# Inference-Aware Neural Optimisation (INFERNO)

**ML-INFN** Meeting

Pablo de Castro, Tommaso Dorigo, Lukas Layer, Giles Strong









# **Problem Statement**

Our data is high dimensional. How do we choose the best summary statistics for the inference of the parameter of interest?

**Classical approach in HEP** 

- Signal vs. background classification with simulated samples by minimizing Binary Cross-Entropy of Neural Networks to approximate the optimal Likelihood Ratio
- Inference with histograms by building binned non-parametric
   Poisson counts likelihood



# **Problem Statement**

### Our data is high dimensional. How do we choose the best summary statistics for the inference of the parameter of interest?

2.0

0.0+-20

30

40

50

s parameter of interest

60

70

80

no nuisance effect

with nuisance effect

- Problem: simulations are imperfect → lack of knowledge for inference accounted by additional nuisance parameters
- additional nuisance parameters
  Neglecting nuisance parameters in the training process leads to larger measurement uncertainties
- Upper limit of ML usefulness in LHC analyses



# INFERNO

### Overview

- Inference-Aware Neural Optimization developed by Pablo de Castro and Tommaso Dorigo: P. de Castro Manzano, T. Dorigo, "INFERNO: Inference-Aware Neural Optimization", Comp. Phys. Commun. 244 (2019) 170; <u>https://doi.org/10.1016/j.cpc.2019.06.007</u>
- Objective: optimize statistical inference in the presence of nuisance parameters by learning non-linear summary statistics with neural networks via minimization of an inferencemotivated loss
- Code written in TF 1.x (<u>https://github.com/pablodecm/paper-inferno</u>) → key elements reproduced in TF 2.x (<u>https://github.com/llayer/inferno</u>)

#### Summary statistics with INFERNO

- "Standard HEP summary statistics" histograms are not differentiable → softmax as differentiable approximation
- Nuisance parameters approximated by changes of mixture *ŝ* coefficients, translations of a subset of features, or conditional density ratio re-weighting → restricted to likelihood-free inference problems



$$\hat{s}_i(D; \boldsymbol{\phi}) = \sum_{x \in D} \frac{e^{f_i(\boldsymbol{x}; \boldsymbol{\phi})/\tau}}{\sum_{j=0}^b e^{f_j(\boldsymbol{x}; \boldsymbol{\phi})/\tau}}$$



#### Algorithm 1 Inference-Aware Neural Optimisation.

Input 1: differentiable simulator or variational approximation  $g(\theta)$ . Input 2: initial parameter values  $\theta_s$ . Input 3: parameter of interest  $\omega_0 = \theta_k$ . Output: learned summary statistic  $s(D; \phi)$ .



#### Algorithm 1 Inference-Aware Neural Optimisation.

Input 1: differentiable simulator or variational approximation  $g(\theta)$ .

Input 2: initial parameter values  $\theta_s$ .

*Input 3:* parameter of interest  $\omega_0 = \theta_k$ .

*Output:* learned summary statistic  $s(D; \phi)$ .

1: **for** i = 1 to *N* **do** 

2: Sample a representative mini-batch  $G_s$  from  $g(\boldsymbol{\theta}_s)$ .



#### Algorithm 1 Inference-Aware Neural Optimisation.

Input 1: differentiable simulator or variational approximation  $g(\theta)$ .

Input 2: initial parameter values  $\theta_s$ .

*Input 3*: parameter of interest  $\omega_0 = \theta_k$ .

*Output:* learned summary statistic  $s(D; \phi)$ .

- 1: **for** i = 1 to *N* **do**
- 2: Sample a representative mini-batch  $G_s$  from  $g(\boldsymbol{\theta}_s)$ .
- 3: Compute differentiable summary statistic  $\hat{s}(G_s; \phi)$ .



#### Algorithm 1 Inference-Aware Neural Optimisation.

Input 1: differentiable simulator or variational approximation  $g(\theta)$ .

Input 2: initial parameter values  $\theta_s$ .

*Input 3:* parameter of interest  $\omega_0 = \theta_k$ .

