# **Unsupervised Anomaly Detection with CATHODE**

14th Annual Meeting of the Helmholtz Alliance "Physics at the Terascale"

Gregor Kasieczka, <u>Tobias Quadfasel</u>, Manuel Sommerhalder (Institute of Experimental Physics, University of Hamburg)

23.11.2021





Bundesministerium für Bildung und Forschung



**Motivation** 

Majority of searches for new physics rely heavily on both signal and SM background models

Impossible to cover all models/phase space regions with a dedicated search  $\rightarrow$  Need **model-independent** searches

Test scenario:

→Dijet anomaly search X→YZ, Y & Z decaying hadronically
 →Look for resonant new physics with anomalous jet substructure





- Benchmark: LHC Olympics 2020 challenge R&D dataset (arxiv:2101.08320)
- Background: 1M simulated QCD multijet events
- \* Signal: 100k W'→YZ events where Y→qq and Z→qq
- $*~m_{W'}=3.5\,{
  m TeV},\,m_{
  m Y}=500\,{
  m GeV},\,m_{
  m Z}=100\,{
  m GeV}$
- \* Input: 4 variables
  - Lower jet mass m<sub>j1</sub>
  - \*~ mass difference  ${igt \Delta} m_{1,2}$
  - \* Jet subjettiness ratios  $au_{21,j1}$  and  $au_{21,j2}$

#### Benchmark Dataset



arXiv:2001.04990

Given distributions of signal  $p_S(\mathbf{x})$  and background  $p_B(\mathbf{x})$  in some set of variables  $\mathbf{x}$ , Neyman-Pearson-Lemma:  $\rightarrow$  best test based on likelihood ratio

Problem: Signal is buried under large amount of background

 $\rightarrow$ We can't estimate  $p_{S}(\mathbf{x})$  directly

The best we can do: Estimate  $p_{S+B}(\mathbf{x}|SR)$  in "Signal Region" and  $p_B(\mathbf{x}|SB)$  from region without signal ("Sidebands")

 $\rightarrow$  Conditional variable containing resonance:  $m_{jj}$ 

We need to take LR in SR:  $\rightarrow$ Interpolate  $p_B(\mathbf{x}|\text{SB})$  into SR  $\rightarrow$ Construct estimate of LR:  $\frac{p_{S+B}(\mathbf{x}|\text{SR})}{p_B(\mathbf{x}|\text{SR})}$ 

#### **General Principle**



Get estimate of LR using: →**classification** →**density estimation** (Normalizing Flows)

# » Normalizing Flows

Introduction

\* Flows based on random variable transformation

\* 
$$f: U \to X; \ p(x) = p(u) \left| \frac{df(u)}{du} \right|$$

- \* Learn invertible mapping f from latent variables  $m{u}$  to data  $m{x}$
- \* Flow: stack many invertible transformations  $f_i$ :  $f=f_k\circ ...\circ f_2\circ f_1$

$$p(\mathbf{x}) = p(f^{-1}(\mathbf{x})) \prod_{i} \left| \det \left( \frac{\partial f_{i}^{-1}}{\partial \mathbf{x}} \right) \right| = p(\mathbf{u}) \prod_{i} \left| \det \left( \frac{\partial f_{i}}{\partial \mathbf{u}} \right) \right|^{-1}$$



#### **Classification Without Labels (CWoLa)**



- + Simple classification task
  - Basic DNN architecture
  - Highly dependent on correlations between **x** and *m* →Variables **x** need to be hand-picked



 Direct estimation of conditional densities

**Classification & Density Estimation** 

**Anomaly Detection with Density** 

**Estimation (ANODE)** 

- + Easy to interpolate p<sub>bg</sub>
- Robust against correlations between **x** and *m*
  - →Arbitrary choice of **x**
  - Computationally intense
    - Estimation of signal contribution difficult

#### DOI:10.1007/JHEP10(2017)174

#### DOI:10.1103/PhysRevD.101.075042

[5/9]

#### Classifying Anomalies THrough Outer Density Estimation (CATHODE)

- + Only one density estimator needed
- + Due to interpolation: robust against correlations between *x* and *m* →Arbitrary choice of *x*
- + Final tagger based on simple classification task
- + No density estimator for signal contribution needed
- Computationally intense

#### arxiv:2109.00546



CATHODE

[6/9]

