



Technical Report

RISE Germany Internship Unfolding FACT Data

Jacob Bieker, Mathis Börner, Kai
Brügge, Maximilian Nöthe

09/2017



Part of the work on this technical report has been supported by Deutsche Forschungsgemeinschaft (DFG) within the Collaborative Research Center SFB 876 "Providing Information by Resource-Constrained Analysis", project C3.

Speaker: Prof. Dr. Katharina Morik
Address: TU Dortmund University
Joseph-von-Fraunhofer-Str. 23
D-44227 Dortmund
Web: <http://sfb876.tu-dortmund.de>

1 Introduction

In this report the results from a 10 week internship are presented. The goal of the internship was to apply different unfolding approaches to conduct measurements of energy spectra from data acquired by FACT, the First G-APD Cherenkov Telescope. FACT is the first operational telescope of its kind, employing a camera equipped with silicon photo multipliers (G-APD aka SiPM) to primarily detect gamma rays. Improving the unfolding method can help with better interpretation of the data and more accurate physics results without the need for new equipment or more observations. The approaches tested during this internship range from simplistic matrix inversion to an improvement over of the previous standard (TRUEE [7]).

2 Unfolding

The goal of unfolding is to determine which input distribution $f(y)$ of a desired quantity y leads to the measured distribution $g(x)$ of the observable x . For continuous distribution this can be written as an integral equation

$$g(x) = \int A(x, y) f(y) dy. \quad (1)$$

In this equation $A(x, y)$ is the so-called detector response function and returns the conditional probability to measure a value x given a value for y . In most applications discrete distributions f and g are used and (1) transforms to

$$\vec{g} = \mathbf{A}\vec{f}. \quad (2)$$

The discretization for f is often motivated by the physics that is being investigated and is usually a simple equidistant binning in y or a spline representation of the distribution (e.g. in TRUEE). For an equidistant binning \vec{f} contains the number of events in the different bins and for a spline representation it the coefficient vector for the splines. In (2) the former detector response function $A(x, y)$ is replaced by the detector response matrix \mathbf{A} . Analog to $A(x, y)$ the detector response matrix contains the conditional probabilities to measure an event in a bin i of the observable distribution \vec{g} given the event is in bin j of the desired distribution \vec{f} .

As shown later the challenges of solving (2) is driven by the detector response matrix \mathbf{A} . The characteristics of this matrix mainly depends on the binning chosen for the observable space. In subsection 2.1 the different combinations of unfolding and binning approaches are presented and used to illustrate the evolution from the most basic approach to the the new approach that improves upon the current standard (TRUEE) used for FACT data [8].

2.1 Different Unfolding Approaches

In this section the different approaches and the challenges associated with them are explained.

2.1.1 Inverse Matrix

The simplest way to unfold data is to invert the detector response matrix. The inverted matrix can then be multiplied with the measured signal \vec{g} to unfold back to the real spectrum \vec{f} using $\vec{f} = \mathbf{A}^{-1}\vec{g}$. Depending on the uncertainty in the measured distribution, this will lead to large fluctuations in the solution. The cause for those fluctuations can be shown by transforming (2) into a basis in which \mathbf{A} is a diagonal matrix [1] ($\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^\top$):

$$\vec{g} = \mathbf{A}\vec{f} \rightarrow \mathbf{U}^\top\vec{g} = \mathbf{D}\mathbf{U}^\top\vec{f} \rightarrow \vec{c} = \mathbf{D}\vec{b}$$

Matrix \mathbf{D} is a diagonal matrix with \mathbf{A} 's eigenvalues on the diagonal¹. Therefore, the unfolding via inverting the matrix in this space is simply:

$$b_i = \frac{1}{\lambda_i}c_i. \quad (3)$$

In (3) the eigenvalues of matrix \mathbf{A} are inverted and therefore small eigenvalues become large factors for their corresponding coefficients c_i . With increasing i the coefficients become less and less significant. This is what shows up as large fluctuations in the solution and is sometimes referenced as *fitting to noise*. This characteristic of the unfolding emerges from ambiguous connections between a bin in the observable space and a bin of the sought-after quantity. The so-called condition number κ of matrix \mathbf{A} is a useful measure of the uncertainty of the measurement and how strong the solution will be dominated by noise. It is defined by the ratio between the biggest and the smallest eigenvalues $\kappa = \lambda_0/\lambda_N$.

To tackle the problem of fluctuating solutions multiple options are available: Increase statistics to get more significant coefficients, ignore small eigenvalues (*regularization*²) or find a way to increase the eigenvalues of matrix \mathbf{A} . In most analysis all available data is used, so using regularization or increasing the condition of matrix \mathbf{A} are the only feasible options.

2.1.2 Singular Value Decomposition

Singular Value Decomposition (SVD) [6] improves on inverting the matrix by allowing for \mathbf{A} 's that are not square. Without that restriction, more information can be utilized, generally leading to a better condition³ of matrix \mathbf{A} .

2.1.3 Poissonian Likelihood

It is known that the underlying distribution in each observable bin g_i is a Poisson distribution. This knowledge is not taken into account by neither the matrix inversion nor

¹We assume that the eigenvalues are ordered by their size

²As shown later for this unfolding no regularization is needed, it is not covered in full detail in this report.

³For rectangular matrices the condition is defined with the singular values instead of the eigenvalues.

by the SVD approach. This can lead to biased solutions when many bins contain small (≤ 10) numbers of events.

The poissonian distribution can be taken into account when the problem is formulated as a likelihood function:

$$\mathcal{L}(\vec{g}|\vec{f}) = \prod_{i=1}^m \left(\frac{\hat{g}_i^{g_i}}{g_i!} \exp(-\hat{g}_i) \right) \cdot \underbrace{\frac{1}{\sqrt{(2\pi)^n \det(\tau \mathbb{1})}} \exp\left(-\frac{1}{2} \tau \vec{f}^\top \mathbf{C}^\top \mathbb{1} \mathbf{C} \vec{f}\right)}_{\text{Tikhonov Regularization}}. \quad (4)$$

To find the maximum of the likelihood function (4) two different methods were used.

