

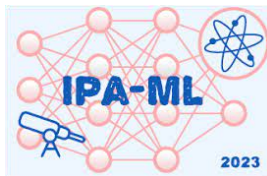
Analysis techniques in High Energy Physics

Florian Eble

22/03/2023

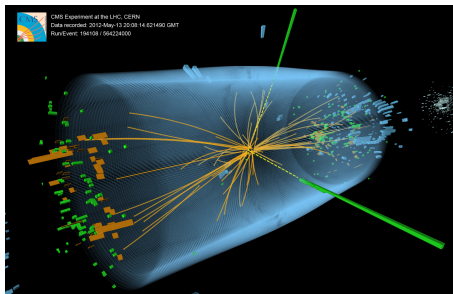
IPA ML Workshop

ETH zürich



IPA

- Very wide topic!
- Will focus on some common HEP problems
- Will discuss how ML can help us solve them
- Hopefully, the problems discussed in this talk are general enough and applicable to other disciplines!

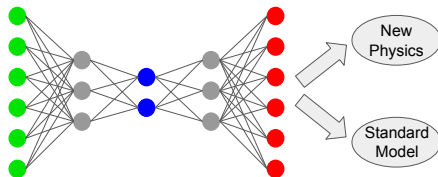
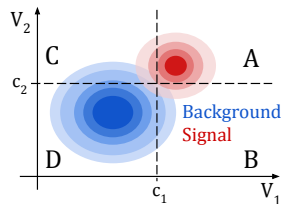
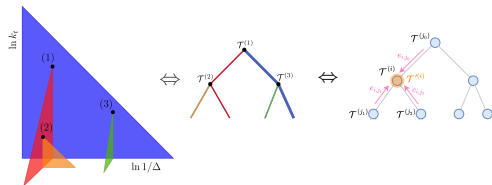


Incorporating physics inductive bias in ML models

Background estimation

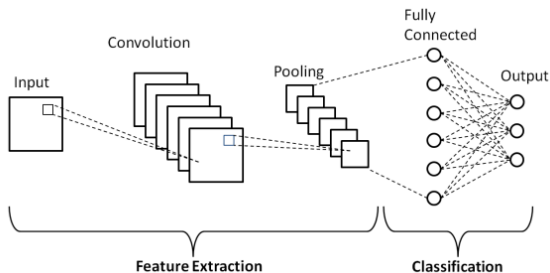
Training decorrelated models

Searching for new physics via anomaly detection



- 1 Introducing inductive bias
- 2 Background estimation in HEP
- 3 Training decorrelated models
- 4 Searching for new physics

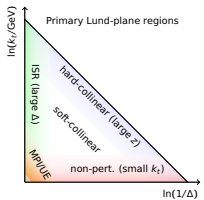
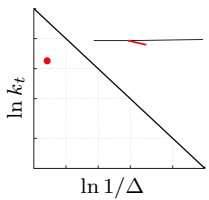
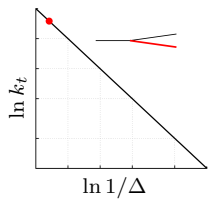
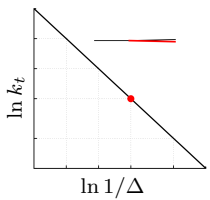
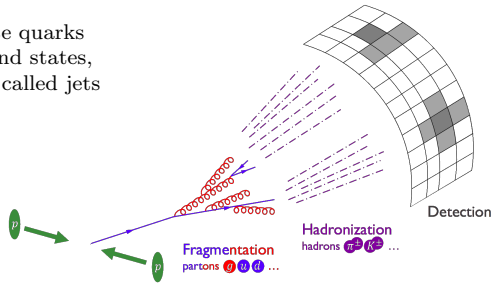
- The **universal approximation theorem** states every continuous function can be approximated by a neural network (NN)
- However, designing architectures exploiting specificities of a problem is often a necessity for a successful learning!
 - Introducing **inductive bias** in NN
- *E.g.* Convolutional Neural Network (CNN) architectures make use of translational invariance of images by implementing dedicated convolutional layer



Same output!

