No classification without representation

Gregor Kasieczka Email: <u>gregor.kasieczka@uni-hamburg.de</u> Twitter: <u>@GregorKasieczka</u> Mastodon: <u>@gregor_k@sciencemastodon.com</u> Al and Quantum in Fundamental Physics

CLUSTER OF EXCELLENCE

QUANTUM UNIVERSE





CDCS





*

Bundesministerium für Bildung und Forschung

CENTER FOR DATA AND COMPUTING

Partnership of Universität Hamburg and DESY

Introductions



State of the field

- Large data volumes and complex structures: Long history of advanced analysis techniques in particle physics
- Machine learning (ML) quickly becoming omnipresent in all aspects



THE USE OF NEURAL NETWORKS IN HIGH ENERGY PHYSICS*

BRUCE DENBY Fermi National Accelerator Laboratory M.S. 318 Batavia, Illinois 60510 U.S.A. denby@fnal.bitnet

ABSTRACT

In the past few years a wide variety of applications of neural networks to pattern recognition in experimental high energy physics has appeared. The neural network solutions are in general of high quality, and, in a number of cases, are superior to those obtained using 'traditional' methods. But neural networks are of particular interest in high energy physics for another reason as well: much of the pattern recognition must be performed online, i.e., in a few microseconds or less. The inherent parallelism of neural network algorithms, and the ability to implement them as very fast hardware devices, may make them an ideal technology for this application.

3

Micro-Intro: Particle Physics

- Particle physics: study smallest constituents of matter
- Standard Model: incredible scientific achievement, describes 3/4 fundamental forces
- Mathematical, quantum theoretical understanding of matter at the smallest scales





- Experimental evidence (e.g. dark matter) & theoretical considerations: Standard Model is not sufficient, need *new physics*
- Comprehensive program at Large Hadron Collider (LHC)
- Experimental data is complemented by large volumes of high quality simulations (synthetic data)

Micro-Intro: Data

- Particle collisions with ~1 MB/ event happen at a rate of 40 MHz
- Distill to ~1 kHz via lossy, irreversible filtering algorithms (Trigger)
- Very heterogenous data: lowlevel readouts in ~100M channels; can condense to O(10) high-level features
- One collision/event = "one image"

sample i.i.d. from underlying physics distribution (e.g. the Standard Model + potential new physics)



Experimental particle physics



Micro-Intro: Machine Learning

- Rephrase task as a minimisation problem..
- ..and "simply" solve:

$$\theta^* = \operatorname{argmin}_{\theta} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \left[\mathcal{L}(f_{\theta}(\mathbf{x}), \mathbf{x}) \right]^{\theta}$$

- Modern ML: function f is a deep neural network & minimisation carried out via gradient descent
- Devil in the details:
 - How to map physics objective to loss function L
 - How to structure f to make maximum use of physics knowledge
 - How learn in a robust way from minimum amount of data



hidden layer

hidden layer 2

hidden layer 3

- Detailed domain knowledge mathematical structures, useful observables, symmetries — in particle physics
- How to include in neural networks to improve performance and/or data and resource efficiency?
- Will cover three different areas & approaches

- Detailed domain knowledge mathematical structures, useful observables, symmetries — in particle physics
- How to include in neural networks to improve performance and/or data and resource efficiency?
- Will cover three different areas & approaches



Select physically relevant features for supervised classification

- Detailed domain knowledge mathematical structures, useful observables, symmetries — in particle physics
- How to include in neural networks to improve performance and/or data and resource efficiency?
- Will cover three different areas & approaches



Select physically relevant features for supervised classification



Add permutation symmetry to generative networks

- Detailed domain knowledge mathematical structures, useful observables, symmetries — in particle physics
- How to include in neural networks to improve performance and/or data and resource efficiency?
- Will cover three different areas & approaches



Select physically relevant features for supervised classification



Add permutation symmetry to generative networks



See how these issues affect the current frontier of anomaly detection

Particle taggers



Concrete Task

- Distinguish jets initiated by a top quarks from jets from other particles
 - Binary classification task
- Use simulation as synthetic training data: perfect class labels available
 - (Leads to domain shift when applied to collider data)





- 1.2M training examples (*jets*),
 400k each for testing and validation
- Each example: Up to 200 particles with 3 features/particle
 (2D position on detector surface+ energy)
- Metrics: AUC: area under curve and R₃₀:1/FPR @ TPR=0.3

GK, Plehn, et al 1902.09914



- Many ways to encode symmetries in network architecture to improve tagging performance: Dedicated talk later today
- Instead, a way to automatically find best features to describe a jet



Motivation

- Advantage of few high-level features:

 -easy to understand and calibrate
 -cheap to evaluate
- Advantage of complex architecture and low-level features: performance
- Can we combine both?



