HISS Summer School

Intro and Overview

Prof. Gregor Kasieczka gregor.kasieczka@uni-hamburg.de

Different fields, unified by digitisation challenges:

- Immense data volumes with high complexity
- Weak signals overshadowed by much larger background processes
- Large data-throughput requiring custom solutions



The planned Square Kilometre Array (SKA) telescope will produce one exabyte of raw data/day (~10 PB after compression)

Different fields, unified by digitisation challenges:

- Immense data volumes with high complexity
- Weak signals overshadowed by much larger background processes
- Large data-throughput requiring custom solutions



At the High Luminosity LHC, overlapping collisions mean that thousands of trajectories need to be reconstructed in parallel

Different fields, unified by digitisation challenges:

- Immense data volumes with high complexity
- Weak signals overshadowed by much larger background processes
- Large data-throughput requiring custom solutions



Rare interesting processes only occur in 1 of 10 billion events in particle physics experiments

Different fields, unified by digitisation challenges:

- Immense data volumes with high complexity
- Weak signals overshadowed by much larger background processes
- Large data-throughput requiring custom solutions



Trigger systems filter 40 million particle collisions/second at experiments at the Large Hadron Collider (LHC) at CERN

State of AI in Physics



See also https://iml-wg.github.io/HEPML-LivingReview/

Plan

- Lecture 1:
 - Basics of Neural Networks
 - Classification example: Top Tagging
- Lecture 2:
 - Generative Models
 - Anomaly Detection

Resources

Deep Learning

An MIT Press book

Ian Goodfellow and Yoshua Bengio and Aaron Courville

Free online: https://www.deeplearningbook.org/



https://www.worldscientific.com/ worldscibooks/ 10.1142/12294#t=aboutBook (available in UHH library)

Lecture Slides



Prof. Gregor Kasieczka gregor.kasieczka@uni-hamburg.de

Learning like a machine

Terminology

Artificial Intelligence (AI) General term

Machine Learning (ML)

Also includes e.g. boosted decision trees (BDTs), shallow neural networks, ...

Deep Learning (DL)

Neural networks with many layers, unprocessed inputs

Terminology

Artificial Intelligence (AI) General term

Machine Learning (ML)

Also includes e.g. boosted decision trees (BDTs), shallow neural networks, ...

Deep Learning (DL)

Neural networks with many layers, unprocessed inputs

Basic idea

- Classical approach:
 - Write a sequence of instructions to solve a specific task
 - e.g.:
 - Tracking algorithm
 - Jet clustering
 - Calculation of physical observables

```
if (trackerTopology.pxfSide(detId) == 1) {
22
           trackingLayer.__side = TrackingLayer::Side::NEG_ENDCAP;
23
        } else if (trackerTopology.pxfSide(detId) == 2) {
24
           trackingLayer._side = TrackingLayer::Side::POS_ENDCAP;
25
        } else {
           throw cms::Exception("FastSimulation/Tracking")
27
               << "Tracker hit for seeding in FPix seems neither on positive nor on negative disk side: "</p>
28
               << trackerTopology.print(detId).c_str();
30
        trackingLayer._layerNumber = trackerTopology.pxfDisk(detId);
      }
       //TIB
       else if (subdet == StripSubdetector::TIB) {
34
        trackingLayer._subDet = TrackingLayer::Det::TIB;
        trackingLayer._side = TrackingLayer::Side::BARREL;
         trackingLayer._layerNumber = trackerTopology.tibLayer(detId);
      }
38
      //TID
       else if (subdet == StripSubdetector::TID) {
40
        trackingLayer._subDet = TrackingLayer::Det::TID;
41
        if (trackerTopology.tidSide(detId) == 1) {
42
           trackingLayer._side = TrackingLayer::Side::NEG_ENDCAP;
43
        } else if (trackerTopology.tidSide(detId) == 2) {
44
           trackingLayer._side = TrackingLayer::Side::POS_ENDCAP;
45
        } else {
46
           throw cms::Exception("FastSimulation/Tracking")
47
               << "Tracker hit for seeding in TID seems neither on positive nor on negative disk side: "
48
               << trackerTopology.print(detId).c_str();
49
50
         trackingLayer._layerNumber = trackerTopology.tidWheel(detId);
51
       3
52
       //T0B
       else if (subdet == StripSubdetector::TOB)
```

Basic idea

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

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

• Will now go step-by-step to understand the underlying ideas, focusing on neural networks.

General Strategy

Loss function \mathcal{L}

• Define an optimisation target (loss function)

Loss function: Supervised

Supervised Learning:

Attempt to infer some target (truth label): classification, regression (often also clustering/inference)

Use training data with known labels (often from Monte Carlo simulation)



observable features such as kinematics, tracks,...

(e.g. true energy)

predicted energy

Regression: Minimize mean squared error:

 $\mathcal{L} = (y - \hat{y})^2$

Classification Tasks

Distinguish a pair of classes (binary) or several (multi-class).

Images



Loss function: Cross entropy

$$\mathcal{L} = -y \log \left(\hat{y} \right) - \left(1 - y \right) \log \left(1 - \hat{y} \right)$$

Classification Tasks

Distinguish a pair of classes (binary) or several (multi-class).