Lund plane representation of a jet

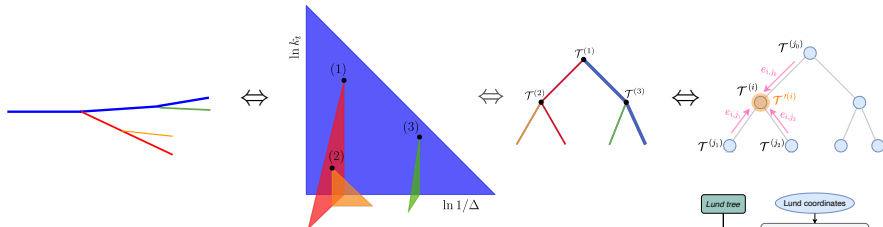
- Because of color confinement, final state quarks and gluons fragment until forming bound states, forming a collimated spray of particles called jets
- The Lund plane¹ is 2-dimensional representation of gluon/quark emission in hadronic showers
- Natural description of the radiation pattern inside of a jet



$$\Delta = \sqrt{(y_a - y_b)^2 + (\phi_a - \phi_b)^2} \quad \text{opening angle of the splitting}$$

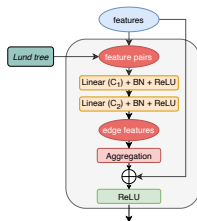
$$k_t = p_t \Delta$$

¹See [arXiv:1807.04758](https://arxiv.org/abs/1807.04758)

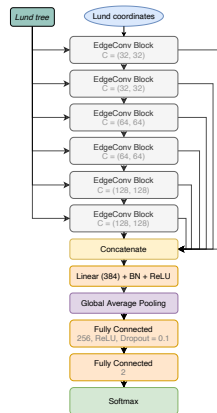


- At each splitting, compute a tuple of kinematic features $\mathcal{T} = (k_t, \Delta, \dots)$, used as node features in the graph
- EdgeConv operation as a fully connected NN using features of connected nodes
- Stack EdgeConv layers to build up the LundNet¹
- State-of-the-art jet tagger with small computational cost (fixed graph via Lund decomposition)

$$\mathcal{T}^{(i)} = \prod_{k=0}^3 \mathbf{h}_{\Theta}(\mathcal{T}^{(i)}, \mathcal{T}^{(j_k)})$$



The EdgeConv equation and EdgeConv block

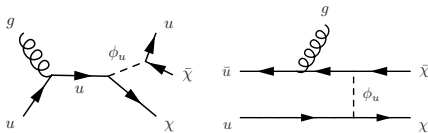


¹See [arXiv:2012.08526](https://arxiv.org/abs/2012.08526)

- 1 Introducing inductive bias
- 2 Background estimation in HEP
- 3 Training decorrelated models
- 4 Searching for new physics

- General problem in HEP: estimate the expected number of background events in signal-enriched region of the phase-space
- *E.g.* searches for WIMPs by looking for an excess of event at high E_T (Missing Transverse Energy)
 - Need to know expected number of Standard Model events in the signal region!
- Example of [arXiv:1712.02345](https://arxiv.org/abs/1712.02345), backgrounds estimated from simulation:

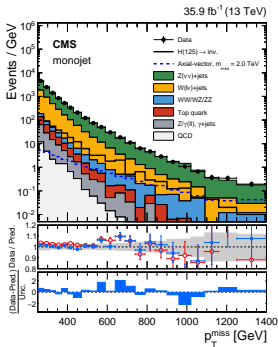
Example diagrams:



Monojet event selection:

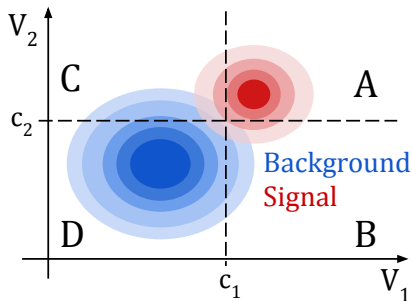
Variable	Selection	Target background
Muon (electron) veto	$p_T > 10 \text{ GeV}, \eta < 2.4(2.5)$	$Z(\ell\ell)+\text{jets}, W(\ell\nu)+\text{jets}$
τ lepton veto	$p_T > 18 \text{ GeV}, \eta < 2.3$	$Z(\ell\ell)+\text{jets}, W(\ell\nu)+\text{jets}$
Photon veto	$p_T > 15 \text{ GeV}, \eta < 2.5$	$\gamma+\text{jets}$
Bottom jet veto	$\text{CSV}v2 < 0.8484, p_T > 15 \text{ GeV}, \eta < 2.4$	Top quark
p_T^{miss}	$> 250 \text{ GeV}$	QCD, top quark, $Z(\ell\ell)+\text{jets}$
$\Delta\phi(\vec{p}_T^{\text{jet}}, \vec{p}_T^{\text{miss}})$	> 0.5 radians	QCD
Leading AK4 jet p_T and η	$> 100 \text{ GeV}$ and $ \eta < 2.4$	All

Searching for high E_T :



- Sometimes backgrounds cannot be estimated from simulation, *e.g.* cross-section of processes with large number of hadronic jets is difficult to calculate
- Need to use signal-free regions (control regions) to predict background in the signal region
- The ABCD method is one of the methods for this task
- Let V_1 and V_2 be **2 independent** variables for the background distribution
- Signal region A: $V_1 > c_1$ and $V_2 > c_2$
- Number of events in each region: N_A, N_B, N_C, N_D
- Background estimation in signal region:

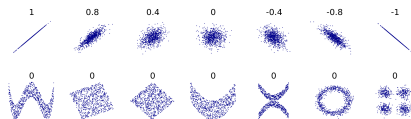
$$N_A^{\text{bkg}} = \frac{N_B}{N_D} N_C$$



- One requirement for the ABCD method is to have **2 independent variables**
- This can be checked using Distance Correlation (DisCo) [1] [2] [3]

- **Pearson correlation** only evaluates **linear correlations**:

$$\rho_{\text{Pearson}}^2(X, Y) = \frac{\text{Cov}^2(X, Y)}{\text{Cov}(X, X)\text{Cov}(Y, Y)}$$

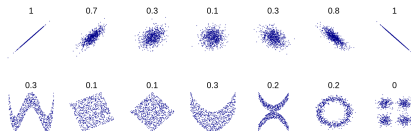


Pearson correlation coefficient

- **Distance correlation (DisCo)** defined using the probability distributions of X and Y , and their joint probability distribution

→ Makes use of all information of the random variables!

$$\text{DisCo}^2(X, Y) = \frac{d\text{Cov}^2(X, Y)}{d\text{Cov}(X, X)d\text{Cov}(Y, Y)}$$



Distance correlation coefficient

If no two clear discriminative and independent variables can be found:

- Can train one NN and use DisCo regularization to force its output to be independent of a discriminative physics observable V

For a batch of examples X :

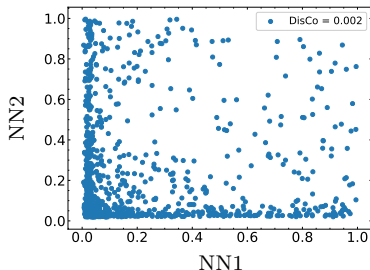
$$L(X) = L_{\text{NN}}(X) + \lambda \cdot \text{DisCo}(\text{NN}(X), V(X)) \quad (1)$$

- Can train two NNs and use DisCo regularization to force them to be independent of each other!

For a batch of examples X :

$$L(X) = L_{\text{NN1}}(X) + L_{\text{NN2}}(X) + \lambda \cdot \text{DisCo}(\text{NN1}(X), \text{NN2}(X)) \quad (2)$$

Where L_{NN} , L_{NN1} , L_{NN2} are the NNs usual loss, *e.g.* for a supervised binary classifier, binary cross-entropy.



