Tools for unbinned unfolding and an application for jet measurements

Ryan Milton AI at the EIC Workshop 2025

In collaboration with: V. Mikuni, T. Lee, M. Arratia, T. Wamorkar, B. Nachman

ΔR



NSF CSSI 2311666



Overview

- Review of binned and unbinned unfolding
- Status of unbinned unfolding tools and application to jet measurements
- H1 unbinned unfolding
- Unbinned unfolding pipeline

Overview

• Review of binned and unbinned unfolding

- Status of unbinned unfolding tools and application to jet measurements
- H1 unbinned unfolding
- Unbinned unfolding pipeline

Overview of unfolding

• Objective: Remove detector distortions from experimental data



Experimental data

Physics information

m = Rt

Measured data

m = Rt

Measured data

Truth (physics) data

Response matrix R_{ij} = Pr(measure in bin i | truth in bin j) Made with simulation data

$$m = Rt$$

Measured data

Truth (physics) data

Response matrix R_{ij} = Pr(measure in bin i | truth in bin j) Made with simulation data







Main idea of binned unfolding is to ~invert the response matrix!



Examples: Singular value decomposition, Iterative Bayesian unfolding

 $t = R^{-1}m$

m = Rt

Main idea of binned unfolding is to ~invert the response matrix!



- How to define optimal binning?
 - Must be chosen before unfolding procedure
 - Binning choices makes it difficult to compare between experiments and to publish data

- How to define optimal binning?
 - Must be chosen before unfolding procedure
 - Binning choices makes it difficult to compare between experiments and to publish data
- Difficult to scale histograms to multiple dimensions by including multiple distributions

- How to define optimal binning?
 - Must be chosen before unfolding procedure
 - Binning choices makes it difficult to compare between experiments and to publish data
- Difficult to scale histograms to multiple dimensions by including multiple distributions
- Integrate over quantities thought to be irrelevant for a specific analysis
 - For other observables, like those motivated by future theoretical insight, may have to repeat analysis from scratch

Unbinned unfolding motivation

- Motivates an unbinned unfolding method using machine learning
- Naturally unbinned and can handle high dimensions







Unbinned unfolding with OmniFold

- Iterative Bayesian Unfolding (IBU) is a popular binned unfolding algorithm
- Each iteration can be broken down into two steps
- In each step, a likelihood ratio is approximated
- Can estimate these two ratios using classifiers instead
- Use these likelihood ratios as unfolding weights!



Applications of OmniFold

OmniFold has been applied to many experimental analyses!

PHYSICAL REVIEW LETTERS 128, 132002 (2022)

Measurement of Lepton-Jet Correlation in Deep-Inelastic Scattering with the H1 Detector Using Machine Learning for Unfolding

Measurement of event shapes in minimum bias events from pp collisions at 13 TeV

The CMS Collaboration

PHYSICAL REVIEW LETTERS 133, 261803 (2024)

Simultaneous Unbinned Differential Cross-Section Measurement of Twenty-Four Z + jets Kinematic Observables with the ATLAS Detector

> G. Aad *et al.** (ATLAS Collaboration)

PHYSICAL REVIEW D 108, L031103 (2023)

Multidifferential study of identified charged hadron distributions in Z-tagged jets in proton-proton collisions at $\sqrt{s} = 13$ TeV

> R. Aaij *et al.** (LHCb Collaboration)

Measurement of Collinear Drop jet mass and its correlation with SoftDrop groomed jet substructure observables in $\sqrt{s}=200~{\rm GeV}~pp$ collisions by STAR

YOUQI SONG (WRIGHT LABORATORY, YALE UNIVERSITY)

on behalf of the STAR Collaboration

Overview

- Review of binned and unbinned unfolding
- Status of unbinned unfolding tools and application to jet measurements
- H1 unbinned unfolding
- Unbinned unfolding pipeline

Status of unbinned unfolding tools

- <u>PyPi omnifold</u>:
 - Architectures: Multilayer perceptron (MLP), point-edge transformer (PET)
 - Any number of dimensions
 - May need GPUs to train
- <u>RooUnfold inspired omnifold</u>:
 - Architecture: Boosted decision tree (BDT)
 - Low dimensional unfolding
 - Simple set-up: No data preprocessing, no GPUs needed