*Output:* learned summary statistic  $s(D; \phi)$ .

- 1: **for** i = 1 to *N* **do**
- 2: Sample a representative mini-batch  $G_s$  from  $g(\boldsymbol{\theta}_s)$ .
- 3: Compute differentiable summary statistic  $\hat{s}(G_s; \phi)$ .
- 4: Construct Asimov likelihood  $\mathcal{L}_A(\boldsymbol{\theta}, \boldsymbol{\phi})$ .
- 5: Get information matrix inverse  $I(\boldsymbol{\theta})^{-1} = \boldsymbol{H}_{\boldsymbol{\theta}}^{-1}(\log \mathcal{L}_A(\boldsymbol{\theta}, \boldsymbol{\phi})).$
- 6: Obtain loss  $U = I_{kk}^{-1}(\boldsymbol{\theta}_s)$ .



Algorithm 1 Inference-Aware Neural Optimisation.

Input 1: differentiable simulator or variational approximation  $g(\theta)$ .

Input 2: initial parameter values  $\theta_s$ .

*Input 3:* parameter of interest  $\omega_0 = \theta_k$ .

*Output:* learned summary statistic  $s(D; \phi)$ .

- 1: **for** i = 1 to *N* **do**
- 2: Sample a representative mini-batch  $G_s$  from  $g(\boldsymbol{\theta}_s)$ .
- 3: Compute differentiable summary statistic  $\hat{s}(G_s; \phi)$ .
- 4: Construct Asimov likelihood  $\mathcal{L}_A(\boldsymbol{\theta}, \boldsymbol{\phi})$ .
- 5: Get information matrix inverse  $I(\boldsymbol{\theta})^{-1} = \boldsymbol{H}_{\boldsymbol{\theta}}^{-1}(\log \mathcal{L}_A(\boldsymbol{\theta}, \boldsymbol{\phi})).$
- 6: Obtain loss  $U = I_{kk}^{-1}(\boldsymbol{\theta}_s)$ .
- 7: Update network parameters  $\phi \to \text{SGD}(\nabla_{\phi} U)$ .
- 8: **end for**

## **Application to 3D Synthetic Mixture model**

୫

0

0

, A

°, 0, , ∕,

3.r

2.0

0.00

%

A

0

 $x_1$ 

Taken from https://arxiv.org/pdf/1806.04743.pdf

D

ზ

ծ

*x*<sub>0</sub>

 $\mathbf{x}_{1}$ 

#### Model



- 3D Mixture model: 2D Multivariate
  Normal and 1D Exponential
- Signal distribution is fully specified
- Background distribution depends on the parameters *r* (shift of the 2D Normal mean) and λ (rate of the exponential)

0.° ~ 2.° 3. ~

 $X_2$ 

## Results for the synthetic model



Taken from https://arxiv.org/pdf/1806.04743.pdf

Comparison of inference with INFERNO summary statistics and classical binary cross-entropy summary statistics

- **Confidence intervals** considerably **narrower** with INFERNO than those using BCE and closer to those expected when using the true model
- Improvement over standard classification increases with more nuisance parameters
- The inference-aware summary statistics learnt for θ<sub>s</sub> work well when θ<sub>true</sub> /= θ<sub>s</sub> in the range of variation explored.

# Benchmarks

- Several inference problems are considered based on the 3D benchmark mixture model
- Systematic comparison, shows that **INFERNO** clearly **outperforms any classifier** (even optimal Bayes) when nuisance parameters are relevant

Table 1: Expected uncertainty on the parameter of interest s for each of the inference benchmarks considered using a cross-entropy trained neural network model, INFERNO customised for each problem and the optimal classifier and likelihood based resuls.

