#### Machine Learning for PDF determination

Roy Stegeman Supervisor: Stefano Forte

University of Milan and INFN Milan

University of Milan, 19 December 2022









This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 740006.

### Precision physics at the LHC



SM cross-sections [ATLAS, PRD 87, 112003 (2013)]

- Currently good agreement between data and theory
- Observing deviations requires sub-percent accuracy

## Precision physics at the LHC



SM cross-sections [ATLAS, PRD 87, 112003 (2013)]

- Currently good agreement between data and theory
- Observing deviations requires sub-percent accuracy

PDFs uncertainties is often the dominant source of uncertainty for LHC cross-sections



Projection for HL-LHC [HL/HE-LHC WG2, 2019]7

#### Importance of faithful uncertainties

- This years'  $m_W$  determination highlights the importance of understanding all sources of uncertainty
- Understand apparent discrepancies between measurements
- Important that predictions are accurate



SM parameters:  $m_W$  [LHCb-FIGURE-2022-003]

#### Importance of faithful uncertainties

- This years'  $m_W$  determination highlights the importance of understanding all sources of uncertainty
- Understand apparent discrepancies between measurements
- Important that predictions are accurate
- Important role for PDF uncertainties



SM parameters:  $m_W$  [LHCb-FIGURE-2022-003]

Source	Size [MeV]
Parton distribution functions	9
Theory (excl. PDFs) total	17
Transverse momentum model	11
Angular coefficients	10
QED FSR model	7
Additional electroweak corrections	5
Experimental total	10
Momentum scale and resolution modelling	7
Muon ID, trigger and tracking efficiency	6
Isolation efficiency	4
QCD background	2
Statistical	23
Total	32

SM parameters:  $m_W$  [JHEP01(2022)036]

#### Status of modern PDF sets



[Snowmass (2022), 2203.13923]

- PDF predictions are **consistent** but with **different uncertainties**
- Ingredients of a PDF fit:
  - Data
  - Theory
  - Methodology

Main difference between fitting groups is the methodology

#### Status of modern PDF sets



- PDF predictions are **consistent** but with **different uncertainties**
- Ingredients of a PDF fit:
  - Data
  - Theory
  - Methodology

Main difference between fitting groups is the methodology

[Snowmass (2022), 2203.13923]

"Does the NNPDF methodology produce faithful uncertainties?"

#### Particularly relevant today!

#### Parton distributions need representative sampling

Aurore Courtey, 1-3 Josef Huston,<sup>2,1</sup> Pavel Nadolky,<sup>3,1</sup> Keping Xia,<sup>4,1</sup> Mengahi Yan,<sup>5,4</sup> and C.-P. Yuan,<sup>2,+\*</sup> <sup>1</sup> Institute de Fúsica, Universida Nacional Ardionona de Mésico, Apartado Partal 20-54,<sup>6</sup> 10100 Cudad de Mésico, Mesico <sup>2</sup> Dipartemat of Physics, Swattern Michaelds University, Boal, 77, 7757-1018, USA <sup>4</sup> Pattshuph Particle Physics, Astrophysics, and Cosmology Center, Department of Physics and Astronomy. University of Nacional Cost, 11, 12500, USA <sup>5</sup> School of Physics and State Key Laboratory O Nuclear Physics and Technology, Pesing University, Beijing 100871, China (Date).

In global OCD fits of parton distribution functions (PDFs), a large part of the estimated uncertainty on the PDFs originates from the choices of parametric functional forms and fitting methodology. We argue that these types of uncertainties can be underestimated with common PDF ensembles in high-stake measurements at the Large Hadron Collider and Tevatron. A fruitful approach to quantify these uncertainties is to view them as arising from sampling of allowed PDF solutions in a multidimensional parametric space. This approach applies powerful insights gained in recent statistical studies of large-scale population surveys and quasi-Monte Carlo integration methods. In particular. PDF fits may be affected by the big data paradox, which stipulates that more experimental data do not automatically raise the accuracy of PDFs - close attention to the data quality and sampling of possible PDF solutions is as essential. To test if the sampling of the PDF uncertainty of an experimental observable is truly representative of all acceptable solutions, we introduce a technique ("a hopscotch scan") based on a combination of parameter scans and stochastic sampling. With this technique, we show that the PDF uncertainty on key LHC cross sections at 13 TeV obtained with the public NNPDF4.0 fitting code is larger than the nominal uncertainty obtained with the published NNPDF4.0 Monte-Carlo replica sets. For example, the uncertainties on the charm distribution at a large momentum fraction x and gluon PDF at small x are enlarged. In PDF anomples obtained in the analytic minimization (Hessian) formalism, the tolerance on the PDF uncertainty must be based on sufficiently complete sampling of PDF functional forms and choices of the experiments

