<u>deep ai image editor</u>

Machine Learning in Analysis Techniques for future physics

Hannah Bossi (MIT) Quark Matter 2025 Frankfurt, Germany April 12th, 2025









HOW IS ML CURRENTLY BEING USED AS AN ANALYSIS **TECHNIQUE?**

*Focus on new items since QM 2023 or items presented at QM 2025

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

WHAT DOES THE FUTURE HOLD?

ORIGINS OF ML

 First electronic computers developed with the purpose of performing computations that were difficult and timeconsuming for humans to perform accurately.

[Electronic Numerical Integrator and Computer, 1945]

 Later push for computers to replicate computations humans are naturally good at, i.e. pattern recognition, but on larger datasets.

As a result, created algorithms that imitate how humans learn, i.e. learning from examples in a way that improves over time \rightarrow machine learning

Hannah Bossi (hannah.bossi@cern.ch)

• Fundamental building blocks of machine learning directly mimic the brain!

[Press release Nobel Prize in Physics 2024]

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

Neurons \rightarrow Nodes Synapses \rightarrow Weighted couplings Learning \rightarrow training

 Connections between some neurons (synapses) become stronger and others weaker as we learn.

<u>.O Hebb, "The organization of behavior", 1949]</u>

Collections of nodes give rise to computational collective behavior. - artificial neural networks

[Press release Nobel Prize in Physics 2024]

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

• Hopfield network: binary nodes ($s_i = 0/1$), information read in/read out from all nodes.

ML Analysis Techniques

Massachusetts Institute of

Collections of nodes give rise to computational collective behavior. - artificial neural networks

• Boltzmann machine: Nodes are subdivided into visible nodes and hidden nodes.

[Press release Nobel Prize in Physics 2024]

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

• Hopfield network: binary nodes ($s_i = 0/1$), information read in/read out from all nodes.

John Hopfield

Geoffrey Hinton (Hopfield Network) (Boltzmann Machine)

Hannah Bossi (hannah.bossi@cern.ch)

[Press release Nobel Prize in Physics 2024]

ML Analysis Techniques

Massachusetts Institute of Technology

ML AND PHYSICS

Goal of measurements: To extract relevant physics information from available data!

[CMS observes jet quenching]

physics/experimental constraints (2) perform a statistical analysis on selected data.

Quark-Gluon Plasma (QGP)

Historical approach: (1) make selection using a series of boolean decisions motivated by

What do we do when decision becomes difficult to derive from expert-knowledge alone?

alysis Techniques

Number of HEP-ML Papers by Year

[HEPML-Living Review]

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

 Increase in computing power and collection of large datasets at modern facilities have lead to increased use of ML for physics.

In the coming decades, ML will help usher high energy physics into it's next era! (More on this later)

*Focus on new items since QM 2023 or items being presented at QM 2025

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

ML Analysis Techniques

HOW IS ML CURRENTLY BEING USED AS AN ANALYSIS **TECHNIQUE?**

WHAT DOES THE FUTURE HOLD?

Massachusetts Institute of

GENERATIVE MODELS

CLASSIFICATION

EXPERIMENTAL RESULTS

Radius of circle = # of papers/4

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

ML Analysis Techniques

REGRESSION

UNCERTAINTY QUANTIFICATION

* Note in this scheme papers are allowed to count in more than one category.

SIMULATION

BASED

INFERENCE

🐹 Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

- Higher particle multiplicities, higher degree of complexity.
- Large variance in simulations of heavyion collisions causes high dependence on sample used in training.

ML Analysis Techniques

Massachusetts Institute of

ML IN HEAVY IONS

EXPERIMENTAL RESULTS

CLASSIFICATION

UNCERTAINTY QUANTIFICATION

Radius of circle = 4 X # of papers

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

SIMULATION BASED INFERENCE

- Investigate paper subset within HEP that specifically relates to HIs or is a technique readily employed in HIs.
- Similar distribution between the different subtopics as HEP overall.

* Note experimental results include those performed with TMVA * Note in this scheme papers are allowed to count in more than one category.

FIRST: WHAT CAN ML NOT DO?

Garbage In

Garbage Out

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

SIGNAL/BACKGROUND CLASSIFICATION

 Large amount of experimental applications of machine learning rely on signal vs. background classification.

Hannah Bossi (hannah.bossi@cern.ch)

 Decision trees are commonly used for signal classification •Each node is a classification rule that splits the data into two or more parts.

- In training you determine the proper rules that
- maximizes the information gain and minimize entropy

$E(x) = \sum -p(x)log_2(p(x))$

SIGNAL/BACKGROUND DISCRIMINATION

 Boosted decision trees are used when multiple weaker learners are combined in a series where each additional component seeks to minimize error of previous one.

