

Statistics in Data Analysis

All you ever wanted to know about statistics but never dared to ask

part 7

Paweł Brückman de Renstrom
(pawel.bruckman@ifj.edu.pl)

April 23, 2025

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Variable transformation

(revisited)

Let y be a function of a random variable $y = a(x)$ (or variables) which itself is a random variable.

x is distributed according to p.d.f. $f(x)$ How do we find $g(y)$, the p.d.f. of y ?

$$\Delta P = \int_{x' \in \{x, x+\Delta x\}} f(x') dx' = \int_{y' \in \{a(x), a(x+\Delta x)\}} g(y') dy' \quad (1)$$

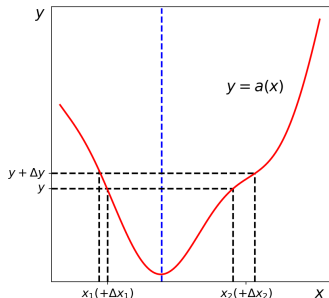
Irregardles of the sign of the derivative we have:

$$\Delta x \rightarrow 0 \Rightarrow f(x)|dx| = g(y)|dy|$$
$$g(y) = f(x(y)) \left| \frac{dx}{dy} \right| = f(x(y)) |a'(x)|^{-1} \quad (2)$$

where $x(y) \equiv a^{-1}(y)$.

If a not single-valued, contributions must be added:

$$g(y) = \sum_i f(x(y))_i \left| \frac{dx}{dy} \right|_i = \sum_i f(x(y))_i |a'(x_i)|^{-1} \quad (3)$$



Variable transformation

(revisited)

In a more general case of multidimensional joint p.d.f. one can define a one-to-one (invertable) transformation. We have:

$$g(y_1, \dots, y_n) = f(x_1, \dots, x_n) |J| \quad (4)$$

$$\text{where } J \equiv \det \begin{pmatrix} \frac{\partial x_1}{\partial y_1} & \dots & \dots \\ \dots & \dots & \dots \\ \dots & \dots & \frac{\partial x_n}{\partial y_n} \end{pmatrix},$$

$$x_i(y_1, \dots, y_n) \equiv a_i^{-1}(y_1, \dots, y_n).$$

The Monte Carlo method

Cumulative function & uniform distribution $[0, 1]$

- Let y be the cumulative of an arbitray p.d.f. $f(x)$:

$$y = F(x) = \int_{-\infty}^x f(x')dx' \quad (5)$$

- $g(y)$, the p.d.f. of y is given by the transformation:

$$g(y) = f(x(y)) \left| \frac{1}{F'(x)} \right| = f(x(y)) \frac{1}{f(x(y))} = \begin{cases} 1 & \text{for } y \in (0, 1) \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

- For any continuous p.d.f. $f(x)$, $y = F(x)$ is distributed according to $[0, 1]$. Hence, p.d.f. of $x = F^{-1}(y)$ will be $f(x)$ if y has a uniform distribution $[0, 1]$.
- $f(x; 0, 1)$ (or simply $[0, 1]$) is commonly used in statistics as the base for random number generators.

The Monte Carlo method

Generating the uniform distribution $[0, 1]$

- In order to generate random variables according to an arbitrary p.d.f. one needs the $[0, 1]$ to start with.
- This task is usually accomplished by a **random number generator** algorithm. More specifically, computers can generate **pseudorandom** numbers (e.g. RANMAR, RANLUX from CernLib or `numpy.random.random()` in Python).
- Pseudorandom sequence behaves has a feel & touch of randomness, but given the initial **seed** remains entirely deterministic. It also has a cycle period (the longer the better, of course).
- There are various algorithms on the market. A common simple implementation generates the sequence n_1, n_2, \dots according to:

$$n_{i+1} = (a \times n_i) \bmod m.$$

The **multiplier** a and the **modulus** m are constants that determine the sequence, in particular its period (e.g. $a = 40692$ and $m = 2147483399$ were used on 32-bit machines giving the period of $\approx 2 \times 10^9$).

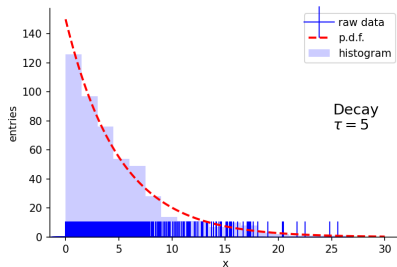
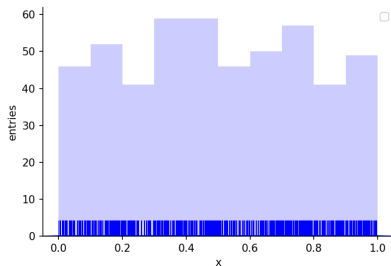
- Much better, more sophisticated algorithms are routinely used nowadays.