- And then perform ABCD background estimation in the (NN, V) plane or $(\text{NN1}, \text{NN2})$ plane
- **ML provides a systematic way of addressing this issue!**

References:

[arXiv:2001.05310](https://arxiv.org/abs/2001.05310)

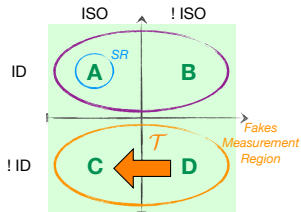
[arXiv:2007.14400](https://arxiv.org/abs/2007.14400)

- Learn transfer factor \mathcal{T} from signal-free control region D to C for N physics variables \mathcal{V}_i

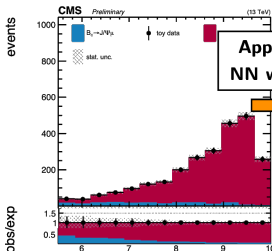
$$\mathcal{T} = \frac{\alpha - \beta \cdot \gamma}{1 - \gamma}, \quad \alpha = \frac{\text{Data}_C}{\text{Data}_D}, \quad \beta = \frac{\text{MC}_C}{\text{MC}_D}, \quad \gamma = \frac{\text{MC}_D}{\text{Data}_D}$$

- Train 3 NNs to learn the mappings α, β, γ between the N -dimensional distributions of the different “regions”
- Apply \mathcal{T} from signal-free region B to signal region A:

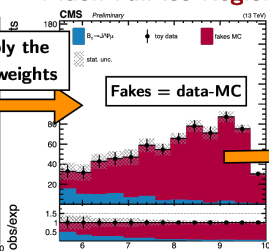
$$\text{Fake}_A = \mathcal{T}(\text{Data}_B) \cdot \text{Data}_B - \mathcal{T}(\text{MC}_B) \cdot \text{MC}_B$$



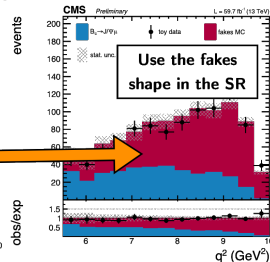
Muon Fail Iso Region



NN Reweighted Muon Fail Iso Region



Signal Region



- 1 Introducing inductive bias
- 2 Background estimation in HEP
- 3 Training decorrelated models
- 4 Searching for new physics

- **Distance correlation regularization** can also be used to train a classifier decorrelated from a feature \mathcal{F}
 - Particularly interesting when this feature \mathcal{F} is then used in the next step of the analysis, for instance fitted to extract signal
- Example: Mass decorrelated classifier for resonant search (bump hunt in a mass variable)

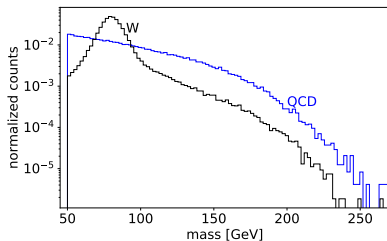


FIG. 1: Invariant mass distribution for the inclusive W and QCD samples.

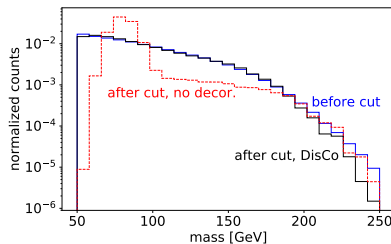
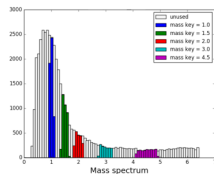
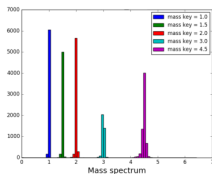
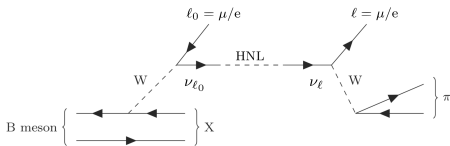


FIG. 4: QCD mass distribution before and after a cut on CNN plus DisCo (W -tagging) with signal efficiency of 50% and JSD $\sim 10^{-3}$.

