

# Robustness and interpretability of machine learning methods applied to LHC data

Steffen Mæland

Western Norway University of Applied Sciences



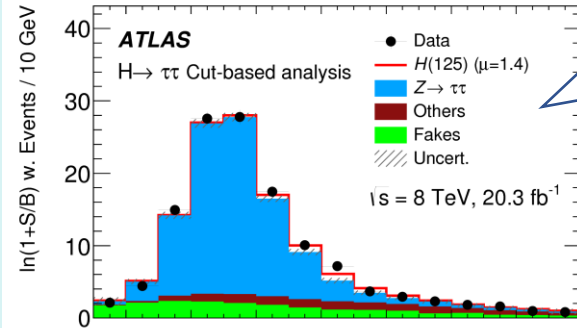
Høgskulen  
på Vestlandet



NORCC

# ML use case: Event classification

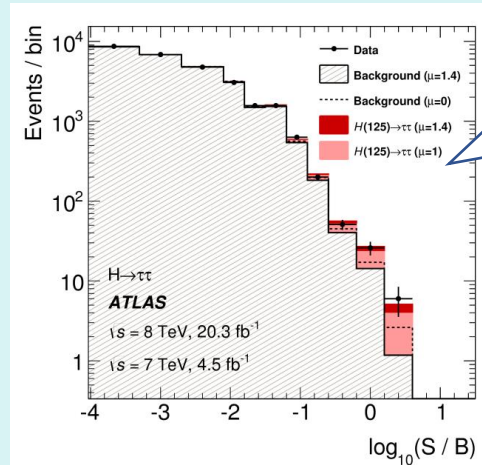
Event  $\longrightarrow$  Observables ( $p, E, \dots$ )  $\longrightarrow$



Carefully  
constructed  
observable

$\longrightarrow$  Fit  $\longrightarrow$  Result

Event  $\longrightarrow$  Observables ( $p, E, \dots$ )  $\longrightarrow$



Machine  
learning

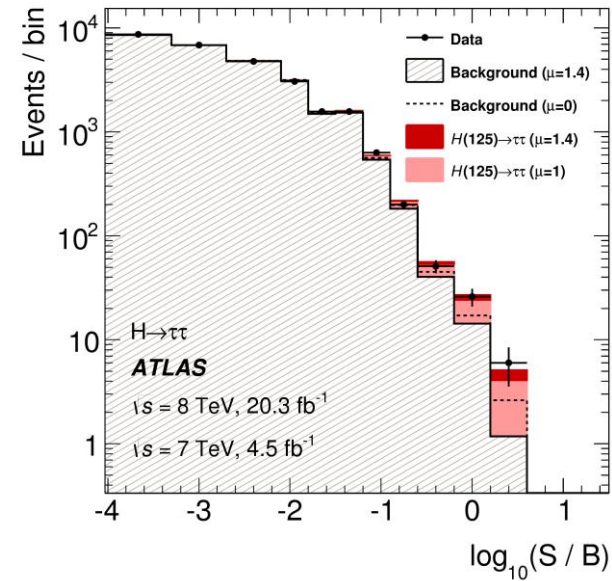
$\longrightarrow$  Fit  $\longrightarrow$  Potentially  
better result



# ML use case: Event classification

Since the transformation

Observables ( $p, E, \dots$ )  $\longrightarrow$

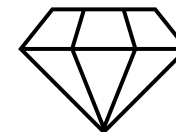
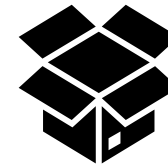
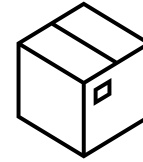


is (highly) non-trivial, verification is typically not straight-forward



# Terms

- *Interpretability:*  
Passive characteristic of a model – to what extent it is understandable by humans
- *Explainability:*  
Active characteristic, involving methods that clarify a model's decision process or internal function
- *Robustness:*  
To what extent a model's prediction is affected by perturbations in the input data





# Explainability approaches

Assume we have an ML model  $f$  with a bunch of parameters  $\theta$ , taking in data  $\mathbf{x}$  and returning predictions  $\mathbf{y}$ :

$$\mathbf{y} = f(\mathbf{x}, \theta)$$

Three options for explaining  $\mathbf{y}$ :

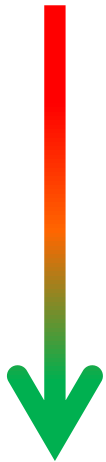
- Replace  $f$  by an interpretable model: *Surrogate model explanations*
- Vary  $\mathbf{x}$  and observe the effect: *Extrinsic explanations*
- Study  $\theta$  (for some given  $\mathbf{x}$ ): *Intrinsic explanations*

# Extrinsic explanations

Model agnostic – study  $y$  for different  $x$

→ Feature importance explanation  
(either on average or for single events)

