# Comparison of Machine Learning Approach to other Commonly Used Unfolding Methods*

Petr Baroň

Joint Laboratory of Optics of Palacký University and Institute of Physics AS CR,
Faculty of Science, Palacky University, 17. listopadu 12, 771 46 Olomouc,
Czech Republic
petr.baron@upol.cz

Unfolding in high energy physics represents the correction of measured spectra in data for the finite detector efficiency, acceptance, and resolution from the detector to particle level.

Compared to other commonly used unfolding methods, recent machine learning approaches provide unfolding on an event-by-event basis using all the information from the collision similarly to the face recognition problem and allows to unfold spectra in a continuous way (independent of binning choice) the simultaneous unfolding of a large number of variables and thus can cover a wider region of the features that effect detector response. This study focuses on a simple comparison of commonly used methods in RooUnfold [1] package to the machine learning package Omnifold [2].

## 1. Introduction

The equation of unfolding can be written as

$$p = \frac{1}{\epsilon} \cdot M^{-1} \cdot \eta \cdot (D - B); \tag{1}$$

where $D$ is the data spectrum from which the background spectrum $B$ is subtracted followed by multiplication of acceptance correction $\eta$ so the main input to the unfolded procedure is prepared. The unfolding is here schematically given by a so-called migration matrix $M^{-1}$ which maps one-to-one events from the detector to particle level. Behind the symbol, $M^{-1}$ one could also imagine not necessarily the algebraic matrix inversion, but rather different unfolding methods, because the treatment of unfolding input differs. However, the aim of this study is not to fully describe all the

---

methods separately, but rather to make a comparison between them. The result of unfolding has to be corrected by the detector efficiency $\epsilon$ to obtain the unfolded "truth" spectrum $p$ ideally close to the truth (particle) level. Figure 1 provides insight to the ingredients on the example of the transverse momentum spectrum of the hadronically decaying top quark in process $pp \to t\bar{t}$ [3].
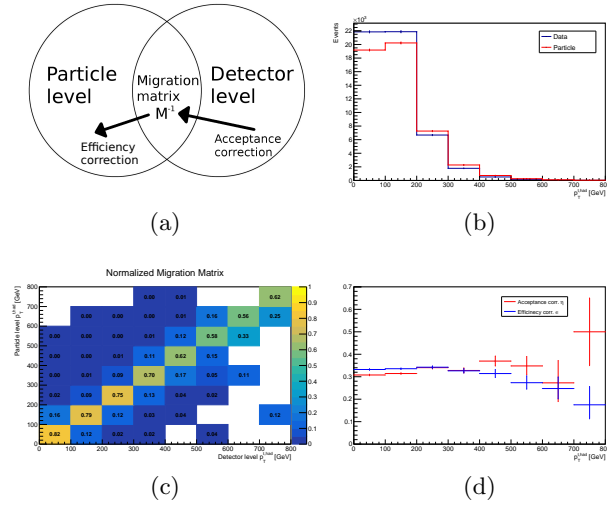


(a)

(b)

(c)

(d)

Fig. 1: Unfolding inputs. **a)** Unfolding procedure diagram **b)** Detector-level (blue) and particle-level (red) spectra. **c)** Migration matrix between particle and detector levels. **d)** Efficiency (blue) and acceptance (red) corrections as a function of the transverse momentum of the hadronically decaying top quark. [4]

## 2. Machine Learning Approach

The machine learning unfolding is similar to e.g. face recognition problem. The idea is to take all the possible information from the detector in a similar way as a photo and train some neural network to classify what process returns such a signature in the detector as the photo of the face is classified to one particular person.