Gradient Descent

The gradient descent (Newton method) approach implemented in this study analytically calculates the negative log likelihood, \mathcal{L} , gradient \vec{h} of the negative log likelihood, and Hessian matrix \mathbf{H} . The Newton method is an iterative method starting with \vec{f}_0 . For each iteration the gradient and Hessian are calculated to determine next position $\vec{f}_{i+1} = \vec{f}_i - \mathbf{H}^{-1} \cdot \vec{h} \cdot \gamma$ ⁴. We run this iteration for a fixed (large) number of steps. The negative log likelihood S , its gradient \vec{h} and Hessian \mathbf{H} used are:

For minimization:

$$S(\vec{g}|\vec{f}) = \sum_{i=1}^m \left((\mathbf{A}\vec{f})_i - g_i \ln((\mathbf{A}\vec{f})_i) \right) + \frac{1}{2} \tau \vec{f}^\top \mathbf{C}' \vec{f}$$

Gradient:

$$\frac{dS}{df_k} = h_k = \sum_{i=1}^m \left(A_{i,k} - \frac{g_i A_{i,k}}{\sum_{j=1}^n A_{i,j} f_j} \right) + \tau \sum_{i=1}^n f_i C'_{k,i}$$

Hesse-Matrix:

$$\frac{d^2 S}{df_k df_l} = H_{k,l} = \sum_{i=1}^m \left(\frac{g_i A_{i,k} A_{i,l}}{\left(\sum_{j=1}^n A_{i,j} f_j \right)^2} \right) + \tau C'_{k,l}$$

The minimum of the negative likelihood provides the most likely spectrum leading to the measured distribution in the observable space. To obtain the uncertainties here the inverse Hessian (parabolic approximation of the likelihood space) is used. This is a reasonable approximation for bins in f with large statistics, but might break down for bins closer to zero, because we can't have negative bin counts and the likelihood space starts to become asymmetric.

⁴The factor γ is used to increase/decrease the step size.

Markov Chain Monte Carlo

The final unfolding method tested during this study is Markov Chain Monte Carlo (MCMC) [5]. The idea of MCMC is to perform a random walk in the likelihood space. For each step of the walk is randomly drawn via rejection sampling based on the likelihood function. This leads to a sample of the posterior probability distribution of the likelihood function.

The sample can be used to determine the best fit and different uncertainty levels. One advantage of the MCMC approach above the gradient descent is that no approximations are made about the likelihood space. This leads to reasonable error contours for small entries in \vec{f} . Another great advantage is that the sample of the likelihood also provides a straightforward way to compare the compatibility between a given spectrum \vec{f}_{Test} and the unfolded result. To calculate the compatibility one simply has to calculate the likelihood value for the test spectrum and compare with the likelihood values for each example in the MCMC sample. This provides a p-value for the hypothesis that the test spectrum and the unfolded spectrum have the same underlying distribution. For all MCMC results in this report the `emcee` [4] package was used.

2.2 Binning Methods

As mentioned in section 2 the binning of the observable space is crucial to achieve a good condition for matrix \mathbf{A} . There are no limitations or constraints on the bins' shapes or in the dimensionality of the binning. TRUEEE for example allows up to three observables. In all these observables the number of bins is set and equidistant binning is used for each observable. The danger of increased dimensionality is that the binning can create a lot of empty or sparsely populated bins. This can lead to two major problems. Bins with few entries give very little constraint on the fit, and too few simulated events in a bin introduce a large uncertainty in the entries of matrix \mathbf{A} due to limited statistics.

2.2.1 Equidistant

One way to deal with the previously mentioned problems of sparsely populated bins is by merging bins in a equidistant binning. There are various ways to do this. Ones used here to test the effectiveness of this were merging the closest bins with low counts until the merged bins contained a set minimum number of events, and merging the bins with the lowest counts together until at least the minimum number of events is contained within them. Example of equidistant binning in two dimensions is shown in Figure 1.

2.2.2 Decision Tree Binning

A more sophisticated way [2] to obtain a binning for the problem utilizes a decision tree [3]. The idea is to train a decision tree classifier to classify events into the bins of \vec{f} . The training of a decision tree is an recursive process. In each step a cut is performed to achieve the best separation of the different classes. For the two new datasets created

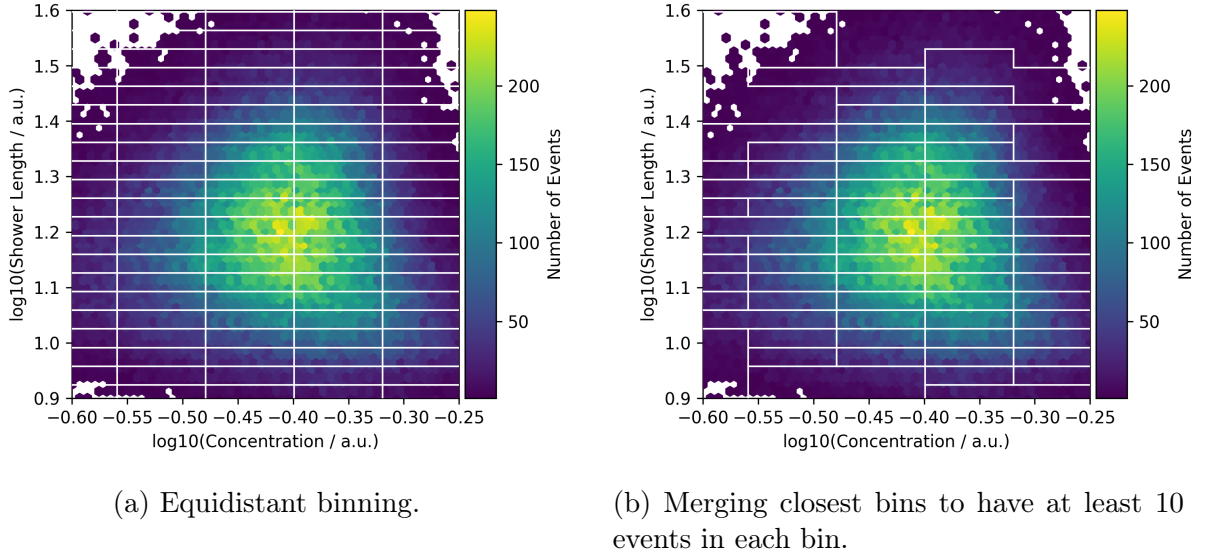


Figure 1: Equidistant without and with merging closest bins until a threshold is reached.

from the cut, a new cut is performed and so on. The recursion stops either when the part of the dataset is purely one class or an external condition is fulfilled. The stopping points are so-called leaves of the tree and the path to the leaves is a sequence of cuts in the observables. Each leaf represents a rectangular space in the observable space and the whole space is covered by the leaves.