Hannah Bossi (hannah.bossi@cern.ch)

Yann Coadou, <u>arXiv: 2206.09645</u>

ML Analysis Techniques

Massachusetts Institute of

HEAVY FLAVOR JET TAGGING **Goal:** identify jets initiated by a heavy-quark Conventional approach: Apply cuts to select jets with displaced decay vertices and large impact parameter tracks.

- ML approach: Learn from low-level features in a supervised approach using BDT or a GNN

CMS-PAS-HIN-24-005

[JINST 13 (2018) 05, P05011]

ML Analysis Techniques

Massachusetts Institute of

Goal: identify je Conventional decay vertices

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

ENABLES THE FIRST OBSERVATION OF B-QUARK DEAD CONE!

upervised approach using BDT or a GNN

CMS-PAS-HIN-24-005

ML Analysis Techniques

Massachusetts Institute of Technology

 $QMMM_{25}$

ENABLES IMPROVED MEASUREMENT OF THE B-JET **CROSS SECTION IN ALICE!!**

Hannah Bossi (<u>hannah.boss</u>i@cern.ch)

Massachusetts Institute of Technology

 $QMMM_{25}$

- Goal: identify jets initiated by a heavy-quark **Conventional approach** SPHENIX decay vertices and large ML approach: Learn from low-level tea Use "long short-term memory (LTSM)" network that allows important information to be retained over long sequences.
- Training variables:
 - Jet variables
 - Jet constitutent variables
 - Event variables

🐹 Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

ML METHOD SHOWS IMPROVED HF JET elect i er track **TAGGING EFFICIENCY OVER TRADITIONAL** APPROACHING approach using BDT or a GNN

ML Analysis Techniques

GMNMM 25

REGRESSION OF JET PT

•Differential measurements of jets are key to understanding jet quenching! •These often involve pushing to large R and/or low $p_{\rm T}$, where background

contribution is difficult to subtract.

By now many methods in which ML can be used to solve this problem! We will discuss one.

See also.....

Shallow NN trained on PYTHIA + HI Background

[PRC 99, 064904 (2019)] ALICE, [PLB 849 (2024) 138412]

- CNN for regression of energy loss fraction [JHEP 03 (2021) 206]
- Interpretable Deep NN

[Phys.Rev.C 108 (2023) 2, L021901]

ML Analysis Techniques

Massachusetts Institute of

- •Use generative AI (unpaired image-to-image translation, cycleGANs) to subtract jet background in an unsupervised way. [UVCGAN: arXiv:2203.02557]
 - Composed of two generator-discriminator pairs w/ cyclic closure (i.e.

 $A \rightarrow B \rightarrow A \sim A$

- One to translate from domain $A \rightarrow B$
- One to translate from domain $B \rightarrow A$ DOMAIN A

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

ML Analysis Techniques

Massachusetts Institute of

CMMM 25

UNFOLDING WITH ML

<u>[PRL 124, 182001 (2020)]</u>

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

Goal of unfolding: Correct measured distributions for smearing.

Conventional Approach: Apply unfolding procedure on a binned distribution and repeat for each observable.

ML-based Approach: Use ML to calculate weighting factors and unfold the phase space all at once, before the choice of binning or observable!

- Tested for the first time on simulation in a HI environment (PYTHIA/HERWIG) + thermal background), similar or better performance to Bayesian unfolding in 3D.
- No explicit background subtraction, built into MultiFold-HI!

[A. Falcao, A. Tackas, to appear on arXiv]

🐹 Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

ML Analysis Techniques

SIMULATION BASED INFERENCE

Massachusetts Institute of

smeareo

elphes

es

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

ML Analysis Techniques

MULTIFOLD ENABLES MULTI-DIFFERENTIAL MEASUREMENTS OF JET SUBSTRUCTURE IN AU+AU.

Conventional Approach: Apply unfolding procedure on a binned distribution and repeat or each observable. TREATMENT OF SIMULATION DEPENDENCE SHOWN

IL-based AppFOR THE FIRST TIME! veighting factors and unfold the phase space all at once, before the choice of binning or observable!

Massachusetts Institute of

HOW IS ML CURRENTLY BEING USED AS AN ANALYSIS **TECHNIQUE?**

*Focus on new items since QM 2023 or items presented at QM 2025

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

WHAT DOES THE FUTURE HOLD?

Very large volumes of will be taken and analyzed in the decades to come - new tools will be increasingly important!

