



# Exploring machine learning methods for unbinned data analysis of Drell-Yan phase space at CMS

---

BSc. David Gutiérrez Menéndez Dr. Fernando Guzmán Martínez  
November 1, 2023

Higher Institute of Technology and Applied Sciences, Havana University

# Introduction

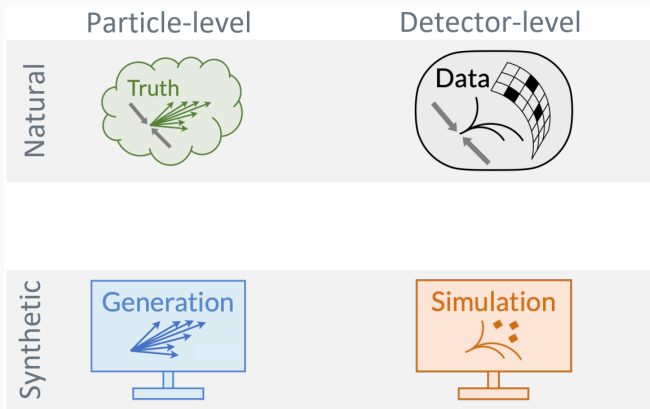
---

## Particle detectors introduce smearing in measurements:

- Finite detector resolution
- Acceptance cuts (geometrical and physical)
- Binning effect in distributions

# Introduction

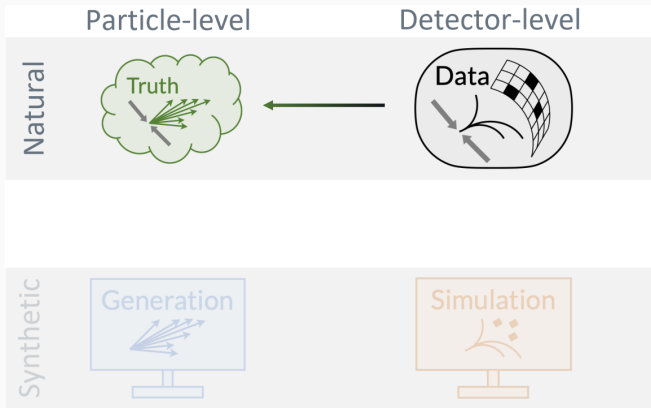
**Goal: Use data (at detector level) and infer true distributions (at particle level).**



**Only then we can reliably compare the measurements against other sources:**

- Semi-Inclusive theoretical predictions
- MC generator parameter tuning
- Other experiments

## Unfolding (deconvolution).



## **Current Unfolding approach**

---

Detector response can be modeled as a linear application.

$$x_{detector} = R * x_{truth} + x_{backg} \quad (1)$$

Where  $R$  represents the detector **response** as  $R_{ij} = P(\text{truth}_i | \text{measure}_j)$ .





## Various approaches exist to unfolding:

- Naive  $R$  inversion

$$x_{truth} = R^{-1} * (x_{detector} - x_{backg}) \quad (2)$$

- Iterative Bayesian Unfolding (IBU): Given a **response** matrix and a prior distribution.

$$x_j^{(n)} = \sum_i P_{n-1}(truth_i | measure_j) * P(measure_j) \quad (3)$$

## Unfolding drawbacks:

- Works for binned data
- 1-dimensional by design (very difficult to extend to higher dimensions)
- Does not account for relations between observables

# Machine Learning approach

---

## Distribution reweight.

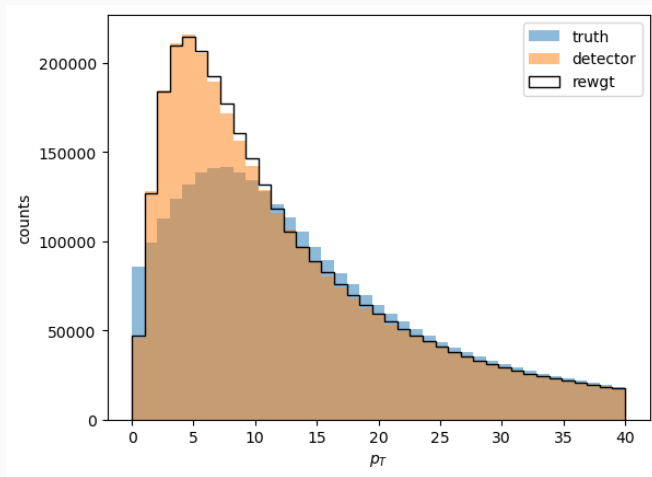
- Using a binary classifier to distinguish between two distributions
- Likelihood ratio can be approximated and used as weight function