This approach has multiple advantages over the equidistant binning. The biggest advantage is that the binning is explicitly optimized for the problem. The decision builds the model to differ between the different bins in \vec{f} . This should directly lead to an optimized condition of the matrix \mathbf{A} .

Another advantage is that the training process can be controlled with external conditions. E.g. the recursion of the training process can be stopped when only k events are left in the dataset. This prevents the binning from having empty/weakly populated bins. A visualization of the binning with two observables is shown in Figure 2

3 Results

The various different methods of unfolding the data are presented below. All methods were applied to simulations for the FACT telescope. The wanted quantity is the energy of the γ -particles. For every unfolding the binning in the energy was equidistant in $\log_{10}(E_{\gamma}/[\text{GeV}])$ between 2.4 and 4.2.

Used Data

When using the FACT Monte Carlo data, the initial dataset consisted of 391 904 pure signal events with 253 different attributes. In this study, only simulated data without

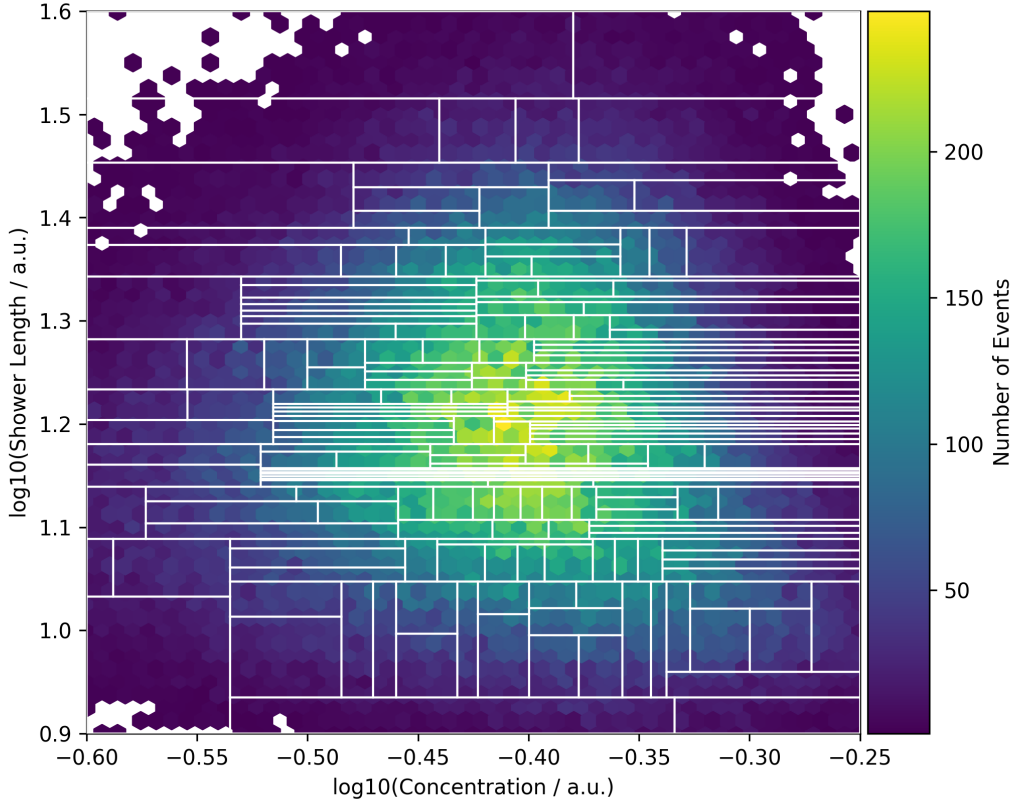


Figure 2: Equidistant without and with merging closest bins until a threshold is reached.

background was used. In addition, only 18 of the attributes were used to train the decision tree classifier for the binning. To simulate an actual measurement run, the events are randomly split into three groups. One group of roughly 5000 events is used as the testing data. All the other data is split into twenty percent for training the decision tree classifier, and eighty percent to build the detector response matrix. For all observable binnings it is ensured to have at least 10 test events in each bin.

3.1 SVD Unfolding

As discussed in subsection 2.1.1 the condition of detector response matrix \mathbf{A} provides deep insight into the characteristics of the unfolding. Figure 3 shows the singular values and the corresponding condition number for five different binnings. The worst condition ($\kappa = 90\,223.0$) is obtained with an equidistant binning in one observable with low correlation to the energy E_γ . Using an observable with a known high correlation to the energy the condition is improved by more than two orders of magnitude ($\kappa = 169.0$). When two observables (visualization Figure 1a) are used the condition can be further improved ($\kappa = 84.6$). The full potential of the decision tree-based binning can be utilized when all 18 observables are used ($\kappa = 23.0$). Using the same two observables as used in the equidistant binning gives a condition number of $\kappa = 67.0$ for the tree based binning.