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

Contact: Karen McNulty Walsh, (631) 344-8350, or Peter Genzer, (631) 344-3174

Brookhaven's Computing Center Reaches 300 Petabytes of Stored Data

Largest compilation of nuclear and particle physics data in the U.S., all easily accessible – with plans for much more

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

"RHIC's newest detector, sPHENIX, with a readout rate of 15,000 events per second, is projected to more than double the data we

Total expected output from sPHENIX ~565 Petabytes DRAMATIC IMPROVEMENTS IN READOUT RATES MAKING SPHENIX A LARGE FRACTION OF TOTAL DATA STORED!

ML Analysis Techniques

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

be increasingly in SAME TRENDS TRUE AT THE LHC!

EVENT FLTERNG

- Data volume is increasing at a fast rate, need solutions for limited computing resources.
 - If we took all raw data, would easily exceed storage capabilities.

ML Analysis Techniques

Massachusetts Institute of Technology

ſĺſĨŧ₂₅₉

EVENT FLTERNG

- Data volume is increasing at a fast rate, need solutions for limited computing resources.
 - If we took all raw data, would easily exceed storage capabilities.
- Perform fast selection/rejection of data with ML integrated into the firmware (FPGAs)
 - Use high level synthesis packages ex: <u>hls4ml</u>

CMS L1 Trigger [CMS-TDR-021] sPHENIX HF Trigger [JINST 19 C02066]

ATLAS Fake Track Rejection in Event Filter [ATLAS-TDR-029-ADD-1]

LHCb track reconstruction for HLT system [See website here]

ML Analysis Techniques

Massachusetts Institute of Technology

GMNN

FULL EVENT SIMULATION

- Full event simulations are highly complex and can be prohibitively computationally expensive (especially in HIs)

ullet

• Ex: using denoising diffusion probabilistic model (DDPM) and Generative Adversarial Networks (GAN) to generate heavy-ion events in sPHENIX designed to mimic HIJING

Time per event					
40 mins					
1.34s					
0.42 ms					

Tradeoff between time and accuracy!

[Phys. Rev. C 110, 034912 (2024)]

PARTIAL EVENT SIMULATION

that with machine learning.

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

Can also take components of simulation that are computationally expensive and replace

(If Needed) Block 4: Detector sim

ML Analysis Techniques

Massachusetts **Institute** of

• Goal: learn data-driven fragmentation function by estimating likelihood ratios for each string break.

 Use event weights derived from base simulation in order to train two neural networks that are used to predict string break weights.

ABLE TO REPRODUCE FRAGMENTATION FUNCTION FROM DATA!

ML Analysis Techniques

Massachusetts Institute of Technology

lace

HOW CAN WE MAKE ML-BASED **APPLICATIONS REPRODUCIBLE?**

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

ML Analysis Techniques

HOW DO WE CONSTRUCT MORE INTERPRETABLE **MODELS?**

DO WE NEED TO STANDARDIZE ML **APPLICATIONS ACROSS EXPERIMENTS?**

Massachusetts Institute of

Bias is not inherently bad!

- Just like any analysis technique, you need to study and quantify the bias on your results. • **Should** correspond to a systematic uncertainty
- No standard way to do this for physics applications, but there has been progress!

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

Bias is not inherently bad!

- Just like any analysis technique, you need to study and quantify the bias on your results. • Should correspond to a systematic uncertainty
- No standard way to do this for physics applications, but there has been progress!

Example case: jet $p_{\rm T}$ **reconstruction**

 Initial ALICE approach is to train on five different fragmentation models, taking difference in final result as systematic uncertainty

Bias is not inherently bad!

- Just like any analysis technique, you need to study and quantify the bias on your results. • Should correspond to a systematic uncertainty
- No standard way to do this for physics applications, but there has been progress!
- If the bias is too large, traditional techniques or ML w/ reduced complexity is preferable.

Example case: jet $p_{\rm T}$ **reconstruction**

 Indications that corresponding systematic uncertainty may be larger at lower $p_{\rm T}$

[arXiv:2412.15440]

Bias is not inherently bad!

- Just like any analysis technique, you need to study and quantify the bias on your results. • Should correspond to a systematic uncertainty
- No standard way to do this for physics applications, but there has been progress!

Example case: jet $p_{\rm T}$ **reconstruction**

- As simulations improve, so will our treatment of the systematic uncertainties.
- •Training on LBT model appears to better reproduce results from traditional approaches than when trained on PYTHIA.
- Applied to LBT + thermal background.

Massachusetts

Institute of

CONCLUSIONS

- We are taking more data and making more complex measurements than ever before!
- analysis pipeline!
 - Many great examples at this conference!

