

MACHINE AND DEEP LEARNING APPLICATIONS IN <u>HIGH ENERGY PHYSICS</u> ATLAS AND CMS

Terzo incontro di fisica con ioni pesanti alle alte energie (IFIPAE) 2021

Federica Legger

ONCE UPON A TIME

- TMVA Toolkit for Multivariate Data Analysis, arXiv:physics/0703039
- TMVA fully integrated in ROOT in <u>2013</u>
- Mainly for classification and regression tasks
 - boosted decision trees, support vector machines, cellular automata, multilayer perceptrons, ...

Multivariate searches for single top quark production with the DØ detector, <u>Phys.Rev.D75:092007</u> (2007)



2

PARTICLE IDENTIFICATION: TAU ID

 Hadronically decaying tau leptons vs jet of hadrons with Boosted Decision Trees (BDT)



RARE PROCESSES: DIRECT STAU PRODUCTION

• **BDT** with low and high level variables



Search for the electroweak production of supersymmetric particles in $\sqrt{s=8}$ TeV pp collisions with the ATLAS detector, <u>Phys. Rev. D 93, 052002</u> (2016)

GOING DOWN THE DEEP ROAD - B TAGGING AT CMS

DeepCSV:

- fully-connected layers
- Multi classification



DeepJet:

- Convolutional layers learn compact feature representation (automatic feature engineering)
- RNN extract information from each set of features



HIGGS (DOUBLE B) TAGGING



DEEP LEARNING FOR SIMULATION

- Computing demands increase nonlinearly with increasing pileup
- LHC Run 2: full detector simulation (Geant4) took ~40% of grid CPU resources for CMS & ATLAS
- Calorimeter simulation most CPU intensive





ATLAS FASTCALOGAN



G	50 (Input latent Space), 50, 100, 200, NVoxel (pid and η dependent)
D	NVoxel, NVoxel, NVoxel, 1
Activation function	ReLU (in all layers)
Optimiser	Adam [21]
Learning Rate	10^{-4}
β	0.5 Fast simulation of the ATLAS calorimeter
Batchsize	¹²⁸ system with Generative Adversarial Networks
Training ratio (D/G)	5 ATL COET DUB 2020 00C (2020)
Gradient penalty λ	$\frac{A1L-SUF1-PUB-2020-006}{2020}$



NEW TRENDS IN RECO Ideal applications for graph neural networks: Hit clouds in Calorimeters: point cloud of Ο energy deposits Filter likely. Filter, convert to adiacent triplets Tracking Ο Raw hit data doublets embedded Jet tagging Ο Train/classify Train/classify doublets in GNN triplets in GNN End-to-end reconstruction of multiple Apply cut for seeds particles simultaneously È Gluon-jet image \rightarrow 3 channels: track

DBSCAN for

track labels

THE FUTURE

- Run Anomaly detection in the trigger
 - Variational autoencoders for new physics mining at the Large Hadron Collider, <u>J. High Energ. Phys. 2019, 36</u> (2019)
- Improve unfolding with invertible networks: detector

 high level
 variables
 - Invertible networks or partons to detector and back again, <u>SciPost Phys. 9</u>, <u>074</u> (2020)
- Use attention to mitigate combinatorics in ttbar events: Network output should be invariant under permutations of the input jet order
 - SPANet: Generalized Permutationless Set Assignment for Particle Physics using Symmetry Preserving Attention, <u>arXiv:2106.03898</u> (2021)

DEEP LEARNING ON FPGAS: HLS4ML

- Tool to deploy NNs to FPGA
 - reads as input models trained on standard DL libraries
 - implements common ingredients (layers, activation functions, etc)
- Uses HLS softwares to provide a firmware implementation of a given network
 - Pruning
 - Quantization

Distance-Weighted Graph Neural Networks on FPGAs for Real-Time Particle Reconstruction in High Energy Physics, <u>arXiv:2008.03601</u>



OPENFORBC



- **Open For Better Computing**
 - Project funded by 2021 INFN Ο Research4Innovation (R4I) call
 - Promote use of GPUs for scientific Ο applications
- Effortless GPU partitioning for hardware from different vendors in Linux KVM

Hardware choice:

Nvidia V100 32GB GPU



DETECTOR DEVELOPMENT: RSD

- RSDs (Resistive AC-Coupled Silicon Detectors): silicon sensors based on LGAD (Low-Gain Avalanche Diode)
 - Signal is seen over several pixels
- Multi-Output regression (MR) and Multi-layer Perceptron (MLP) models using various amplitudes as input to predict hit position

Spatial resolution for sensors with different geometries as a function of the interpad



First application of machine learning algorithms to the position reconstruction in Resistive Silicon Detectors, <u>JINST 16 P03019 (2021)</u>

SMART INFRASTRUCTURE

- Inspired from S.M.A.R.T. (Self-Monitoring, Analysis and Reporting Technology) for hard drives
- <u>Predictive vs reactive</u> maintenance for complex infrastructure
- Can be detectors, computing centers, factories, IOT





- Targets **reduction of operational costs** of distributed computing infrastructure through smart automation
 - Use case: Worldwide LHC Computing Grid (WLCG)
 - Metrics: reduction of number of tickets, number of operators, time to solve, user satisfaction
- Exploits anomaly detection in time series, natural language processing (NLP) and clusterization techniques
- Bonus: Increase resource utilisation efficiency => increase uptime, less resources wasted => "Green" development

https://operational-intelligence.web.cern.ch/

Operational Intelligence for Distributed Computing Systems for Exascale Science, <u>EPJ Web Conf., 245 (2020) 03017</u>