Event Topologies







Loss function: Cross entropy

$$\mathcal{L} = -y \log \left(\hat{y} \right) - (1 - y) \log \left(1 - \hat{y} \right)$$

Cross Entropy

- Example: Binary classification
 - 2 Classes with labels: $y_k = 0, 1$
 - Build the NN such that we the output is a probability: $\hat{y}(x_k) \in [0,1]$
- -> Bernoulli Trial with probabilities p and q = 1 p



Jacob Bernoulli 1654-1705

• Likelihood:

$$L = \prod_{k}^{N} [y_k \cdot \hat{y}(x_k) + (1 - y_k) \cdot (1 - \hat{y}(x_k))]$$

• We want to minimize the average negative log-likelihood:

$$-\frac{1}{N}\log(L) = -\frac{1}{N}\sum_{k}^{N} [y_k \cdot \log(\hat{y}(x_k)) + (1 - y_k) \cdot \log(1 - \hat{y}(x_k))]$$

Loss function: Unsupervised

Unsupervised Learning:

No target, learn the probability distribution (directly from data)

Can use for sampling, anomaly detection, unfolding, ...



Learn to predict: $\hat{p}(\mathbf{x}) = f_{\theta}(\mathbf{x})$

*There also exists a number of other less-than-supervised approaches (weakly supervised learning, semisupervised learning, ...) Not so important for now.

True probablity density

 $p(\mathbf{x})$

Distribution learning: Maximise likelihood (minimize log-likelihood):

(either directly or with approximations)

 $\mathcal{L} = -\log\left(\hat{p}(\mathbf{x})\right)$

General Strategy

Loss function $\mathcal{L} \checkmark$ Neural network f_{θ} Parameters θ Opt. Parameters θ^* Data \mathbf{x} Data distribution $p(\mathbf{x})$

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

General Strategy

Loss function \mathcal{L} Neural network f_{θ} Parameters θ

- Define an optimisation target (loss function)
- Choose a non-linear, expressive, parametrised function (e.g. a neural network)

Neural networks

We need an expressive function (universal approximator) with tuneable parameters and useful implicit biases

$$x_1' = f(W_1 1 \cdot x_1 + W_1 2 \cdot x_2)$$



Neural Networks



Complexity



6 weights



300 weights

stage	output	ResNet-50		ResNeXt-50 (32×4d)			
conv1	112×112	7×7, 64, stride	e 2	7×7 , 64, stride 2			
conv2	56×56	3×3 max pool, st	ride 2	3×3 max pool, stride 2			
		1×1, 64		[1×1, 128]			
		3×3, 64	×3	3×3, 128, C=32	$\times 3$		
		1×1, 256		1×1,256			
conv3	28×28	1×1, 128		[1×1, 256]			
		3×3, 128	×4	3×3, 256, C=32	$\times 4$		
		1×1, 512		1×1,512			
conv4	14×14	1×1,256]	[1×1, 512]			
		3×3, 256	×6	3×3, 512, <i>C</i> =32	×6		
		1×1, 1024		1×1,1024			
conv5	7×7	[1×1,512]	1×1, 1024			
		3×3, 512	×3	3×3, 1024, <i>C</i> =32	×3		
		1×1,2048		1×1, 2048			
	1×1	global average	pool	global average pool			
		1000-d fc, softr	nax	1000-d fc, softmax			
# params.		25.5 ×10 ⁶		25.0 ×10 ⁶			
FLOPs		4.1 ×10 ⁹		4.2 ×10 ⁹			

Deep Learning: Complex network + low level inputs

25 million weights: "

2016 state of the art for image classification

175 billion weights: 2020 GPT-3 text model (GPT-4 ~**1.8 trillion** weights)

Model Name	$n_{\rm params}$	n_{layers}	$d_{ m model}$	$n_{ m heads}$	$d_{ m head}$	Batch Size	Learning Rate
GPT-3 Small	125M	12	768	12	64	0.5M	$6.0 imes 10^{-4}$
GPT-3 Medium	350M	24	1024	16	64	0.5M	$3.0 imes 10^{-4}$
GPT-3 Large	760M	24	1536	16	96	0.5M	2.5×10^{-4}
GPT-3 XL	1.3B	24	2048	24	128	1M	2.0×10^{-4}
GPT-3 2.7B	2.7B	32	2560	32	80	1M	1.6×10^{-4}
GPT-3 6.7B	6.7B	32	4096	32	128	2M	1.2×10^{-4}
GPT-3 13B	13.0B	40	5140	40	128	2M	1.0×10^{-4}
GPT-3 175B or "GPT-3"	175.0B	96	12288	96	128	3.2M	$0.6 imes 10^{-4}$

General Strategy

Loss function $\mathcal{L} \checkmark$ Neural network $f_{\theta} \checkmark$ Parameters $\theta \checkmark$ Opt. Parameters θ^* Data \mathbf{x} Data distribution $p(\mathbf{x})$

- Define an optimisation target (loss function)
- Choose a non-linear, expressive, parametrised function (e.g. a neural network)
- Use training data to optimise parameters, then apply to new examples

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

How do networks learn?

