

Response Matrix Estimation in Unfolding Differential Cross Sections

Mikael Kuusela

Department of Statistics and Data Science
Carnegie Mellon University

Quark Confinement and the Hadron Spectrum 2024
Cairns, Queensland, Australia

August 22, 2024

*Joint work with Richard Zhu (CMU), Larry Wasserman (CMU)
and Andrea Marini (CERN)*

The unfolding problem

- Any differential cross section measurement is affected by the finite resolution of the particle detectors
 - This causes the observed spectrum of events to be “smeared” or “blurred” with respect to the true one
- The *unfolding problem* is to estimate the true spectrum using the smeared observations
- Ill-posed inverse problem with many methodological challenges

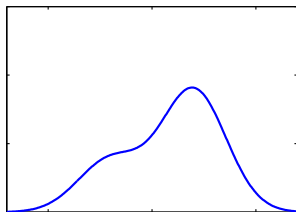


Figure: Smeared spectrum

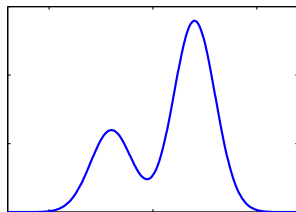
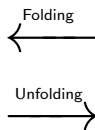


Figure: True spectrum

Problem formulation

- Let f be the true, particle-level spectrum and g the smeared, detector-level spectrum
 - Denote the true space by T and the smeared space by S (both taken to be intervals on the real line for simplicity)
 - Mathematically f and g are the intensity functions of the underlying Poisson point process
- The two spectra are related by

$$g(s) = \int_T k(s, t) f(t) dt,$$

where the smearing kernel k represents the response of the detector and is given by

$$k(s, t) = p(Y = s | X = t, X \text{ observed}) P(X \text{ observed} | X = t),$$

where X is a true event and Y the corresponding smeared event

Task: Infer the true spectrum f given smeared observations from g

Discretization

- Problem usually discretized using histograms (splines are also sometimes used)
- Let $\{T_i\}_{i=1}^p$ and $\{S_i\}_{i=1}^n$ be binnings of the true space T and the smeared space S
- Smeared histogram $\mathbf{y} = [y_1, \dots, y_n]^T$ with mean

$$\boldsymbol{\mu} = \left[\int_{S_1} g(s) ds, \dots, \int_{S_n} g(s) ds \right]^T$$

- Quantity of interest:

$$\boldsymbol{\lambda} = \left[\int_{T_1} f(t) dt, \dots, \int_{T_p} f(t) dt \right]^T$$

- The mean histograms are related by $\boldsymbol{\mu} = \mathbf{K}\boldsymbol{\lambda}$, where the elements of the *response matrix* \mathbf{K} are given by

$$K_{i,j} = \frac{\int_{S_i} \int_{T_j} k(s, t) f(t) dt ds}{\int_{T_j} f(t) dt} = P(\text{smeared event in bin } i \mid \text{true event in bin } j)$$

- The discretized statistical model becomes

$$\mathbf{y} \sim \text{Poisson}(\mathbf{K}\boldsymbol{\lambda})$$

and we wish to make inferences about $\boldsymbol{\lambda}$ under this model

Regularized unfolding

- When the number of true bins p is large, the response matrix \mathbf{K} is severely ill-conditioned
- The unfolded histogram λ is therefore typically estimated using a *regularized* estimator
 - Main idea: **bias** \uparrow , **variance** $\downarrow \Rightarrow$ **MSE** \downarrow
- Two main approaches:

- 1 Tikhonov regularization (e.g., SVD by Höcker and Kartvelishvili (1996) and TUnfold by Schmitt (2012)):

$$\min_{\lambda \in \mathbb{R}^p} (\mathbf{y} - \mathbf{K}\lambda)^T \mathbf{C}^{-1} (\mathbf{y} - \mathbf{K}\lambda) + \delta P(\lambda)$$

with

$$P_{\text{SVD}}(\lambda) = \left\| \mathbf{L} \begin{bmatrix} \lambda_1 / \lambda_1^{\text{MC}} \\ \lambda_2 / \lambda_2^{\text{MC}} \\ \vdots \\ \lambda_p / \lambda_p^{\text{MC}} \end{bmatrix} \right\|^2 \quad \text{or} \quad P_{\text{TUnfold}}(\lambda) = \|\mathbf{L}(\lambda - \lambda^{\text{MC}})\|^2,$$

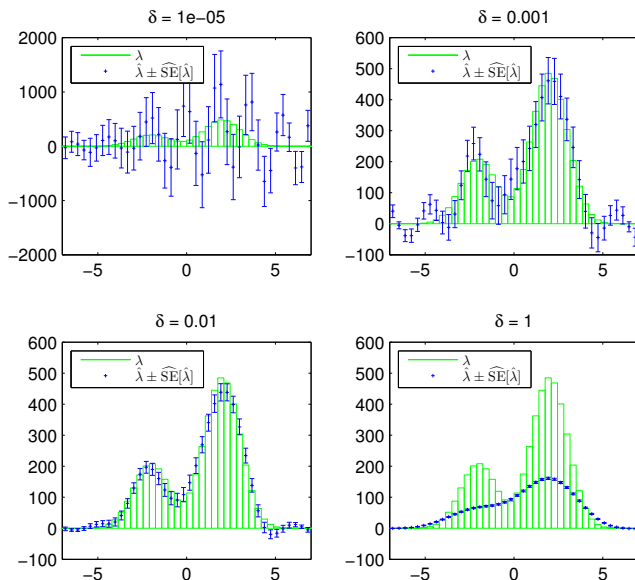
where \mathbf{L} is usually the discretized second derivative (other choices also possible)

- 2 Expectation-maximization iteration with early stopping (D'Agostini, 1995):

$$\lambda_j^{(t+1)} = \frac{\lambda_j^{(t)}}{\sum_{i=1}^n K_{i,j}} \sum_{i=1}^n \frac{K_{i,j} y_i}{\sum_{k=1}^p K_{i,k} \lambda_k^{(t)}}, \quad \text{with } \lambda^{(0)} = \lambda^{\text{MC}}$$

- These methods typically regularize by creating a bias toward a MC ansatz λ^{MC}

Tikhonov regularization, $P(\lambda) = \|\lambda\|^2$, varying δ



D'Agostini demo, $k = 0$

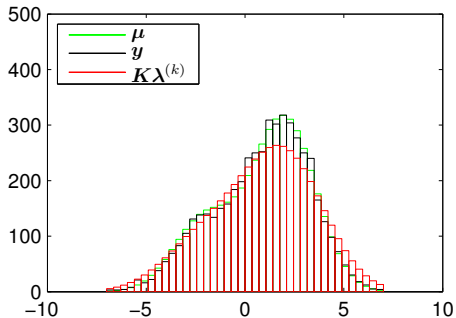


Figure: Smearing histogram

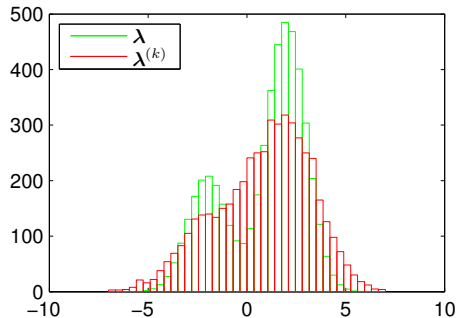


Figure: True histogram

D'Agostini demo, $k = 100$

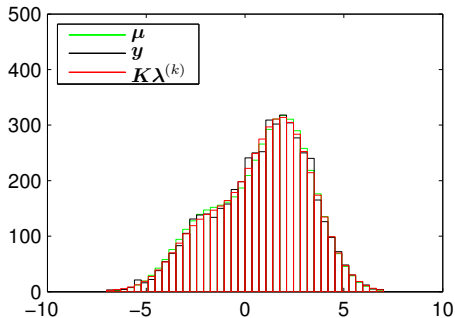


Figure: Smearing histogram

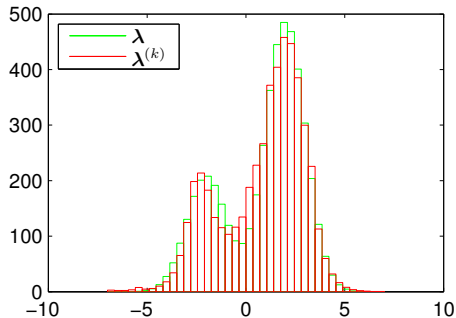


Figure: True histogram

D'Agostini demo, $k = 10000$

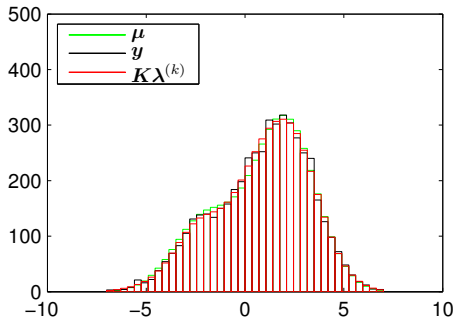


Figure: Smearing histogram

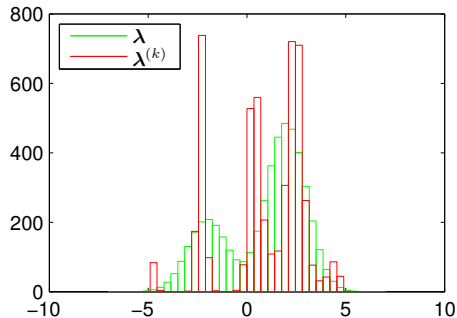


Figure: True histogram

D'Agostini demo, $k = 100000$

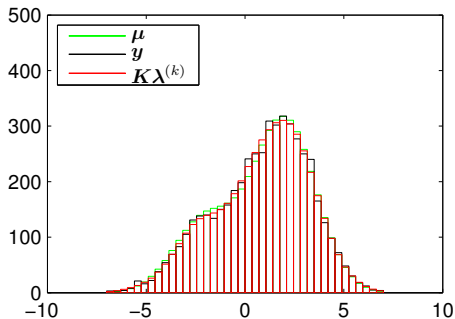


Figure: Smearred histogram

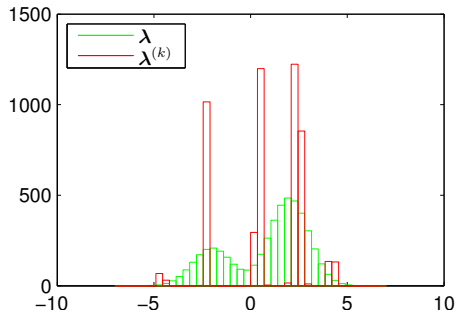


Figure: True histogram

Response matrix estimation

An aspect of the unfolding problem that has received relatively less attention is that in practice we do not know the response kernel $k(s, t)$ or the response matrix \mathbf{K}

Instead, we are given paired MC samples (Y_i, X_i) from

$$p^{\text{MC}}(Y = s, X = t) = p^{\text{MC}}(Y = s | X = t)p^{\text{MC}}(X = t)$$

These samples are then used to produce an estimator $\hat{\mathbf{K}}$ of \mathbf{K} which is used in the chosen unfolding method as if it were the true matrix \mathbf{K}

For example, Tikhonov regularization for known \mathbf{K} is

$$\hat{\lambda} = (\mathbf{K}^T \mathbf{K} + \delta \mathbf{I})^{-1} \mathbf{K}^T \mathbf{y},$$

but in practice we use

$$\hat{\lambda} = (\hat{\mathbf{K}}^T \hat{\mathbf{K}} + \delta \mathbf{I})^{-1} \hat{\mathbf{K}}^T \mathbf{y}$$

This raises the following questions:

- 1 Which estimator $\hat{\mathbf{K}}$ should one use?
- 2 How does the estimated matrix affect the unfolded solution?