May 2022 [2205.10444]



Edinburgh 2022/19 TIF-UNIMI-2022-21

#### Response to "Parton distributions need representative sampling"

#### The NNPDF Collaboration:

Richard D. Ball,<sup>1</sup> Juan Cruz-Martinez,<sup>2</sup> Luigi Del Debbio,<sup>1</sup> Stefano Forte,<sup>3</sup> Zahari Kassabov,<sup>4</sup> Emanuele R. Nocera,<sup>5</sup>Juan Rojo,<sup>6,7</sup> Roy Stegeman,<sup>1</sup> Maria Ubiali<sup>4</sup>

<sup>14</sup>The Higgs Centre for Theoretical Physics, University of Eduburgh, JCMB, R.R. Magdiel M.E. diaburgh, EH9 3.27, Scotland <sup>2</sup> CERN, Theoretical Physics Department, CH-1211 Geneva 23, Switzerland <sup>3</sup> MJ I.A.D. Dipartimento di Fision, Università di Milano and INFN, Sesione di Milano, Yao Celoria 16, F-20137 Milano, Italy <sup>4</sup>DAMTP, University of Cambridge, Wilkerforce Road, Cambridge, CB OWA, Unital Kingdom <sup>5</sup> Dipartimento di Fisica, Università degli Studi di Tornio and INFN, Scione di Tornio, Via Petro Ginvin I, L10123 Tornio, Italy <sup>9</sup>Department of Physics and Astronomy, Vrije Universiteit, NL-1081 HV Ansterlann <sup>9</sup>Nikhof Theory Group, Science Davi 105, 1098 CA materiana, The Verkerlandes

#### November 2022 [2211.12961]

## Outline

#### NNPDF4.0

- Data
- Methodology
- Phenomenology

#### PDF correlations

- Correlation between different sets of PDFs
- Combination of PDF sets

#### Future ML developments

- A data-based parametrization
- An overfitting metric

2109.02653 and 2109.02671

PDF correlations

## PDFs from data

- Evaluating LHC cross-sections:
  - $\sigma_{ab} = \sum_{ab} f_a \otimes f_b \otimes \hat{\sigma}_{ab}$
  - PDF <u>f</u><sub>a</sub> of flavor a (non-perturbative, from data)
  - hard-scattering matrix element  $\hat{\sigma}_{ab}$ (perturbative QCD)
  - ${\, \bullet \,} \otimes$  denotes a convolution over momentum fraction x
- PDFs  $f_a$  depends only on x and  $Q^2$
- Other kinematic variables in  $\hat{\sigma}$



PDF correlations

## PDFs from data

- Evaluating LHC cross-sections:
  - $\sigma_{ab} = \sum_{ab} f_a \otimes f_b \otimes \hat{\sigma}_{ab}$
  - PDF *f<sub>a</sub>* of flavor *a* (non-perturbative, from data)
  - hard-scattering matrix element  $\hat{\sigma}_{ab}$  (perturbative QCD)
  - ullet  $\otimes$  denotes a convolution over momentum fraction x
- PDFs  $f_a$  depends only on x and  $Q^2$
- Other kinematic variables in  $\hat{\sigma}$

#### The problem

- Given a dataset D, determine p(f|D) in the space of PDFs  $f:[0,1]\to \mathbb{R}$
- PDFs are multivariate probability distributions in infinite dimensional space
- However data is discrete



PDF correlations

## Data in NNPDF4.0



9/37

PDF correlations

### Data in NNPDF4.0



PDF correlations

<sup>=</sup>uture ML developments 000000

#### The NNPDF approach: probabilities in a space of function

Data is fully defined by central values  $\mu_i$  and covariance matrix  $\mathrm{cov}_{ij}$ 

