### Patrick L.S. Connor (IEXP & CDCS)

CDCS CENTER FOR DATA AND COMPUTING IN NATURAL SCIENCES



CMS Experiment at the LHC, CERN Data recorded: 2016-Sep-27 14:40:45.336640 GMT Run / Event / LS: 281707 / 1353407816 / 851

# CDCS Opening Symposium, 27-April, 2022 Precision measurements in High Energy Physics: correlated and uncorrelated uncertainties







## **Motivation (1)**

- Kepler & Brahe: planets' orbits are ellipses
- Newton: perihelion precession
- Le Verrier: residual discrepancy of Mercury's orbit of 42.980±0.001" per century
- **Einstein**: tiny (but significant!) gap between experimental observations and prediction explained by General Relativity



### Uncertainties tell you how much you can trust your data

## **Motivation (2)**

CDCS CENTER FOR DATA AND COMPUTING IN NATURAL SCIENCES



p<sub>Tiet</sub> [GeV]

### Uncertainties of such spectra require a sophisticated description.

- Bin-to-bin fully correlated, partially correlated, and fully decorrelated from bin to bin.
- Typically O(10) different sources of systematic uncertainties.
- Direct impact on extraction of fundamental physics quantities such as the  $\alpha_s$ .

CDL1

## **Goal & Outline**





CDI 1

Discuss tools for more or less broad application:

- 1) DAS: Das Analysis System
  - Jet analysis in context of the CMS Collaboration

2) RAN: Refinement Adversarial Networks

- By-pass computing-demanding simulation of detector with Geant4
- **3) Step**: Smoothness Tests based on Expansion of Polynomials
  - Investigate the quality of the statistical description of a differential measurement
- 4) Teddy: Treat Entries ouf of the Diagonal DecentlY
  - Extract high-level distributions from histograms with complex uncertainty schemes

## **Typical analysis strategy**

#### CDCS CENTER FOR DATA AND COMPUTING IN NATURAL SCIENCES



CDL1

### **DAS** (with CMS colleagues)

- Optimised for **debugging**
  - Event loop based

CDL1

- **Factorising** all steps
- Accounts natively for systematic variations at event level
  - •• Vectors of weights and factors
- Currently dedicated to CMS jet analysis exclusively
  - Modular design makes it easily extendable...
  - Essentially depends on demand and person power

Patrick L.S. Connor

#### CDCS CENTER FOR IN NATURA







### **RAN** (with TUHH and IEXP colleagues)

#### CDCS CENTER FOR DATA AND COMPUTING IN NATURAL SCIENCES

- Replace full simulation based on Geant4 by fast simulation based on higherlevel phenomenological models (e.g. Delphes)
- **Refine** obtained simulation with adversarial networks
- Interface with DAS ...?



### Promising approach to improve the model dependence of the data reduction!

### CDL1

## **STEP (1)** (with Radek Zlebcik)

• Typical jet measurements span over several orders of magnitude.

 $\frac{\mathrm{d}^2\sigma}{\mathrm{d}p_{\mathrm{T}}\,\mathrm{d}y} = \frac{1}{\mathcal{L}} \frac{N_{\mathrm{jets}}^{\mathrm{eff}}}{\Delta p_{\mathrm{T}}\,\Delta y}$ 

- Even *counting* jets is not trivial, while we typically target %-level precision.
- Residual artifacts in the spectrum may render the data difficult to fit.
- Issues may often be spotted only after the data have been published.
- $\rightarrow$  find build a function with same shape



CDL1

Patrick L.S. Connor

p\_JET (GeV/c)

- Fit with and divide by an ad hoc, (nearly) agnostic, smooth function  $f_n(p_{\rm T}) = \exp\left(\sum_{i=0}^n b_i T_i \left(2 \frac{\log p_{\rm T}/\log p_{\rm T}^{\rm min}}{\log p_{\rm T}^{\rm max}/\log p_{\rm T}^{\rm min}} - 1\right)\right) \qquad \qquad \chi_n^2 = \min_{b_i \le n} \left[\left(\mathbf{x} - \mathbf{y}_{b_i \le n}\right)^{\sf T} \mathbf{V}^{-1} \left(\mathbf{x} - \mathbf{y}_{b_i \le n}\right)\right]$
- Track down possible artifacts introduced in the data reduction (e.g. trigger)



Patrick L.S. Connor

## **STEP (3)** (with Radek Zlebcik)

- Iterative fit procedure with early stopping criterion based on cross validation.
- Two sets of 10k replicas (fake data generated from statistical properties of the original spectrum) are produced:
  - training replicas are used to run the fit;
  - validation replicas are used to determine the level of overfitting.
- Not more than 10% of the validation replicas should have better  $\chi^2$  than the training replicas.



2111.09968v2



- Apply **transformation** on distributions with complex correlations
  - ~ e.g. ratio of inclusive 2- and 3-jet cross sections
    - Many systematic effects cancel  $\rightarrow$  stronger sensitivity to  $\alpha_s$
- Normalise 2D distribution in bins
  - $_{\star \star}~$  e.g. jet substructure variable  $\lambda$  in bins of  $p_{_T}$ 
    - Factour out irrelevant physics effects
- Extract a fraction, etc. etc.
  → generic approach based on MC techniques



 $\frac{1}{N(p_{\rm T})} \frac{{\rm d}^2 \sigma}{{\rm d} p_{\rm T} {\rm d} \lambda} \quad \text{where} \quad N(p_{\rm T}) = \frac{{\rm d} \sigma}{{\rm d} p_{\rm T}}$ 



CDL1

Patrick L.S. Connor

11

## **TEDDY (2)** (with Radek Zlebcik)

#### CDCS CENTER FOR DATA AND COMPUTING IN NATURAL SCIENCES

- We want to apply y = f(x)
  - ~ x has covariance matrix V

$$\boldsymbol{\theta}_n = \mathbf{f} \left( \mathbf{x} - \mathbf{R}^{-1} \boldsymbol{\delta'}_n \right) \quad ext{with} \quad \delta'_{n,i} \sim \mathcal{N} \left( 0, \sqrt{\max(0, k_i)} \right)$$

• The covariance matrix of y is simply obtained by MC integration:

$$\mathbf{W} = \left(\frac{1}{N}\sum_{n=1}^{N}\boldsymbol{\theta}_{n} \cdot \boldsymbol{\theta}_{n}^{\mathsf{T}}\right) - \frac{1}{N^{2}}\left(\sum_{n=1}^{N}\boldsymbol{\theta}_{n}\right) \cdot \left(\sum_{n=1}^{N}\boldsymbol{\theta}_{n}\right)^{\mathsf{T}}$$

 $\rightarrow$  illustrated with partial normalisation of H1 dijet measurement



Patrick L.S. Connor

CDL1

- Statistical and systematic uncertainties are crucial in measurements.
- When designing the analysis software, they should be accounted for from the very beginning → DAS
- Ongoing project with adversarial networks to help produce very large data sets and better cover model uncertainties.
- Generic tools to deal with bin-to-bin partial correlations
  → STEP & TEDDY.

**Thanks for your attention!** 

CDL1