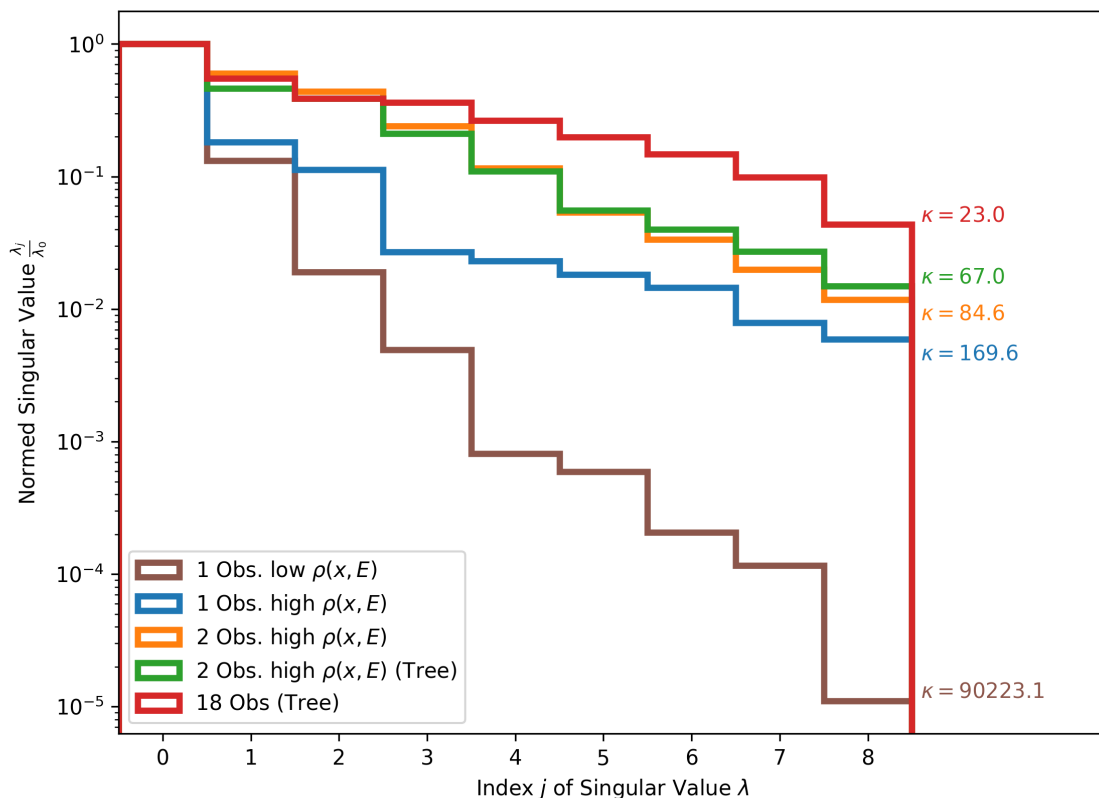


Figure 3: Normed singular values and conditions for different binnings.

In Figure 4 the results of the SVD unfolding for the binnings are shown. The differences of the condition can directly be seen in the fluctuations of the results. Only the unfolding using the tree binning with all observables shows no obvious fluctuations. This is an indication that such a binning removes or reduces greatly the regularization needed to unfold FACT data.

3.2 Poissonian Likelihood

Figure 5 shows the path of the gradient descent approach. The algorithm is capable of finding the maximum of the likelihood function. Looking at the likelihood sample from the MCMC in the background of Figure 5 indicates that for the very first bin a parabolic approximation of the likelihood space might not be justified. The distribution is asymmetric because it is close to zero.

In the figures 6, 7, 8, 9 the MCMC samples and best fits for the same and for different binning as shown in Figure 3 are depicted.

Similar to the SVD unfolding, all binnings, except the tree binning using all 18 observables, show no clear maximum in the likelihood space and a large variance in the sample for all bins.