- Generate  $N_{\text{rep}}$  Monte Carlo data "replicas"  $\hat{\mu}_i$  such that as  $N_{\text{rep}} \to \infty$   $\mu_i = \frac{1}{N_{\text{rep}}} \sum_{i=1}^{N_{\text{rep}}} \hat{\mu}_i$  $\operatorname{cov}_{ij} = \operatorname{cov}[\hat{\mu}_i, \hat{\mu}_j]$
- Perform a PDF fit to each replica
- $\begin{array}{l} \textcircled{\textbf{O}} \quad \text{Compute observables } X \text{ and their uncertainties} \\ \langle X\left[f\right] \rangle = \frac{1}{N_{\mathrm{rep}}} \sum_{r=1}^{N} X\left[f^{(r)}\right] \\ \text{Var}\left[X\left[f\right]\right] = \frac{1}{N_{\mathrm{rep}}} \sum_{r=1}^{N_{\mathrm{rep}}} \left(X\left[f^{(r)}\right] \langle X\left[f\right] \rangle\right)^{2} \end{array}$



PDF correlations

<sup>=</sup>uture ML developments 000000

#### The NNPDF approach: probabilities in a space of function

Data is fully defined by central values  $\mu_i$  and covariance matrix  $\mathrm{cov}_{ij}$ 

- Generate  $N_{\text{rep}}$  Monte Carlo data "replicas"  $\hat{\mu}_i$  such that as  $N_{\text{rep}} \to \infty$   $\mu_i = \frac{1}{N_{\text{rep}}} \sum_{i=1}^{N_{\text{rep}}} \hat{\mu}_i$  $\operatorname{cov}_{ij} = \operatorname{cov}[\hat{\mu}_i, \hat{\mu}_j]$
- Perform a PDF fit to each replica
- $\begin{aligned} & \textbf{O} \quad \text{Compute observables } X \text{ and their uncertainties} \\ & \langle X\left[f\right] \rangle = \frac{1}{N_{\mathrm{rep}}} \sum_{r=1}^{N} X\left[f^{(r)}\right] \\ & \text{Var}\left[X\left[f\right]\right] = \frac{1}{N_{\mathrm{rep}}} \sum_{r=1}^{N_{\mathrm{rep}}} \left(X\left[f^{(r)}\right] \langle X\left[f\right] \rangle\right)^2 \end{aligned}$





PDF correlations

## The NNPDF approach: an importance sampling

- Importance sampling produces Gaussian posterior distribution
- All replicas equally likely



PDF correlations

## The NNPDF approach: an importance sampling

- Importance sampling produces Gaussian posterior distribution
- All replicas equally likely



How did we obtain this distribution?

PDF correlations

#### Parametrization



Physical constraints:

- PDF positivity [JHEP 11 (2020)]
- Integrability of nonsinglet distributions (Gottfried sum rules)
- $\bullet\,$  Train by minimizing  $\chi^2$  loss function comparing data to prediction

 $f_i(x,Q_0) = x^{-\alpha_i} (1-x)^{\beta_i} \operatorname{NN}_i(x)$ 

PDF correlations

## Training the neural network

- Optimize using a gradient descent algorithm
- NN should generalize the underlying law, but if trained to long noise is fitted

Cross-validation

- Divide data into training and validation
- Minimize training  $\chi^2$
- Stop if validation  $\chi^2$  no longer improves





PDF correlations 000000000

## Verify the importance sampling assumption



- All PDF replicas are fitted equally well to their data replica
- Thus outliers correspond to unlikely data replicas

PDF correlations

## Automated model selection

NNPDF aims to minimize sources of bias in the PDF:

- $\bullet$  Functional form  $\rightarrow$  Neural Network
- Model parameters  $\rightarrow$  ?