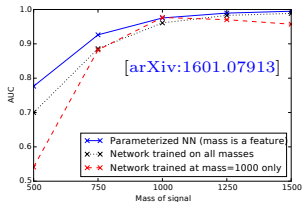
[arXiv:2001.05310]

- What if the signal is unknown?
E.g. searching for resonance at unknown mass
- Can use **parametric neural network** (pNN)
- Train a NN on a mixture of signals with an additional input feature: the true value of the signal parameter p
- For background examples:
 - If p not meaningful for background, use random value, following signal distribution
 - Else, *e.g.* p directly translates as a physics observable A , use $p = A$

pNN used at ETH in search for long-lived Heavy Neutral Leptons:

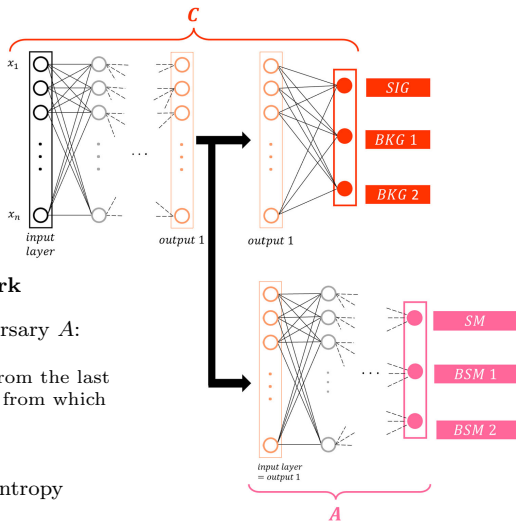


- pNNs achieve same performance on hypothesis p_0 as a single NN trained only on signal hypothesis p_0
- pNNs better interpolate between the different signal hypotheses than a plain NN trained on a mixture of signals



Training physics model decorrelated classifier (2)

- What if the model parameters to decorrelate against are not directly related to a physics observable?
- *E.g.* varying anomalous coupling of the Higgs boson
 - impacts many observables in a non-trivial way and signal shape for template fit!



- Can use **adversarial neural network**
- Train both a classifier C and an adversary A :
 - C classifies signal vs background
 - A takes the latent representation from the last hidden layer of C and tries to find from which process it comes:

$$L = L_C - \alpha L_A$$

where $L_{C/A}$ is the categorical cross-entropy
 α is an hyper-parameter

- If the adversary A cannot figure out from which process the event is, then the output of C is model-independent!

[Eur. Phys. J. C 82, 921 (2022)]

- ① Introducing inductive bias
- ② Background estimation in HEP
- ③ Training decorrelated models
- ④ Searching for new physics

An autoencoder (AE) is composed of:

- an encoder NN f
- a “symmetric” decoder NN g

The AE network is trained to learn to reconstruct the input examples it is given.

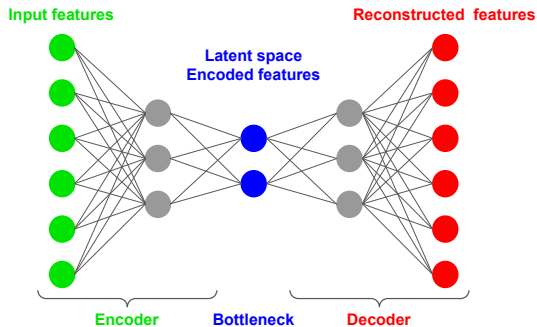
Loss for an example x :

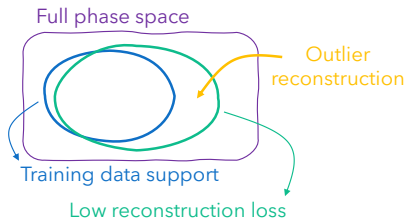
$$L(x) = \|g(f(x)) - x\|$$

where $\|\cdot\|$ is a distance

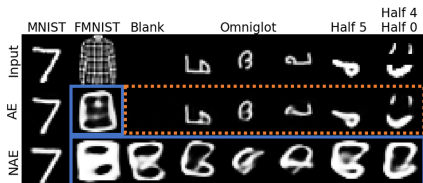
The aim of an AE for anomaly detection is to reconstruct with low error only the examples it is trained on but not others!