$$w_{A \rightarrow B} = \frac{p_B(x)}{p_A(x)} \approx \frac{f(x)}{1 - f(x)} \quad (4)$$

- Distributions are binned for representation but weights are estimated per event



# Machine Learning approach



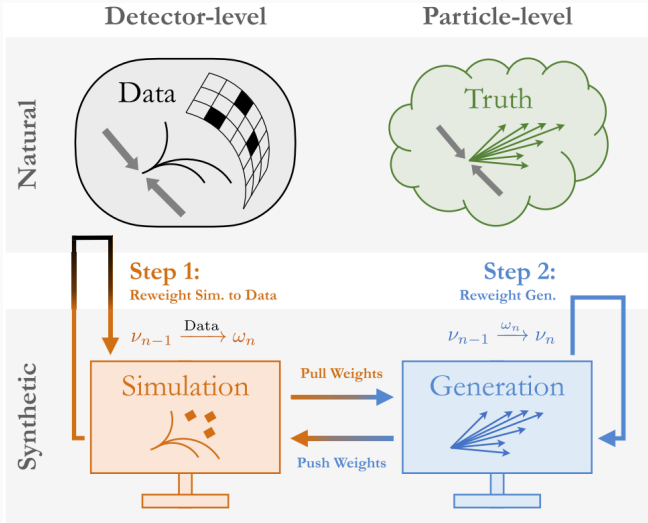
**Figure 1:** Transverse momentum distribution.  $pp \rightarrow DY7TeV$ .

## Advantages of ML reweight.

- Used for non-trivial distributions
- Easily extended to multiple dimensions
- Learns complex relations between observables

# Machine Learning approach

## Iterative Neural Network Reweight: OmniFold.



## Iterative Neural Network Reweight: OmniFold.

- Continuous generalization of IBU
- Inherits all the advantages from ML reweight
- Unfolds experimental data using multiple dimensions (possibly full phase space)





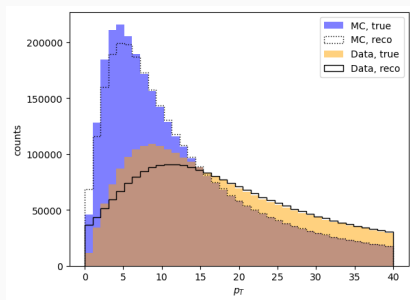
**Work in progress**

---

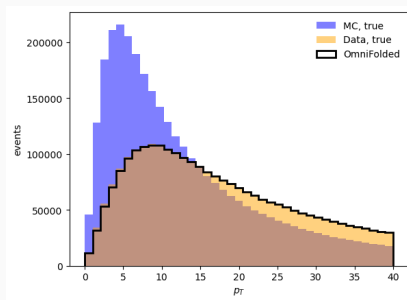
## Drell-Yan phenomenology analysis with CMS data.

- Explore HEP analysis viability with **Python**
- Unlock ML ecosystem and methods
- Validation of NanoAODRun1 format for CMS data
- Improve accuracy of inferred distributions

## Drell-Yan phenomenology analysis with CMS data.



**Figure 2:** Synthetic and Natural distributions.



**Figure 3:** Result of reweighted distribution.

## Future work.

- Full phase space unfolding
- Unbinned and data driven analysis
- Evaluate statistical and systematical uncertainties
- Search for more **robust** observables where current theoretical predictions struggle



# Exploring machine learning methods for unbinned data analysis of Drell-Yan phase space at CMS

---

BSc. David Gutiérrez Menéndez Dr. Fernando Guzmán Martínez  
November 1, 2023

Higher Institute of Technology and Applied Sciences, Havana University