PDF correlations

#### Automated model selection

NNPDF aims to minimize sources of bias in the PDF:

- $\bullet$  Functional form  $\rightarrow$  Neural Network
- $\bullet \ \ \mathsf{Model} \ \mathsf{parameters} \ \to \ \mathbf{Hyperoptimization}$

Scan over thousands of hyperparameter combinations and select the best one

**k-fold cross-validation**: used to define the reward function based on a **test dataset** 

Objective function:  $L = mean(\chi_1^2, \chi_3^2, \chi_2^2, \dots, \chi_k^2)$ 

Final step requires human input



PDF correlations

## High-precision: gluon

$$\mathcal{L}_{ij}\left(M_X, y, \sqrt{s}\right) = \frac{1}{s} \sum_{i,j} f_i\left(\frac{M_X e^y}{\sqrt{s}}, M_X\right) f_j\left(\frac{M_X e^{-y}}{\sqrt{s}}, M_X\right)$$



PDF correlations

## High-precision: singlet

$$\mathcal{L}_{ij}\left(M_X, y, \sqrt{s}\right) = \frac{1}{s} \sum_{i,j} f_i\left(\frac{M_X e^y}{\sqrt{s}}, M_X\right) f_j\left(\frac{M_X e^{-y}}{\sqrt{s}}, M_X\right)$$



Typical uncertainties in the data region:

- $\bullet$  singlet: from  $\sim 3\%$  to  $\sim 1\%$
- $\bullet$  nonsinglet: from  $\sim 5\%$  to  $\sim 2\%$

PDF correlations

#### Impact of the new data



Individual datasets have a limited impact, but collectively they result in:

- Moderate reduction of PDF uncertainties
- Shifts in central value at the one-sigma level

PDF correlations

## Impact of the new fitting methodology



- Significant reduction of PDF uncertainties
- Good agreement between the central values

PDF uncertainties are validated using closure tests and future tests Validation tests successful for both NNPDF4.0 and NNPDF3.1

PDF correlations 000000000

## The open-source NNPDF code

- The full NNPDF code has been made public along with user friendly documentation
- This includes: fitting, hyperoptimization, theory, data processing, visualization
- It is possible to reproduce all results of NNPDF4.0 and more!



Eur.Phys.J.C 81 (2021) 10, 958
https://github.com/NNPDF/nnpdf
https://docs.nnpdf.science

2110.08274

### PDF correlations: definition and application

#### Definition:

Covariance:

 $\operatorname{Cov} [f_a, f_b](x, x') = E \left[ f_a \left( x, Q_0^2 \right) f_b \left( x', Q_0^2 \right) \right] - E \left[ f_a \left( x, Q_0^2 \right) \right] E \left[ f_b \left( x', Q_0^2 \right) \right]$ 

#### Correlation:

 $\rho\left[f_a, f_b\right](x, x') = \frac{\operatorname{Cov}[f_a, f_b](x, x')}{\sqrt{\operatorname{Var}[f_a](x)\operatorname{Var}[f_a^q](x')}}$ 

In the MC approach this can be calculated using  $E\left[f_a\left(x,Q_0^2\right)f_b\left(x',Q_0^2\right)\right] = \frac{1}{N}\sum_{r=1}^N f_a^{(r)}\left(x,Q_0^2\right)f_b^{(r)}\left(x',Q_0^2\right)$ 

Correlation **induced by** data, theory (e.g. sumrules), and methodology (e.g. preprocessing).



ttH

#### Correlations between different PDF sets

General cross-covariance between PDFs in different sets, e.g. NNPDF4.0 and MSHT20:

$$\operatorname{Cov}\left[f_{a}^{\operatorname{NNPDF}}, f_{b}^{\operatorname{MSHT}}\right]\left(x, x'\right) = E\left[f_{a}^{\operatorname{NNPDF}}\left(x, Q_{0}^{2}\right) f_{b}^{\operatorname{MSHT}}\left(x', Q_{0}^{2}\right)\right] - E\left[f_{a}^{\operatorname{NNPDF}}\left(x, Q_{0}^{2}\right)\right] E\left[f_{b}^{\operatorname{MSHT}}\left(x', Q_{0}^{2}\right)\right] = E\left[f_{b}^{\operatorname{MSHT}}\left(x', Q_{0}^{2}\right)\right]$$

Special cases of cross-correlation:

- F-correlation same PDF, different flavor  $ho\left[f_{a}^{\mathrm{NNPDF}},f_{b}^{\mathrm{NNPDF}}
  ight]$
- S-correlation different PDF, same flavor  $\rho \left[ f_a^{\text{NNPDF}}, f_a^{\text{MSHT}} \right]$

The same replica must be used when calculating covariance

If  $f^{(r)}$  and  $f^{(r')}$  are uncorrelated, covariance vanishes:  $E[f_a f_b] = \frac{1}{N} \sum_{r=1}^{N} f_a^{(r)} f_b^{(r')} = [f_a][f_b]$ 

Problem: What is the meaning of "same replica" across PDF sets?

#### Correlations between different PDF sets

General cross-covariance between PDFs in different sets, e.g. NNPDF4.0 and MSHT20:

$$\operatorname{Cov}\left[f_{a}^{\operatorname{NNPDF}}, f_{b}^{\operatorname{MSHT}}\right]\left(x, x'\right) = E\left[f_{a}^{\operatorname{NNPDF}}\left(x, Q_{0}^{2}\right)f_{b}^{\operatorname{MSHT}}\left(x', Q_{0}^{2}\right)\right] - E\left[f_{a}^{\operatorname{NNPDF}}\left(x, Q_{0}^{2}\right)\right]E\left[f_{b}^{\operatorname{MSHT}}\left(x', Q_{0}^{2}\right)\right] = E\left[f_{b}^{\operatorname{MSHT}}\left(x', Q_{0}^{2}\right)\right] =$$

Special cases of cross-correlation:

- F-correlation same PDF, different flavor  $\rho \left[ f_a^{\text{NNPDF}}, f_b^{\text{NNPDF}} \right]$
- S-correlation different PDF, same flavor  $\rho\left[f_{a}^{\mathrm{NNPDF}}, f_{a}^{\mathrm{MSHT}}\right]$

The same replica must be used when calculating covariance

If  $f^{(r)}$  and  $f^{(r')}$  are uncorrelated, covariance vanishes:  $E[f_a f_b] = \frac{1}{N} \sum_{r=1}^{N} f_a^{(r)} f_b^{(r')} = [f_a][f_b]$ 

Problem: What is the meaning of "same replica" across PDF sets?

**Possible solution**: PDF replicas fitted to the same data replica

- Fit  $f^{\text{MSHT}(r)}$  and  $f^{\text{NNPDF}(r)}$  to the same data replica r
- Calculate covariance as  $E [f_a f_b] = \frac{1}{N} \sum_{r=1}^{N} f_a^{\text{NNPDF}(r)} f_b^{\text{MSHT}(r)}$

#### Uncorrelated methodological aspects

The data replica does not uniquely determine the PDF replica



#### Uncorrelated methodological aspects

The data replica does not uniquely determine the PDF replica



PDF depends on uncorrelated methodological aspects:

- initialization of the neural network
- preprocessing exponents
- training/validation mask
- ...

#### Data-induced correlation

Let us distinguish

- data replicas r
- $\bullet$  methodological replicas r'

$$\left|\frac{1}{N}\sum_{r=1}^{N}f_{a}^{(r,r')}f_{b}^{(r,r'')} - E\left[f_{a}\right]E\left[f_{b}\right]\right| \leq \left|\frac{1}{NM}\sum_{r=1}^{N}\sum_{r'=1}^{M}f_{a}^{(r,r')}f_{b}^{(r,r')} - E\left[f_{a}\right]E\left[f_{b}\right]\right|$$

Only data-induced contributions are calculated



#### Correlated methodological replicas

- Can easily be done for parametric components such as preprocessing or architecture
- Noticeable impact of preprocessing, negligible for architecture



Non-trivial for non-parametric aspects of the methodology

<sup>=</sup>uture ML developments 000000

#### Data-induced self-correlation

- Correlation between two determinations of the same PDF set (e.g. NNPDF4.0)
- The deviation from 100% correlation is due to uncorrelated aspects of the methodology



