ANOMALY DETECTION AS A VALUABLE TOOL FOR UNCOVERING UNEXPECTED PHENOMENA AT COLLIDERS

LIP

Simão Silva Cardoso Supervisors: Nuno Castro, Rute Pedro e Miguel Peixoto



Standard Model of Elementary Particles



MOTIVATION

The Standard Model is the most successful theory of particle physics to date, but there are many open questions. Research on this topic can be assisted by using anomaly detection methods on colliders' data, finding possible signals that might hint at new physics.

SIMULATED DATASETS

Monte Carlo events were simulated for all processes in the background. The signals used are examples.



Resonant Dark Matter particles (S1)



Dark Matter production by two Higgs doublets (S2)

Heavy Vectorial Triplet (S4)

Gluino pair (S3)

Methods used

Supervised:

- Semi-supervised:

 - Autoencoder

31 features used - some examples: MET, 4-momentum jets and large jets, HT (scalar momentum sum),...

• Deep Neural Networks Isolation Forest (w/o contamination) Variational Autoencoder (w/o contamination)

SUPERVISED LEARNING - NEURAL NETWORKS

It's an algorithm whose architecture comprises layers with many neurons, leading to an output. Labels are used for training.

$$\mathbf{f}_{l}(\mathbf{Z}) = \mathbf{g}_{l}(\mathbf{W}_{l}\mathbf{Z} + \mathbf{b}_{l})$$

min_{w,b} $\frac{1}{N} \sum_{i}^{N} [y_{i} \log_{2}[NN(\mathbf{x}_{i}, \mathbf{W}, \mathbf{b})] +$

 $(1 - y_i)\log_2[1 - NN(\mathbf{x}_i, \mathbf{W}, \mathbf{b})]]$



Input

NEURAL NETWORKS -RESULTS





- - signals.
- - labelled.

WHY PREFERING SEMI-SUPERVISED TO **SUPERVISED METHODS?**

 Neural Network's performance might degrade when applied to other

• Real datasets aren't

Semi-supervised methods

are independent and offer

great performance.

SEMI-SUPERVISED LEARNING -ISOLATION FOREST

Isolation Forest creates tree structures that can represent recursive partitioning during learning. The path length of the trees averaged over a forest of such random trees is used to measure the anomaly score.



ISOLATION FOREST -RESULTS





ISOLATION FOREST - RESULTS (CONTAMINATED)





SEMI-SUPERVISED LEARNING -AUTOENCODER

The Autoencoder is a deep architecture that learns to compress (encode) and then decompress (decode) data through a bottleneck called the latent space.



 $loss = min_{W_{i}}$

$$\mathbf{b} \frac{1}{N} \sum_{i}^{N} \|AE(\mathbf{x}_{i}, \mathbf{W}, \mathbf{b}) - \mathbf{x}_{i}\|^{2}$$



AUTOENCODER -RESULTS



SEMI-SUPERVISED LEARNING -VARIATIONAL AUTOENCODER

The Variational Autoencoder (VAE) is an AE whose training is regularised to avoid overfitting and ensure that the latent space has suitable properties that enable the generative process.



$loss = MSE(\mathbf{x}, \hat{\mathbf{x}}) + \beta KL[N(\mu_x, \sigma_x), N(0, 1)]$

 $z = \mu_x + \sigma_x \odot \zeta$

VARIATIONAL AUTOENCODER -RESULTS

14

ROC curve - Variational Autoencoder

VARIATIONAL AUTOENCODER -RESULTS (CONTAMINATED)

CORRELATION BETWEEN METHODS

Correlation Between AD algorithms on S2

CORRELATION BETWEEN METHODS

Correlation Between AD algorithms on S4

AUC SCORES OF THE METHODS

CONCLUSION

Distinct algorithms for anomaly detection are highly effective in isolating diverse types of BSM physics. Furthermore, these algorithms can complement each other in unsupervised searches for new physics, making them potential tools in particle physics research.