We need a basis

- Energy Flow Polynomials (EFPs) form a basis of jet substructure
 - Nodes: energy fractions
 - Edges: angular distances
- Depending on order considered, too many (e.g 7k) to efficiently train NN (many features work if there is structure, not so much for EFPs)

$$\bullet_{j} \iff \sum_{i_{j}=1}^{M} z_{i_{j}}, \qquad k \longrightarrow \ell \iff \theta_{i_{k}i_{\ell}}$$
e.g.
$$\bullet = \sum_{i_{1}=1}^{M} \sum_{i_{2}=1}^{M} \sum_{i_{3}=1}^{M} \sum_{i_{4}=1}^{M} z_{i_{1}}z_{i_{2}}z_{i_{3}}z_{i_{4}}\theta_{i_{1}i_{2}}\theta_{i_{2}i_{3}}\theta_{i_{2}i_{4}}^{2}\theta_{i_{3}i_{4}}.$$

Look for optimal feature set

Solution: Iterative feature selection



Das, GK, Shih 2212.00046; Faucett, Thaler, Whiteson, 2010.11998

Algorithm



Algorithm



Example

Example

Algorithm

Das, GK, Shih 2212.00046;

 DiscoFFS find relevant features quicker than alternative feature selection methods

- DiscoFFS find relevant features quicker than alternative feature selection methods
- New best trade-off between performance and complexity

- DiscoFFS find relevant features quicker than alternative feature selection methods
- New best trade-off between performance and complexity
- Features stable under re-training

 $\mathrm{EFP}_{G}^{(\kappa,\beta)} = \sum_{i_{1}\in J} \cdots \sum_{i_{N}\in J} z_{i_{1}}^{(\kappa)} \cdots z_{i_{N}}^{(\kappa)} \prod_{(m,\ell)\in G} \theta_{i_{m}i_{\ell}}^{(\beta)}.$; chromatic number c

Fast simulation

Simulation is crucial to connect experimental data with theory predictions

Simulation is crucial to connect experimental data with theory predictions but computationally very expensive

Use generative models trained on initial data to augment statistics

Overview of generative architectures

Simulation is crucial to connect experimental data with theory predictions but computationally very expensive

Use generative models trained on initial data to augment statistics

Want to achieve particle showers in calorimeters

CALICE AHCal testbeam

ILD Detector

Illustration of particle shower in a sampling calorimeter

One data example

Simulation is crucial to connect experimental data with theory predictions but computationally very expensive

Use generative models trained on initial data to augment statistics

Want to achieve particle showers in calorimeters

But start with a simpler task (Claudius will talk about Calorimeters) Kansal et al 2106.11535;

JetNet dateset: Pythia final-state particles from a jet. Either up to 30 or 150 particles/jet

Point cloud generation

- In both cases (calorimeter showers, JetNet):
 - Each jet/shower is a point cloud with 30 to O(100k) points (particles)
 - Each point has a 2D/3D position in space and other associated quantities
- This poses a limit to using fixed structures (e.g. convolutions)
 - How to simulate point clouds directly?
 - Specifically: how to add permutation symmetry?

Deep Sets

Theorem 2 A function f(X) operating on a set X having elements from a countable universe, is a valid set function, i.e., **invariant** to the permutation of instances in X, iff it can be decomposed in the form $\rho\left(\sum_{x \in X} \phi(x)\right)$, for suitable transformations ϕ and ρ .

EPiC

Equivariant Point Cloud interaction (EPiC) block: Similar to deep sets, but with additional global information exchange.

Still permutation equivariant

Buhmann, GK, Thaler 2301.08128;

EPiC GAN

(b) Discriminator

Buhmann, GK, Thaler 2301.08128;

Can use to build

for classical GAN

architecture, but

fully respecting

discriminator blocks

permutation symmetry.

generator and

Buhmann, GK, Thaler 2301.08128;

Where to go from here?

- Extend to larger point clouds
- Increase fidelity
- GANs work, what about other architectures (VAEs, flows, diffusion)?

Finding new physics

Motivation

- Theoretical and experimental reasons to expect new physics beyond the Standard Model
- However, so far only negative results in direct (model driven) searches
- Two discovery strategies:
 - Model-specific
 - Model independent
- Machine learning plays a key role in both — focus on anomaly detection now

Selection of observed limits at 95% C.L. (theory uncertainties are not included). Probe up to the quoted mass limit for light LSPs unless stated otherwise. The quantities ΔM and x represent the absolute mass difference between the primary sparticle and the LSP, and the difference between the intermediate sparticle and the LSP relative to ΔM , respectively, unless indicated otherwise.