Search for new physics in HEP with AE is based on learning SM physics and flag new physics as anomalous!





- Outlier reconstruction happens when the network assigns low reconstruction error to out-of-distribution (OOD) examples
- OOD reconstruction not suppressed during training in plain AE
- Sometimes phrased as “OOD examples need to be more ‘complex’ to not be reconstructed”



Outlier reconstruction example: AE and NAE trained on MNIST, other inputs are outliers.

- **Normalized autoencoder¹(NAE)** features a mechanism to suppress OOD reconstruction!

¹NAE first introduced in [arXiv:2105.05735](https://arxiv.org/abs/2105.05735) and used in HEP in [arXiv:2206.14225](https://arxiv.org/abs/2206.14225)

- **Ensure that low reconstruction error phase-space matches that of training data**
- *i.e.* OOD examples are constrained to have high reconstruction error
- The model probability p_θ is defined from the reconstruction error E_θ via the Boltzmann distribution:

$$p_\theta(x) = \frac{1}{\Omega_\theta} \exp(-E_\theta(x))$$

- The loss is designed to learn $p_\theta = p_{\text{data}}$:

$$\mathbb{E}_{x \sim p_{\text{data}}} [L_\theta(x)] = \mathbb{E}_{x \sim p_{\text{data}}} [E_\theta(x)] - \mathbb{E}_{x' \sim p_\theta} [E_\theta(x')]$$

positive energy E_+ negative energy E_-

- Positive energy is the reconstruction error of the training examples
- Need to sample from the model to get the “negative samples” x' and compute E_-
 → Monte Carlo Markov Chain (MCMC) employed

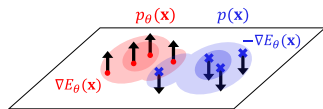


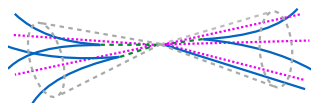
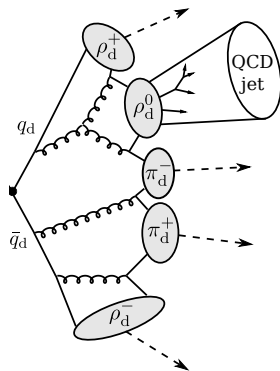
Figure 2. An illustration of the energy gradients in Eq. (7). The red and blue shades represent the model and the data density, respectively. The gradient update following Eq. (7) increases the energy of samples from $p_\theta(x)$ (the red dots) and decreases the energy of training data (the blue crosses).

→ Low energy examples have high probability

- Semi-visible jets (SVJ) are new physics signatures arising from theories where dark matter is made of dark quarks and a dark QCD force, very similar its SM counterpart
- Dark quarks hadronize to form dark hadrons, a fraction of which promptly decays to SM quarks which hadronize in the SM sector
- SVJs are jets made of visible SM hadrons with different substructure than SM QCD jets
- Currently developing NAE using substructure variables and a fully connected NN
- Loss function:

$$L = \log(\cosh(E_+ - E_-)) + \lambda_+ E_+^2$$

- First term to suppress OOD reconstruction
- Second term to learn training examples reconstruction



SM hadrons

Stable dark hadrons

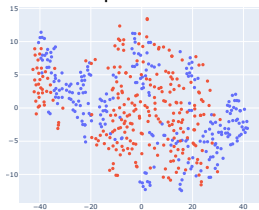
Visualizing positive and negative samples:

- Energy Mover's Distance (EMD)
- t-distributed Stochastic Neighbor Embedding (t-SNE) plots

→ Check suppression of OOD reco, e.g. “that the reco loss is high outside the training manifold”