Machine learning has led to new physics insights and can be used throughout the whole

- O
- Will be crucial at future facilities such as FAIR, HL-LHC, and the EIC!
- Stay tuned for many interesting results in the future!!

<u>deep ai image editor</u>

Thank you!

Special thanks to Lee Barnby, Xuan Li, Changwhan Choi, and the MIT Heavy lon group for useful discussions and feedback!

deep ai image editor

Backup

111

example, with auto-encoders.

ML Analysis Techniques

Massachusetts Institute of

WHAT IS AI/ML?

Artificial Intelligence: Programs with the ability to acquire and apply knowledge and skills.

ARTIFICIAL INTELLIGENCE

MACHINE LEARNING

Machine Learning: computational algorithms that imitate how humans learn, i.e. learning from examples in a way that improves over time

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

ML Analysis Techniques

Ex: Chatbots (humans give rules)

ACCELERATOR COMPLEX

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

MACHINE LEARNING CAN BE USED THROUGHOUT THE ANALYSIS PIPELINE!

ML Analysis Techniques

Massachusetts Institute of

~ Given an answer ~ "White Box" ML ~ Underlying physics

Helpful in understanding uncertainties or shortcomings of models!

Proof of concept identifying the AP splitting function exists [PLB 829 (2022) 137055]

Kannah Bossi (<u>hannah.bossi@cern.ch</u>)

- "Data"-based learning complements simulation-based inference.
 - ~ Domain knowledge
 - ~ "Black Box" ML
 - ~ Answer

- THIS IS A LONG TERM EFFORT!

HOW DOES THE MACHINE LEARN?

SUPERVISED LEARNING

Algorithm learns from a labeled set of "true values".

UNSUPERVISED LEARNING

Algorithm finds structure in the data without knowing the desired outcome.

Driven by the Task Analogy: Taking a test

Driven by the Data Analogy: Clustering

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

REINFORCEMENT LEARNING

Algorithm learns in a reward based system to determine a series of actions.

ML Analysis Techniques

Massachusetts Institute of Technology

STUDIES WITH NN JET PT RECONSTRUCTION

Hannah Bossi (hannah.bossi@cern.ch)

ML Analysis Techniques

- See offset (bias) in $\delta p_{\rm T}$ when ML is trained in PYTHIA vs. LBT.
- Crucial for applications in data to correct for this bias in an unfolding procedure.
 - Apply same model on your data and the response matrix.

Massachusetts Institute of Technology

Use supervised learning on jet images with a CNN to perform the regression task of predicting the energy loss ratio in HI collisions (hybrid model).

 Very useful to separate and study quenched vs. unquenched jets as well as extracting the initial energy of the jet. (Ideal probe of selection bias!)

[JHEP 2021, 206 (2021)]

Shows good performance!

🐹 Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

Weights: $w(x) = p_0(x)/p_1(x)$ Ok for 1D

 $\approx f(x)/(1 - f(x)) \frac{(\text{Andreassen and Nachman})}{(\text{PRD 101, 091901 (2020)})}$

where f(x) is a neural network and trained with the binary crossentropy loss function

> to distinguish jets coming from <u>data</u> vs from simulation

Unfolding \rightarrow Reweighting histograms \rightarrow Classification \rightarrow Neural network

Slide from Yougi Song

Extract splitting function from the network in white-box ML.

Done with a GAN split into two components.

1. Time independent learns the z, ϕ

2. Time dependent learns the θ

Was able to reproduce AP splitting function.

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

Massachusetts

Institute of

EVENT CLASSIFICATION AT THE EIC

- Study the effectiveness of ML-based classifiers to
 - Identify the flavor of the jet
 - Identify the underlying hard process of the collision
- Additionally study the effectiveness of different ways of representing information
 - Particle Flow Networks [JHEP 01 (2019) 121]

$$F(p_1, \dots, p_N) = F\left(\sum_{i=1} \Phi(p_i)\right) \quad p_i = (z_i, \eta_i, \phi_i, \text{PII})$$

Energy Flow Polynomials [JHEP 04 (2018) 013]

$$\operatorname{EFP}_{G} = \sum_{i_{1}} \cdots \sum_{i_{V}} z_{i_{1}} \cdots z_{i_{V}} \prod_{(k,l) \in E} \theta_{i_{k}i_{l}}$$

Indications that ML-based methods will have an improved performance over traditional techniques! See also event classification with large language models, [arXiv:2404.05752]

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

<u>[JHEP 03 (2023) 085]</u>

better balance of precision (purity) and recall (efficiency)!