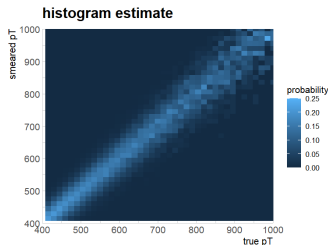
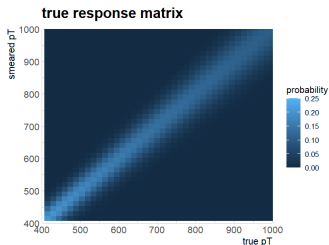
Binned estimator

In most cases, \mathbf{K} is understood as a 2D (conditional) histogram whose entries are estimated as follows:

$$\hat{K}_{i,j} = \frac{\# \text{ events originating from bin } j \text{ that have been recorded in bin } i}{\# \text{ events originating from bin } j}$$

Assuming f^{MC} is correct, this gives an unbiased estimator of \mathbf{K}

However, this estimator can be very noisy, especially in the tails of steeply falling spectra where there are few MC events available:



Estimating the response kernel

We should be able to do better by borrowing strength from nearby bins

We propose to do this by estimating the response kernel

$k(s, t) = p(Y = s | X = t)$ based on the MC samples (Y_i, X_i)

Doing this is a well-studied problem in statistics called nonparametric *conditional density estimation* (CDE)

- Lots of existing methods and software available to produce the estimator

Once we have the estimator $\hat{k}(s, t)$, we can use it to obtain an estimator of the response matrix elements:

$$\hat{K}_{i,j} = \frac{\int_{S_i} \int_{T_j} \hat{k}(s, t) f^{\text{MC}}(t) dt ds}{\int_{T_j} f^{\text{MC}}(t) dt}$$

Note also that if we discretize $f(t)$ using a basis expansion (such as splines), we also need $\hat{k}(s, t)$ to obtain an estimate of the discretized forward model

Estimating the response kernel

We consider the following nonparametric estimators of the response kernel:

1 Kernel CDE:

$$\begin{aligned}\hat{p}(y|x) &= \arg \min_a \sum_{i=1}^n (K_{h_2}(y - Y_i) - a)^2 K_{h_1}(x - X_i) \\ &= \sum_{i=1}^n w_i(x) K_{h_2}(y - Y_i), \text{ where } w_i(x) = \frac{K_{h_1}(x - X_i)}{\sum_{j=1}^n K_{h_1}(x - X_j)}\end{aligned}$$

2 Local linear CDE:

$$\begin{aligned}(\hat{a}, \hat{b}) &= \arg \min_{a,b} \sum_{i=1}^n (K_{h_2}(y - Y_i) - a - b(X_i - x))^2 K_{h_1}(x - X_i) \\ \hat{p}(y|x) &= \hat{a}\end{aligned}$$

- 3 **Kernel CDE with local bandwidths:** Make the bandwidths h_1 and h_2 functions of x by estimating them within some window of size $\delta(x)$ around x

$$\hat{p}(y|x) = \sum_{i:|x-X_i|<\delta(x)} w_i(x) K_{h_2(x)}(y - Y_i),$$

where

$$w_i(x) = \frac{K_{h_1(x)}(x - X_i)}{\sum_{j:|x-X_j|<\delta(x)} K_{h_1(x)}(x - X_j)}$$

Estimating the response kernel

- ④ **Location-scale model:** Assume the following model for the smeared observations

$$Y = \mu(X) + \sigma(X)\varepsilon,$$

where ε has some distribution p_ε with mean 0 and variance 1

It follows that

$$p(y|x) = \frac{1}{\sigma(x)} p_\varepsilon \left(\frac{y - \mu(x)}{\sigma(x)} \right)$$

We can estimate $\mu(x)$ and $\sigma(x)$ using nonparametric regression and p_ε using KDE on the standardized observations $(y_i - \hat{\mu}(x_i))/\hat{\sigma}(x_i)$

The estimated response kernel is then

$$\hat{p}(y|x) = \frac{1}{\hat{\sigma}(x)} \hat{p}_\varepsilon \left(\frac{y - \hat{\mu}(x)}{\hat{\sigma}(x)} \right)$$

Simulation setup

The following simulation study is designed to mimic the unfolding of inclusive jet p_{\perp} spectrum at the LHC

The particle-level spectrum is

$$f(p_{\perp}) = LN_0 \left(\frac{p_{\perp}}{\text{GeV}} \right)^{-\alpha} \left(1 - \frac{2}{\sqrt{s}} p_{\perp} \right)^{\beta} e^{-\gamma/p_{\perp}}, \quad 0 < p_{\perp} \leq \frac{\sqrt{s}}{2},$$

where $L, N_0, \alpha, \beta, \gamma, \sqrt{s}$ are parameters set to mimic conditions at the LHC

The response kernel is $k(p'_{\perp}, p_{\perp}) = N(p'_{\perp} - p_{\perp} | 0, \sigma(p_{\perp})^2)$ with

$$\left(\frac{\sigma(p_{\perp})}{p_{\perp}} \right)^2 = \left(\frac{C_1}{\sqrt{p_{\perp}}} \right)^2 + \left(\frac{C_2}{p_{\perp}} \right)^2 + C_3^2$$

so this is a heteroscedastic deconvolution problem

The problem is discretized using $n = p = 40$ bins over [400 GeV, 1000 GeV]

The resulting response matrix \mathbf{K} is severely ill-conditioned but not singular

Comparison of estimated response matrices

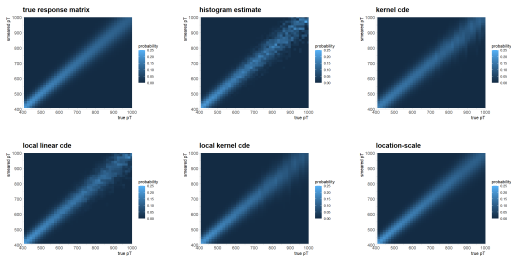


Figure: Estimated response matrices

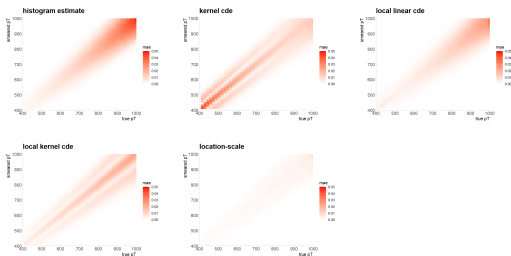


Figure: Mean absolute errors

Now let's see how the different response matrix estimators affect the quality of unfolded point estimators

We will consider regularized unfolding using

- 1 Tikhonov regularization
- 2 D'Agostini iteration

Tikhonov regularization

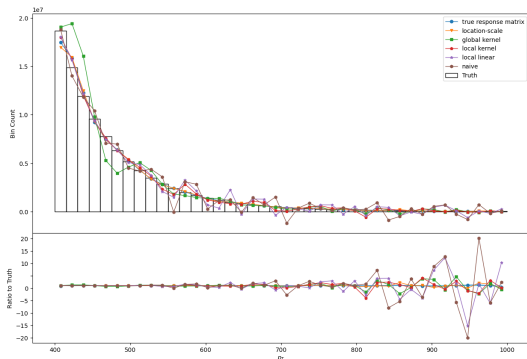


Figure: Tikhonov regularization solutions with $\delta = 10^{-10}$

With moderate regularization, things behave as expected: better response matrix estimators tend to give better unfolded histograms

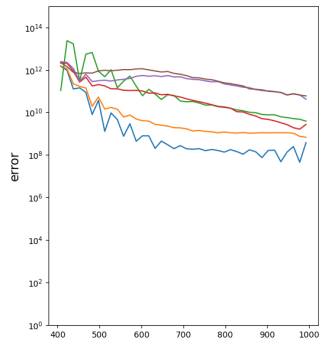


Figure: Mean squared error for Tikhonov regularization with $\delta = 10^{-10}$

Tikhonov regularization

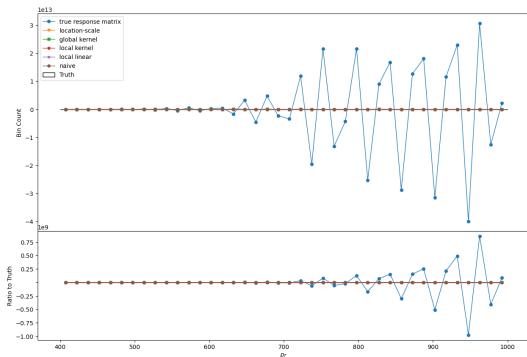


Figure: Tikhonov regularization solutions with $\delta = 0$

With no regularization (i.e., matrix inversion), we see something unexpected: the estimated response matrices give better unfolded histograms than the actual true response matrix

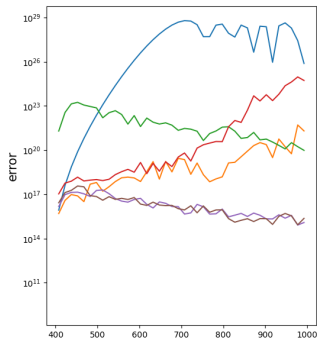


Figure: Mean squared error for Tikhonov regularization with $\delta = 0$

Implicit regularization

The explanation is that the estimated response matrices implicitly perform regularization (an ill-conditioned matrix with some additive random noise becomes well-conditioned with high probability (Tao and Vu, 2007))

Table: The median condition numbers for the estimated response matrices over $M = 1000$ simulations

Estimation method	Median condition number
True	$1.7 \cdot 10^{17}$
Kernel	$3.9 \cdot 10^7$
Local linear	$6.7 \cdot 10^3$
Local kernel	$1.5 \cdot 10^8$
Location-scale	$3.9 \cdot 10^4$
Naive histogram	$2.6 \cdot 10^3$

D'Agostini iteration

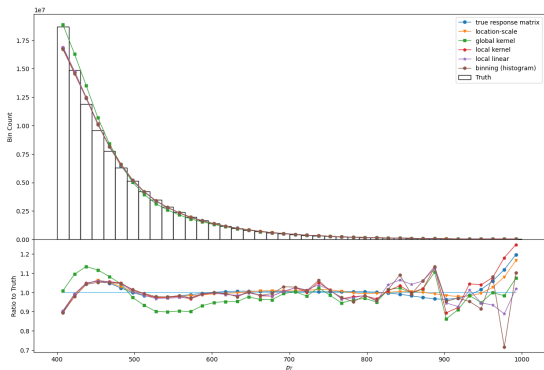


Figure: D'Agostini solutions with 5 iterations

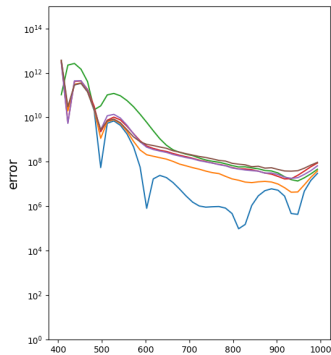


Figure: Mean squared error for D'Agostini with 5 iterations

With the D'Agostini iteration, most of the estimators behave similarly for a small number of iterations with the true matrix providing the best solution, as expected