→ Good anomaly detection: low reco error of training examples (SM physics) AND suppression of OOD (BSM physics) reco (low EMD, overlap in t-SNE plots)

Epoch 500

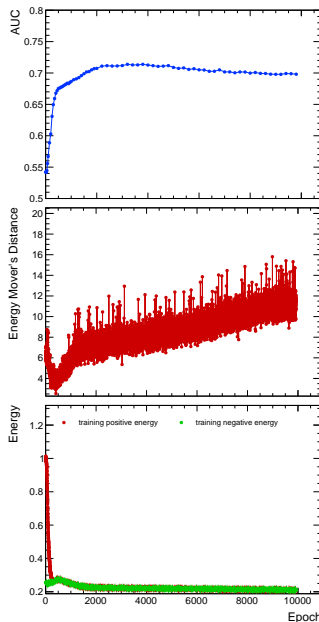


Epoch 10000



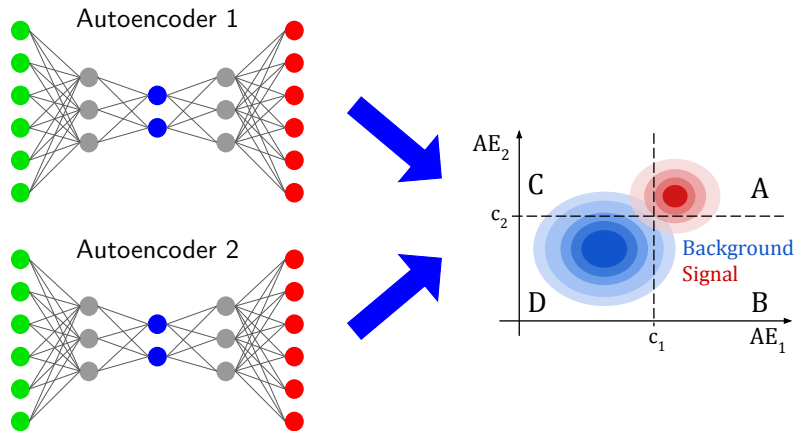
labels

- Positive samples
- Negative samples



The idea is¹:

- to train 2 autoencoders, decorrelated from each other using DisCo regularization
- such that the new physics enriched region is the high loss region of the AEs
- to perform ABCD background estimation using the losses of the two AEs



¹[arXiv:2111.06417](https://arxiv.org/abs/2111.06417)

- ML provides tools to address several HEP problems:
 - Background estimation
 - Building decorrelated classifiers with respect to signal hypotheses or a physics observable
- ... not mentioning building classifiers for jet tagging, searches or precision measurement!
- Many exciting developments to search for new physics with unsupervised learning!
- Still ongoing developments to incorporate physics knowledge into new ML models and improve interpretability

Backup

The **Pearson correlation** only evaluates **linear correlations**:

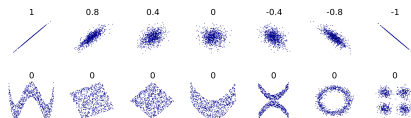
$$\rho_{\text{Pearson}}^2(X, Y) = \frac{\text{Cov}^2(X, Y)}{\text{Cov}(X, X)\text{Cov}(Y, Y)} \quad (3)$$

The **Distance correlation (DisCo)** makes use of all information of the random variables:

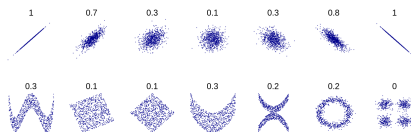
$$\text{dCov}^2(X, Y) = \int d^p s d^q t |f_{X,Y}(s, t) - f_X(s)f_Y(t)|^2 w(s, t)$$

where f_X (resp. f_Y) is the characteristic function of X (resp. Y), $f_{X,Y}$ is the joint characteristic function of X and Y .
 $f_{X,Y} = f_X f_Y$ iff X and Y are **independent**.

$$\text{DisCo}^2(X, Y) = \frac{\text{dCov}^2(X, Y)}{\text{dCov}(X, X)\text{dCov}(Y, Y)} \quad (4)$$