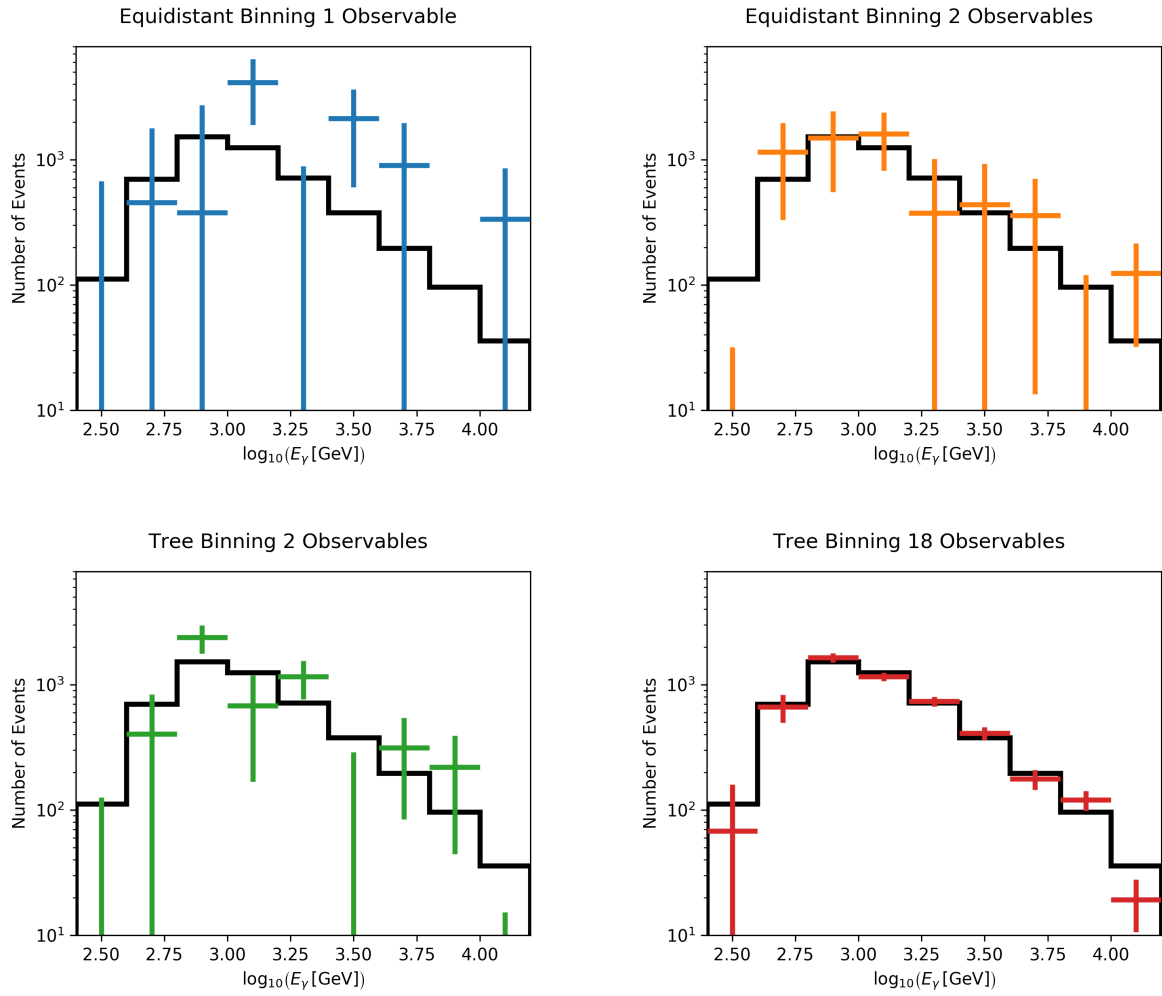


Figure 4: Unfolding results from SVD unfolding for the different detector response matrices obtained with the different binning approaches.

4 Conclusion

In the first chapters the benefits and drawbacks of four different unfolding approaches were discussed. In the application of those approaches it clearly showed that using a MCMC sampler with poissonian likelihood for the unfolding provides the best results.

The difference between the SVD and the gradient descent unfolding (both with 18 observable tree binning) were the uncertainties on the unfolded spectra. The first bin in the unfolding illustrated why an approach that doesn't rely on a parabolic approximation to obtain the uncertainties should be favored. Even the best fit in the bin was ≈ 80 results close to 0 were not excluded. For results close to 0 the parabolic approximation can be expected to break down.

Another insight gained from the studies is that using the decision tree based binning approach (subsubsection 2.2.2) with all 18 observables in the training improves the condition of the problem up to a point where no regularization might be needed. This is a promising finding, but needs to be validated in more extensive studies. A first check would be to repeat the unfolding of simulated events multiple times ($\mathcal{O}(1000)$) and investigate if the result is unbiased. To do that for each unfolding, the p-value for the distribution of the simulation has to be calculated as described in subsection 3.2. For an unbiased solution the p-values should have a uniform distribution.

After showing that the unfolding with is unbiased the next logical step is to apply the approach on real data.

References

- [1] Volker Blobel. “An unfolding method for high energy physics experiments”. In: *arXiv preprint hep-ex/0208022* (2002).
- [2] M. Börner et al. “Measurement/Simulation Mismatches and Multivariate Data Discretization in the Machine Learning Era”. In: *ADASS XXVII*. Ed. by TBD. Vol. TBD. ASP Conf. Ser. San Francisco: ASP, 2018, TBD.
- [3] L. Breiman et al. *Classification and Regression Trees*. Monterey, CA: Wadsworth and Brooks, 1984.
- [4] Daniel Foreman-Mackey et al. “emcee: the MCMC hammer”. In: *Publications of the Astronomical Society of the Pacific* 125.925 (2013), p. 306.
- [5] Charles J Geyer. “Markov chain Monte Carlo maximum likelihood”. In: (1991).
- [6] Andreas Hoecker and Vakhtang Kartvelishvili. “SVD approach to data unfolding”. In: *arXiv preprint hep-ph/9509307* (1995).
- [7] Natalie Milke et al. “Solving inverse problems with the unfolding program TRUEE: Examples in astroparticle physics”. In: *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* 697 (2013), pp. 133–147.

- [8] Maximilian Nöthe et al. “FACT Performance of the First Cherenkov Telescope Observing with SiPMs”. In: *Proceedings of the 35th International Cosmic Ray Conference*. PoS(ICRC2017)791. Proceedings of Science, 2017. URL: <https://pos.sissa.it/301/791/pdf>.