- Test unfolding quality using closure test
 - Use two Monte Carlo simulations so we know what the truth looks like!
- Use Pythia as Monte Carlo (Gen/Sim), Herwig as pseudodata ("Truth"/"Data") to get unfolded jet observables

- Test unfolding quality using closure test
 - Use two Monte Carlo simulations so we know what the truth looks like!
- Use Pythia as Monte Carlo (Gen/Sim), Herwig as pseudodata ("Truth"/"Data") to get unfolded jet observables

Method	ML type	Input	Unfolding type
OmniFold	PET	All particles	Unbinned

- Test unfolding quality using closure test
 - Use two Monte Carlo simulations so we know what the truth looks like!
- Use Pythia as Monte Carlo (Gen/Sim), Herwig as pseudodata ("Truth"/"Data") to get unfolded jet observables

Method	ML type	Input	Unfolding type
OmniFold	PET	All particles	Unbinned
MultiFold	BDT, MLP (DNN)	All jet observables	Unbinned

- Test unfolding quality using closure test
 - Use two Monte Carlo simulations so we know what the truth looks like!
- Use Pythia as Monte Carlo (Gen/Sim), Herwig as pseudodata ("Truth"/"Data") to get unfolded jet observables

Method	ML type	Input	Unfolding type
OmniFold	PET	All particles	Unbinned
MultiFold	BDT, MLP (DNN)	All jet observables	Unbinned
UniFold	BDT	Individual jet observables	Unbinned

- Test unfolding quality using closure test
 - Use two Monte Carlo simulations so we know what the truth looks like!
- Use Pythia as Monte Carlo (Gen/Sim), Herwig as pseudodata ("Truth"/"Data") to get unfolded jet observables

Method	ML type	Input	Unfolding type
OmniFold	PET	All particles	Unbinned
MultiFold	BDT, MLP (DNN)	All jet observables	Unbinned
UniFold	BDT	Individual jet observables	Unbinned
IBU	N/A	Individual jet observables	Binned

Performance on jet observables

- Trying to match the "Truth" distributions by applying weights to Gen
- Other jet observables in paper: Width, Soft Drop mass, N-subjettiness ratio, groomed jet momentum fraction





27

Overview

- Review of binned and unbinned unfolding
- Status of unbinned unfolding tools and application to jet measurements
- H1 unbinned unfolding
- Unbinned unfolding pipeline

Application to H1

(H1 slides courtesy of Vinicius Mikuni)

Using **228 pb⁻¹** of data collected by the **H1 Experiment** during **2006** and **2007** at **318 GeV center-of-mass energy**





Q² = - q² y = Pq / pk

P: incoming proton 4-vector
k: incoming electron 4-vector
q=k-k' : 4-momentum transfer

Reconstructed hadrons using combined detector information: **energy flow algorithm**

Experimental setup



Fiducial Phase space definition:

- 0.2 < y < 0.7
- $Q^2 > 150 \text{ GeV}^2$

Particle selection:

- p_T > 0.1 GeV
- $-1 < \eta_{lab} < 2.75$
- Charge information used if $\eta_{\rm lab}$ < 2

Reco Phase space definition:

- 0.08 < y < 0.7
- $Q^2 > 150 \text{ GeV}^2$
- p_{τ} miss < 10 GeV,
- 45 < em/p₇ < 65

Particle selection:

- $p_{\rm T} > 0.1 \ {\rm GeV} \\ -1 < \eta_{\rm lab} < 2.75$
- Pass reco selection: Red -> Orange: 77%
- Pass fiducial selection: Red -> Blue: 58%
- Pass fiducial and reco selection: Blue -> Orange: 96%
- Don't pass fiducial but pass reco: Red -> Orange (without blue): 50%



Red box used during unfolding, but only fiducial results shown

Pretraining





Pre-training:

- We use a smaller version of the OmniLearn model¹
- We have around 20M simulated samples in Djangoh and Rapgap, but only around 200k data events
- We can improve the reweighting quality by pre-training the model first
- Train the model to classify Rapgap from
 Djangoh and use that as the initialization for the rest of the unfolding