ANOMALY DETECTION IN THE CERN CLOUD

- Aims to identify problematic nodes in the CERN cloud
- Metrics are encoded as images or vectors according to model



Anomaly detection in the CERN cloud infrastructure, EPJ Web Conf 251, 02011 (2021)

SUMMARY

- HEP is using MVA methods (aka ML) since > 20 years
- DL entered the scene with jet tagging
 - Now successfully used in analysis (S/B, jet to parton assignment,...)
 - Anomaly detection, attention, GANs, ... for Run 3 and beyond
- At the LHC we are resource-limited at L1, HLT and offline
 - DL may be a way to save resources and extend physics reach
 - **Sparse data**: traditional NN (CNN, RNN) may work but at a cost
 - Custom edge computing: inference needs to run everywhere (FPGA, custom chips, grid)
 - Real time: inference within 1 μs (trigger boundary)
- Many more applications: detector, computing, ...

OUTLOOK

- Physics applications of ML and DL in a wide range of domains, and growing
- Many challenges ahead:
 - Keep the pace with AI research
 - Nowadays mainly driven by industry, science should not stand behind!
 - Foster the use of common tools/technologies
 - Exploit heterogeneous hardware
 - Deploy to production

A PROPOSAL FOR THE DARTMOUTH SUMMER RESEARCH PROJECT ON ARTIFICIAL INTELLIGENCE

> J. McCarthy, Dartmouth College M. L. Minsky, Harvard University N. Rochester, I. B. M. Corporation C. E. Shannon, Bell Telephone Laboratories

BACKUP



• Recurrent Neural Network (RNN)



- Long short term memory(LSTM)
- Gated Recurrent Unit (GRU)



Convolutional Neural Networks (CNN)



- **Convolutional** layer: two functions produce a third that describes how the shape of one is changed by the other
- **pooling** layer: reduce dimensionality



TRANSFORMERS (2017)

- All you need is attention
- Self-attention: query, key, value:
 - the output is a weighted sum of the values, where the weight assigned to each value is determined by the dot-product of the query with all the keys:

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

arXiv:1706.03762 [cs.CL]





ATTENTION



GENERATIVE ADVERSARIAL NETWORK - GAN (2014)



- Double network: generator net and discriminator net
 - generator produces samples close to training samples
 - discriminator differentiates samples from generator and training set
 - training until discriminator can no longer distinguish



GRAPH NEURAL NETWORK (GNN)

Node predictionEdge prediction

Graph prediction

Node **Message Passing** State update Message Message Vote Embedding Passing & Passing & State State update update $s_{1}^{(1)}$ $s_{1}^{(t)}$ y_1



GNN FOR TRACKING

 Graphs can capture the sparsity, manifold, relational structures of physics data
Doublets:



ARTIFICIAL GENERAL INTELLIGENCE

• Common sense:

- Current systems may be easily fooled by just slight changes in the input data (for example image taken from another viewpoint)
- Embed coordinate systems, whole-part relationship (capsules)

• Abstract concepts:

- Current models may be able to distinguish between a jet and a tau, but do not know what a particle is
- Creativity:
 - Current models highly specialised and engineered to solve specific problems

[Murray Shanahan, Geoff Hinton]

- Supervised learning needs many labeled data
- Reinforced learning:
 - Not practical to train in real world (when no simulation is available)
 - takes longer than an average human for a machine to learn a new task
- **Self supervised learning:** Predict everything from everything else learn representations, rather than learning specific tasks
 - Very large networks trained with large amount of data
 - Fill_ing the bl_anks Word2Vec, Transformer architecture for NLP
 - Not (yet) so successful for continuous problems (image, video)

CONSCIOUSNESS PRIOR

- Current deep learning:
 - **System 1:** fast, unconscious task solving
- Future deep learning:
 - **System 2:** slow, conscious task solving like reasoning, planning
- How?
 - Learn by predicting in <u>abstract space</u>
 - Learn representations (low dimensional vector), derived using <u>attention</u> from a high dimensional vector
 - The prior: the factor graph (joint distribution between a set of variables) is <u>sparse</u> $\Pi_{f_i}(S_i)$

Yoshua Bengio, arXiv:1709.08568 [cs.LG]

$$P(S) = \frac{\prod_j f_j(S_j)}{Z}$$

INFORMATION BOTTLENECK

- Hidden layers represent a Markov chain of topologically distinct representations
 - Information about the inputs decreases along the hidden layers
 - $\circ \quad I(X, h_1) > ... > I(X, h_i) > I(X, h_i+1)$
- In the first epochs, the network is trained to fully represent the input data; then, it learns to forget the irrelevant details by compressing the representation of the input

arXiv:1503.02406 [cs.LG]



A INITIAL STATE: Neurons in Layer 1 encode everything about the input data, including all information about its label. Neurons in the highest layers are in a nearly random state bearing little to no relationship to the data or its label.

B FITTING PHASE: As deep learning begins, neurons in higher layers gain information about the input and get better at fitting labels to it.

C PHASE CHANGE: The layers suddenly shift gears and start to "forget" information about the input.

D COMPRESSION PHASE: Higher layers compress their representation of the input data, keeping what is most relevant to the output label. They get better at predicting the label.

E FINAL STATE: The last layer achieves an optimal balance of accuracy and compression, retaining only what is needed to predict the label.