In both cases, machine learning needs the truth information to train the neural network, e.g. the face on the photo belonging to a particular person or the detector signature belonging to the process $pp \to t\bar{t}$. The model describing the physics is still necessary. Authors of Omnifold [2] call the truth process and truth detector signature natural and the model is called
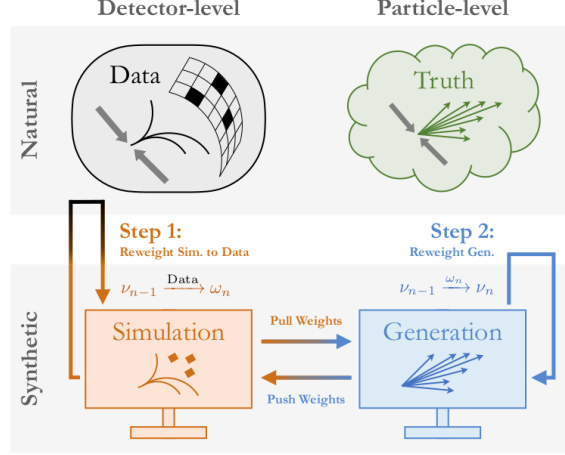
synthetic, see Figure 2.



Fig. 2: An illustration of the OmniFold method applied to a set of synthetic and natural data. As a first step, starting from prior weights $\nu_0$, the detector-level synthetic data ("simulation") is reweighted to match the detector-level natural data (simply "data"). These weights $\omega_1$ are pulled back to induce weights on the particle-level synthetic data ("generation"). As a second step, the initial generation is reweighted to match the new weighted generation. The resulting weights $\nu_1$ are pushed forward to induce a new simulation, and the process is iterated [2].

The authors extended the idea of Iterative Bayes unfolding to continuous form and with machine learning concept enabled to perform unfolding event-by-event so the trained network returns a set of weights for each event or a function which can be applied to measured data. The detailed description of the algorithm is in Appendix of [2].

## 3. Performing Closure Test

The closure test of the unfolding method is to unfold not the measured data, but rather the generated particle level spectrum. If the method is consistent their ratio should be close to unity. To avoid other systematical uncertainties from efficiency and acceptance corrections the events are chosen only from the overlap of the particle and detector level phase spaces, see the intersection of the circles in Figure 1.

As the process of study, the process of top quark pair production in proton-proton collisions $pp \rightarrow t\bar{t}$ in $\ell+$ jets channel was chosen, see Figure 3, simulated using MadGraph [5] software with a generation of events using

Pythia8 [6] with the detector-level and the pseudo data simulated using Delphes [7] with ATLAS detector card. The basic selection and cuts were applied to obtain spectra with a similar shape to those measured at Large Hadron Collider (LHC) in the real ATLAS experiment, although comparison of unfolding methods could be performed with an arbitrary measured process.



Fig. 3: **a)** Final state diagram of the process $pp \rightarrow t\bar{t}$ **b)** Top quark pair production branching ratios.

The input data set was divided into two statistically independent sets, so the classical unfolding methods were using migration matrix build from one set and the input to unfolding procedure from the other set. The same exclusive sets were used for training the neural network and to perform the machine-learning unfolding.

## 4. Results

The following spectra of interest were chosen: the transverse-momentum, mass, energy, and pseudo-rapidity of the hadronically, the leptonically decaying top quark, and also the $t\bar{t}$ system. The distribution $\phi$ was omitted due to its flat shape. In total $4 \times 3 = 12$ binned spectra were unfolded classically using the RooUnfold package with Bayes, SVD, and Ids methods [1]. The variables were used event-by-event in the neural network with 100 epochs to later perform simultaneous unfolding.

Even though results from the machine learning approach could be shown as continuous spectra, for comparison purposes particular fine binning was chosen.

Particle-level spectra as the input into the unfolding procedure were chosen to perform the closure test. Thus the ratio between input particle-level spectrum and unfolded spectrum should be ideally close to unity as it was already discussed in Section 3.

The metric of comparison is $\chi^2$ divided by the number of degrees of freedom NDF which is equal to the number of bins in the spectrum.