The Monte Carlo method

Using transformation for random number generation

- If the Eq. 5 can be solved for x ($x = F^{-1}(y)$) the transformation offers the most effective (efficient) way to generate random numbers.
- A simple example comes from the exponential decay:

$$y = F(x) = \int_0^{x(y)} \frac{1}{\tau} e^{-x'/\tau} dx' \implies x(y) = -\tau \log(1 - y) \quad (7)$$

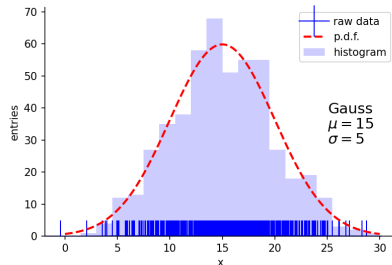
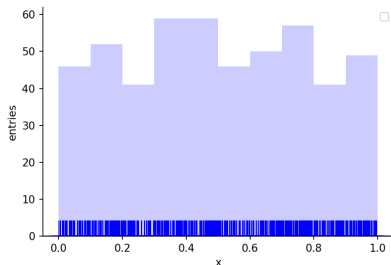


The Monte Carlo method

Using transformation for random number generation

- If the Eq. 5 can be solved for x ($x = F^{-1}(y)$) the transformation offers the most effective (efficient) way to generate random numbers.
- Numerically, it can also be done for a Gaussian:

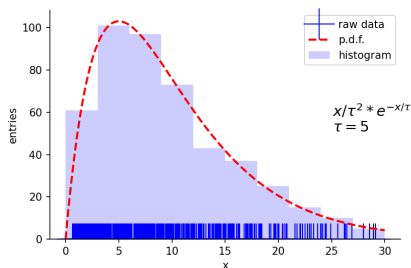
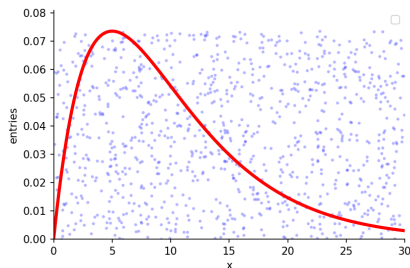
$$y = F(x) = \int_{-\infty}^{x(y)} G(x'; \mu, \sigma) dx' \quad \Rightarrow \quad x(y) = \sqrt{2}\sigma \operatorname{erf}^{-1}(2y - 1) + \mu \quad (8)$$



The Monte Carlo method

Acceptance-rejection for random number generation

- In more complex cases the acceptance-rejection method proves handy (also generally less efficient).
 - 1) Generate pairs of random numbers satisfying: $x = x_{\min} + r_1(x_{\max} - x_{\min})$, $u = r_2 f_{\max}$.
 - 2) If $u < f(x)$, then accept, otherwise reject.



Applications of the Monte Carlo method

- Various tests of estimators (we have seen multiple examples).
- Numerical integration. The accuracy improves as $1/\sqrt{N}$. Numerical integration using trapezoidal rule follows $1/N^{2/d}$ where d is the domain dimension. In 1D it is much faster but for $d > 4$ Monte Carlo comes into the game!
- Determination of p.d.f. of a function of random variables in more general case when $y(x_1, \dots, x_n)$ when p.d.f.'s of x_i are known.
- More advanced examples come e.g. from particle physics when experimental data are simulated using **event generators** and subsequently **detector simulation** programs. In both cases randomness of the physics process is the paramount requirement.

From measured to true distributions

Unfolding - general concepts

In experimental physics (more generally, science) we usually take measurement and compare them to some theoretical models. What we observe (measure) in an experiment is never perfect. There are two basic ways around this problem:

- 1 Generate Monte Carlo events according to our model and add all experimental effects using simulation. Then, we can compare experimental data directly to the Monte Carlo prediction,
- 2 Given the observed data construct *estimator* for the underlying true distribution - called **unfolding**.

The first option is usually easier, but does not allow for comparison outside of the experimental context (e.g. at a later time) and does not allow for direct comparison or combination of different experiments. Unfolding, although more involved, gives considerably larger flexibility of interpretation.