- Backpropagation + Gradient descent
- Important: Loss function needs to be differentiable
 - (Or find a differentiable approximation)
- Pass input (x₁, x₂, ...) to networks
- From output calculate loss function
 Find gradient of loss function with respect to weights
- Use gradient to find new weights

$$\begin{split} \theta_{t+1} &= \theta_t - \eta \frac{\partial \mathcal{L}}{\partial \theta_t} = \theta_t - \eta \nabla \mathcal{L} \\ & \text{Learning rate} \end{split}$$

• In practice, handled by optimise algorithm (e.g. Adam)

How do networks learn?



Batches and Epochs

- The loss depends on O(1k-1b) parameters
- To minimize it these parameters need to be adjusted
- For large networks the loss landscape can get very complicated with many local minima
- Better convergence
 - Multiple passes (epochs) over training data
 - Split training data in batches with a minimization call and model update after each batch







Learning Rate

- Many popular minimization algorithms have a tunable parameter called Learning rate α
- Should be chosen such that the training smoothly converges



https://www.jeremyjordan.me/nn-learning-rate/

Adam

$$egin{aligned} m_w^{(t+1)} &\leftarrow eta_1 m_w^{(t)} + (1 - eta_1) \,
abla_w L^{(t)} \ v_w^{(t+1)} &\leftarrow eta_2 v_w^{(t)} + (1 - eta_2) \left(
abla_w L^{(t)}
ight)^2 \end{aligned}$$

Maintain history of previous changes: momentum

Exponential moving average of gradient and gradient squares

$$\hat{m}_w = rac{m_w^{(t+1)}}{1-eta_1^t} \ \hat{v}_w = rac{v_w^{(t+1)}}{1-eta_2^t}$$

$$w^{(t+1)} \leftarrow w^{(t)} - \eta rac{\hat{m}_w}{\sqrt{\hat{v}_w} + \epsilon}$$

https://arxiv.org/abs/1412.6980

The loss curve

- Once the training is done it is time to analyze the training process and the performance of the model
- For the training process it is a good idea to look at the loss as a function of the number of epochs – for the training data and validation data
- Ideally the training converges with similar loss values for the validation and training data



Overfitting

- If the model has too many parameters it can overfit the data
- This leads to a worse performance on data with different statistical fluctuations than the training data
- -> the training and validation losses diverge



Overfitting

- Also an issue in classification tasks (which are after all a form of regression of discrete values)
- Overfitting can be interpreted as remembering the training samples
- We can see that the performance would change if the dots move around even a bit



Overfitting - Mitigation

- There are multiple ways to mitigate overfitting such as dropout, regularization or early stopping
- Dropout:
 - Add layers that randomly discard a fraction of inputs in training phase
- Early stopping:
 - Monitor training and validation loss during the training process and stop when they diverge
- Regularisation:
 - Add penalty for large weight values in loss function



Side note: Overfitting

Reconciling modern machine-learning practice and the classical bias– variance trade-off (PNAS, Belkin et al)



- Recent observation in computer science
 - Extremely high capacity does not guarantee overfitting
- Even more recent:
 - Grokking: Sudden generalisation of a trained model
- Active research into foundations, does not stop us from using it

General Strategy

Loss function $\mathcal{L}\checkmark$ Neural network $f_{\theta}\checkmark$ Parameters $\theta\checkmark$ Opt. Parameters $\theta^*\checkmark$ Data **x** Data distribution $p(\mathbf{x})$

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

High-level/no structure: Fully Connected



Regular grid: Convolution



Sequence/Time Series: Recurrent

Point cloud: Sets & Graphs



General Strategy

Loss function $\mathcal{L} \checkmark$ Neural network $f_{\theta} \checkmark$ Parameters $\theta \checkmark$ Opt. Parameters $\theta^* \checkmark$ Data $\mathbf{x} \checkmark$ Data distribution $p(\mathbf{x}) \checkmark$

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

Summary so far

- Key idea of deep learning:
 - Define a optimisation target and use an expressive function (neural network) to minimise it using gradient descent
- Let's look at a typical problem in particle physics...

Jet Tagging



A jet is a collimated shower of particles in the detector



We want to know which particle produced a jet



How to build ML algorithms for complex, heterogenous data?



1810.05165; **GK** et al 2312.00123



* techniques are also applied to collider data

GK, Plehn, et al 1902.09914

Landscape Dataset

- Open dataset for the development of better tagging algorithms for particle physics
- 2 million simulated examples
- Input: momentum sorted list of 200 particles/jet with 3 features/particle (px, py, pz)
- Perfect class labels: top jet or light quark/gluon jet
- Supervised learning problem

