



Machine learning introduction

Rui Zhang (张 瑞) <u>rui.zhang@cern.ch</u> University of Wisconsin-Madison, Wisconsin 03 Aug, 2024

Books and references



- 台湾大学李宏毅视频课程(<u>https://github.com/Fafa-DL/</u> Lhy_Machine_Learning)
- 深度学习论文精读(<u>https://github.com/mli/paper-reading</u>)

A typical LHC Physics Analysis Workflow



ML is an old friend of HEP



ML is an old friend of HEP



R. Zhang

ML is an old friend of HEP





R. Zhang

03.08.2024

ML is not a magic

It's built upon linear algebra and information theory



Neural network is a function that maps input to output; "universal approximation theorem"

Learning procedure is to compress the input to output.

$$y_1 = f_1(x)$$

$$y_2 = f_2(x)$$

$$\vdots$$

$$y_n = f_n(x)$$

Which function is close to truth?

Need to quantify "similarity" between y_i and

- Also known as "loss" => min(Loss(y_i , y_{truth}))

Information theory offers measures for quantifying similarity

- Entropy: disorder of 1 PDF
- Divergence: disorder between 2 PDFs

ML is not a magic: capacity and regularisation

$$y = f_W(x)$$



- f_W has too few free parameters → underfit
- f_W has too many free parameters → overfit
- Avoid underfit/overfit? → check on the test sample (generalisation error)



- Model's effective capacity can be affected by some factors:
 - Optimisation is very difficult and may not find the global optimum
 - Adding regularisation can limit the capacity
- An example regularisation: L2-norm $J(w) = \text{MeanSquaredError} + \lambda ||w||_2^2$

$$||w||_2 = \sqrt{\sum_{k=1}^n |w_k|^2}$$

Large weight are suppressed

ML is not a magic: distance between distributions

- In information theory, learning that an unlikely event has occurred is more informative than learning that a likely event has occurred.
 - $p(x_1) < p(x_2) \Rightarrow I(x_1) > I(x_2)$
 - I(x) >= 0
 - If x_1 , x_2 are independent, ie, $p(x_1) \cdot p(x_2) = p(x_1 + x_2)$, $I(x_1 + x_2) = I(x_1) + I(x_2)$.
- Shannon entropy: $H(p) = E_{x \sim p}[I(x)] = -E_{x \sim p}[\log p(x)]$
 - Shannon entropy of a distribution is the expected amount of information in an event drawn from that distribution; when p is continuous, it is differential entropy
 - Distributions that are nearly deterministic (where the outcome is nearly certain) have low entropy; distributions that are closer to uniform have high entropy.



 $I = -\log p$

ML is not a magic: distance between distributions

Divergence is a measure of statistical distance between two distributions.

Most popular one: Kullback-Leibler (KL) divergence: $D_{\text{KL}}(P||Q) = \mathbb{E}_{x \sim P} \left[\log \frac{P(x)}{Q(x)} \right]$

Information needed to send a message containing symbols drawn from P, when we use a code that was designed to minimise the length of messages drawn from probability distribution Q.

Asymmetric between p and q

 $q^* = \operatorname{argmin}_q D_{\mathrm{KL}}(p \| q)$

Wish all events in p will be found in q

q

Wish all events in q will be found in p

 $q^* = \operatorname{argmin}_{q} D_{\mathrm{KL}}(q \| p)$

Jensen-Shannon Divergence (JSD)

Wasserstein Distance (Earth Mover's Distance)

Total Variation Distance (TV Distance)

Bhattacharyya Distance

Hellinger Distance

f-Devergence

Rényi divergence

R. Zhang

ML is a Tool for HEP

• Step 1: how to represent data • Step 2: how to train the model





ML is a Tool for HEP

• Step 1: how to represent data

• Step 2: how to train the model



Step 1: how to represent data?



1.1 Input has fixed length

 Decide in advance variable list for training, then train a deep neural network / BDT



A typical signal extraction using NN



CMS tau ID deep network

1.1 Perceptron / Deep neural network

 Perceptron was invented in 1943 by Warren McCulloch and Walter Pitts, sharing the very similar concepts with the modern neural network



- Step 1: multiply input with weights and sum over
- Step 2: apply activation function on the output



reference

1.1 Why deep (MLP)? Example of XOR

XOR truth table



Multiple layers increase non-linearity of the function represented by the NN

R. Zhang

1.1 Gradient decent





细长的二次函数类似一个长峡谷。 梯度下降把时间浪费于在峡谷壁反 复下 降,因为它们是最陡峭的特征

x 是 f 某个横截面 的局部极大点,却 是另一个横截面的局部极小点

1.2 Input as sequences

- In some situations, fixed length is not suitable
 - e.g. Jets contain a variable number of particles
 - Recurrent Neural Networks shows great performance for Natural Language Processing tasks
 - Information across the entire sequence can be accumulated and used