D'Agostini iteration

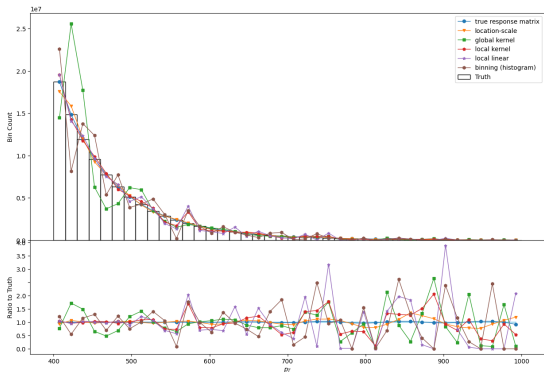


Figure: D'Agostini solutions with 5000 iterations

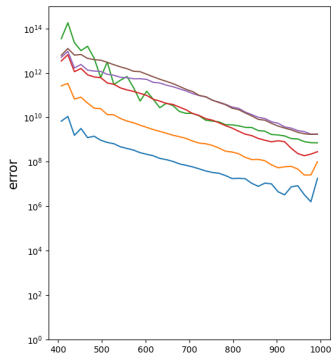


Figure: Mean squared error for D'Agostini with 5000 iterations

For a large number of iterations, differences emerge with better response matrix estimators providing overall better unfolded histograms

Here the true matrix always provides the best solution in contrast with Tikhonov with a vanishing regularization strength

Outlook: Unfolding with machine learning

A major development in the past few years: using machine learning to help solve the unfolding problem

Two main approaches:

- 1 **OmniFold** (Andreassen et al., 2020): iteratively reweight particle-level MC events using classifier-based density ratios
- 2 **Generative unfolding** (Bellagente et al., 2020): Train a generative model to sample from $p(X = t|Y = s)$; iterate to reduce dependence on $p^{\text{MC}}(X = t)$

Benefits of ML-based unfolding:

- Does not rely on binning
- Provides event-level unfolded results
- Can handle (moderately) high-dimensional phase spaces
- Does not need a separate estimate of the response kernel $k(s, t)$

There are many open questions regarding the type of regularization ML-based unfolding imposes on the unfolded solution

However, the fact that these methods don't need a plug-in estimate of the response kernel seems like a potential way to simplify that aspect of the unfolding problem

Discussion and conclusions

- Unfolding is a complex data analysis task with many statistical challenges
- Any binned unfolding method needs the response matrix \mathbf{K} as an input
 - In practice, a plug-in estimate $\hat{\mathbf{K}}$ is used
- We have introduced several new ways of estimating \mathbf{K} using conditional density estimation
- We found that there are non-trivial interactions between the estimate of \mathbf{K} and the unfolding method used
 - In particular, we found that noisy estimates of \mathbf{K} can implicitly regularize the problem in the absence of other regularization
 - Potential workaround: ML-based unfolding inverts the smeared data without needing a plug-in estimate of the forward operator
- Several other challenges related to the response matrix:
 - Dependence on the Monte Carlo ansatz (wide-bin bias)
 - Dependence on nuisance parameters
 - Other systematic uncertainties (e.g., confounding variables)
 - Error propagation to the unfolded solution
 - ...

- T. Auye. Unfolding algorithms and tests using RooUnfold. In H. B. Prosper and L. Lyons, editors, *Proceedings of the PHYSTAT 2011 Workshop on Statistical Issues Related to Discovery Claims in Search Experiments and Unfolding*, CERN-2011-006, pages 313–318, CERN, Geneva, Switzerland, 17–20 January 2011.
- A. Andreassen, P. T. Komiske, E. M. Metodiev, B. Nachman, and J. Thaler, OmniFold: A method to simultaneously unfold all observables, *Physical Review Letters*, 124: 182001, 2020. doi: 10.1103/PhysRevLett.124.182001.
- P. Batlle, P. Patil, M. Stanley, H. Owhadi, and M. Kuusela, Optimization-based frequentist confidence intervals for functionals in constrained inverse problems: Resolving the Burrus conjecture, Preprint arXiv:2310.02461 [math.ST], 2024.
- M. Bellagente, A. Butter, G. Kasieczka, T. Plehn, A. Rousselot, R. Winterhalder, L. Ardizzone, and U. Köthe, Invertible networks or partons to detector and back again, *SciPost Phys.*, 9:074, 2020. doi: 10.21468/SciPostPhys.9.5.074.
- L. Brenner, R. Balasubramanian, C. Burgard, W. Verkerke, G. Cowan, P. Verschuur, and V. Croft, Comparison of unfolding methods using RooFitUnfold, *International Journal of Modern Physics A*, 35(24):2050145, 2020.
- G. D’Agostini, A multidimensional unfolding method based on Bayes’ theorem, *Nuclear Instruments and Methods A*, 362:487–498, 1995.

References II

- A. P. Dempster, N. M. Laird, and D. B. Rubin, Maximum likelihood from incomplete data via the EM algorithm, *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1):1–38, 1977.
- P. J. Green and B. W. Silverman. *Nonparametric Regression and Generalized Linear Models: A Roughness Penalty Approach*. Chapman & Hall, 1994.
- P. C. Hansen, Analysis of discrete ill-posed problems by means of the L-curve, *SIAM Review*, 34(4):561–580, 1992.
- A. Höcker and V. Kartvelishvili, SVD approach to data unfolding, *Nuclear Instruments and Methods in Physics Research A*, 372:469–481, 1996.
- M. Kuusela. *Uncertainty quantification in unfolding elementary particle spectra at the Large Hadron Collider*. PhD thesis, EPFL, 2016. Available online at: <https://infoscience.epfl.ch/record/220015>.
- M. Kuusela and V. M. Panaretos, Statistical unfolding of elementary particle spectra: Empirical Bayes estimation and bias-corrected uncertainty quantification, *The Annals of Applied Statistics*, 9(3):1671–1705, 2015.
- K. Lange and R. Carson, EM reconstruction algorithms for emission and transmission tomography, *Journal of Computer Assisted Tomography*, 8(2):306–316, 1984.

References III

- L. B. Lucy, An iterative technique for the rectification of observed distributions, *Astronomical Journal*, 79(6):745–754, 1974.
- W. H. Richardson, Bayesian-based iterative method of image restoration, *Journal of the Optical Society of America*, 62(1):55–59, 1972.
- S. Schmitt, TUnfold, an algorithm for correcting migration effects in high energy physics, *Journal of Instrumentation*, 7:T10003, 2012.
- L. A. Shepp and Y. Vardi, Maximum likelihood reconstruction for emission tomography, *IEEE Transactions on Medical Imaging*, 1(2):113–122, 1982.
- M. Stanley, P. Patil, and M. Kuusela, Uncertainty quantification for wide-bin unfolding: one-at-a-time strict bounds and prior-optimized confidence intervals, *Journal of Instrumentation*, 17(10):P10013, 2022.
- M. Stanley, P. Batlle, P. Patil, H. Owhadi, and M. Kuusela, Confidence intervals for functionals in constrained inverse problems via data-adaptive sampling-based calibration, In preparation, 2024.
- P. B. Stark, Inference in infinite-dimensional inverse problems: Discretization and duality, *Journal of Geophysical Research*, 97(B10):14055–14082, 1992.

- T. Tao and V. H. Vu. The condition number of a randomly perturbed matrix. In *Proceedings of the Thirty-Ninth Annual ACM Symposium on Theory of Computing*, pages 248–255. Association for Computing Machinery, 2007.
- Y. Vardi, L. A. Shepp, and L. Kaufman, A statistical model for positron emission tomography, *Journal of the American Statistical Association*, 80(389):8–20, 1985.
- E. Veklerov and J. Llacer, Stopping rule for the MLE algorithm based on statistical hypothesis testing, *IEEE Transactions on Medical Imaging*, 6(4):313–319, 1987.
- I. Volobouev. On the expectation-maximization unfolding with smoothing. arXiv:1408.6500v2 [physics.data-an], 2015.

Backup

Uncertainty quantification in unfolding

- Let's assume that we are interested in some linear functional $\theta = \mathbf{h}^T \boldsymbol{\lambda}$ of $\boldsymbol{\lambda}$ (or potentially some collection of functionals)
 - For example, $\theta = \mathbf{e}_i^T \boldsymbol{\lambda} = i$ th unfolded bin or
 $\theta =$ average of several unfolded bins or $\theta = \mathbf{1}^T \boldsymbol{\lambda} =$ sum of all unfolded bins
- We can use $\hat{\theta} = \mathbf{h}^T \hat{\boldsymbol{\lambda}}$ as a natural point estimator of θ
- For uncertainty quantification, our goal is to find a random interval $[\underline{\theta}(\mathbf{y}), \bar{\theta}(\mathbf{y})]$ with *coverage probability* $1 - \alpha$:

$$P(\theta \in [\underline{\theta}(\mathbf{y}), \bar{\theta}(\mathbf{y})]) \approx 1 - \alpha$$

- Most implementations construct the interval based on the variance of $\hat{\theta}$:

$$[\underline{\theta}, \bar{\theta}] = \left[\hat{\theta} - z_{1-\alpha/2} \sqrt{\text{var}(\hat{\theta})}, \hat{\theta} + z_{1-\alpha/2} \sqrt{\text{var}(\hat{\theta})} \right]$$

- But: These intervals may suffer from significant undercoverage because they ignore the **regularization bias**

An alternative approach to explicit regularization that has become increasingly popular in LHC data analysis is to simply use very few unfolded bins (i.e., use small p)

⇒ **Regularization using wide bins**

Intuition: The detector should not be able to recover features smaller than its intrinsic resolution so should chose

$$\text{bin size} \gtrsim \text{detector resolution}$$

This intuition is sound but the typical implementation is problematic

Wide-bin unfolding

The response matrix elements are:

$$K_{i,j} = \frac{\int_{S_i} \int_{T_j} k(s, t) f(t) dt ds}{\int_{T_j} f(t) dt}$$

This depends on the unknown intensity function f (specifically, the shape of f inside the true bins T_j)

To get around this, $K_{i,j}$ is approximated based on a MC ansatz f^{MC} :

$$K_{i,j}^{\text{MC}} = \frac{\int_{S_i} \int_{T_j} k(s, t) f^{\text{MC}}(t) dt ds}{\int_{T_j} f^{\text{MC}}(t) dt}$$

This means that unfolding is performed using an approximate matrix \mathbf{K}^{MC} instead of the true matrix \mathbf{K}