The Machine Learning Landscape of Top Taggers	
 G. Kasieczka (ed)¹, T. Plehn (ed)², A. But M. Fairbairn⁶, D. A. Faroughy⁵, W. Fedo P. T. Komiske¹⁰, S. Leiss¹, A. Lister⁷, S B. Nachman,^{12,13}, K. Nordström^{14,15}, J. Pea J. M. Thompse 	 ter², K. Cranmer³, D. Debnath⁴, B. M. Dillon⁵ rko⁷, C. Gay⁷, L. Gouskos⁸, J. F. Kamenik^{5,9}, Macaluso^{3,4}, E. M. Metodiev¹⁰, L. Moore¹¹, rkes⁷, H. Qu⁸, Y. Rath¹⁶, M. Rieger¹⁶, D. Shih⁷ on², and S. Varma⁶
 Institut für Experimentalphy Institut für Theoretische Phys Center for Cosmology and Particle Phy NHECT, Dept. of Physics and Astronom Jozef Stefan Instit Theoretical Particle Physics and Cosmo Department of Physics and Astronomy Bepartment of Physics, Universi Faculty of Mathematics and Physics, 10 Center for Theoretical I 11 CP3, Universitéxx Catholique of 12 Physics Division, Lawrence Berk Simons Inst. for the Theory of Comp 14 National Institute for Subatomic PI 15 LPTHE, CNRS & Sort 16 III. Physics Institute A, R 	sik, Universität Hamburg, Germany sik, Universität Heidelberg, Germany vsics and Center for Data Science, NYU, USA my, Rutgers, The State University of NJ, USA tute, Ljubljana, Slovenia ology, King's College London, United Kingdom , The University of British Columbia, Canada sity of California, Santa Barbara, USA University of Ljubljana, Ljubljana, Slovenia Physics, MIT, Cambridge, USA de Louvain, Louvain-la-Neuve, Belgium eley National Laboratory, Berkeley, USA uting, University of California, Berkeley, USA hysics (NIKHEF), Amsterdam, Netherlands ponne Université, Paris, France WTH Aachen University, Germany
gregor.kasieczł plehn@un	æ@uni-hamburg.de i-heidelberg.de
Inter	24, 2010

Abstract

arXiv:1902.09914v3 [hep-ph] 23 Jul 2019

Based on the established task of identifying boosted, hadronically decaying top quarks, we compare a wide range of modern machine learning approaches. Unlike most established methods they rely on low-level input, for instance calorimeter output. While their network architectures are vastly different, their performance is comparatively similar. In general, we find that these new approaches are extremely powerful and great fun.



(*Subset of methods, graph from V Breso at ML4Jets)



Qu, Gouskos 1902.08570; Quret al 2202.03772; Gong et al 2201.08187; Shlomi et al 2007.13681; Bogatskiy et al 2211:00454;



Qu, Gouskos 1902.08570; Qu et al 2202.03772; Gong et al 2201.08187; Shlomi et al 2007.13681; Bogatskiy et al 2211.00454; Brehmer et al 2411.00446





Important issues on application side

- Domain shift
- Calibrateable
- Compute cost

Modern tagging algorithms are widely used in searches for new particles

End of Part I.

Bonus Material

Symmetries



Reminder: Space-time described by Lorentz group (rotations, boosts)

$$(t,x,y,z)\mapsto t^2-x^2-y^2-z^2$$

Equivariance: L-GATr $(\Lambda(x)) = \Lambda(L-GATr(x))$

0

0

Multivector:

"Grades" C with well defined behaviour under Lorentz transformations

$$x = x^{S} 1 + x^{V}_{\mu} \gamma^{\mu} + x^{B}_{\mu\nu} \sigma^{\mu\nu} + x^{A}_{\mu} \gamma^{\mu} \gamma^{5} + x^{P} \gamma^{5}$$
Scalar Bi-Vector Pseudo-Scalar
tz Vector Axial-Vector
$$\sqrt[p]{e^{-\binom{1}{0} \frac{1}{0} \frac{0}{0} \frac{0}{0} - 1}} \sqrt[p]{e^{-\binom{1}{0} \frac{1}{0} \frac{0}{0} \frac{0}{0} - 1}} \sqrt[p]{e^{-\binom{1}{0} \frac{1}{0} \frac{0}{0} \frac{0}{0} - 1}} \sqrt[p]{e^{-\binom{1}{0} \frac{0}{0} \frac{0}{0} - 1}}} \sqrt[p]{e^{-\binom{1}{0} \frac{0}{0} \frac{0}{0} - 1}} \sqrt[p]{e^{-\binom{1}{0} \frac{0}{0} - 1}}} \sqrt[p]{e^{-\binom{1}{0} \frac{0}{0} - 1}}} \sqrt[p]{e^{-\binom{1}{0} \frac{0}{0} - 1}}} \sqrt[p]{e^{-\binom{1}{0} \frac{0}{0} - 1}} \sqrt[p]{e^{-\binom{1}{0} \frac{0}{0} - 1}}} \sqrt[p]{e^{$$