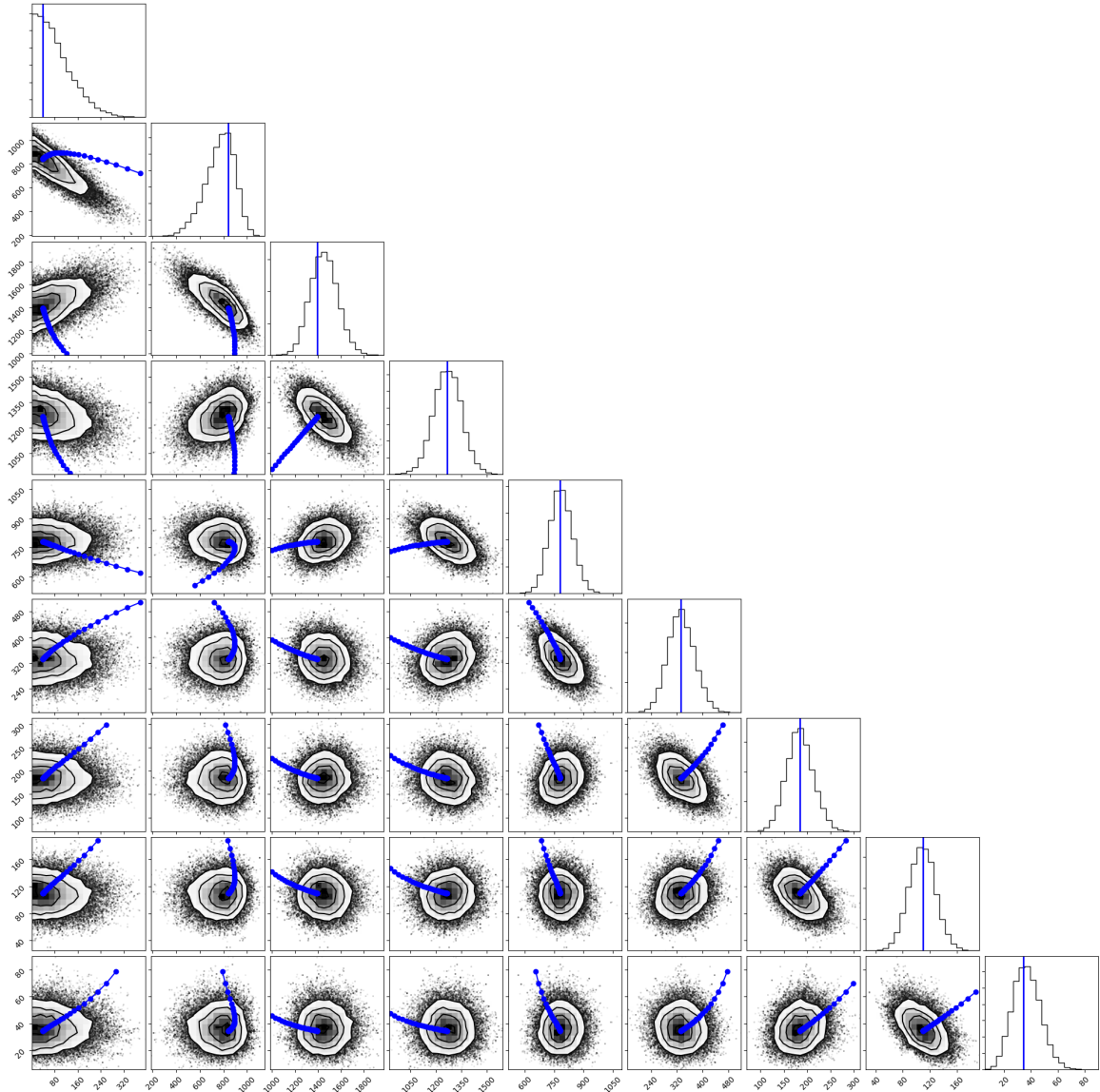


Figure 5: Visualization of the path in the gradient descent minimization. The blue points of every 10-th step of the 500 iterations. In the background the likelihood sample obtained with the MCMC approach is shown.

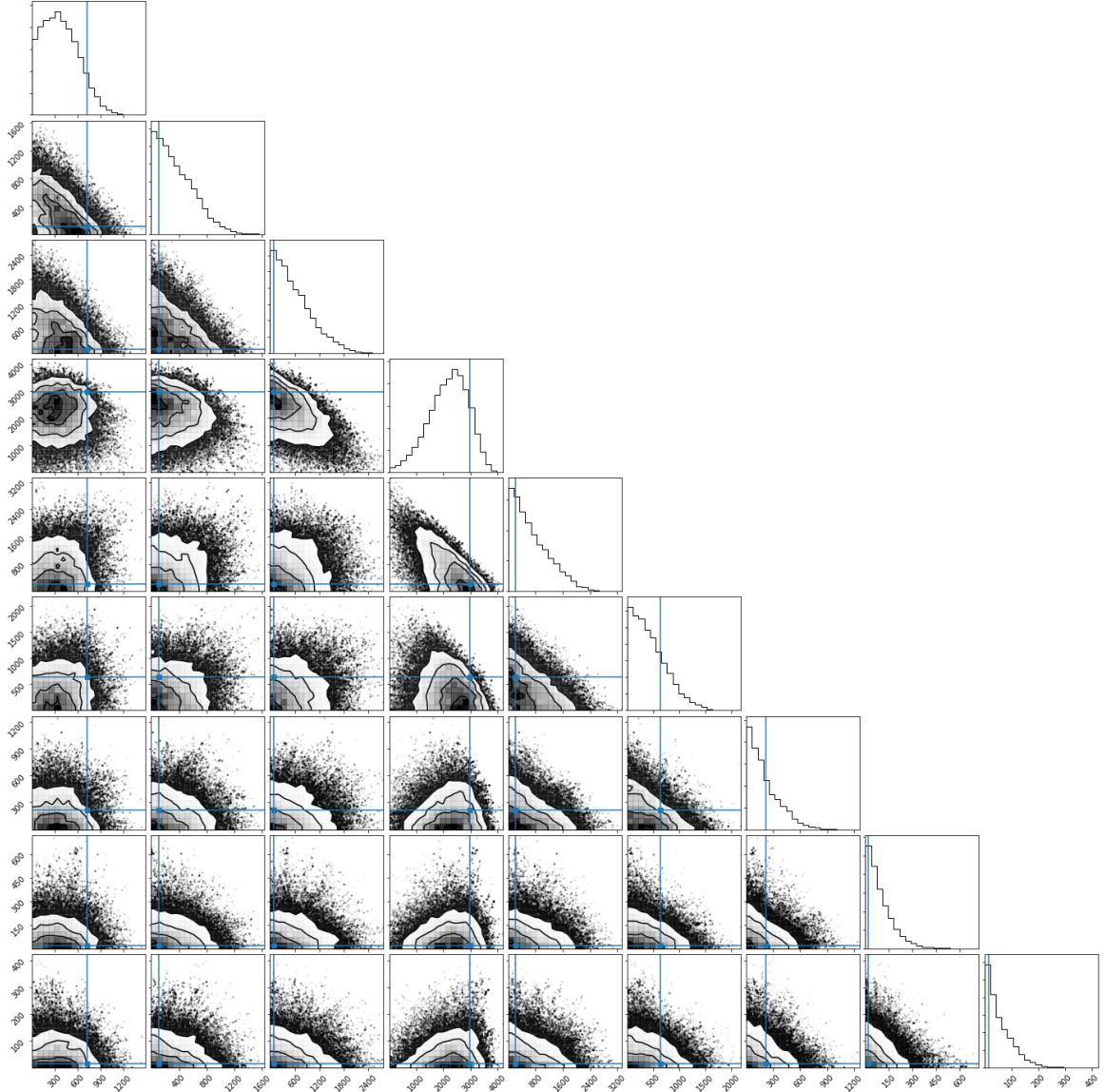


Figure 6: Result of the MCMC unfolding using equidistant binning with a single observable. The plots on the diagonal show the distribution of the samples for the different bins in \vec{f} . The non-diagonal plots show the two dimensional distribution of two bins. All axes show number of events in bin. The marked points are the best fit.