(a) $p_T$       (b) $p_T$       (c) $p_T$

(d) Mass       (e) Mass       (f) Mass

(g) Energy       (h) Energy       (i) Energy

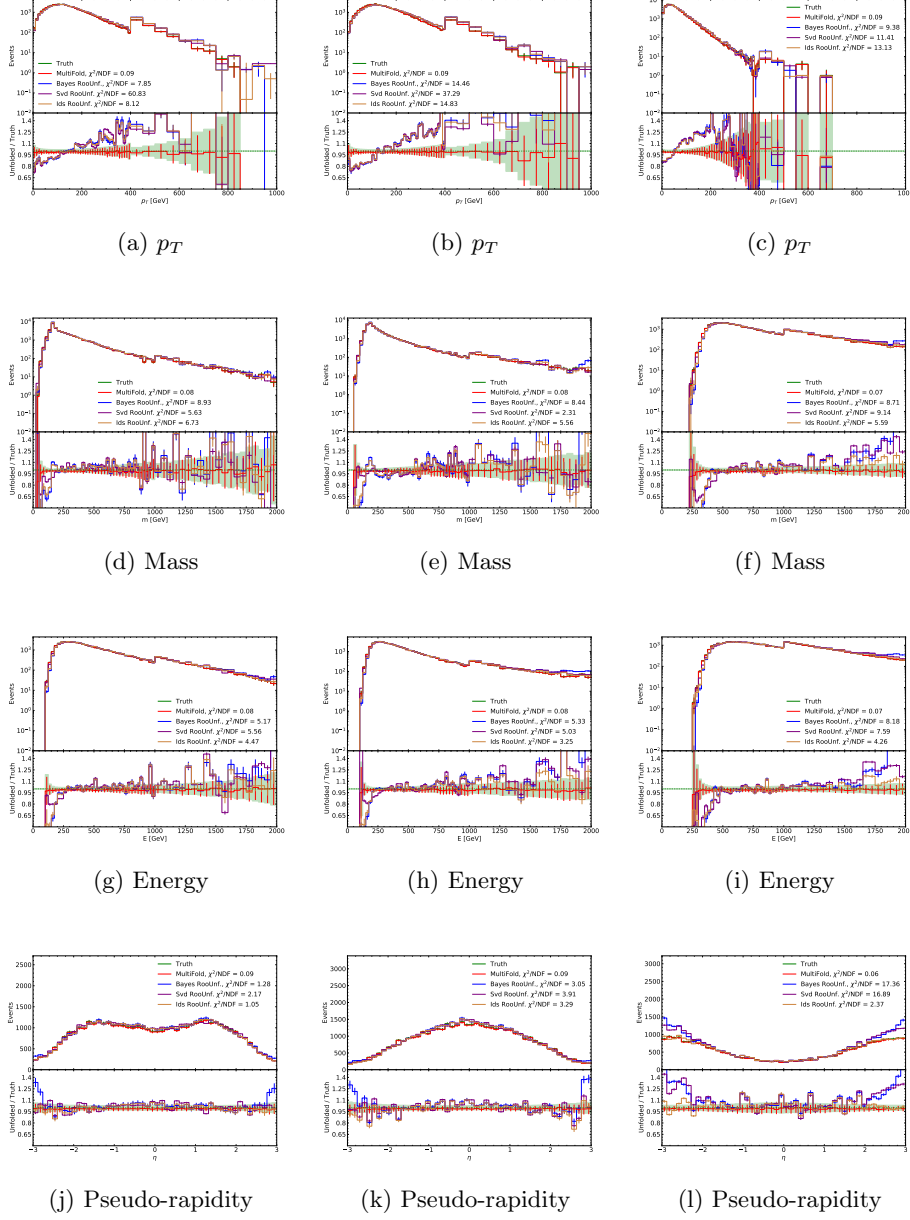(j) Pseudo-rapidity     (k) Pseudo-rapidity     (l) Pseudo-rapidity

Fig. 4: Closure test of the spectra of the hadronically **- left** leptonically the **- middle** decaying top quark and of the $t\bar{t}$ system the **- right** of the $p_T$, mass, energy and pseudo-rapidity with the lower pads showing the ratio of the unfolded to the truth spectrum.

|  | Tr. momentum | Mass | Energy | Pseudo-rapidity |
|---|---|---|---|---|
| Hadronic top | 0.09 | 0.08 | 0.08 | 0.09 |
| Leptonic top | 0.09 | 0.08 | 0.08 | 0.09 |
| $t\bar{t}$ system | 0.09 | 0.07 | 0.07 | 0.06 |