 $\gamma^2 =$

Sum over grades **Symmetries** Transformer L-GAIr Layer type Linear(x) Attention(q, k, v)_{ic} $\sum_{j=1}^{n_t} \operatorname{Softmax}_j \left(\sum_{c'=1}^{n_c} \frac{q_{ic'}k_{jc'}}{\sqrt{n_c}} \right) v_{jc}$ $\sum_{j=1}^{n_t} \operatorname{Softmax}_j \left(\sum_{c'=1}^{n_c} \frac{\langle q_{ic'}, k_{jc'} \rangle}{\sqrt{16n_c}} \right) v_{jc}$ LayerNorm(x) $x \left[\frac{1}{n_c} \sum_{c=1}^{n_c} x_c^2 + \epsilon \right]^{-1/2}$ $x \left[\frac{1}{n_c} \sum_{c=1}^{n_c} \sum_{k=0}^{4} \left| \left\langle \langle x_c \rangle_k, \langle x_c \rangle_k \right\rangle \right| + \epsilon \right]^{-1/2}$ GELU($\langle x \rangle_0$)x GELU(x)Activation(x) GP(x, y)xу Output Probabilities Softmax Inner product Linear Norm of Add & No (requires some work but Feed Forward inner product can compute Algebra-Add & Norm by grade Multi-Head rules in advances) Feed Attention N× Forward Add & Norn N× $x_\mu x^\mu = \eta_{\mu u} x^\mu x^ u = (ct)^2 - \mathbf{x} \cdot \mathbf{x} \stackrel{\mathrm{def}}{=} (c au)^2$ Add & No Masked Multi-Head Multi-Hea Attention Attention $\eta^{\mu u} = \eta_{\mu u} = egin{pmatrix} 1 & 0 & 0 & 0 \ 0 & -1 & 0 & 0 \ 0 & 0 & -1 & 0 \ \end{pmatrix}.$ Positional Positional Encoding Encoding Output Input Embedding Embedding Inputs Outputs 2411.00446 (shifted right)

Classification

$$-\frac{1}{N}\log(L) = -\frac{1}{N}\sum_{k}^{N} [y_k \cdot \log(\hat{y}(x_k)) + (1 - y_k) \cdot \log(1 - \hat{y}(x_k))]$$

• Note:

$$y_k \cdot \log(\hat{y}(x_k)) + (1 - y_k) \cdot \log(1 - \hat{y}(x_k)) = \sum_i q_i \log(p_i)$$

Now remember entropy

entropy = $\int p(x) \log (p(x)) dx$

cross-entropy = $\int q(x) \log (p(x)) dx$

Gibbs inequality: Cross entropy is minimal if p=q

For multiple classes: Categorical cross entropy

Multiple classes

- For more than 2 classes it is usually a good idea to use **One-Hot-Encoding**
 - -> For m classes the output is a m-dimensional vector with components

 $\hat{y}_c(x_k) \in [0,1]$

- One neuron per class in output layer
- All outputs should sum to one
- Loss function: Extension of cross entropy to multiple classes: Categorical cross entropy



https://medium.com/@michaeldelsole/

Activation functions for Classification

- Need activations in the output layer that allow probability interpretation:
 - $\hat{y}(x_k) \in [0,1]$
 - $\sum_{classes} \hat{y}(x_k) = 1$
- Binary classification:
 - Only 1 output node needed
 - Activation: Sigmoid

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



- Multiclassification:
 - One output node per class i
 - Softmax activation:

$$\sigma(x) = \frac{e^{x_i}}{\sum_j^{classes} e^{x_j}}$$

Binary Classification Outputs

- The NN outputs a score for belonging to class "signal" (high value) or class "background" (low value)
- Often you need to make a decision based on this score -> Choose a threshold
- Key quantities:
 - True Positives (TP)
 - False Positives (FP)
 - True Negatives (TN)
 - False Negatives (FN)



- Purity/Precision = TP/(TP+FP): How pure is the sample predicted as signal?
- True Positive Rate/Sensitivity/Efficiency/Recall = TP/(TP+FN): What fraction of the signal is classified correctly?
- False Positive Rate = FP/(FP +TN)
- Accuracy = (TP+TN)/(TP+FP+TN+FN): What fraction of all data is classified correctly?

Choosing a threshold

- Which threshold to choose?
- Depends on your problem:
 - Think about what questions you want to answer
 - The cost of being wrong
- Do you need a high purity? ("We need to be sure that all animals we tag as harmless are actually harmless!")
- Do you need a high Efficiency? ("We need to detect as many positive Covid cases as possible")



Multiclassification

- In multiclassification we have one output neuron for each class
- If we use softmax as final activation function, the values of each output neuron sum to 1 and can be interpreted as probability
- Common approach: Classify a sample according to the neuron with the highest output value



Receiver Operator Characteristic (ROC)

- Overall performance of a classifier is usually visualized using the ROC curve:
 - Scan all possible thresholds and evaluate True Positive Rate and False Positive Rate
 - Plot one as a function of the other
- Area under the curve (AUC) often quoted as performance metric
- ROCs can aide in choosing a classifier and threshold
- Note: Sometimes ROCs of different classifiers cross each other -> they are better at different TPR regions





Confusion matrices

- With a chosen threshold we can plot confusion (migration) matrices
- Shows migration from one class to another
- Very usefull for multiclassification:
 - Which background class looks most signal like?
 - Which classes are difficult to separate?
- Can be normalized to:
 - True classes (efficiency matrix)
 - Predicted class (purity matrix)





HISS Summer School

Machine Learning: Anomalies and Generative Models

Prof. Gregor Kasieczka gregor.kasieczka@uni-hamburg.de

Plan

- Lecture 1:
 - Basics of Neural Networks
 - Classification example: Top Tagging
- Lecture 2:
 - Anomaly Detection
 - Generative Models

Plan

- Lecture 1:
 - Basics of Neural Networks
 - Classification example: Top Tagging
- Lecture 2:
 - Anomaly Detection
 - Generative Models





(Rubbia and van der Meer)







Divide into 25 categories according to purity, relative amount of signal using simple ML





Weight events according to excepted purity of category



Nobel price for Higgs mechanism 2013 (Higgs and Englert)


How to discover a new particle?