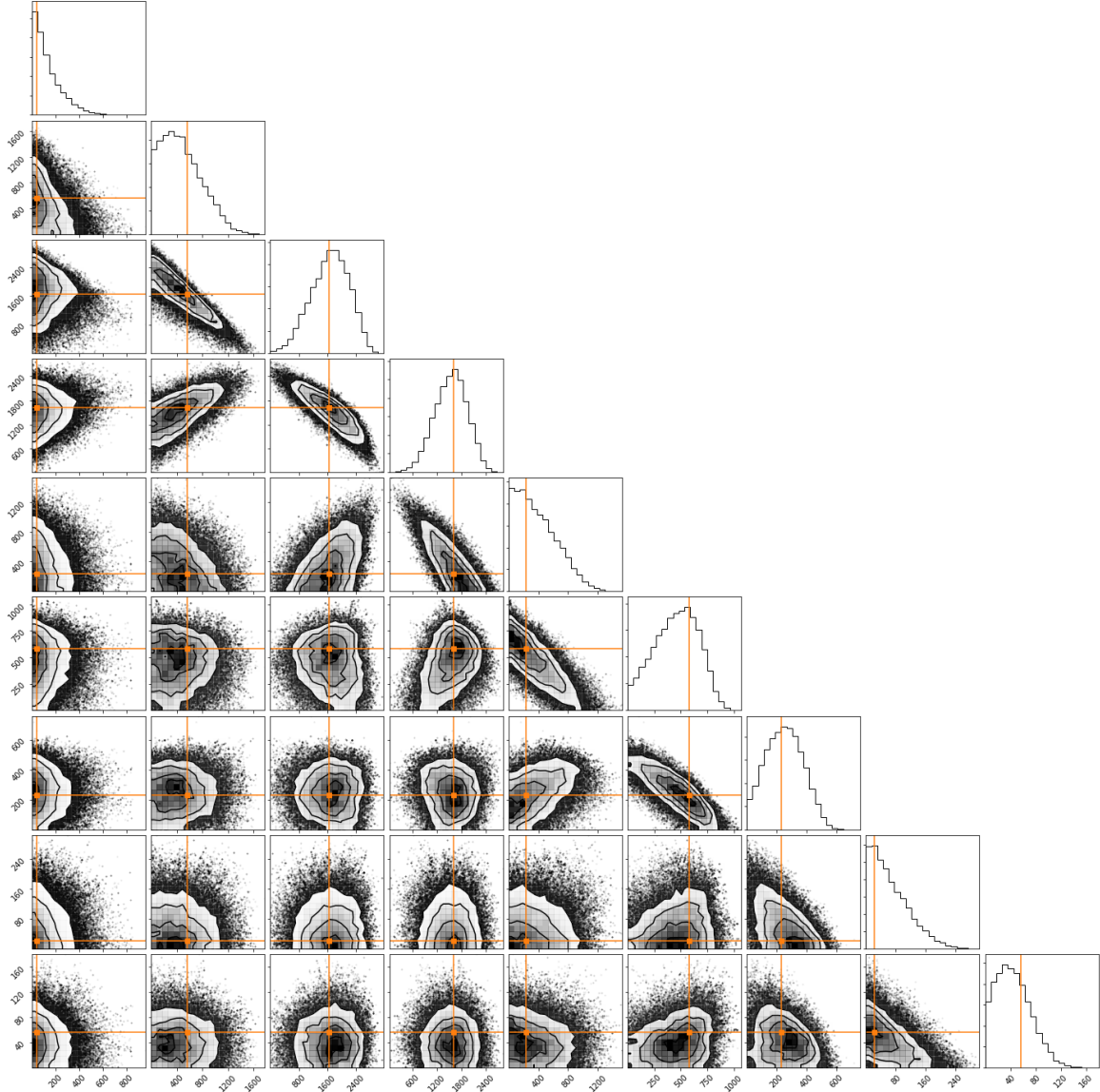


Figure 7: Result of the MCMC unfolding using equidistant binning with two observables. The plots on the diagonal show the distribution of the samples for the different bins in \vec{f} . The non-diagonal plots show the two dimensional distribution of two bins. All axes show number of events in bin. The marked points are the best fit.