	Benchmark 0	Benchmark 1	Benchmark 2	Benchmark 3	Benchmark 4
NN classifier	$14.99^{+0.02}_{-0.00}$	$18.94_{-0.05}^{+0.11}$	$23.94^{+0.52}_{-0.17}$	$21.54_{-0.05}^{+0.27}$	$26.71_{-0.11}^{+0.56}$
INFERNO 0	$15.51\substack{+0.09 \\ -0.02}$	$18.34_{-0.51}^{+5.17}$	$23.24_{-1.22}^{+6.54}$	$21.38^{+3.15}_{-0.69}$	$26.38^{+7.63}_{-1.36}$
INFERNO 1	$15.80^{+0.14}_{-0.04}$	$16.79\substack{+0.17 \\ -0.05}$	$21.41^{+2.00}_{-0.53}$	$20.29^{+1.20}_{-0.39}$	$24.26^{+2.35}_{-0.71}$
INFERNO 2	$15.71_{-0.04}^{+0.15}$	$16.87^{+0.19}_{-0.06}$	$16.95\substack{+0.18 \\ -0.04}$	$16.88^{+0.17}_{-0.03}$	$18.67\substack{+0.25\\-0.05}$
INFERNO 3	$15.70^{+0.21}_{-0.04}$	$16.91\substack{+0.20\\-0.05}$	$16.97\substack{+0.21\\-0.04}$	$16.89\substack{+0.18 \\ -0.03}$	$18.69^{+0.27}_{-0.04}$
INFERNO 4	$15.71_{-0.06}^{+0.32}$	$16.89^{+0.30}_{-0.07}$	$16.95\substack{+0.38\\-0.05}$	$16.88^{+0.40}_{-0.05}$	$18.68\substack{+0.58\\-0.07}$
Optimal classifier	14.97	19.12	24.93	22.13	27.98
Analytical likelihood	14.71	15.52	15.65	15.62	16.89

Taken from https://arxiv.org/pdf/1806.04743.pdf

# Future of INFERNO

### Applications

- Application of INFERNO to a real world CMS analysis
- Optimizations of detector designs
- Development of a common package across experiments



- Training and optimization needs powerful GPUs
- Data preprocessing and ML will benefit from Spark / Kubernetes clusters
- Ideally move entirely to cloud computing

# Summary

- INFERNO shows promising performance on toy problems
- Several advanced applications of INFERNO are planned / have started
- Need for powerful hardware

BACKUP

## Top pair cross section with CMS open data

#### TOP-11-004 - Top pair cross section in tau+jets $\sigma(tt) = 152 \pm 12$ (stat.) $\pm 32$ (syst.) $\pm 3$ (lum.) pb Events / 0.1 10<sup>4</sup> CMS vs=7 TeV, 3.9 fb Data stat.+syst. uncertainty tt τ<sub>+</sub>+iets single top W/Z + iets 10<sup>3</sup> background 00000 10<sup>2</sup> 10 Data/MC 1.5 0.5 9<u>5</u> -0.5 0.5 1.5 iet D<sub>NN</sub> https://arxiv.org/abs/1301.5755

- Replication of full TOP-11-004 with CMS open data including systematics
- Relevant 2011 legacy samples available for Data and MC
- Analysis dominated by systematic uncertainties and uses a ANN with TMVA
- Simple analysis, but can be made more challenging by analyzing more decay channels

# Status: implemented basic framework to analyze CMS open data with full systematic uncertainties

## Recent related work on differentiable analysis

### Work by S. Wunsch et. al (CMS)

- Optimal statistical inference in the presence of systematic uncertainties using neural network optimization based on binned Poisson likelihoods with nuisance parameters (https://arxiv.org/ abs/2003.07186)
- Idea: construct binned Poisson Likelihood by using histograms and approximate the gradient of the bin function by the derivative of a Gaussian → avoids the softmax approximation

### Work by N. Simpson, L. Heinrich et. al (ATLAS)

 Idea: optimize neural networks with differentiable profile likelihood in the loss function via fixed-point differentiation → allows to directly optimize CLs limits



https://github.com/pyhf/neos

Library *neos* based on autodiff library **JAX** and library **FAX** for differentiating fixed point problems in JAX

### HSF activity for differentiable analysis

- New HSF activity for differential analysis related to various HEP problems across experiments
- Discussion channel for the community on gitter