangle - \left\langle \hat{f}(y)\right\rangle \left\langle \hat{f}(y')\right\rangle$$

and with

$$f(y) = \sum_{k=0}^m \hat{eta}_k \ b_k(y)$$

one finds

$$egin{aligned} C(y,y') &= \sum_{k,l=0}^m ig\langle \hat{eta}_k \hat{eta}_l ig
angle \ b_k(y) b_l(y') - \sum_{k,l=0}^m ig\langle \hat{eta} ig
angle ig\langle \hat{eta}_l ig
angle \ b_k(y) b_l(y') \ &= \sum_{k,l=0}^m C_{kl}(\hat{eta}) b_k(y) b_l(y') = \sum_{k,l=0}^m rac{\delta_{kl}}{\lambda_k} ig
brace b_k(y) b_l(y') = \sum_{k=0}^m rac{1}{\lambda_k} \ b_k(y) b_l(y') \end{aligned}$$