Pearson correlation coefficient



Distance correlation coefficient

Energy-based models (EMBs)

- EMBs are models where the probability is defined through the Boltzmann distribution
- Let θ denote the model parameters
- The model probability p_θ is defined from the energy E_θ

$$p_\theta(x) = \frac{1}{\Omega_\theta} \exp(-E_\theta(x)/T) \quad (5)$$

where the normalization constant Ω_θ is

$$\Omega_\theta = \int \exp(-E_\theta(x)/T) dx \quad (6)$$

- The EBM loss for a training example x is the negative log-likelihood:

$$L_\theta(x) = -\log p_\theta(x) = E_\theta(x)/T + \log \Omega_\theta \quad (7)$$

- The gradient of the EBM loss is thus:

$$\nabla_\theta L_\theta(x) = \nabla_\theta E_\theta(x) - \mathbb{E}_{x' \sim p_\theta} [\nabla_\theta E_\theta(x')] \quad (8)$$

- The expectation value over the training dataset, with probability p_{data} is:

$$\mathbb{E}_{x \sim p_{\text{data}}} [\nabla_\theta L_\theta(x)] = \mathbb{E}_{x \sim p_{\text{data}}} [\nabla_\theta E_\theta(x)] - \mathbb{E}_{x' \sim p_\theta} [\nabla_\theta E_\theta(x')] \quad (9)$$

Loss

$$\mathbb{E}_{x \sim p_{\text{data}}} [L_{\theta}(x)] = \mathbb{E}_{x \sim p_{\text{data}}} [E_{\theta}(x)] - \mathbb{E}_{x' \sim p_{\theta}} [E_{\theta}(x')] = E_{+} - E_{-}$$

positive energy negative energy

Positive energy

- Simply the reconstruction error over the training dataset
- Take SM jets and compute the reconstruction error!

Negative energy

- Reconstruction error of the “negative samples” x' from the probability distribution p_{θ}
 - Need to sample from the model to get the “negative samples”
- Monte Carlo Markov Chain (MCMC) employed

MCMC

- Start from an initial point x'_0
- Run n Langevin MCMC steps:

$$x'_{i+1} = x'_i - \lambda_i \nabla_x E_{\theta}(x'_i) + \sigma_i \epsilon \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (10)$$

drift diffusion

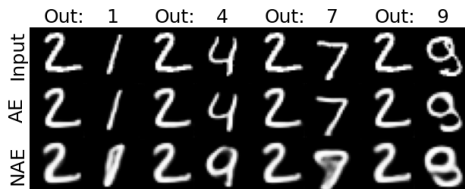
- Repeat with several points $x'^{(j)}$, the negative samples are the $x_n'^{(j)}$

Table 1. MNIST hold-out class detection AUC scores. The values in parentheses denote the standard error of mean after 10 training runs.

HOLD-OUT: 0	1	2	3	4	5	6	7	8	9	AVG	
NAE-OMI	.989 _(.002)	.919 _(.013)	.992 _(.001)	.949 _(.004)	.949 _(.005)	.978 _(.003)	.938 _(.004)	.975 _(.024)	.929 _(.004)	.934 _(.005)	.955
AE	.819	.131	.843	.734	.661	.755	.844	.542	.902	.537	.677

Signal	NAE	
	AUC	$\epsilon_B^{-1}(\epsilon_S = 0.2)$
top (AE)	0.875	68
top (NAE)	0.91	80
QCD (AE)	0.579	12
QCD (NAE)	0.89	350

AUC score for top tagging (2 first rows) and QCD tagging (2 last rows) for AE and NAE. The AE is a pre-training phase of the NAE.



Reconstruction examples in MNIST hold-out class detection for AE (middle row) and NAE (bottom row). Each pair of column is a different training for a different hold-out class.

- The NAE brings huge improvement compared to the plain AE on image classification task
- NAE achieves symmetric tagging, not only tagging of more complex objects!
- State-of-the-art anomaly detection on images