1.2 Recurrent Neural Networks



1.2 Recurrent Neural Networks



1.2 Recurrent Neural Networks



1.2 RNN application



R. Zhang

1.3 Input as images

Jets can be viewed as images

• So as electrons



Electron classes



1.3 Convolutional neural network

 Convolutional neural network shows great performance for computer vision tasks

- Nice features: "sparse interactions": kernel's dim is much smaller than image's dim
- "parameter sharing": kernel is shared by different patches of the image
- "equivariance" if f(g(x)) = g(f(x)), then f and g are equivariance



1.3 CNN applications

E_T^{miss} reconstruction

ATL-PHYS-PUB-2019-028



Hybrid: DNN + RNN + CNN application

Particle and vertex based DNN: Deeplet



1.4 Input as sets

Sequence (and also image) implies certain ordering

- Lack of permutation invariance $f(x_1, x_2) \neq f(x_2, x_1)$
- Deepset [Manzil et al]
 - for any permutation $\pi : f(\{x_1, ..., x_M\}) = f(\{x_{\pi(1)}, ..., x_{\pi(M)}\})$



1.5 Input as graphs (including point cloud)

Graph is also a natural way to represent LHC data



1.5 Graphs neural networks



- **V** Vertex (or node) attributes e.g., node identity, number of neighbors
- E Edge (or link) attributes and directions e.g., edge identity, edge weight
- **U** Global (or master node) attributes e.g., number of nodes, longest path



Architecture schematic for Message Passing layer. The first step "prepares" a message composed of information from an edge and it's connected nodes and then "passes" the message to the node.

A single layer of a simple GNN. A graph is the input, and each component (V,E,U) gets updated by a MLP to produce a new graph. Each function subscript indicates a separate function for a different graph attribute at the n-th layer of a GNN model.

A Gentle Introduction to Graph Neural Networks

R. Zhang

1.5 Self-attention and transformer

First proposed for natural language processing (NLP)
Later find good performance in computer visual

R. Zhang

Machine learning introductory lecture

Output

1.5 Self-attention and transformer

R. Zhang

1.5 Transformer vs GNN

Transformers are Graph Neural Networks

1.5 Transformer vs RNN

1.5 Transformer vs CNN

1.5 GNN applications

R. Zhang

Machine learning introductory lecture

Step 2: set up the learning task

Step 2: set up the learning task

Weakly-supervised = noisy labels

Unsupervised

Semi-supervised = partial labels

R. Zhang

Unsupervised—fast simulation (than Geant4)

R. Zhang

Machine learning introductory lecture

Fast simulation: Generative models

A generator is a function that maps random numbers to structure.

R. Zhang

Machine learning introductory lecture

Example—Integrated into real experiment

ATLAS fast simulation includes a GAN at intermediate energies for hadrons

COMPUT SOFTW BIG SCI 6, 7 (2022)

R. Zhang

Machine learning introductory lecture

Example—Integrated into real experiment

FastCaloGAN has been expanded from Run 2 to Run 3

	Inner Detector	Calorimeters				Muon Spectrometer
Electrons Photons	Geant	FastCaloGAN V2 <i>E_{kin}</i> < 8 GeV && η < 2.4, Except [0.9< η <1.1, 1.35< η <1.5]		FastCaloSim V2 <i>E_{kin}</i> > 16 GeV && η < 2.4, All <i>E_{kin}</i> && [0.9< η <1.1, 1.35< η <1.5, η >2.4]		
Charged Pions Kaons		Geant4 Pions: E _{kin} < 200 MeV Other hadrons: E _{kin} < 400 MeV	FastCaloSim V2 <i>E_{kin}</i> < 4 GeV && η < 1.4, <i>E_{kin}</i> < 1 GeV && η < 3.15		FastCaloGAN V2 $E_{kin} > 8 \text{ GeV } \& \eta < 1.4,$ $E_{kin} > 2 \text{ GeV } \& 1.4 < \eta < 3.15,$ All $E_{kin} \& \eta > 3.15$	Muon
Baryons			FastCaloGAN V2		+ Geant4	
Muons					Geant4 <u>arXiv</u>	: <u>2404.06335</u>

R. Zhang

Marrying generative techniques (in R&D)

R. Zhang

Fast Calorimeter Simulation Challenge 2022

View on GitHub

Welcome to the home of the first-ever Fast Calorimeter Simulation Challenge!

The purpose of this challenge is to spur the development and benchmarking of fast and high-fidelity calorimeter shower generation using deep learning methods. Currently, generating calorimeter showers of interacting particles (electrons, photons, pions, ...) using GEANT4 is a major computational bottleneck at the LHC, and it is forecast to overwhelm the computing budget of the LHC experiments in the near future. Therefore there is an urgent need to develop GEANT4 emulators that are both fast (computationally lightweight) and accurate. The LHC collaborations have been developing fast simulation methods for some time, and the hope of this challenge is to directly compare new deep learning approaches on common benchmarks. It is expected that participants will make use of cutting-edge techniques in generative modeling with deep learning, e.g. GANs, VAEs and normalizing flows.