When p is small, one can typically unfold simply using the unregularized generalized least-squares estimator

$$\hat{\lambda}^{\text{MC}} = ((\mathbf{K}^{\text{MC}})^T \mathbf{C}^{-1} \mathbf{K}^{\text{MC}})^{-1} (\mathbf{K}^{\text{MC}})^T \mathbf{C}^{-1} \mathbf{y}$$

But this is biased because $\mathbf{K}^{\text{MC}} \neq \mathbf{K} \Rightarrow$ [Wide-bin bias](#)

Wide-bins-via-fine-bins unfolding

Because of the wide-bin bias, variability intervals based on $\hat{\lambda}^{\text{MC}}$ will undercover

We could try to inflate the intervals by an amount corresponding to the bias, but this bias is very difficult to estimate and quantify

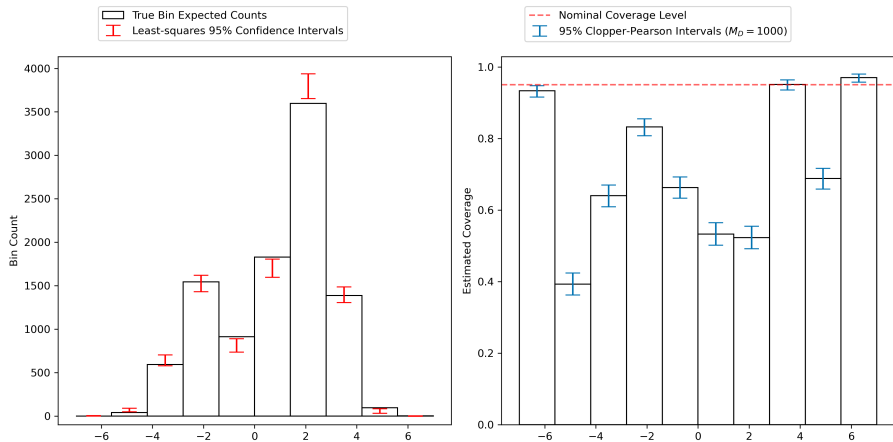
Alternative idea (Stanley et al., 2022):

The wide-bin bias gets reduced the smaller the bins in the true space

So we can *first unfold with fine bins (and no regularization) and then aggregate into wide bins, keeping track of the bin-to-bin correlations in the error propagation*

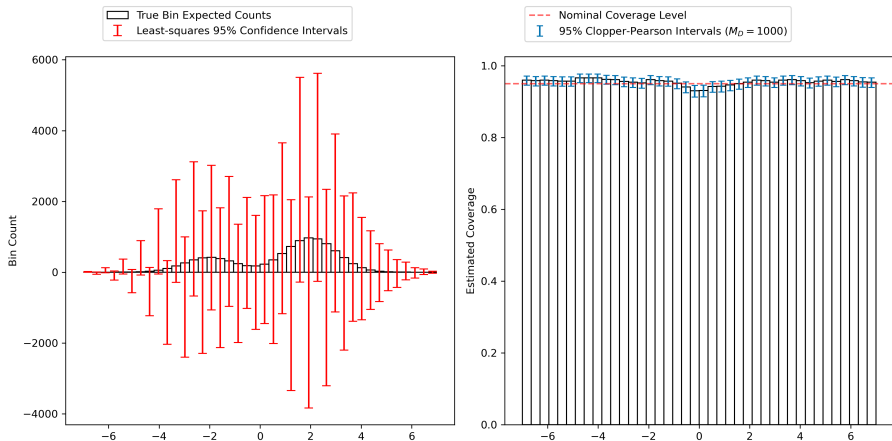
This [wide-bins-via-fine-bins unfolding](#) approach provides reasonably sized unfolded confidence intervals that do not suffer from regularization bias and have minimal wide-bin bias

Wide bins, standard approach, misspecified MC



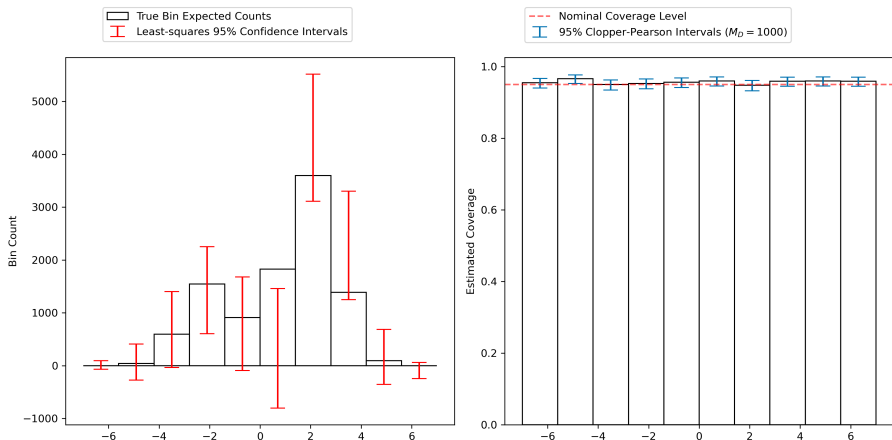
Intervals undercover because they ignore the wide-bin bias caused by the misspecified f^{MC}

Fine bins, standard approach, misspecified MC



With narrow bins, there is less dependence on f^{MC} so coverage is improved, but the intervals are very wide
⇒ Let's aggregate these into wide bins

Wide bins via fine bins, misspecified MC



With the same misspecified f^{MC} , wide-bins-via-fine-bins unfolding gives both correct coverage and reasonably sized intervals

Handling constraints and rank-deficient matrices

The previous example shows that the wide-bins-via-fine-bins approach can circumvent both the regularization bias and the wide-bin bias

But the simple approach based on the least-squares variability intervals has two important limitations:

- It cannot easily impose constraints (such as positivity) on the solution
- It cannot handle column-rank-deficient response matrices \mathbf{K} (such as when $\#$ of true bins $>$ $\#$ of smeared bins)

Handling constraints and rank-deficient matrices

In Stanley et al. (2022), we developed two new methods that can incorporate constraints and handle rank-deficient matrices:

- One-at-a-time strict bounds (OSB) intervals
- Prior-optimized (PO) intervals

The OSB intervals are a modification of the simultaneous strict bounds (SSB) intervals of Stark (1992) with the intervals designed to provide binwise coverage instead of simultaneous coverage

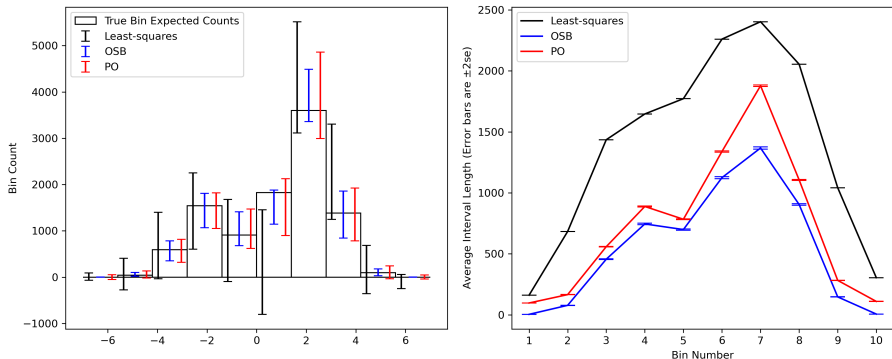
The PO intervals are decision-theoretic intervals where the interval length is optimized using a prior subject to a constraint on correct coverage¹

Both intervals have correct empirical coverage in most scenarios; PO also has a rigorous proof of coverage; details in Stanley et al. (2022)

¹Importantly, finite-sample frequentist coverage is guaranteed even for misspecified priors, but the interval length might be suboptimal in those cases.

Wide bins via fine bins, with positivity constraint

The interval lengths can be reduced by imposing a positivity constraint on the solution:



All of the above intervals have correct empirical coverage

Test inversion confidence intervals for unfolding

The OSB intervals are closely related to the inversion of the following test with respect to θ (Batlle et al., 2024):

$$H_0 : \boldsymbol{\lambda} \in \Phi_\theta \cap \mathcal{C} \quad \text{versus} \quad H_1 : \boldsymbol{\lambda} \in \mathcal{C} \setminus \Phi_\theta,$$

where $\Phi_\theta = \{\boldsymbol{\lambda} : \mathbf{h}^\top \boldsymbol{\lambda} = \theta\}$ and \mathcal{C} is a constrained set of solutions

In fact, they are equivalent to the inversion of the likelihood ratio test

$$\Lambda(\theta) = \frac{\sup_{\boldsymbol{\lambda} \in \Phi_\theta \cap \mathcal{C}} L(\boldsymbol{\lambda})}{\sup_{\boldsymbol{\lambda} \in \mathcal{C}} L(\boldsymbol{\lambda})}$$

assuming that the null distribution of $-2 \log \Lambda(\theta)$ is χ_1^2

In the presence of constraints (i.e., $\mathcal{C} \subsetneq \mathbb{R}^p$), this is only approximately true

We are currently finalizing a manuscript (Stanley et al., 2024) showing how to calibrate this test for high-dimensional $\boldsymbol{\lambda}$ using sampling and quantile regression

Test inversion confidence intervals for unfolding

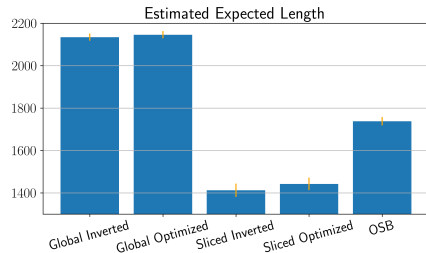
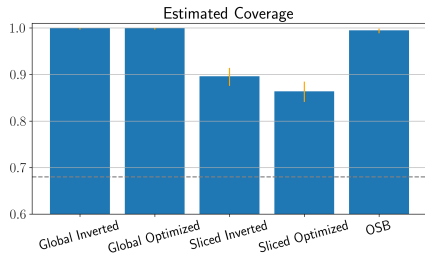


Figure: Test inversion intervals maintain nominal coverage (left panel) but are substantially shorter than the OSB intervals (right panel)

Test inversion confidence intervals for unfolding

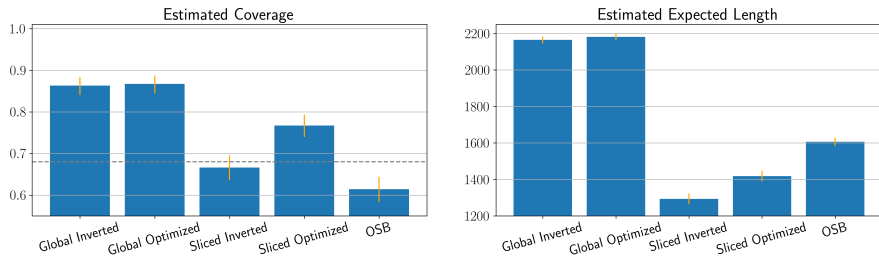


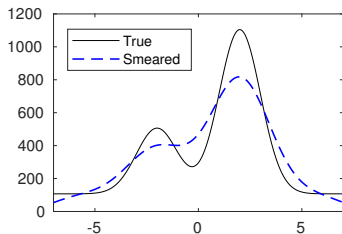
Figure: Test inversion intervals for an adversarial particle-level spectrum