#### Data-induced self-correlation

- Correlation between two determinations of the same PDF set (e.g. NNPDF4.0)
- The deviation from 100% correlation is due to uncorrelated aspects of the methodology
- Calculate correlation between different methodologies
- "weakest link" NNPDF3.1 $\circ$ NNPDF4.0  $\approx$  NNPDF3.1 $\circ$ NNPDF3.1

Higher correlation indicates a more efficient methodology



PDF4I HC combination

PDF correlations

Future ML developments 000000

IOP Publishing	Journal of Physics G: Nuclear and Particle Physic
----------------	---

Major Report

The PDF4LHC21 combination of global PDF fits for the LHC Run III\*

- All PDFs assumed to have equal probability
- Monte Carlo combination of 300 replicas from each of the underlying PDF sets



Uncertainty of the Monte Carlo combination:  $Var[f_a^{comb}] = \frac{1}{2} \left( Var[f_a^{MSHT}] + Var[f_a^{NNPDF}] \right) + \frac{1}{4} (E \left[ f_a^{MSHT} \right] - E \left[ f_a^{NNPDF} \right] )^2$ combined uncertainty bigger than average uncertainty if central values disagree

## Correlated PDF combination?

Idea:

- Combine PDF determinations as independent observations
- Correlated combination produces a weighted average

Assuming two sets of same variance,  $Var[f_a^{NNPDF}] = Var[f_a^{MSHT}]$ :

$$\operatorname{Var}[f_a^{\text{comb}}] = \frac{1}{2} \left( 1 + \rho[f_a^{\text{MSHT}}, f_a^{\text{NNPDF}}] \right) \operatorname{Var}[f_a^{\text{NNPDF}}]$$

## Correlated PDF combination?

Idea:

- Combine PDF determinations as independent observations
- Correlated combination produces a weighted average

Assuming two sets of same variance,  $Var[f_a^{NNPDF}] = Var[f_a^{MSHT}]$ :

$$\operatorname{Var}[f_a^{\text{comb}}] = \frac{1}{2} \left( 1 + \rho[f_a^{\text{MSHT}}, f_a^{\text{NNPDF}}] \right) \operatorname{Var}[f_a^{\text{NNPDF}}]$$

#### Problems:

- Does not account for difference in central values
- How to compute correlations reliably? I.e. methodological components

Underestimated data-induced correlations leads to **underestimated uncertainty of combination** 



## Correlated PDF combination?

Idea:

- Combine PDF determinations as independent observations
- Correlated combination produces a weighted average

Assuming two sets of same variance,  $Var[f_a^{NNPDF}] = Var[f_a^{MSHT}]$ :

$$\operatorname{Var}[f_a^{\text{comb}}] = \frac{1}{2} \left( 1 + \rho[f_a^{\text{MSHT}}, f_a^{\text{NNPDF}}] \right) \operatorname{Var}[f_a^{\text{NNPDF}}]$$

#### Problems:

- Does not account for difference in central values
- How to compute correlations reliably? I.e. methodological components

Underestimated data-induced correlations leads to **underestimated uncertainty of combination** 

# What if we combine many repeated determinations of the same PDF set?



## Future ML developments

2111.02954 and 2211.12961

PDF correlations

## Preprocessing

- PDF model:  $f_a = A_a x^{\alpha_a} (1-x)^{\beta_a} NN_a(x, \log x)$
- Exponents  $\alpha_a$  and  $\beta_a$  sampled randomly per replica
- Range determined through self-consistent iterative procedure
- $x \text{ and } \log x \text{ inputs}$



- Need to iterate
- Hierarchy in input scale. A source of bias?
- Sampled exponents add noise during hyperopt



#### Avoid potential bias:

- generate flat distribution: effective cumulative distribution function (eCDF)
- Add interpolation
- Rerun hyperopt
- $\Rightarrow$  Only single input
- $\Rightarrow$  No preprocessing needed



Avoid potential bias:

- generate flat distribution: effective cumulative distribution function (eCDF)
- Add interpolation
- Rerun hyperopt
- $\Rightarrow$  Only single input
- $\Rightarrow$  No preprocessing needed

Good agreement between NNPDF4.0 and feature scaling!