See also LHCb NN to identify calo hits [Int. J. Mod. Phys. A 30, 1530022 (2015)], ATLAS Electron PID w/ CNN [ATL-PHYS-PUB-2023-001] CMS Deep NN to identify hadronic τ -lepton decays [JINST 17 (2022) P07023]

 Use NN trained on a specific particle type to predict a certainty value that is then compared to a pre-set threshold. Decide threshold based on efficiency/purity tradeoff. Takes into account particles from different sub-detectors (here TPC, TOF, TRD of ALICE), robust against missing data.

	Droton					Koor	`
7	Regression	97.38 ± 0.40	93.67 ± 0.38		Regression	91.17 ± 01.00	81.78
05	Proposed	97.80 ± 0.44	93.86 ± 0.27		Proposed	91.55 ± 0.71	83.6
7	Mean	97.85 ± 0.41	93.34 ± 0.32		Mean	90.83 ± 01.71	82.3
1	Ensemble	97.16 ± 0.46	93.74 ± 0.30		Ensemble	91.18 ± 02.00	82.72
1	Standard	99.40 ± 0.01	59.72 ± 0.03		Standard	92.87 ± 0.01	60.3
	Model	Precision	Recall		Model	Precision	Reca

PIULUII

naun

• When comparing the standard method to the proposed method, proposed method has

TRACK RECONSTRUCTION AT THE HL-LHC

- •Data volume and reconstruction will also be a problem for the HL-LHC
 - the power of the multiplicity.

Standard approach: Kalman Filter used to locate hits in charged particle trajectory

ML-based approach: Use ML tools to speed this up such as...

- Recurrent Neural Network [arXiv:2212.02348]
- Convolutional neural network [See Here]

Reconstructing charged particle trajectory is computationally expensive - increases with

SIGNAL/BACKGROUND DISCRIMINATION

Traditional Techniques

- production.
- Trained in a supervised manner with <u>EvtGen</u>
- 50% increase in signal significance with ML!

With **BDT**

[PRL 124, 172301 (2020)]

- Boosted Decision Tree implemented in ROOT TMVA to optimize signal for Λ_c baryon

The performance worsens for Pb-Pb, due to the large UE.

Quark and gluon discrimination is a difficult and ongoing effort in HIs! Future: Apply these methods to data in pp and Pb—Pb!

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

Kannah Bossi (<u>hannah.bossi@cern.ch</u>)

ML Analysis Techniques

Massachusetts Institute of Technology

GENERATIVE ADVERSARIAL NETWORKS (GANS)

Two networks compete with one another in a game.

The generative network seeks to fool the discriminative network.

The discriminative network seeks to find the real sample from the generated samples.

Indirect training \rightarrow generative network never sees the true distribution!

INTRO TO RANDOM FOREST

Random forests are composed of decision trees. Decision trees are a set of rules organized in a tree structure.

Each node is a rule which subdivides the dataset into two or more parts (think 20 questions).

Output of the random forest is a combination of the output of each of the decision trees.

In training, the algorithm sets up the rules of each decision tree.

NEURAL NETWORKS

Flow of information happens between nodes.

A weight is associated with each input to a given node.

The output of each node is a function of the weighted inputs. The output of a node j, is generally written something like

$$O_j = \sum_{i=0}^{N-1} w_{ij}O_i$$

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

In training we seek to learn the set of weights which minimize the total error of the network.

CONVOLUTIONAL NEURAL NETWORKS (CNNS)

Input Layer

Convolution Layer

Key component of a CNN is the **convolution layer**, which (with the help of a filter) will determine if a feature/pattern is present.

🐹 Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

AUTO-ENCODERS

Simple task: NN architecture trained to copy inputs to outputs!

Encoder takes the input and dramatically reduces its complexity via a NN.

Decoder takes the encoded data and reconstructs outputs like the data.

Does not require labeled data as input!

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

3 anomaly!

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

Anomaly detections: If you fail to reconstruct data in the decoding step you have an

DIFFERENT ALGORITHMS FOR DIFFERENT PROBLEMS!

(Shallow or Deep) Neural Networks \rightarrow *Great for* making predictions!

Convolutional Neural Networks (CNNs)→ *Great for* image processing!

Random Forest (Decision Trees)

Hannah Bossi (<u>hannah.bossi@cern.ch</u>)

ML Analysis Techniques

Generative Network "Real" Sample **Generated Sample** Update Network **Discriminative Network** Update Network Binary Classification: Is the sample real or fake? **Generative Adversarial** Networks (GANs) \rightarrow

Powerful tool for generating samples!