From measured to true distributions

Unfolding - general concepts

Let us take a measurement of a spectrum with some features corresponding to measured properties (resonances). We have various *experimental effects* which impact the observed spectrum:

- resolution, which causes migration between bins: the distribution is “smeared out”, peaks broadened,
- efficiency, which causes some events go undetected,
- background, which causes extra events due to spurious processes.

$$f_{\text{meas}}(x) = \int R(x|y) f_{\text{true}}(y) dy \quad (9)$$

Here we consider histogrammed (discretized) data $\mathbf{n} = (n_1, \dots, n_N)$

$$\nu_i = E[n_i] = \sum_{j=1}^M R_{ij} \mu_j + \beta_i, \quad \mu_j = \mu_{\text{tot}} \underbrace{\int_{\text{bin } j} f_{\text{true}}(y) dy}_{p_j}, \quad i = 1, \dots, N \quad (10)$$

From measured to true distributions

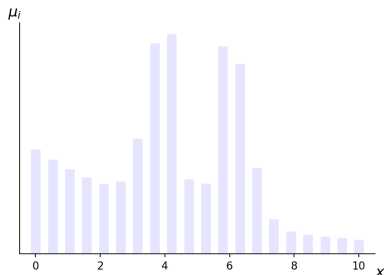
Unfolding - general concepts

$$R_{ij} = P(\text{observed in bin } i \mid \text{true in bin } j). \quad (11)$$

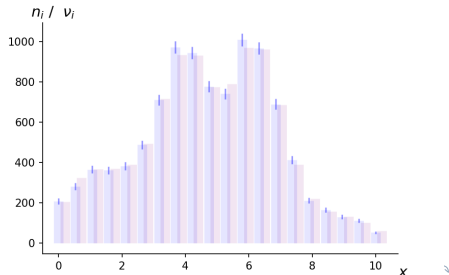
Note that:

$$\sum_{i=1}^N R_{ij} = \varepsilon_j \quad \leftarrow \text{efficiency} \quad (12)$$

“true” distribution histogram



“observed” events histogrammed



From measured to true distributions

Unfolding - general concepts

In the following we shall assume no background or known background contribution which can be subtracted prior to unfolding.

We have:

$$\boldsymbol{\nu} = R\boldsymbol{\mu} + \boldsymbol{\beta} \quad (13)$$

And assume it can be solved (inverted):

$$\boldsymbol{\mu} = R^{-1}\boldsymbol{\nu} \quad (14)$$

Suppose data are independent Poisson:

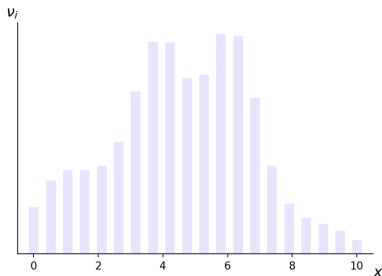
$$P(n_i; \nu_i) = \frac{\nu_i^{n_i}}{n_i!} e^{-\nu_i} \quad \implies \quad \log L(\boldsymbol{\mu}) = \sum_{i=1}^N (n_i \ln \nu_i - \nu_i) \quad \implies \quad \hat{\boldsymbol{\nu}} = \mathbf{n}$$
$$\longrightarrow \hat{\boldsymbol{\mu}} = R^{-1}\mathbf{n} \quad (15)$$

From measured to true distributions

Unfolding - putting it to work

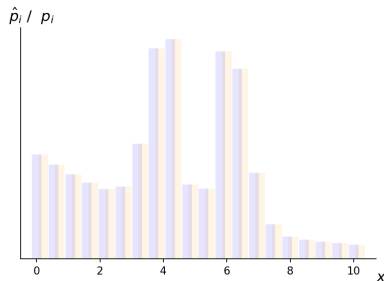
Let us take a finite Gaussian resolution ($\sigma = 0.6$) and an efficiency which linearly grows from 0.5 at 0 to 1 at 10.