1: **V. Mikuni**, B. Nachman, arXiv:2404.16091



Systematic uncertainties

Systematic uncertainties included in the results

- HFS energy scale: +- 1%
- HFS azimuthal angle: +- 20 mrad
- **Lepton energy:** +- 0.5%
- Lepton azimuthal angle: +- 1 mrad
- Model uncertainty: differences in unfolded results between Djangoh and Rapgap
- Non-closure uncertainty: Differences between the expected and obtained values of the closure test
- Statistical uncertainty: Standard deviation of 50 bootstrap samples with replacement (will increase to 100 when done running)

Uncertainties not yet added

QED uncertainty: Use the variation of measured quantities when radiation is turned off in the simulation

Closure test



- Use Djangoh as the **pseudo-data** and unfold Rapgap
- Analysis is currently **preliminary** and **blinded**: Uncertainties shown use data information, but the plots will show only the closure results
- Although we **unfold everything**, we make a small number of observables public
 - Observables that reproduce previous results
 - Observables difficult to unfold classically

	Lab $k_{\rm T}$	Breit $k_{\rm T}$	Centauro
$\Delta \phi^{ m jet}$	Y	Ν	Ν
$\ln(\lambda_1^1)$	Y	Ν	Ν
$p_{ m T}^{ m jet}$	Y	Y	Y
z ^{jēt}	Y	Y	Y

Jet Selections

- Jet algorithm: kT and Centauro both with R = 1
- Jet pT cuts:
 - Lab frame kT jets: $p_{T} > 10 \text{ GeV}$
 - Breit frame kT jets: $\dot{p}_{T} > 5 \text{ GeV}$
 - Centauro jets: z > 0.2 (no pT cut)



Breit Frame provides a natural frame to study ep collisions, where the struck quark forms a jet opposite from the proton beam: useful for jet and TMD studies

Centauro jet algorithm is a longitudinally invariant method designed for DIS studies

Δz and Energy-Energy Correlator are also measured in the Breit frame.

Cluster jets using kT algorithm with radius of 1.0 Notice the analysis is **blind**: Plot shows the closure agreement, but uncertainties are calculated based on the data









Cluster jets using kT algorithm with radius of 1.0 Notice the analysis is **blind**: Plot shows the closure agreement, but uncertainties are calculated based on the data





Cluster jets using kT algorithm with radius of 1.0 Notice the analysis is **blind**: Plot shows the closure agreement, but uncertainties are calculated based on the data





$$z_{\text{jet}} = \frac{P \cdot p_{\text{jet}}}{P \cdot q} \quad \xrightarrow{\text{Breit}}_{\text{frame}}$$



Cluster jets using kT algorithm with radius of 1.0 Notice the analysis is **blind**: Plot shows the closure agreement, but uncertainties are calculated based on the data



Energy-Energy Correlator (only defined) in Breit frame



$$\operatorname{EEC}_{\mathrm{DIS}} = \sum_{a} \int \frac{\mathrm{d}\sigma_{ep \to e+a+X}}{\sigma} \, z_a \, \delta(\cos\theta_{ap} - \cos\theta) \,,$$

$$z_a \equiv \frac{P \cdot p_q}{P \cdot (\sum_i p_i)},$$

Energy-energy correlators in Deep Inelastic Scattering: Li et al. Phys.Rev.D 103 (2021) 9, 094005

Overview

- Review of binned and unbinned unfolding
- Status of unbinned unfolding tools and application to jet measurements
- H1 unbinned unfolding
- Unbinned unfolding pipeline

Binned inference

- To go from measurements to matching with theory, need to do inference
- Typically, a functional form is assumed and parameters are fit
 - Computationally expensive and depends on binning



Conventional inference

ML-based inference

- New approach: Train a model to output probability of data given parameters
- Use this model with unbinned unfolded data!



Summary and outlook

- Unfolding removes detector distortions from data
- Binned unfolding is limited, motivating an ML unbinned unfolding approach
 - With ML: Maintain correlations, high dimensional unfolding, no reliance on bin definitions
- We have developed powerful and easy to use tools to do unbinned unfolding
- BDT OmniFold is in a fork of RooUnfold and plans for including it into RooUnfold are ongoing
- Can apply OmniFold to H1 experimental analysis with good closure
- Can create a fully unbinned pipeline with unbinned inference

Thank you!