Table 1: Results of the closure test, $\chi^2$/NDF between the truth and unfolded spectra using the OmniFold method.

|  | Tr. momentum | Mass | Energy | Pseudo-rapidity |
|---|---|---|---|---|
| Hadronic top | 7.85 | 8.93 | 5.17 | 1.28 |
| Leptonic top | 14.46 | 8.44 | 5.33 | 3.05 |
| $t\bar{t}$ system | 9.38 | 8.71 | 8.18 | 17.36 |

Table 2: Results of the closure test, $\chi^2$/NDF between the truth and unfolded spectra using the Bayes RooUnfold method.

|  | Tr. momentum | Mass | Energy | Pseudo-rapidity |
|---|---|---|---|---|
| Hadronic top | 60.83 | 5.63 | 5.56 | 2.17 |
| Leptonic top | 37.29 | 2.31 | 5.03 | 3.91 |
| $t\bar{t}$ system | 11.41 | 9.14 | 7.59 | 16.89 |

Table 3: Results of the closure test, $\chi^2$/NDF between the truth and unfolded spectra using the Svd RooUnfold method.

|  | Tr. momentum | Mass | Energy | Pseudo-rapidity |
|---|---|---|---|---|
| Hadronic top | 8.12 | 6.73 | 4.47 | 1.05 |
| Leptonic top | 14.83 | 5.56 | 3.25 | 3.29 |
| $t\bar{t}$ system | 13.13 | 5.59 | 4.26 | 2.37 |

Table 4: Results of the closure test, $\chi^2$/NDF between the truth and unfolded spectra using the Ids RooUnfold method.

## 5. Conclusion

Table 1 summarizes the $\chi^2$/NDF results and proves that machine learning approach at this particular study performed the best results compared to values in Tables 2, 3 and 4.

The study aimed to demonstrate the possible potential of the machine learning methods on four typical spectra used in high-energy physics. A slight disadvantage of the machine learning method might be its initial CPU time needed to train the neural network. Although more complex tests of OmniFold unfolding need to be performed in the future, the $\chi^2/\mathrm{NDF}$ presents promising results.

## ACKNOWLEDGEMENTS

## REFERENCES

[1] Tim Adye, in Proceedings of the PHYSTAT 2011 Workshop on Statistical Issues Related to Discovery Claims in Search Experiments and Unfolding, CERN, Geneva, Switzerland, 17–20 January 2011, edited by H.B. Prosper and L. Lyons, CERN–2011–006, pp. 313–318.

[2] A. Andreassen, P. T. Komiske, E. M. Metodiev, B. Nachman and J. Thaler, "OmniFold: A Method to Simultaneously Unfold All Observables," Phys. Rev. Lett. **124** (2020) no.18, 182001 doi:10.1103/PhysRevLett.124.182001 [arXiv:1911.09107 [hep-ph]].

[3] P. Baroň, J. Kvita, "Extending the Fully Bayesian Unfolding with Regularization Using a Combined Sampling Method," Symmetry, Volume 12, 2020, doi:10.3390/sym12122100

[4] P. Baroň, "Fully Bayesian Unfolding in High-energy Physics," Acta Phys. Polon. B **51** (2020) no.6, 1241-1250 doi:10.5506/APhysPolB.51.1241

[5] J. Alwall *et al.*, "The automated computation of tree-level and next-to-leading order differential cross-sections, and their matching to parton shower simulations," JHEP **1407** (2014) 079 doi:10.1007/JHEP07(2014)079 [arXiv:1405.0301 [hep-ph]].

[6] T. Sjöstrand *et al.*, "An Introduction to PYTHIA 8.2," Comput. Phys. Commun. **191** (2015) 159 doi:10.1016/j.cpc.2015.01.024 [arXiv:1410.3012 [hep-ph]].

[7] J. de Favereau *et al.* [DELPHES 3 Collaboration], "DELPHES 3, A modular framework for fast simulation of a generic collider experiment," JHEP **1402** (2014) 057 doi:10.1007/JHEP02(2014)057 [arXiv:1307.6346 [hep-ex]].