How to discover a new particle?



Overview of CMS EXO results



Anomaly Detection

What is an anomaly?



Point anomaly

- Outliers: Datapoints far away from regular distribution
- Examples:
 - Detector malfunctions
 - Background-free search













And now?



Group anomaly



Color (Mass)

Group anomaly

- Individual examples not anomalous, Count but interesting collective behaviour a.u. Examples: New physics searches, e.g. resonances Excess in time series SB SR SB m $p_{\text{data}}(x|m \in SB)$ $p_{\text{data}}(x|m \in SB)$ $p_{\text{data}}(x|m \in SR)$ $= p_{bg}(x|m \in SB)$ $= p_{bg}(x|m \in SB)$
 - Color (Mass)

Group anomalies

Revisit assumptions



Assumptions, revisited







Define a signal cut-out



And interpolate the background







Resonant Anomalies



Consider resonant anomalies: fully data-based construction of anomaly detection score

Resonant Anomalies



CWola



$$L_{M_1/M_2} = \frac{p_{M_1}}{p_{M_2}} = \frac{f_1 \, p_S + (1 - f_1) \, p_B}{f_2 \, p_S + (1 - f_2) \, p_B} = \frac{f_1 \, L_{S/B} + (1 - f_1)}{f_2 \, L_{S/B} + (1 - f_2)}, \quad \text{Problems?}$$

Nachman et al 1708.02949



Aside: Generative Models



Given a set of examples {x_i} drawn from a distribution $p(\mathbf{x})$, learn a mapping f_{Θ} so that $f_{\Theta}(q(\mathbf{z})) = p(\mathbf{x})$ for a random distribution q(**z**)

More later









Aside: CASE

Available on the CERN CDS information server

CMS PAS EXO-22-026

CMS Physics Analysis Summary

Contact: cms-pag-conveners-exotica@cern.ch

2024/03/20

Model-agnostic search for dijet resonances with anomalous jet substructure in proton-proton collisions at $\sqrt{s} = 13$ TeV

The CMS Collaboration

Abstract

This note introduces a model-agnostic search for new physics in the dijet final state. Other than the requirement of a narrow dijet resonance with a mass in the range of 1800-6000 GeV, minimal additional assumptions are placed on the signal hypothesis. Search regions are obtained by utilizing multivariate machine learning methods to select jets with anomalous substructure. A collection of complementary anomaly detection methods – based on unsupervised, weakly-supervised and semi-supervised algorithms – are used in order to maximize the sensitivity to unknown new physics signatures. These algorithms are applied to data corresponding to an integrated luminosity of 138 fb⁻¹, recorded in the years 2016 to 2018 by the CMS experiment at the LHC, at a centre-of-mass energy of 13 TeV. No significant excesses above background expectation are seen, and exclusion limits are derived on the production cross section of benchmark signal models varying in resonance mass, jet mass and jet substructure. Many of these signatures have not previously been searched for at the LHC, making the limits reported on the corresponding benchmark models the first ever and the most stringent to date.

- New result by the CMS collaboration
- Full Run 2 dataset
- 6 anomaly detectors in parallel



What's Beyond?



What's Beyond?



More features

- Better generative models / classifiers
- Other topologies
- Non-resonant features

Outliers

Outliers



What if the events are even weirder?

Autoencoders



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



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

Autoencoders

• Potential issue of models that only use pBackground (e.g. autoencoders)



Example: Triggering Outliers

Learn-compression/decompression on signal free sample and use as anomaly score Now testing in CMS Level 1 trigger



https://indico.cern.ch/event/1283970/ contributions/5554350/attachments/ 2720710/4727877/axol1tl_fastml.pdf






Comments on Anomalies

- Systematic anomaly searches as fail-safe strategy to signature based searches
- Trade off between sensitivity and breadth

• Weak supervision as useful technique beyond anomalies

Plan

- Lecture 1:
 - Basics of Neural Networks
 - Classification example: Top Tagging
- Lecture 2:
 - Anomaly Detection
 - Generative Models

Motivation

Have: input examples (collision events, detector readouts, ...)



Want: more data

Specifically: new data similar to the input, but not exact copies

How to encode in neural net?