# » CATHODE

- Most important performance measure: significance improvement characteristic (SIC)
- \* "Supervised" training using signal/background labels
   →overall upper performance limit
- \* "Idealized AD": distinguish actual sample vs. background-only sample from signal region →upper limit for unsupervised anomaly search
- \* CATHODE shows highest SIC amongst non-idealized anomaly taggers
- \* Performance reaches idealized AD limit
- Significance improvement about factor 14



arxiv:2109.00546

#### » CATHODE

#### Robustness Against Correlations

- Study impact of correlations between **x** and m<sub>jj</sub>
- Introduce artificial correlations
- \* Add 10% of corresponding  $m_{jj}$  value to  $m_{j1}$  and  $\Delta m_{1,2}$
- \* All methods suffer performance
   →"Smearing" of variables
- \* CWoLa performance completely breaks down
- CATHODE retains good performance, similar to idealized AD



arxiv:2109.00546

#### » Summary

- Investigation of dijet resonances with anomalous jet substructure using normalizing flow-based density estimation & classification
- Introduction of new method: CATHODE that combines the advantages of purely density estimation-based (ANODE) and classification-based (CWoLa) approaches
- \* CATHODE outperforms all other non-idealized anomaly detectors
- \* Performance similar to idealized anomaly detector
- $\ast\,$  Robust against correlations between features **x** and conditional variable  $m_{jj}$
- \* Future studies
  - \* Other datasets/topologies
  - \* Study sensitivity for different types of anomalies (e.g. very broad resonances)
  - \* Studies using more (low-level) features



- Comparison for different amounts of signal injected
- CATHODE outperforms other AD methods significantly down to a S/B as low as 0.3%
- \* CATHODE achieves similar performance as idealized anomaly detector
- Below 0.2% S/B: even idealized AD cannot raise significance above 3σ
   →Too limited number of data points in the signal region

#### S/√B 2.14 1.35 1.02 0 68 0 51 0 34 0 17 17.5 15.0 12.5 10.0 7.5 \_\_\_\_\_ and the second 2.5 CATHODE CWoLa 0.0 Supervised ANODE Idealized AD -2.5

0.30 S/B (%)

0.40

arxiv:2109.00546

Achieved Significance

Maximum

0.60

Performance

0.20 0.15 0.10 0.05

#### » CATHODE

- \* Flow trained for 100 epochs
- Model ensembling: pick 10 epochs with lowest validation loss
- \* Draw  $m_{jj}$  values from a KDE in signal region
- \* Use these values as conditional and sample from density estimator
   →Interpolation into SR
- \* We can oversample to produce more samples than we have in data
- Background densities are modelled well by flow



# » Normalizing Flows

#### Introduction

\* Flows based on random variable transformation

\* 
$$f: U \to X; \ p(x) = p(u) \left| \frac{df(u)}{du} \right|$$

- \* Multivariate case: scale with Jacobian determinant
- \* Flow: repeat (invertible) transformation to get complex distributions

$$p(\mathbf{x}) = p(f^{-1}(\mathbf{x})) \prod_{i} \left| \det \left( \frac{\partial f_{i}^{-1}}{\partial \mathbf{x}} \right) \right| = p(\mathbf{u}) \prod_{i} \left| \det \left( \frac{\partial f_{i}}{\partial \mathbf{u}} \right) \right|^{-1}$$



#### CASE (CMS Anomaly SEarch):

Explore different techniques, all involving ML models trained directly on data!

Currently investigated methods:

- \* (Variational) autoencoders
- Weak supervision (CWoLa hunting, Tag N' Train)
- Quasi-anomalous knowledge (QUAK)
- Density estimation & classification (CATHODE)

General procedure:





# » Autoregressive Flows

Autoregressive property:

 $p(\mathbf{x}) = \prod_i p(\mathbf{x}_i | \mathbf{x}_{1:i-1})$ 

Conditional densities depend on trainable parameters:

$$p(\mathbf{x}_i | \mathbf{x}_{1:i-1}) = \mathcal{N}(\mathbf{x}_i | \mu_i, (\exp \alpha_i)^2)$$
$$\mu_i = f_{\mu_i}(\mathbf{x}_{1:i-1})$$
$$\alpha_i = f_{\alpha_i}(\mathbf{x}_{1:i-1})$$

→Earlier variables must not depend on later variables

 $\rightarrow$ Solution: stack transformations into a normalizing flow, change ordering of the  $x_i$  after each transformation



Autoregressive property →Jacobian is upper triangular

$$\left|\det\left(\frac{\partial f}{\partial \boldsymbol{u}}\right)\right| = \exp\left(\sum_{i} \alpha_{i}\right)$$
[5/7]

# » Masked Autoencoder for Distribution Estimation (MADE)

- DNN architecture to implement a single
   *f<sub>i</sub>* in autoregressive flows (1502.03509)
- \* Compute lpha and  $\mu$  in one forward pass
- \* Outputs  $\alpha_j$  and  $\mu_j$  only connected to inputs  $\{x_1, ..., x_{j-1}\} \rightarrow$ autoregressive property
- No connection dropped for conditional input m<sub>jj</sub>



Architecture

# » Anomaly Detection with Density Estimation (ANODE)

#### Architecture

- Stack MADE networks to build "Masked Autoregressive Flow" (MAF)
- \* Learn tranformation  $\mathbf{u} = f^{-1}(\mathbf{x})$  from input features  $\mathbf{x}$  to  $\mathbf{u} \sim \mathcal{N}(0, \mathbb{I})$
- Compute p(x) with normalizing flow from p(u)
- \* Minimize NLL loss  $\mathcal{L} = -log(p(\textbf{\textit{x}}))$
- Architecture used for both density estimators

