

## Binning high-dimensional classifier output for HEP analyses through a clustering algorithm

9th May 2023



Deutsche

Forschungsgemeinschaft



Bundesministerium für Bildung und Forschung

GEFÖRDERT VOM

Svenja Diekmann, Niclas Eich, Martin Erdmann

CHEP 2023





## Introduction

Analysis Context:

search for VH at







### **Gluon Fusion**





## **VH-Analysis Process-classification**

- Events are classified by process
- Highest score serves as fit-variable









# Clustering with K-Means [1/2]

- The DNN delivers a *Nd*-representation of our data
- Correlated to **physics processes**, which will be fitted
- Clustering can create bins in this Nd space without suffering under the curse of dimensionality









## Clustering with K-Means [2/2]

### K-Means Clustering Algorithm:

- random initialisation of K cluster centers
- eventwise assignment to closest cluster (L2 norm)
- iterative update of cluster centers by mean of assigned events
- K-Means already used in this analysis [1]





CHEP 2023



## Approach Comparison

### Standard Binning

- N Histograms
- Each has *m* bins
- Need to be tuned independently





### **Clustering Binning**

- 1 Histogram
- K bins
- Number of tuning parameters very limited





## **Clustered distribution**

Clusters sorted by absolute signal yield		
All bins contain a reasonable amount of		-
events		106-
Some clusters are enriched		$10^{5}$
in specific backgrounds		104-
High purity bins towards	ries	10 <sup>3</sup>
right side	Intr	
		$10^{2}$
		101
		$10^{0} =$

 $10^{-1}$ 

 $\frac{s}{\sqrt{B}} \frac{10^{-1}}{10^{-1}}$ 

 $10^{-2}$ 









## Sensitivity Improvement

### **Benchmark results:**

- evaluation of 95% C.L. limit
  - sensitivity increase for high
  - number of clusters
  - improvement compared to
    - the standard method (~ 400 bins):
    - statistical uncertainties
    - background normalisation

2.00 1.75 (normalised) .50 .25 Limit 1.00 0.75 C.L 0.50 0.50 0.50 0.25

 $0.00^{L}_{0}$ 











## Bias test - MC dependency

Clusters are determined on MC events → Might introduce bias!

Testing-Strategy:

- Dataset split into two equal sets
  - Training
  - Test
- Clustering performed on training set
- Evaluated on sets separately
- Almost no difference
- → Clustering is robust towards MC dataset introduced bias







## Top 25 clusters

•	Clusters a	e determ	nined in	8d-space
---	------------	----------	----------	----------

<b>→</b>	Cluster	center	coordinate	is a	a 8d	vector
----------	---------	--------	------------	------	------	--------

 $10^{6}$ 

Try to get some insight how the clusters are		10 <sup>5</sup> -
constructed!	S	104-
	Entrie	10 <sup>3</sup> -
		10 <sup>2</sup> -
		101-
		10 <sup>0</sup> -
		10-1-
	ς	10 <sup>0</sup> -
	$\frac{S}{\sqrt{B}}$	$10^{-1}$
		$10^{-2}$









## **Cluster Visualisation - Top 1 Cluster**

Visualisation: 8d cluster center coordinate plotted as Radar-Plot

- #1 cluster has the largest signal contribution
- Main coordinate consist of backgrounds
  - $-t\overline{t}$
  - -VV(V)









## Cluster Visualisation - Top 4/5 Cluster

- High signal prediction confidence
- Better signal-background ratio  ${ } \bullet$









# Summary

Binning High-Dimensional Classifier Output through Clustering:

- **Agnostic** binning algorithm in high-dimensional space
- No loss of information
- Improvement of analysis sensitivity
- Radar plots visualise cluster center coordinates
- Robust towards MC training dataset









## Resources

[1] Evidence for associated production of a Higgs boson with a top quark pair in final states with electrons, muons, and hadronically decaying  $\tau$  leptons at  $\sqrt{s}$  = 13 TeV, CMS Collaboration, JHEP 08 (2018) 066,