Uses:

- Simulations
- In-situ background estimation
- Surrogate models

•

This happens in the experiment



This is what we want to know

Simulation is crucial to connect experimental data with theory predictions

This happens in the experiment



This is what we want to know

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



2020 Computing Model -CPU: 2030: Baseline

ATLAS Preliminary



This happens in the experiment



This is what we want to know

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

Use AI to improve efficiency of simulation codes or learn surrogates

Calorimeter Simulation

Interaction of particles with multi-layer detectors to determine their initial energy (and type)

Measurement of energy, position, (and time) of secondary particle hits



Calorimeter Simulation

Interaction of particles with multi-layer detectors to determine their initial energy (and type)

Measurement of energy, position, (and time) of secondary particle hits

Represent data as -fixed grid (3d matrix of detector elements, value=energy) -point cloud (set of hits, each hit is a 3d vector with position+energy)



HEP Strategy

1. Use classical simulation or data as input



(slow)

2. Train generative surrogate



3. Oversample





Paganini, Oliveira, Nachman 1705.02355; Butter, Diefenbacher, **GK**, et al 2008.06545;

Overview





1406.2661 lilianweng.github.io



Training objective: Binary cross entropy

$$\begin{split} \min_{G} \max_{D} V(D,G) &= \mathbb{E}_{\boldsymbol{x} \sim p_{\mathsf{data}}(\boldsymbol{x})} [\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z})))] \\ \uparrow & \uparrow \\ & \uparrow \\ & \mathsf{True examples} \\ \end{split}$$
 Fake examples







Training objective: Binary cross entropy

 $\min_{G} \max_{D} V(D,G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})} [\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z})))]$

At (Nash) equilibrium: Generator produces realistic examples Discriminator is maximally confused



Training objective: Binary cross entropy

 $\min_{G} \max_{D} V(D,G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})} [\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z})))]$

For generation: Sample from Generator Discard Discriminator

Comments on GANs

Architecture:

• Low complexity, fast and adaptable

Learning:

- Unstable training
- Matching of generator/discriminator (vanishing gradients)
- Mode collapse
- Loss function not interpretable

Maturity:

• Well established, many variants and extensions



Mode collapse

Variational Autoencoders

Autoencoder



Two networks Encoder: data \rightarrow latent space Decoder: latent space \rightarrow data

Autoencoder



Two networks Encoder: data \rightarrow latent space Decoder: latent space \rightarrow data

Training objective: L = Minimise input/output difference



Autoencoder



Two networks Encoder: data \rightarrow latent space Decoder: latent space \rightarrow data

 $L = (x - f(g(x)))^2$

Training objective: Minimise input/output difference

Uses:

Dimension reduction Denoising Anomaly detection Generation?

Variational Autoencoder



 $f(x) = (\mu, \sigma)$

Variational Autoencoder (VAE): Split latent space

Variational Autoencoder



Variational Autoencoder (VAE): Split latent space Sample before decoder

$$f(x) = (\mu, \sigma)$$
$$z = \text{Gaussian}(\mu, \sigma)$$

$$x' = g(z)$$

Variational Autoencoder



Variational Autoencoder (VAE):

Split latent space Sample before decoder Penalty so mean/std are close to unit Gaussian

$$J(w) = (\mu, 0)$$

 $f(r) = (\mu \sigma)$

$$z = \text{Gaussian}(\mu, \sigma)$$

x' = g(z)

$$L = (x - g(z))^{2} + \sigma^{2} + \mu^{2} - \log(\sigma) - 1$$

(Calculate KL-divergence
between Gaussians)

VAE Example



Loss terms

Latent space of MNIST VAE



towardsdatascience.com

Comments on VAEs

Architecture:

- Low complexity, fast and adaptable
- Target: Maximise lower bound on likelihood

Learning:

- Stable training
- Average prediction → blurrier output
- Interpretable latent space



VAE

DCGAN

Maturity:

• Well established, many variants and extensions

Normalising Flows



In auto-encoders, the decoder learns to 'undo' the encoder

Can we make this exact?



Learn a diffeomorphism between data and latent-space



Learn a diffeomorphism between data and latent-space

Bijective, invertable



Learn a diffeomorphism between data and latent-space

Bijective, invertable

Learn likelihood of data

Take into account Jacobian determinant to evaluate probability density



Easy-to-calculate Jacobean

Coupling flows



Coupling layers: Not the most expressive, but useful for illustration/understanding

Coupling flows



Simple (e.g. dense) neural networks
Coupling flows



Generative models



Invertible Easy-to-calculate Jacobian probability density

Calculating Jacobian determinant



$$\begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} \xrightarrow{f_1} \begin{pmatrix} \mathbf{z}_1 \\ \mathbf{x}_2 \end{pmatrix} \xrightarrow{f_2} \begin{pmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{pmatrix} \text{ with } \begin{aligned} \mathbf{x}_1 \xrightarrow{f_1} \mathbf{z}_1 &= \mathbf{x}_1 \odot \exp(s_2(\mathbf{x}_2)) + t_2(\mathbf{x}_2) \\ \mathbf{x}_2 \xrightarrow{f_1} \mathbf{x}_2. \end{aligned}$$

$$\mathbf{J_1} = \begin{pmatrix} \frac{\partial \mathbf{z}_1}{\partial \mathbf{x}_1} & \frac{\partial \mathbf{z}_1}{\partial \mathbf{x}_2} \\ \frac{\partial \mathbf{x}_2}{\partial \mathbf{x}_1} & \frac{\partial \mathbf{x}_2}{\partial \mathbf{x}_2} \end{pmatrix} = \begin{pmatrix} \operatorname{diag}(\exp(s_2(\mathbf{x}_2))) & \frac{\partial \mathbf{z}_1}{\partial \mathbf{x}_2} \\ 0 & 1 \end{pmatrix}$$

Triangular by construction

$$\det \mathbf{J_1} = \prod \exp(s_2(\mathbf{x}_2)) = \exp\left(\sum s_2(\mathbf{x}_2)\right)$$

Composition



Composition of bijective functions remains bijective

Chain rule: Jacobian determinant of composition is product of determinants

How to train NF?