This challenge is modeled after two previous, highly successful data challenges in HEP – the top tagging community challenge and the LHC Olympics 2020 anomaly detection challenge.

Datasets

The challenge offers three datasets, ranging in difficulty from "easy" to "medium" to "hard". The difficulty is set by the dimensionality of the calorimeter showers (the number layers and the number of voxels in each layer).

<u>Link</u>

Un-/weakly/semi supervised—anomaly detection

Why anomaly detection?

Typical Searches

- Looking for a specific,
 physics motivated signal
- Maximum sensitivity (using supervised learning e.g.
 BDT) for a specific model
- Not very useful for other signal models

Anomaly Detection

- Model agnostic/ independent search
- Looking for deviations
 from background only
- Less sensitive to any specific model, but can look for multiple different models

Two types of anomaly detection

Outlier Detection

o Searching for unique or unexpected events

 In HEP, this is the tails of distributions or uncovered phase space

Overdensity detection

Analogous to the traditional bump hunting

[1805.02664, 1806.02350, 1902.02634, 1912.12155, 2001.05001, 2001.04990, 2012.11638, 2106.10164, 2109.00546, 2202.00686, 2203.09470, 2208.05484, 2210.14924, 2212.11285, 2305.04646, 2305.15179, 2306.03933, 2307.11157, 2309.12918, 2310.06897, 2310.13057,]

Inspired by this presentation

Outlier Detection in experiments (ATLAS)

- Full event level anomaly detection
- Searched in 9 invariant masses including di-jet, di-b-jet, with three anomaly regions => demonstrating high efficiency in the search arXiv: 2307.01612

Overdensity Detection in experiments (ATLAS)

- Performed on di-fatjet resonant search
- Network is learning difference between Prob(b) and Prob(s+b)

Phys. Rev. Lett. 125 (2020) 131801

Summary

- ML and HEP: an enduring partnership
 - ML has been a longstanding companion in HEP in various stages of the data analysis pipeline
- ML as a Toolset for HEP
 - ML serves as a valuable assistant, maximising the exploration of costly collision data
 - Choosing ML architectures based on the data structures to optimise efficiency
 - Evolution towards unsupervised and semi-supervised learning on more generative tasks
- Future directions
 - Expanding training data to refine ML models
 - Delving into lower level features to uncover hidden patterns
 - Incorporating physics knowledge for a deeper contextual understanding

ML continues to unlock breakthroughs within the realm of HEP.

References

A Living Review of Machine Learning for Particle Physics

- <u>https://iml-wg.github.io/HEPML-LivingReview/</u>
- Updated summary on arXiv available submission in machine learning in HEP
- Neural Networks, Types, and Functional Programming
 - <u>http://colah.github.io/posts/2015-09-NN-Types-FP/</u>
 - Deep learning introduction in 10 min
- (New) Machine learning chapter in the particle data group book:
 - <u>https://pdg.lbl.gov/2023/reviews/rpp2022-rev-machine-learning.pdf</u>

GN1 vs GN2

Updated Attention Mechanism

GN2 follows more closely the *transformer* architecture [1706.03762]

R. Zhang

Generative model: GAN

Generative model: VAE

A summary blog

Variational Autoencoders: A pair of networks where one embed the data into a latent space with a given prior and the other decode back to the data space.

Generative model: NF

A summary blog Normalising flow: invertible transformations to map a simple distribution to a complex one. Exact likelihood computation Slow sampling High generative capacity Inverse Flow \mathbf{x}' Flow-based models: \mathbf{X} \mathbf{z} $f^{-1}(\mathbf{z})$ $f(\mathbf{x})$ Invertible transform of distributions Diefenbacher et al. 2023 JINST 18 P10017 Pang et al, arXiv:2308.11700 Krause et al, PhysRevD.107.113003 Buckley et al, arXiv:2305.11934

Generative model: Diffusion model

A summary blog

R. Zhang

Non-uniform ATLAS geometry

Decorrelation

Caution Part I

35

How can we learn a classifier that does not sculpt a bump in the background?

Enforcing Independence

Train e.g. a neural network with a **custom loss functional** $\mathcal{L}[f(x)] = \sum_{i \in s} L_{\text{classifier}}(f(x_i), 1) + \sum_{i \in b} L_{\text{classifier}}(f(x_i), 0) + \lambda \sum_{i \in b} L_{\text{decor}}(f(x_i), m_i)$

Recent proposals:

Adversaries: L_{decor} is the loss of **a 2nd NN** (adversary) that tries to learn *m* from f(x).

Distance Correlation: L_{decor} is **distance correlation** (generalizes Pearson correlation) between *m* and *f*(*x*).

Mode Decorrelation: L_{decor} is small when the **CDF** of f(x) is the same across different values of *m*.

Nachman, Overview of Machine Jet Image Learning for Particle Physics

R. Zhang