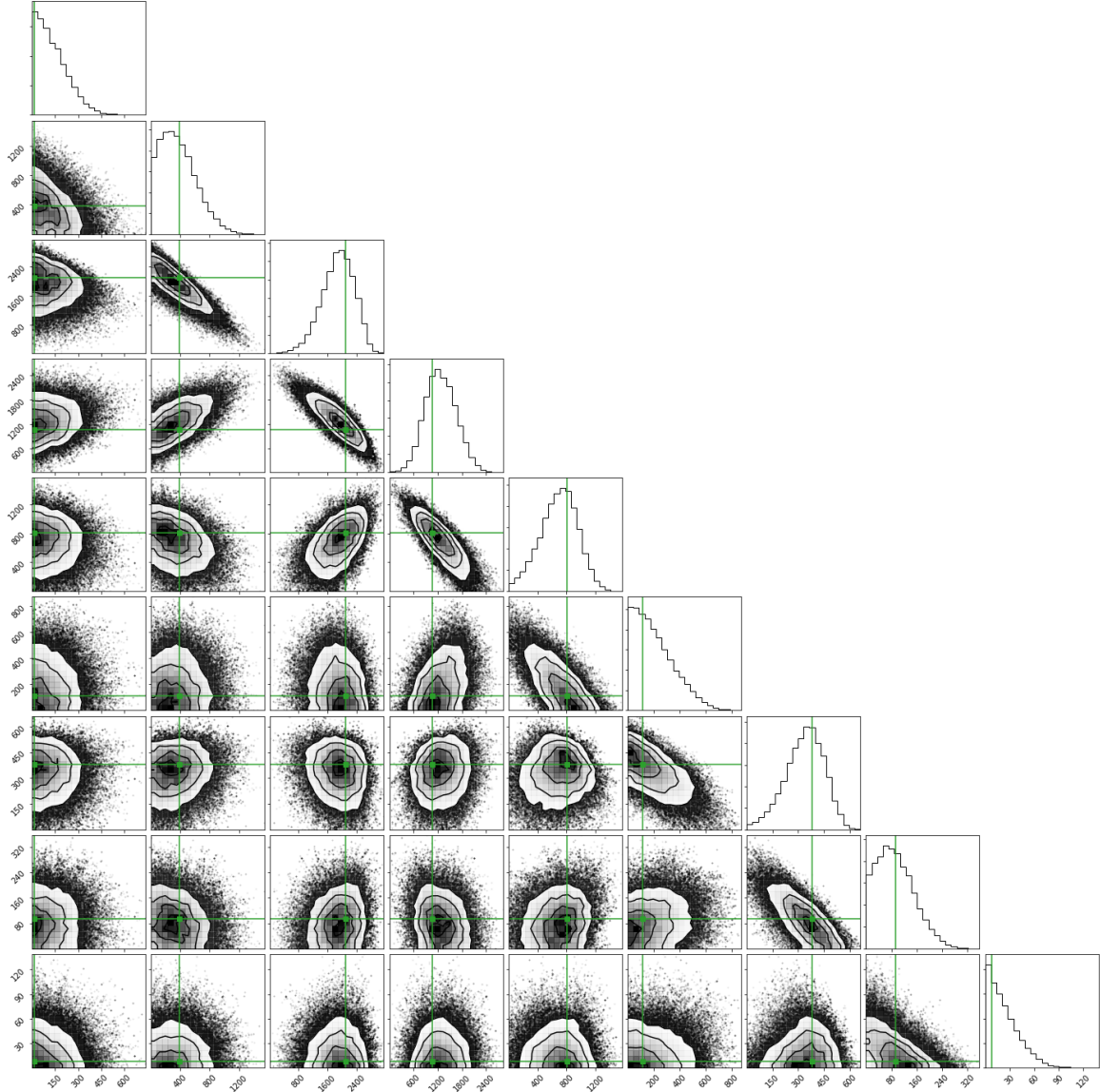


Figure 8: Result of the MCMC unfolding using decision tree based binning with two observables. The plots on the diagonal show the distribution of the samples for the different bins in \vec{f} . The non-diagonal plots show the two dimensional distribution of two bins. All axes show number of events in bin. The marked points are the best fit.

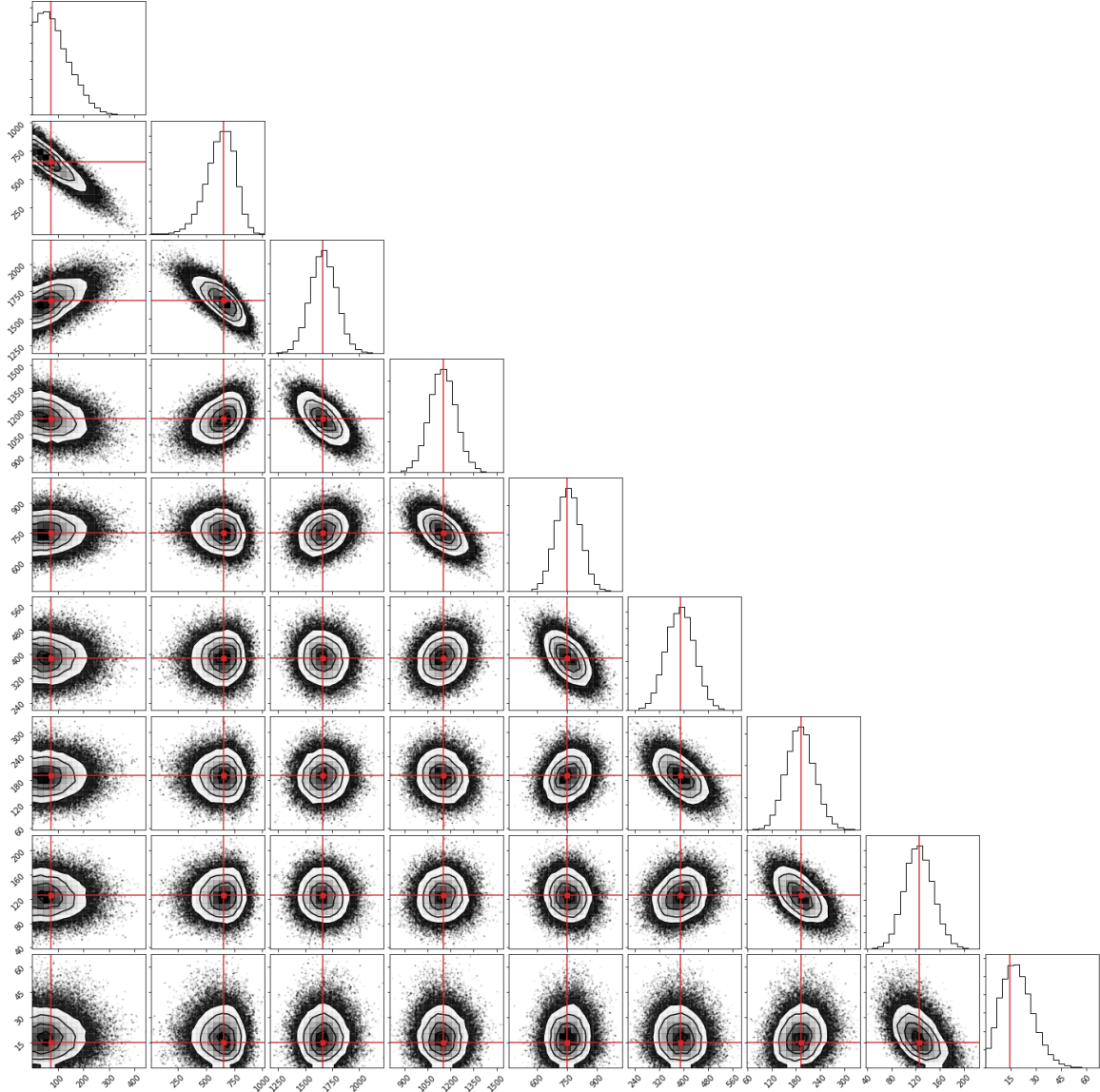


Figure 9: Result of the MCMC unfolding using decision tree based binning with all (18) observables. The plots on the diagonal show the distribution of the samples for the different bins in \vec{f} . The non-diagonal plots show the two dimensional distribution of two bins. All axes show number of events in bin. The marked points are the best fit.