In fact, if we approximate the Poisson noise using a Gaussian and use an affine estimator $\hat{\lambda}$ (e.g., Tikhonov-type estimators), then the coverage of the variability intervals can be written down in closed form (Kuusela, 2016):

$$\mathbb{P}(\theta \in [\underline{\theta}, \bar{\theta}]) = \Phi\left(\frac{\text{bias}(\hat{\theta})}{\sqrt{\text{var}(\hat{\theta})}} + z_{1-\alpha/2}\right) - \Phi\left(\frac{\text{bias}(\hat{\theta})}{\sqrt{\text{var}(\hat{\theta})}} - z_{1-\alpha/2}\right)$$

These intervals have coverage $1 - \alpha$ if and only if $\text{bias}(\hat{\theta}) = 0$; otherwise coverage $< 1 - \alpha$ and symmetric w.r.t. the sign of $\text{bias}(\hat{\theta})$

Simulation setup



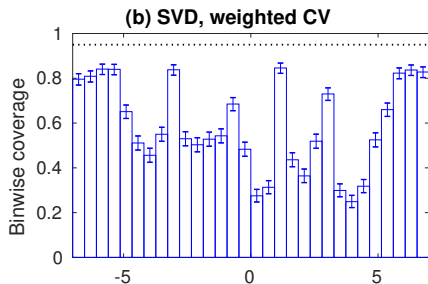
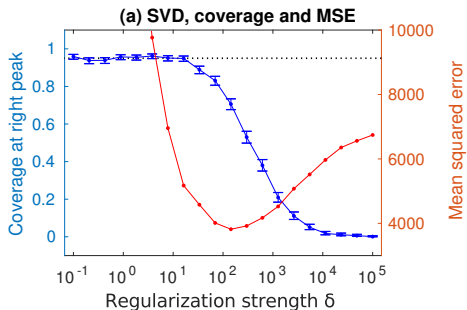
$$f(t) = \lambda_{\text{tot}} \left\{ \pi_1 \mathcal{N}(t|-2, 1) + \pi_2 \mathcal{N}(t|2, 1) + \pi_3 \frac{1}{|T|} \right\}$$

$$g(s) = \int_T \mathcal{N}(s-t|0, 1) f(t) dt$$

$$f^{\text{MC}}(t) = \lambda_{\text{tot}} \left\{ \pi_1 \mathcal{N}(t|-2, 1.1^2) + \pi_2 \mathcal{N}(t|2, 0.9^2) + \pi_3 \frac{1}{|T|} \right\}$$

[Or slight variations of this setup.]

Undercoverage in unfolding

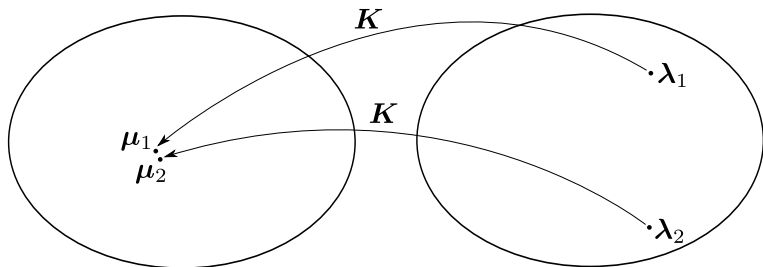


Coverage in SVD unfolding: as a function of the regularization strength (left) and for cross-validated regularization strength (right)

- The optimal point estimator in terms of the MSE has a sizeable regularization bias
- As a result, the unfolded variability intervals have substantial undercoverage
- Similar conclusions hold for other common methods (D'Agostini, TUnfold,...)

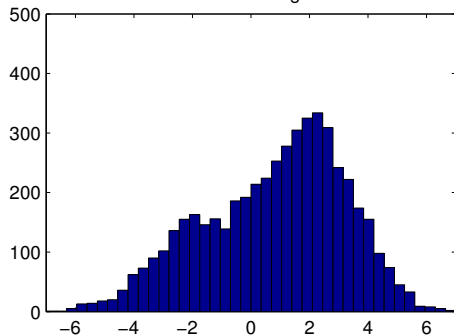
Unfolding is an ill-posed inverse problem

- The main challenge in unfolding is that \mathbf{K} is an ill-conditioned matrix
- When the linear system $\boldsymbol{\mu} = \mathbf{K}\boldsymbol{\lambda}$ is ill-conditioned, true histograms $\boldsymbol{\lambda}_1$ and $\boldsymbol{\lambda}_2$ that are very different can map into smeared histograms $\boldsymbol{\mu}_1$ and $\boldsymbol{\mu}_2$ that are very similar
- As a result, distinguishing between $\boldsymbol{\lambda}_1$ and $\boldsymbol{\lambda}_2$ based on noisy data in the $\boldsymbol{\mu}$ -space is very difficult

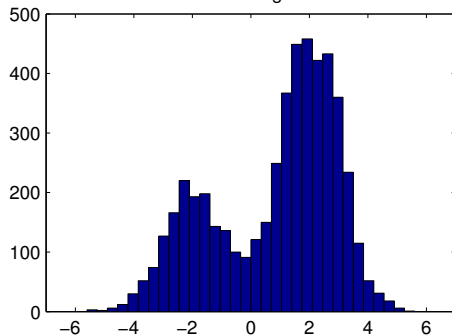


Demonstration of ill-posedness

Smeared histogram

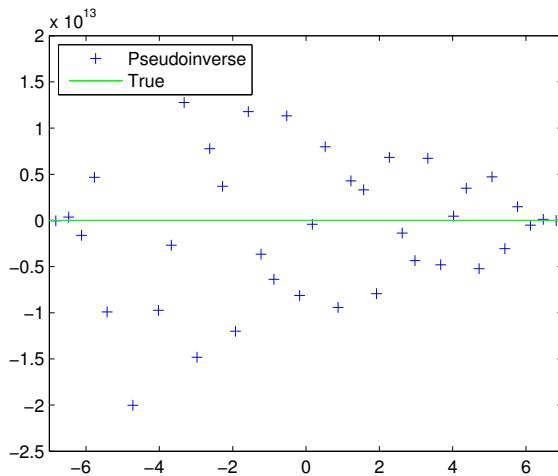


True histogram

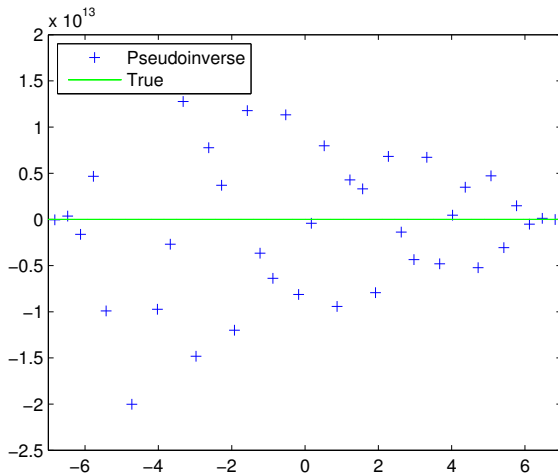


$$\mu = K\lambda, \quad \mathbf{y} \sim \text{Poisson}(\mu) \quad \xRightarrow{??} \quad \hat{\lambda} = K^{-1}\mathbf{y}$$

Demonstration of ill-posedness



Demonstration of ill-posedness



$$\text{MSE}(\hat{\theta}) = \mathbb{E}((\hat{\theta} - \theta)^2) = [\text{bias}(\hat{\theta})]^2 + \text{var}(\hat{\theta})$$

Regularization: bias \uparrow , variance $\downarrow \Rightarrow$ MSE \downarrow

Two main approaches to regularization:

① Explicit penalty term

- Tikhonov regularization / SVD unfolding / TUnfold (Höcker and Kartvelishvili, 1996; Schmitt, 2012)

② Early stopping of an iterative algorithm

- EM iteration with early stopping / D'Agostini iteration (D'Agostini, 1995; Richardson, 1972; Lucy, 1974; Shepp and Vardi, 1982; Lange and Carson, 1984; Vardi et al., 1985)

Tikhonov regularization

- Tikhonov regularization estimates λ by solving:

$$\min_{\lambda \in \mathbb{R}^p} (\mathbf{y} - \mathbf{K}\lambda)^T \mathbf{C}^{-1}(\mathbf{y} - \mathbf{K}\lambda) + \delta P(\lambda)$$

- The first term is a Gaussian approximation to the Poisson log-likelihood
- The second term penalizes physically implausible solutions
- Common penalty terms:
 - **Norm**: $P(\lambda) = \|\lambda\|^2$
 - **Curvature**: $P(\lambda) = \|\mathbf{L}\lambda\|^2$, where \mathbf{L} is a discretized 2nd derivative operator
 - **SVD unfolding** (Höcker and Kartvelishvili, 1996):

$$P(\lambda) = \left\| \mathbf{L} \begin{bmatrix} \lambda_1 / \lambda_1^{\text{MC}} \\ \lambda_2 / \lambda_2^{\text{MC}} \\ \vdots \\ \lambda_p / \lambda_p^{\text{MC}} \end{bmatrix} \right\|^2,$$

where λ^{MC} is a MC prediction for λ

- **TUnfold**² (Schmitt, 2012): $P(\lambda) = \|\mathbf{L}(\lambda - \lambda^{\text{MC}})\|^2$

²TUnfold implements also more general penalty terms

- Starting from some initial guess $\boldsymbol{\lambda}^{(0)} > \mathbf{0}$, iterate

$$\lambda_j^{(k+1)} = \frac{\lambda_j^{(k)}}{\sum_{i=1}^n K_{i,j}} \sum_{i=1}^n \frac{K_{i,j} y_i}{\sum_{l=1}^p K_{i,l} \lambda_l^{(k)}}$$

- Regularization by stopping the iteration before convergence:
 - $\hat{\boldsymbol{\lambda}} = \boldsymbol{\lambda}^{(K)}$ for some small number of iterations K
 - This will bias the solution towards $\boldsymbol{\lambda}^{(0)}$
 - Regularization strength controlled by the choice of K
- RooUnfold (Adey, 2011) defaults to $\boldsymbol{\lambda}^{(0)} = \boldsymbol{\lambda}^{\text{MC}}$
 - It used to be not possible to change this but recent versions of RooUnfold include an undocumented method `SetPriors` for changing the initial guess

D'Agostini iteration

$$\lambda_j^{(k+1)} = \frac{\lambda_j^{(k)}}{\sum_{i=1}^n K_{i,j}} \sum_{i=1}^n \frac{K_{i,j} y_i}{\sum_{l=1}^p K_{i,l} \lambda_l^{(k)}}$$