“smeared” distribution histogram



$$\nu_i = \sum_{j=1}^M R_{ij} \mu_j$$

“unfolded” histogram



$$\hat{\mu}_i = \sum_{j=1}^M R_{ij}^{-1} \nu_j$$

From measured to true distributions

Unfolding - covariance matrix

For Poisson data the ML estimators are unbiased:

$$E[\hat{\boldsymbol{\mu}}] = R^{-1}E[\mathbf{n}] = \boldsymbol{\mu} \quad (16)$$

What is the covariance matrix of the unfolded histogram ($\hat{\boldsymbol{\mu}}$)?

$$U_{ij} = \text{cov}[\hat{\mu}_i, \hat{\mu}_j] = \sum_{k,l=1}^N (R^{-1})_{ik} (R^{-1})_{jl} \text{cov}[n_k, n_l] \quad (17)$$

or in short: $U = R^{-1}V(R^{-1})^T$

Recall the RCF bound:

$$(U^{-1})_{kl} = -E \left[\frac{\partial^2 \log L}{\partial \mu_k \partial \mu_l} \right] = \sum_{i=1}^N \frac{R_{ik} R_{il}}{\nu_i} \quad (18)$$
$$\implies U_{ij} = \sum_{k,l=1}^N (R^{-1})_{ik} (R^{-1})_{jl} \nu_k, \quad \underbrace{(\text{cov}[n_k, n_l] = \delta_{kl} \nu_k)}_{\text{independent Poisson}}.$$

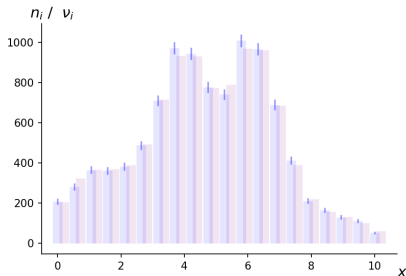
Unfolding realises the minimal variance among all unbiased estimators!

From measured to true distributions

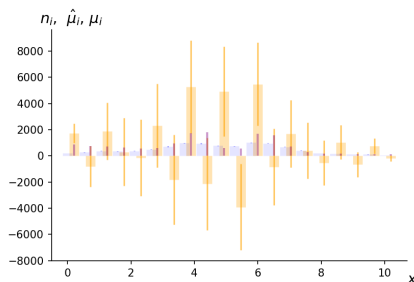
Unfolding - real life

Take the same Gaussian resolution ($\sigma = 0.6$) and an efficiency which linearly grows from 0.5 at 0 to 1 at 10 and generate 10,000 random events:

real data histogram



unfolded data histogram



Now, apply the same unfolding \rightarrow

$$\hat{\mu}_i = \sum_{j=1}^M R_{ij}^{-1} n_j$$

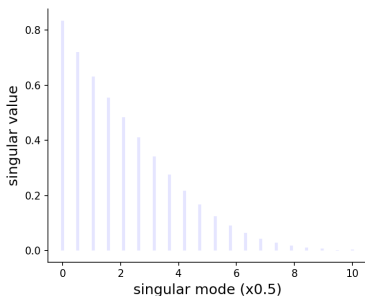
Huston, we have a problem!



From measured to true distributions

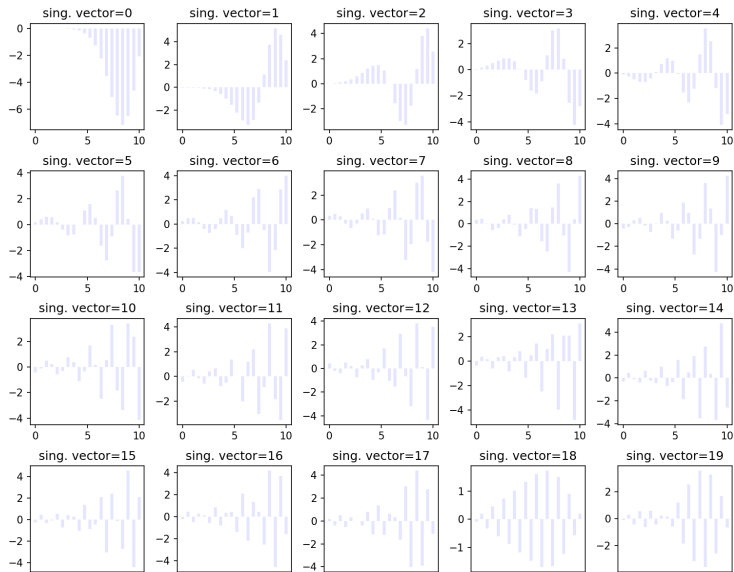
What has happened?

- 1 The perfect “unfolding” for the ν is assured by construction.
- 2 Data n seem very close but...
there are Poisson fluctuations w.r.t. expectation ν .
- 3 Smearing is larger than the width of one bin. Information about fine structure at this level is lost. This corresponds to the fact that solutions corresponding to the highest “frequency” modes are ill-defined.
- 4 Mathematically, such modes have very small singular values.