Theoretical foundation




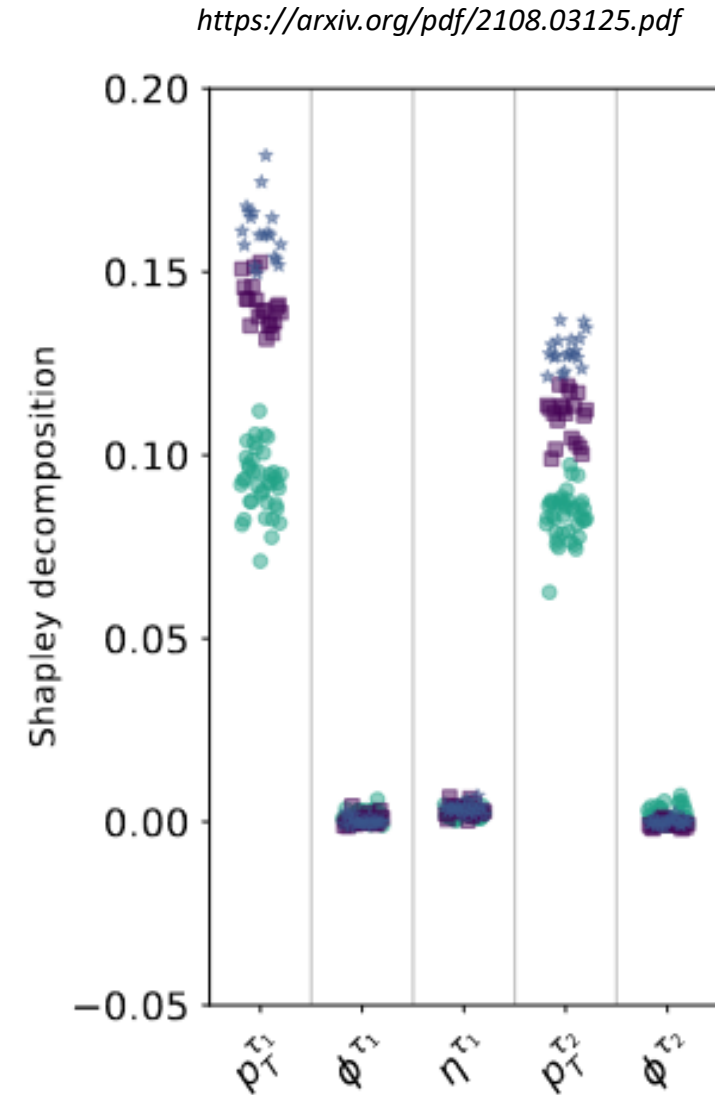
Computational effort



Assumptions

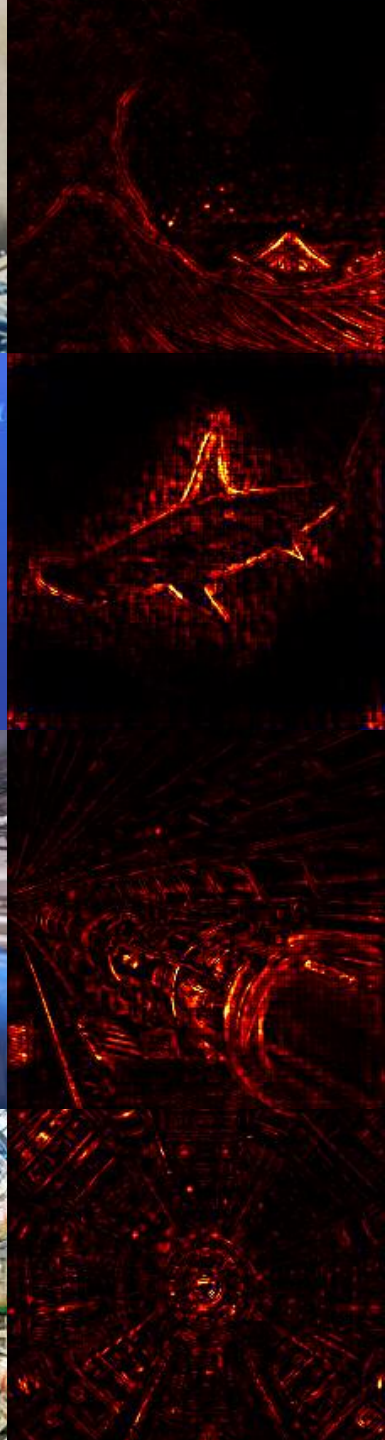
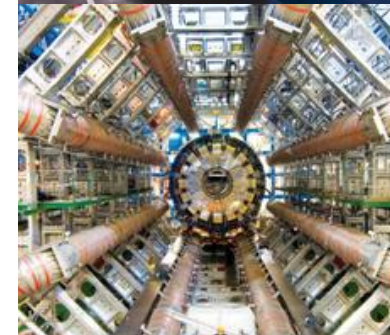


- Randomise feature values  
<https://arxiv.org/abs/1801.01489>
- SHAP   
<https://arxiv.org/abs/1705.07874>
- Shapley values  
<https://arxiv.org/abs/1705.07874>



# Intrinsic explanations

- Model-specific – requires all model parameters
- Use gradients to quantify how a change in input would change the prediction (per event)
- Can be combined with randomisation of feature values