Training objective: Minimise negative log likelihood of data

Sample points from training data

$$\mathcal{L} = -\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[-\frac{1}{2} ||f(\mathbf{x})||_2^2 + \sum s(\mathbf{x}) \right]$$

How to train NF?

Training objective: Minimise negative log likelihood of data



How to train NF?

Training objective: Minimise negative log likelihood of data

$$\begin{aligned} \mathcal{L} &= -\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \left[-\frac{1}{2} || f(\mathbf{x}) \rangle ||_2^2 + \sum s(\mathbf{x}) \right] \\ & \text{Contribution from Jacobian} \\ & \text{determinant} \\ & \text{det } \mathbf{J} = \exp\left(\sum s(\mathbf{x})\right) \\ & -\log(\det \mathbf{J}) = -\sum s(\mathbf{x}) \end{aligned}$$

Animation



Comments on Flows

Only scratched the surface: more constructions available

→ Better generative fidelity
→ Can evaluate likelihood of

data

More complex

→ Slower, choice of fast direction







Core idea: Stepwise transition from pure noise to data

Markov chain





Core idea: Stepwise transition from pure noise to data



https://medium.com/mlearning-ai/enerating-images-withddpms-a-pytorch-implementation-cef5a2ba8cb1



Core idea: Stepwise transition from pure noise to data

z [layers]



CaloCloud, time stamp: t_{99}



with stochastic differential equations (SDEs)

Comments on Generative Models

• Simulations are a core need of fundamental physics but computationally very costly

• Aided by developments in generative models

 Main application: Detector simulation, but also other developments (theory simulation, in-situ backgrounds, ...)



Closing

- Broad developments on ML in Particle Phsycis
- Covered some key topics, but wide range of other applications
- -> https://iml-wg.github.io/HEPML-LivingReview/



Bonus Material



This is added in the training process



transition to any time



This is added in the training process



This is added in the training process





Optimal expected mean (take into account known noise schedule)



and learn to predict noise



Reverse Resulting learning objective
(Noise
$$\rightarrow$$
 data) $L_{simple}(\theta) \coloneqq \mathbb{E}_{t,\mathbf{x}_0,\epsilon} \left[\left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta} (\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \right\|^2 \right]$
Noisy image \mathbf{R}_{t}
Reminder: Forward diffusion to time t $\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}) = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$

Timestep



Core idea: Stepwise transition from pure noise to data

Algorithm 1 Training

- 1: repeat
- 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
- 4: $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- Take gradient descent step on 5:

$$\nabla_{\theta} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta} (\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \right\|^2$$

6: **until** converged



Core idea: Stepwise transition from pure noise to data

Algorithm 1 Training	Algorithm 2 Sampling
1: repeat 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 3: $t \sim \text{Uniform}(\{1, \dots, T\})$ 4: $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$ 5: Take gradient descent step on $\nabla_{\theta} \ \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta} (\sqrt{\overline{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \overline{\alpha}_t} \boldsymbol{\epsilon}, t) \ ^2$ 6: until converged	1: $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$ 2: for $t = T, \dots, 1$ do 3: $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$ if $t > 1$, else $\mathbf{z} = 0$ 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\overline{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 5: end for 6: return \mathbf{x}_0

Forward SDE:



(Correspond to the noise schedule in discrete case)

Forward SDE (data \rightarrow noise) $\mathbf{x}(0)$ $\mathbf{dx} = \mathbf{f}(\mathbf{x}, t) dt + g(t) d\mathbf{w}$ $\mathbf{x}(T)$ $\mathbf{x}(T)$ $\mathbf{x}(0)$ $\mathbf{dx} = [\mathbf{f}(\mathbf{x}, t) - g^2(t) \nabla_{\mathbf{x}} \log p_t(\mathbf{x})] dt + g(t) d\bar{\mathbf{w}}$ $\mathbf{x}(T)$ Reverse SDE (noise \rightarrow data)

Probability density of x(t)

Reverse SDE:
$$d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t) - g(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x})] dt + g(t) d\bar{\mathbf{w}}$$

Score function

Reverse of a diffusion process is also a diffusion



Reverse SDE:
$$d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t) - g(t)^2 \nabla_{\mathbf{x}} \log p_t(\mathbf{x})] dt + g(t) d\bar{\mathbf{w}}$$

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{arg\,min}} \mathbb{E}_t \Big\{ \lambda(t) \mathbb{E}_{\mathbf{x}(0)} \mathbb{E}_{\mathbf{x}(t)|\mathbf{x}(0)} \Big[\left\| \mathbf{s}_{\boldsymbol{\theta}}(\mathbf{x}(t), t) - \nabla_{\mathbf{x}(t)} \log p_{0t}(\mathbf{x}(t) \mid \mathbf{x}(0)) \right\|_2^2 \Big] \Big\}$$

Learn to approximate score function with neural network