From measured to true distributions

The singular vectors (modes):



From measured to true distributions

Linear algebra - reminder

$$\mathcal{M}\mathbf{X} = \mathbf{Y} \implies \mathbf{X} = \mathcal{M}^{-1}\mathbf{Y} \quad (19)$$

Assuming \mathcal{M} a square matrix we can perform *Singular Value Decomposition* (SVD):

$$\mathcal{M} = S^T \mathcal{D} V \quad (20)$$

$$\underbrace{S \mathcal{M} V^T}_{\mathcal{D}} \underbrace{V \mathbf{X}}_{\tilde{\mathbf{X}}} = \underbrace{S \mathbf{Y}}_{\tilde{\mathbf{Y}}} \implies \tilde{\mathbf{X}} = \mathcal{D}^{-1} \tilde{\mathbf{Y}}, \quad (21)$$

where \mathcal{D} is a diagonal matrix with singular values λ_i on the diagonal and S (V) is a orthonormal matrix with eigenvectors \mathbf{S}_i (\mathbf{V}_i) as rows.

$$\mathcal{D} \begin{pmatrix} \lambda_1 & \dots & 0 \\ \dots & \dots & \dots \\ 0 & \dots & \lambda_N \end{pmatrix}, \quad U \propto \sum_i \frac{1}{\lambda_i^2} \quad (22)$$

From measured to true distributions

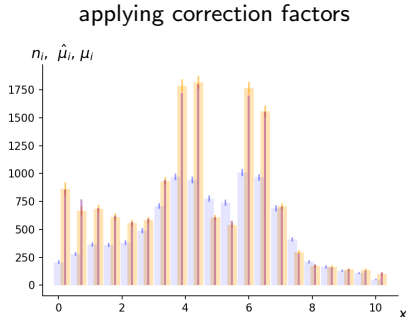
Correction factors

By far, the simplest method of correcting distributions is by applying **correction factors**.

$$\hat{\mu}_i = C_i n_i, \quad \text{with} \quad C_i = \frac{\mu_i^{\text{MC}}}{\nu_i^{\text{MC}}} \quad \leftarrow \text{from MC.} \quad (23)$$

$$U_{ij} = \text{cov}[\hat{\mu}_i, \hat{\mu}_j] = C_i C_j \text{cov}[n_i, n_j] = \delta_{ij} C_i^2 n_i \quad \leftarrow \text{small!}. \quad (24)$$

We get a seemingly ideal response!
but...



From measured to true distributions

Correction factors

Now, real peaks have been shifted by one left.

This method is NOT unbiased. We actually have:

$$b_i = E[\hat{\mu}_i] - \mu_i \longrightarrow b_i = \left(\frac{\mu_i^{\text{MC}}}{\nu_i^{\text{MC}}} - \frac{\mu_i}{\nu_i} \right) \quad (25)$$

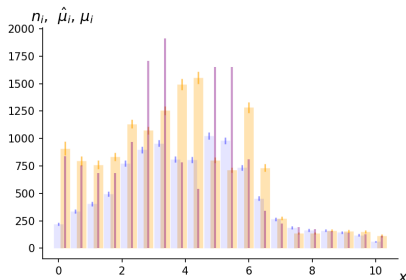
When reality is different from the assumed model, we get it completely wrong ☹️

In other words, correction factor method works **ONLY** if we are sure of the model!

In particular:

- **bad** if we don't know the exact shape of the signal.
- **good** if we want to correct for the known inefficiency.

applying correction factors

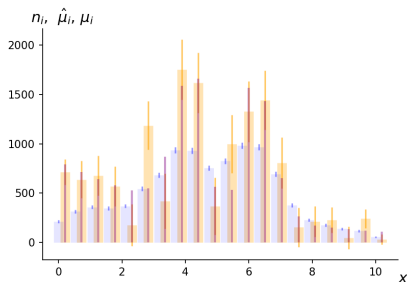


From measured to true distributions

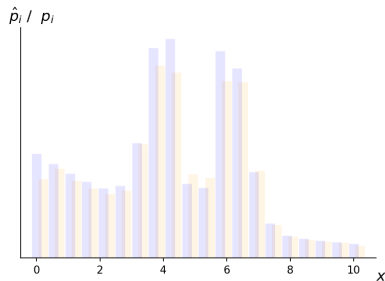
Unfolding - with conditioning

We execute the same unfolding but using the R^{-1} matrix with suppressed *weak modes*. This can be done in various ways. Here a very simple conditioning of the R is used, just to demonstrate the principle.