- This iteration has been discovered in various fields, including optics (Richardson, 1972), astronomy (Lucy, 1974) and tomography (Shepp and Vardi, 1982; Lange and Carson, 1984; Vardi et al., 1985)
- In particle physics, it was popularized by D'Agostini (1995) who called it “Bayesian” unfolding
- **But:** This is in fact an expectation-maximization (EM) iteration (Dempster et al., 1977) for finding the *maximum likelihood estimator* of λ in the Poisson regression problem $\mathbf{y} \sim \text{Poisson}(\mathbf{K}\lambda)$
- As $k \rightarrow \infty$, $\lambda^{(k)} \rightarrow \hat{\lambda}_{\text{MLE}}$ (Vardi et al., 1985)
- *This is a fully frequentist technique for finding the (regularized) MLE*
 - The name “Bayesian” is an unfortunate misnomer

D'Agostini demo, $k = 0$

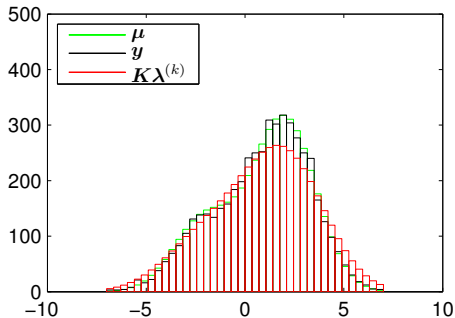


Figure: Smearing histogram

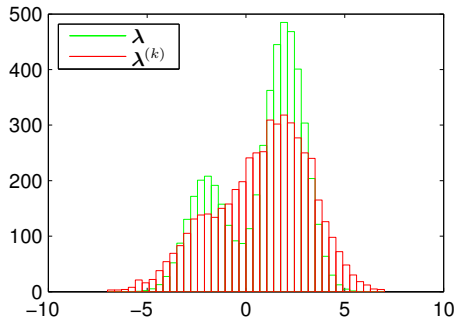


Figure: True histogram

D'Agostini demo, $k = 100$

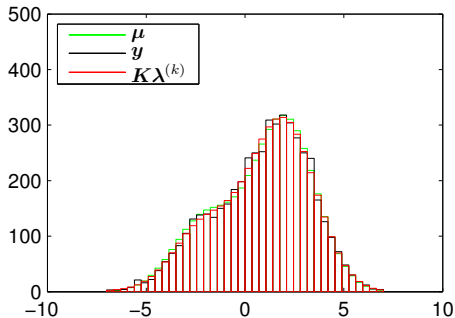


Figure: Smearing histogram

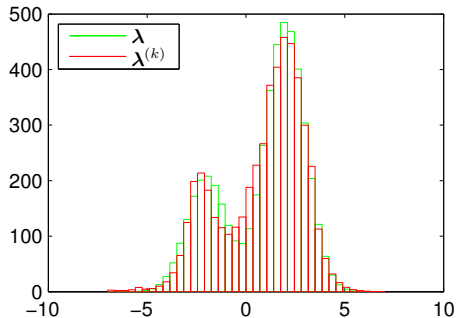


Figure: True histogram

D'Agostini demo, $k = 10000$

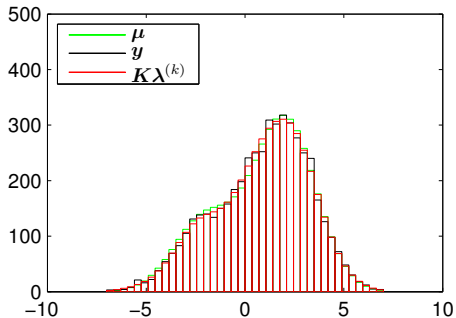


Figure: Smearred histogram

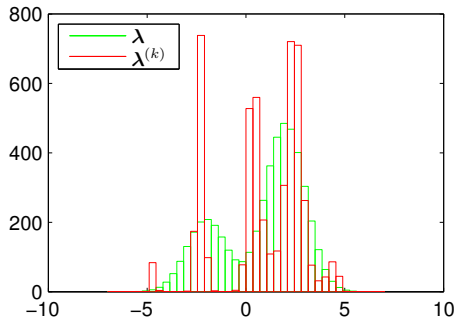


Figure: True histogram

D'Agostini demo, $k = 100000$

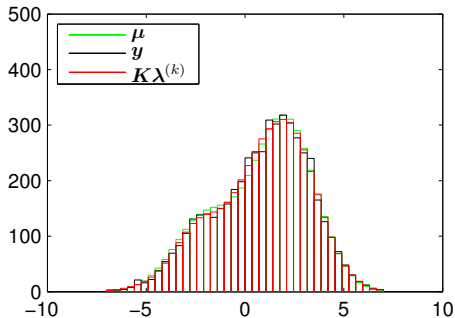


Figure: Smearred histogram

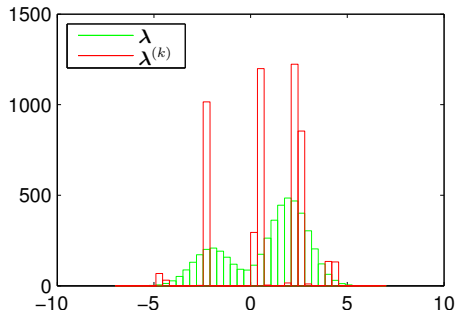
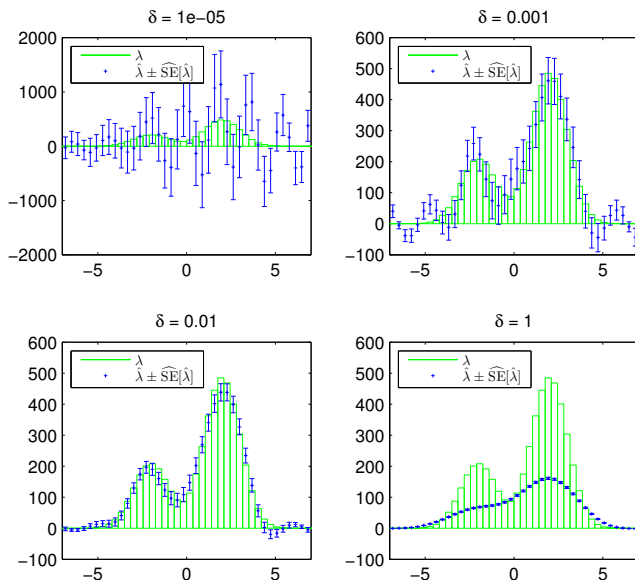


Figure: True histogram

Choice of the regularization strength

- The choice of the regularization strength (δ in Tikhonov, # of iterations in D'Agostini) is a key issue in unfolding
 - Controls the bias-variance trade-off inherent in regularization
 - The solution and especially the uncertainties depend heavily on this choice
- This choice should ideally be done using an objective data-driven criterion
 - In particular, one must not rely on the software defaults for the regularization strength (such as 4 iterations of D'Agostini in RooUnfold)
- Many data-driven methods have been proposed:
 - 1 (Weighted/generalized) cross-validation (e.g., Green and Silverman, 1994)
 - 2 L-curve (Hansen, 1992)
 - 3 Marginal maximum likelihood (MMLE; Kuusela and Panaretos (2015))
 - 4 Goodness-of-fit test in the smeared space (Veklerov and Llacer, 1987)
 - 5 Akaike information criterion (Volobouev, 2015)
 - 6 Minimization of a global correlation coefficient (Schmitt, 2012)
 - 7 Stein's unbiased risk estimate (SURE; new in TUnfold V17.9)
 - 8 Confidence interval coverage (Kuusela, 2016; Brenner et al., 2020)
 - 9 ...
- Limited experience about the relative merits of these in typical unfolding problems

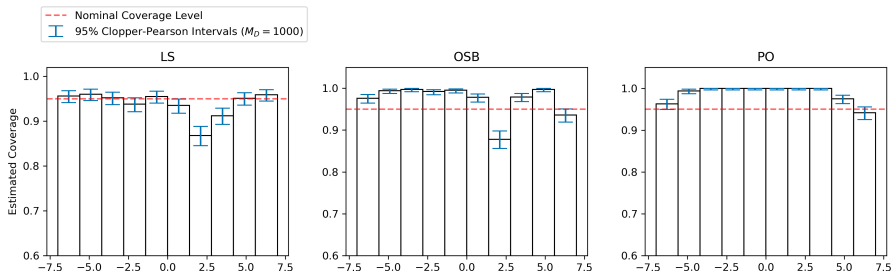
Tikhonov regularization, $P(\lambda) = \|\lambda\|^2$, varying δ



Motivation for the rank-deficient case

However, even with a 40×40 response matrix, the wide-bin bias can be sizeable for heavily misspecified f^{MC}

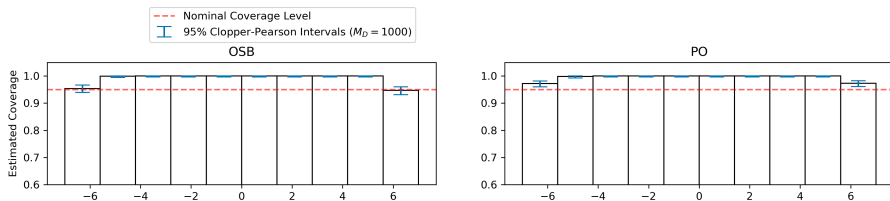
Coverage of the previous three methods for an adversarial f^{MC} :



Wide bins via fine bins, with rank-deficient K

This can be fixed by using an even larger number of true bins, which requires methods that can handle a rank-deficient K

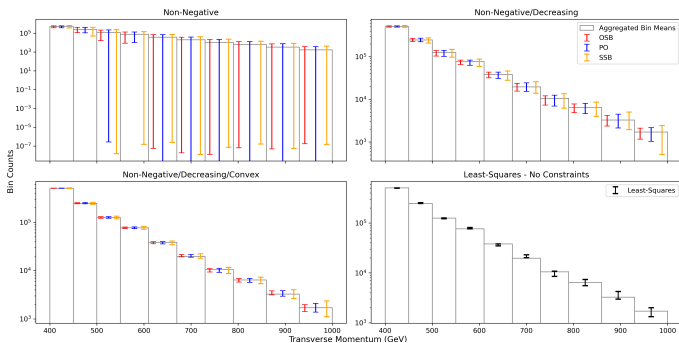
Coverage of the OSB and PO intervals with a 40×80 response matrix:



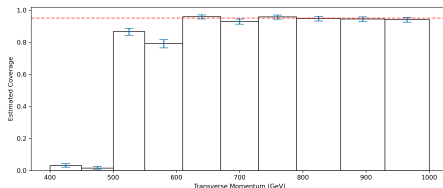
We have additionally found that:

- The interval width of both methods flattens out as the number of true bins is further increased
- The PO interval width has little sensitivity to the choice of the prior

Application to unfolding a steeply falling spectrum

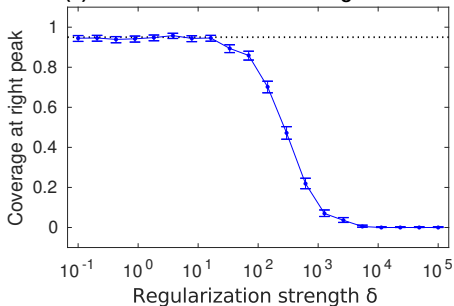


The OSB, PO and SSB intervals based on a 30×60 response matrix all have at least 95% coverage, while the least-squares intervals with a 30×10 matrix do not cover:

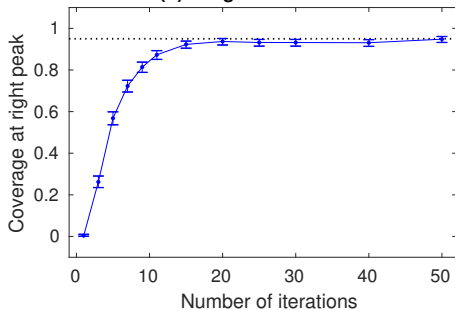


Coverage as a function of regularization strength

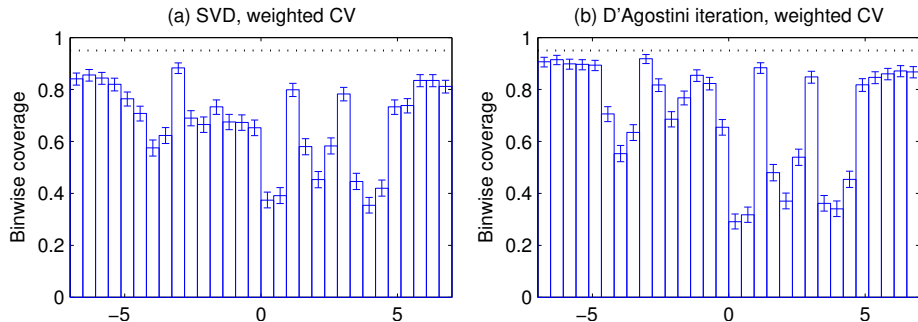
(a) SVD variant of Tikhonov regularization



(b) D'Agostini iteration



Undercoverage of existing methods



There is major undercoverage if regularization strength chosen using (weighted) cross-validation; same is true for L-curve and MMLE.

Key point: These methods are designed for optimal point estimation, but:
optimal point estimation \neq optimal uncertainty quantification

Undersmoothed unfolding

- Standard methods for picking the regularization strength choose too much bias from the perspective of the variance-based uncertainties
- One possible solution is to *debias* the estimator, i.e., to adjust the bias-variance trade-off to the direction of less bias and more variance
- The simplest form of debiasing is to reduce δ from the cross-validation / L-curve / MMLE value until the intervals have close-to-nominal coverage
- The challenge is to come up with a data-driven rule for deciding *how much to undersmooth*
- With Lyle Kim, we have implemented the data-driven methods from Kuusela (2016) as an extension of TUnfold
- The code is available at:

<https://github.com/lylejkim/UndersmoothedUnfolding>

- If you're already working with TUnfold, then trying this approach requires adding only one extra line of code to your analysis

Unfolded histograms, $\lambda^{\text{MC}} = 0$

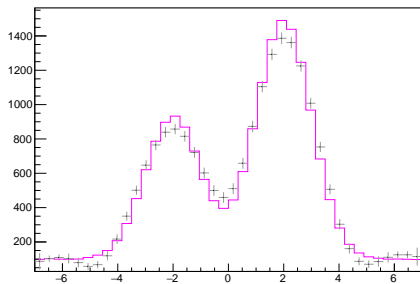


Figure: L-curve, $\tau = \sqrt{\delta} = 0.01186$

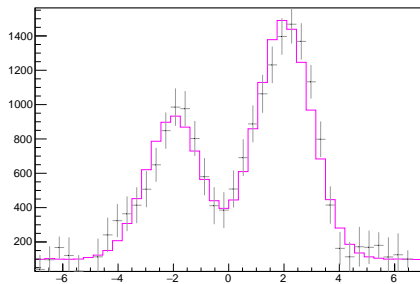


Figure: Undersmoothing, $\tau = \sqrt{\delta} = 0.00177$

Binwise coverage, $\lambda^{\text{MC}} = 0$

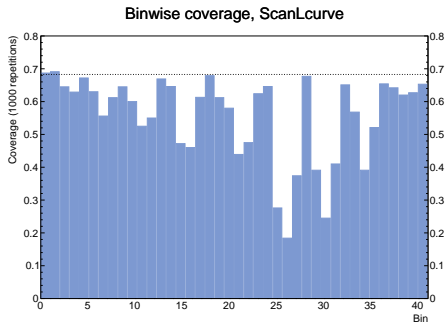


Figure: L-curve

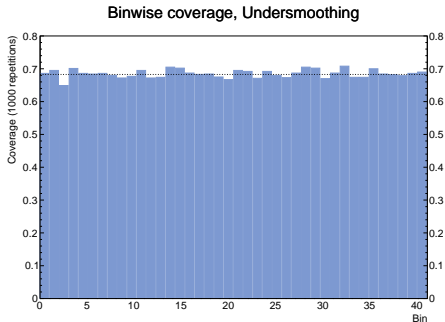


Figure: Undersmoothing

Unregularized unfolding?

- At the end of the day, *any regularization technique makes unverifiable assumptions about the true solution*
 - If these assumptions are not satisfied, the uncertainties will be wrong
 - In the absence of oracle information about the true λ , there does not seem to be any obvious way around this
- So maybe we should reconsider whether explicit regularization is such a good idea to start with?
- Instead of finding a regularized estimator of λ , what if we simply used³ the unregularized matrix inverse $\hat{\lambda} = \mathbf{K}^{-1}\mathbf{y}$?
- This is unbiased ($\mathbb{E}(\hat{\lambda}) = \lambda$) and hence also the corresponding estimator $\hat{\theta} = \mathbf{h}^T \hat{\lambda}$ of the functional $\theta = \mathbf{h}^T \lambda$ is unbiased
- Therefore, by the previous discussion, the resulting variability intervals have correct coverage $1 - \alpha$

³For simplicity, I assume here that $\mathbf{K} \in \mathbb{R}^{n \times p}$ is an invertible square matrix. The case where $n > p$ with \mathbf{K} having full column rank is also easy using the pseudoinverse $\hat{\lambda} = (\mathbf{K}^T \mathbf{K})^{-1} \mathbf{K}^T \mathbf{y}$. The case where \mathbf{K} is column-rank deficient (including when $p > n$) is trickier but probably doable; see <https://indico.cern.ch/event/882374/>.

Implicit regularization

- Of course, when \mathbf{K} is ill-conditioned, the unregularized estimator $\hat{\boldsymbol{\lambda}}$ will have a huge variance
- *But this does not mean that $\hat{\theta} = \mathbf{h}^T \hat{\boldsymbol{\lambda}}$ needs to have a huge variance!*
- The mapping $\hat{\boldsymbol{\lambda}} \mapsto \hat{\theta} = \mathbf{h}^T \hat{\boldsymbol{\lambda}}$ can act as an implicit regularizer resulting in a well-constrained interval $[\underline{\theta}, \bar{\theta}]$ for the functional $\theta = \mathbf{h}^T \boldsymbol{\lambda}$
- This is especially the case when the functional is a smoothing / averaging / aggregation operation
 - For example, inference for aggregated unfolded bins (demo to follow)
- Of course, there are also functionals that are more difficult to constrain (e.g., individual bins $\theta = \mathbf{e}_i^T \boldsymbol{\lambda}$, derivatives,...)
- In those cases, the intervals $[\underline{\theta}, \bar{\theta}]$ are wide—as they should be, since there is simply not enough information in the data \mathbf{y} to constrain these functionals

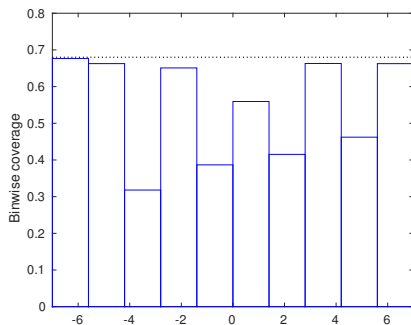
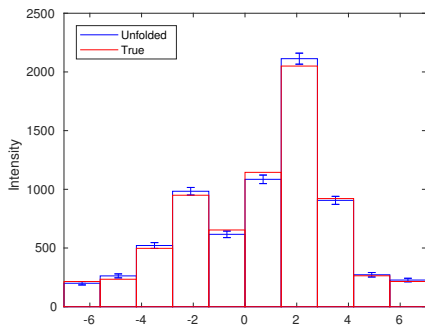
Wide bin unfolding

- One functional we should be able to recover without explicit regularization is the integral of f over a *wide* unfolded bin:

$$H_j[f] = \int_{T_j} f(t) dt, \quad \text{width of } T_j \text{ large}$$

- But one cannot simply arbitrarily increase the particle-level bin size in the conventional approaches, since this increases the MC dependence of \mathbf{K}
- To circumvent this, *it is possible to first unfold with fine bins (without regularization) and then aggregate into wide bins*
- Let's see how this works using a similar deconvolution setup as before

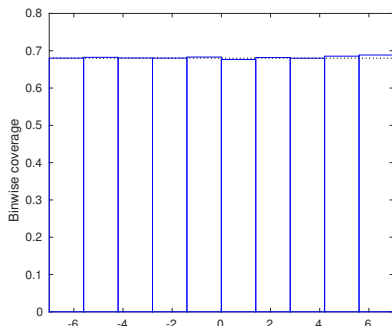
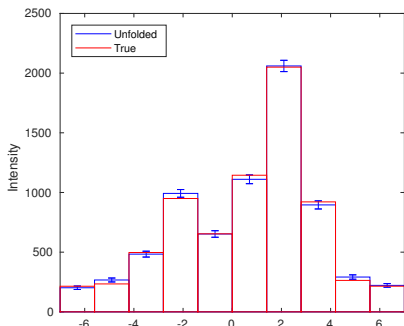
Wide bins, standard approach, perturbed MC



The response matrix $K_{i,j} = \frac{\int_{S_i} \int_{T_j} k(s,t) f^{\text{MC}}(t) dt ds}{\int_{T_j} f^{\text{MC}}(t) dt}$ depends on f^{MC}

\Rightarrow Undercoverage if $f^{\text{MC}} \neq f$

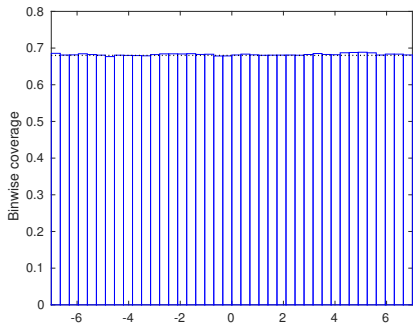
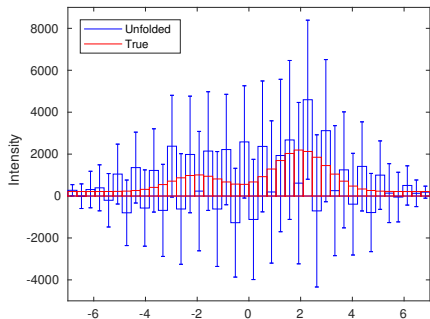
Wide bins, standard approach, correct MC



If $f^{\text{MC}} = f$, coverage is correct

⇒ But this situation is unrealistic because f of course is unknown

Fine bins, standard approach, perturbed MC

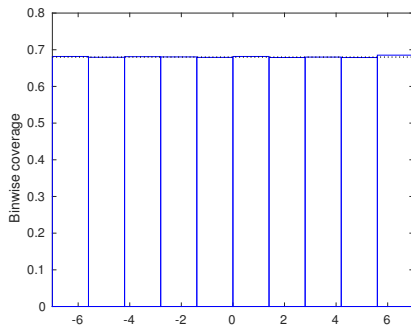
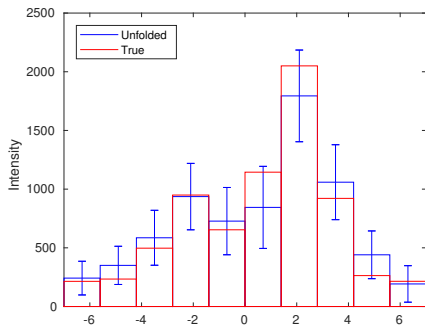


With narrow bins, less dependence on f^{MC} so coverage is correct, but the intervals are very wide⁴

⇒ Let's aggregate these into wide bins, keeping track of the bin-to-bin correlations in the error propagation

⁴More unfolded realizations given in the [backup](#).