Strategies

- Orthogonal strategy to model specific searches:
 - Discover new physics with making minimal assumptions
- Less sensitive to one specific model, broader coverage

ML-assisted global comparison

- Systematically compare simulation to recorded data, look for differences
- Con: Rely on imperfect simulation, maximally background model dependent
- Pro: Sensitive to all types of anomalies

Resonant anomaly detection / Enhanced bump hunts

- Estimate background in-situ from data
- Con: Need to make assumptions about signal shape
- Pro: Data-driven on background model

Need to find a feature in which signal is resonant and background smooth.

No assumptions in other features.

Further generalisation as open issue.

Autoencoders

Autoencoders: Learn-compression/ decompression on signal free sample and use loss as anomaly score

Heimel, GK, et al 1808.08979; Farina et al 1808.08992; ...

Setup

Welcome to the home of the LHC
Olympics 2020!Image: Construction of the LHC
Olympics 2020!Image: Construction

reduces generality — but allows fully data-driven construction of anomaly detection score

Autoencoder Challenges

• Complexity bias

If anomalies are much simpler than backgrounds: L will still be lower, despite never encountered in training. Overcome e.g. with normalising AE

• Change of variables

Autoencoder approximates background density — not stable under change of variables

GK, ...Shih, et al 2209.06225; Le Lan, Dinh 2012.03808; Weber MSc thesis; Finke et al 2104.09051; Dillon et al 2206.14225;

Weak supervision

Metodiev, Nachman, Thaler, 1708.02949; Howe, Nachman 1805.02664

LaCATHODE

- If R(x) is only calculated in signal region, it's extrapolation is not well-defined
- Potential problem for bumphunt if it shapes distributions

LaCATHODE

- If R(x) is only calculated in ^B
 signal region, it's extrapolation is not well-defined
- Potential problem for bumphunt if it shapes distributions
- Instead, train classifier in latent space

LaCATHODE

- If R(x) is only calculated in signal region, it's extrapolation is not well-defined
- Potential problem for bumphunt if it shapes distributions
- Instead, train classifier in latent space to achieve flat distributions

 ΔR dataset (bkg-only training), selecting 1%

Comments on anomaly detection

• As CATHODE approximates a likelihood ratio, it should be robust compared to methods that only use p_{Background} (e.g. autoencoders)

Comments on anomaly detection

- As CATHODE approximates a likelihood ratio, it should be robust compared to methods that only use p_{Background} (e.g. autoencoders)
- However, still can be sensitive to choice of input features
 - Here shown: idealised anomaly detector (perfect DE)

Comments on anomaly detection

- As CATHODE approximates a likelihood ratio, it should be robust compared to methods that only use p_{Background} (e.g. autoencoders)
- However, still can be sensitive to choice of input features
- Need also consider
 - Higher dimensional input data
 - Stable training
 - Holistic 'end-to-end' setups

Closing & future outlook

Conclusions

- Large potential and wide range of applications of machine learning in particle physics
- Key issue of including physics knowledge in ML training:
 - Either by starting from physics observables and constructing efficient selections
 - or by building architectures that are invariant/equivariant under relevant symmetries
- Recent progress in anomaly detection but scaling stable techniques to more complex data representations still unsolved

Thank you!

Backup

Distance Correlation

$$dCov^{2}(\vec{X}, \vec{Y}) := E[\|\vec{X} - \vec{X'}\| \|\vec{Y} - \vec{Y'}\|] + E[\|\vec{X} - \vec{X'}\|] E[\|\vec{Y} - \vec{Y'}\|] - 2E[\|\vec{X} - \vec{X'}\| \|\vec{Y} - \vec{Y''}\|].$$

$$\mathrm{dCor}^2(\vec{X}, \vec{Y}) = \frac{\mathrm{dCov}^2(\vec{X}, \vec{Y})}{\sqrt{\mathrm{dCov}^2(\vec{X}, \vec{X}) \ \mathrm{dCov}^2(\vec{Y}, \vec{Y})}},$$

$$\overline{\mathrm{dCor}}^2(\vec{X}, \vec{Y}) = \mathrm{dCor}^2(\Sigma_X^{-1/2}\vec{X}, \Sigma_Y^{-1/2}\vec{Y}).$$

Das, GK, Shih 2212.00046;