# Robustness

- Again using our ML model  $f$  and a test data point  $\mathbf{x}$ , how robust is the prediction  $\mathbf{y}$  to a perturbation  $\mathbf{x}'$  in the input?  
i.e. is

$$\mathbf{y} = f(\mathbf{x}) \quad \text{equal to} \quad \mathbf{y}' = f(\mathbf{x} + \mathbf{x}')?$$

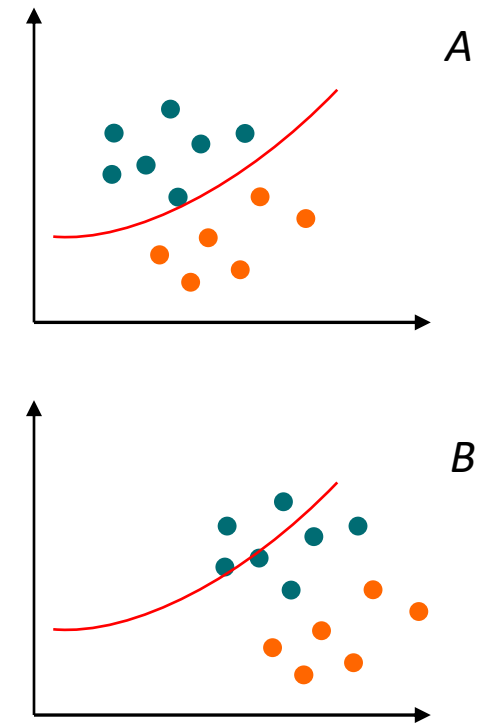
- Types of perturbations:
  - Random noise (  $E(\mathbf{x}') = 0$  )
  - Distribution shifts (  $E(\mathbf{x}') \neq 0$  )
  - Adversarial (  $\mathbf{x}'$  selected so that  $\mathbf{y} \neq \mathbf{y}'$  )





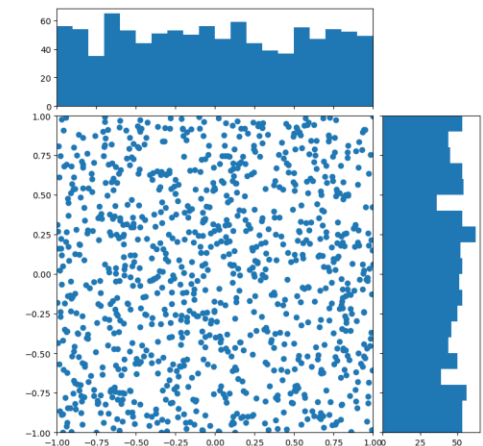
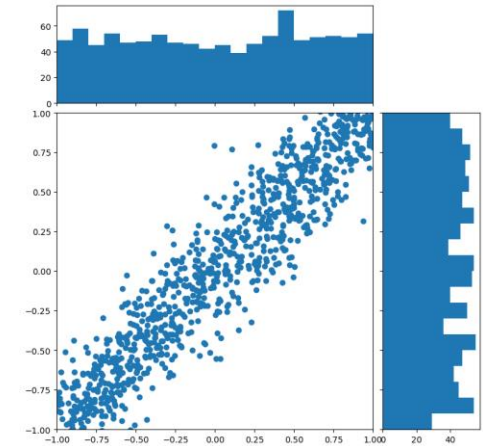
# Robustness under distribution shifts

- Under distribution shifts, feature correlations remain but numerical values are consistently shifted
  - ML methods typically not happy about this
- Mitigated by *domain adaptation*
  - Methods applicable to analysis re-interpretation  
<https://arxiv.org/abs/2207.09293>



# Robustness under random noise

- Can be improved through data augmentation (randomly sampling  $\mathbf{x}'$  during training)
  - Requires augmentation to be realistic
- Common measures of robustness rely on the same sampling, at different noise levels
  - Estimate is only as good as the sampled values
  - Sampling from the marginal distribution leads to unlikely data points if features are correlated
- Realistic sampling gives
  - Data augmentation ✓
  - Feature relevance estimate ✓
  - Robustness estimate (at statistical level) ✓



# Robustness to adversarial examples

- Prediction accuracy will vary in different regions of feature space
- Adversarial attacks exploit this to find and insert the smallest  $x'$  that will change the prediction
- Won't see this in HEP data, *but* method is useful for identifying regions of low robustness

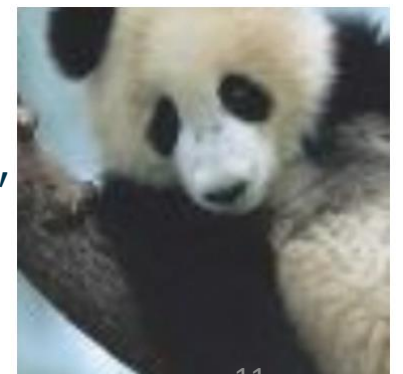
$x$   
 $y = \text{«panda»}$



$x'$



$x + 0.007 \times x'$   
 $y = \text{«gibbon»}$





# Our projects

- Realistic data augmentation for improved robustness

## *Framework for*

- Data augmentation (improves also generalisation)
- Adversarial testing (model diagnostics and verification)
  - Mostly NN specific
- Develop suitable robustness score for HEP ML