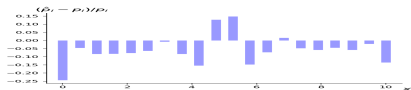
unfolded with conditioned R^{-1}



same applied to ν



Bias can be assessed
from unfolding $\nu \rightarrow$



That's way better!



Regularized unfolding

The problem of unfolding are typically large uncertainties on “high frequency” modes which are unphysical and are driven by the statistical fluctuations.

Regularisation: maximize the following with respect to $\boldsymbol{\mu}$:

$$\Phi(\boldsymbol{\mu}) = \alpha \log L(\boldsymbol{\mu}) + S(\boldsymbol{\mu}) \quad (26)$$

$S(\boldsymbol{\mu})$ regularization function (measure of smoothness),
 α regularization parameter (tradeoff between $\log L$ and S).

In addition require $\sum_i^N \nu_i = \sum_{i,j} R_{ij} \mu_j = n_{\text{tot}}$, i.e. maximize:

$$\phi(\boldsymbol{\mu}, \lambda) = \alpha \log L(\boldsymbol{\mu}) + S(\boldsymbol{\mu}) + \lambda \left[n_{\text{tot}} - \sum_i^N \nu_i \right], \quad (27)$$

where λ is a Lagrange multiplier: $\frac{\partial \phi}{\partial \lambda} = 0 \Rightarrow \sum_i^N \nu_i = n_{\text{tot}}$.

Regularized unfolding

One needs: a) regularization function $S(\boldsymbol{\mu})$, b) a prescription for choosing α .

1 Tikhonov regularization:

$$S[f_{\text{true}}(y)] = - \int \left(\frac{d^k f_{\text{true}}(y)}{dy^k} \right)^2, \quad \text{where } k = 1, 2, \dots \quad (28)$$

E.g. for $k = 2$ and $\log L = -\frac{1}{2}\chi^2$ one gets:

$$\phi(\boldsymbol{\mu}, \lambda) = -\frac{\alpha}{2}\chi^2(\boldsymbol{\mu}) - \sum_{i=2}^{N-1} (-\mu_{i-1} + 2\mu_i - \mu_{i+1})^2, \quad (29)$$

which after setting derivatives equal zero, gives a system of linear equations solvable using e.g. SVD.

2 Entropy: $H = -\sum_{i=1}^N p_i \ln p_i$ (max when all p_i equal).

Use entropy-based regularization function:

$$S(\boldsymbol{\mu}) = H(\boldsymbol{\mu}) = -\sum_{i=1}^N \frac{\mu_i}{\mu_{\text{tot}}} \ln \frac{\mu_i}{\mu_{\text{tot}}}, \quad (30)$$

$\propto \ln(\text{number of ways to arrange } \mu_{\text{tot}} \text{ entries in } M \text{ bins})$

From measured to true distributions

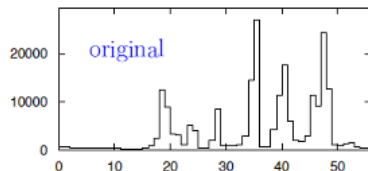
Unfolding - image unblure

The **entropy-based** unfolding of a blurred image.

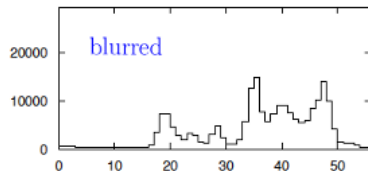
credit: G. Cowan



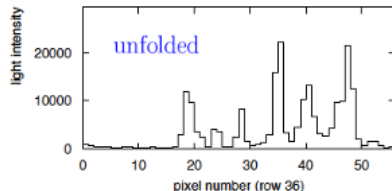
(a)



(b)



(c)



Regularized unfolding

Choosing α parameter

One needs a trade-off between bias and variance:

Common choices include:

1

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (U_{ii} + \hat{b}_i^2), \quad \text{or} \quad \text{MSE}' = \frac{1}{N} \sum_{i=1}^N \frac{U_{ii} + \hat{b}_i^2}{\hat{\mu}_i}.$$

2 Allow for a “reasonable” change of the χ^2 :

$$\Delta\chi^2 = 2\Delta \log L = N$$

3 Require bias be consistent with zero within its uncertainty:

$$\chi_b^2 = \sum_{i=1}^N \frac{\hat{b}_i^2}{W_{ii}} = N, \quad \text{where} \quad W_{ij} = \text{cov}[\hat{b}_i, \hat{b}_j].$$

Note: there is no optimal choice. All depends on the particular analysis context.

Thank you

Back-up