Wide bins via fine bins, perturbed MC



Wide bins via fine bins gives both correct coverage and intervals with reasonable length⁵

⁵More unfolded realizations given in the [backup](#).

Current unfolding methods

- Two main approaches:

- 1 Tikhonov regularization (i.e., SVD by Höcker and Kartvelishvili (1996) and TUnfold by Schmitt (2012)):

$$\min_{\lambda \in \mathbb{R}^p} (\mathbf{y} - \mathbf{K}\lambda)^T \hat{\mathbf{C}}^{-1} (\mathbf{y} - \mathbf{K}\lambda) + \delta P(\lambda)$$

with

$$P_{\text{SVD}}(\lambda) = \left\| \mathbf{L} \begin{bmatrix} \lambda_1 / \lambda_1^{\text{MC}} \\ \lambda_2 / \lambda_2^{\text{MC}} \\ \vdots \\ \lambda_p / \lambda_p^{\text{MC}} \end{bmatrix} \right\|^2 \quad \text{or} \quad P_{\text{TUnfold}}(\lambda) = \|\mathbf{L}(\lambda - \lambda^{\text{MC}})\|^2,$$

where \mathbf{L} is usually the discretized second derivative (also other choices possible)

- 2 Expectation-maximization iteration with early stopping (D'Agostini, 1995):

$$\lambda_j^{(t+1)} = \frac{\lambda_j^{(t)}}{\sum_{i=1}^n K_{i,j}} \sum_{i=1}^n \frac{K_{i,j} y_i}{\sum_{k=1}^p K_{i,k} \lambda_k^{(t)}}, \quad \text{with } \lambda^{(0)} = \lambda^{\text{MC}}$$

- All these methods typically regularize by biasing towards a MC ansatz λ^{MC}
- Regularization strength controlled by the choice of δ in Tikhonov or by the number of iterations in D'Agostini
- Uncertainty quantification: $[\underline{\lambda}_i, \bar{\lambda}_i] = \left[\hat{\lambda}_i - z_{1-\alpha/2} \sqrt{\widehat{\text{var}}(\hat{\lambda}_i)}, \hat{\lambda}_i + z_{1-\alpha/2} \sqrt{\widehat{\text{var}}(\hat{\lambda}_i)} \right]$, with $\widehat{\text{var}}(\hat{\lambda}_i)$ estimated using error propagation or resampling

Coverage as a function of $\tau = \sqrt{\delta}$

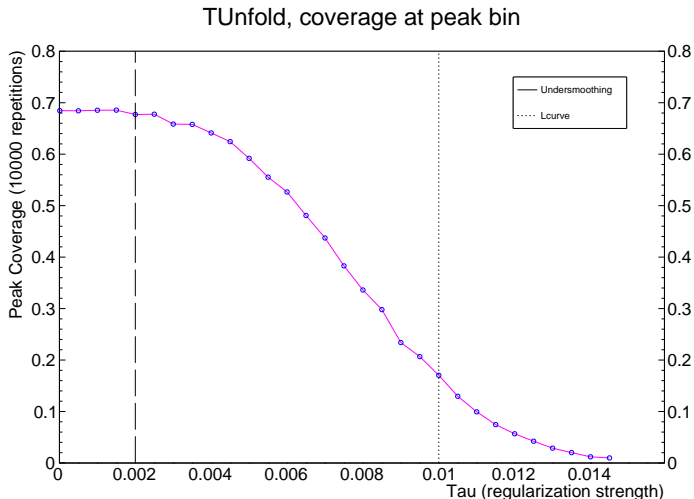
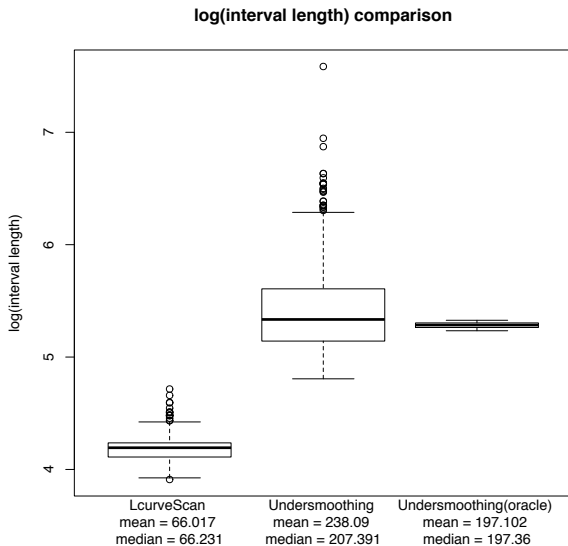
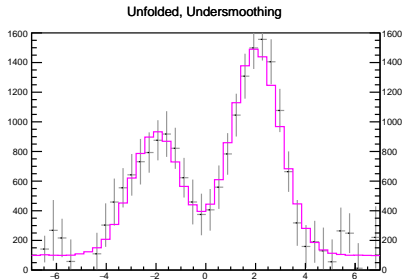
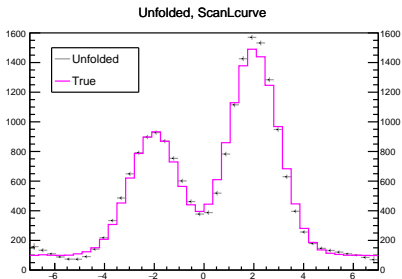
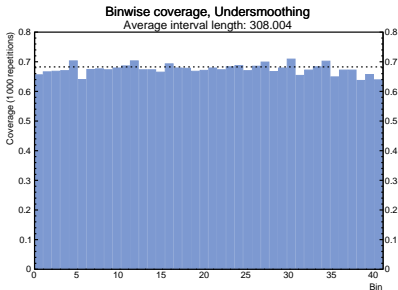
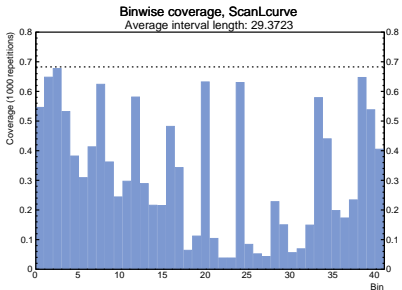


Figure: Coverage at the right peak of a bimodal density

Interval lengths, $\lambda^{\text{MC}} = 0$



Histograms, coverage and interval lengths when $\lambda^{MC} \neq 0$



Coverage study from Kuusela (2016)

Method	Coverage at $t = 0$	Mean length
BC (data)	0.932 (0.915, 0.947)	0.079 (0.077, 0.081)
BC (oracle)	0.937 (0.920, 0.951)	0.064 (0.064, 0.064)
US (data)	0.933 (0.916, 0.948)	0.091 (0.087, 0.095)
US (oracle)	0.949 (0.933, 0.962)	0.070 (0.070, 0.070)
MMLE	0.478 (0.447, 0.509)	0.030 (0.030, 0.030)
MISE	0.359 (0.329, 0.390)	0.028
Unregularized	0.952 (0.937, 0.964)	40316

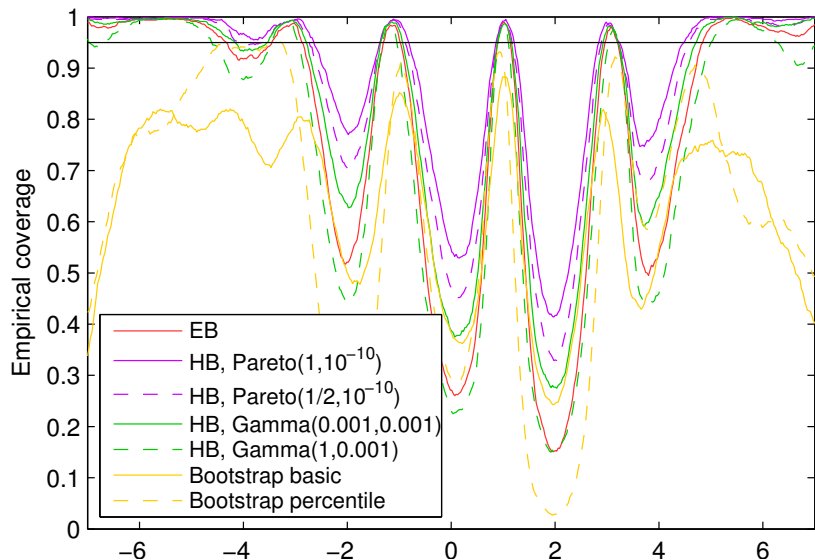
BC = iterative bias-correction

US = undersmoothing

MMLE = choose δ to maximize the marginal likelihood

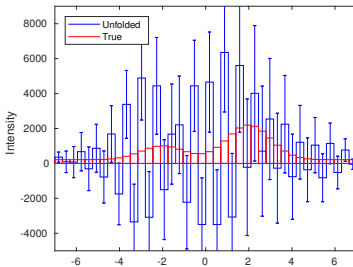
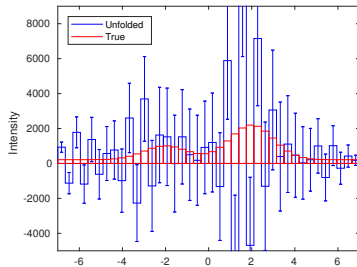
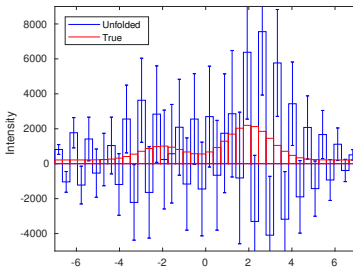
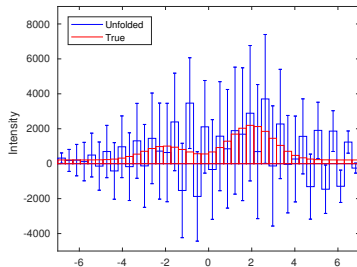
MISE = choose δ to minimize the mean integrated squared error

UQ in inverse problems is challenging



[Kuusela and Panaretos (2015)]

Fine bins, standard approach, perturbed MC, 4 realizations



Wide bins via fine bins, perturbed MC, 4 realizations