## Overfitting in hyperopt

- Statistical fluctuations affecting hyperopt results
- Hyperopt solutions can be overfitted or underfitted



- What does it mean for a PDF to be overfitted?
- More wiggles  $\rightarrow$  overfitting?
- $\bullet~\mbox{More}~\mbox{wiggles} \rightarrow \mbox{better}~\mbox{agreement}~\mbox{to}~\mbox{data}$

#### overfitting metric

Validation data only used for early stopping

- Take set of PDF replicas
- ② Calculated expected validation loss
- Ompare to real validation loss

$$\mathcal{R}_{O} = \chi_{\mathrm{val,r}}^{2} \left[ \mathcal{T} \left[ f^{(r)} \right], \mathcal{D}^{(r)} \right] - \frac{1}{N} \sum_{r'=1}^{N} \chi_{\mathrm{val,r}}^{2} \left[ \mathcal{T} \left[ f^{(r)} \right], \mathcal{D}^{(r')} \right].$$

Negative  $\mathcal{R}_O \rightarrow \text{overfitting}$ 

#### overfitting metric

Validation data only used for early stopping

- Take set of PDF replicas
- ② Calculated expected validation loss
- Ompare to real validation loss

$$\mathcal{R}_{O} = \chi_{\mathrm{val,r}}^{2} \left[ \mathcal{T} \left[ f^{(r)} \right], \mathcal{D}^{(r)} \right] - \frac{1}{N} \sum_{r'=1}^{N} \chi_{\mathrm{val,r}}^{2} \left[ \mathcal{T} \left[ f^{(r)} \right], \mathcal{D}^{(r')} \right].$$

Negative  $\mathcal{R}_O \rightarrow \text{overfitting}$ 



Future ML developments

#### I did not discuss other sources of uncertainty

- Theory and missing higher order uncertainties (Andrea, next days)
- The negligible impact of data inconsistencies (see 2212.07703)
- Future tests and closure tests

#### Possible future directions

- Optimize hyperopt folds
- Overfitting metric in hyperopt
- Parallelization on GPU
- Understanding PDF fitting in a Bayesian framework

More work still to be done!

Future ML developments

#### I did not discuss other sources of uncertainty

- Theory and missing higher order uncertainties (Andrea, next days)
- The negligible impact of data inconsistencies (see 2212.07703)
- Future tests and closure tests

#### Possible future directions

- Optimize hyperopt folds
- Overfitting metric in hyperopt
- Parallelization on GPU
- Understanding PDF fitting in a Bayesian framework

More work still to be done!

# Thank you!

Backup

## Experimental data in NNPDF4.0

- 44 new datasets included
- 323 more data points in NNPDF4.0 than in NNPDF3.1

#### New data is mostly from the LHC RUN II

Data set	Ref.	NNPDF3.1	NNPDF4.0	ABMP16	CT18	MSHT20
ATLAS $W, Z$ 7 TeV ( $\mathcal{L} = 35 \text{ pb}^{-1}$ )	[51]	1	1	1	1	1
ATLAS W, Z 7 TeV $(\mathcal{L} = 4.6 \text{ fb}^{-1})$	[52]	1	1	×	(✔)	1
ATLAS low-mass DY 7 TeV	[53]	1	1	×	(✔)	×
ATLAS high-mass DY 7 TeV	[54]	1	1	×	(✔)	1
ATLAS W 8 TeV	(79)	×	(✔)	×	×	1
ATLAS DY 2D 8 TeV	[78]	×	1	×	×	1
ATLAS high-mass DY 2D 8 TeV	(77)	×	1	×	(✔)	1
ATLAS $\sigma_{W,Z}$ 13 TeV	[81]	×	1	1	×	×
ATLAS W+jet 8 TeV	[93]	×	1	×	×	1
ATLAS Z p <sub>T</sub> 7 TeV	(260)	(✔)	×	×	(✔)	×
ATLAS Z p <sub>T</sub> 8 TeV	[63]	1	1	×	1	1
ATLAS $W + c$ 7 TeV	[83]	×	1	x	(✔)	×
ATLAS $\sigma_{tt}^{1ot}$ 7, 8 TeV	(65)	1	1	1	×	×
ATLAS $\sigma_{tt}^{tot}$ 7, 8 TeV	[261 - 266]	×	×	1	×	×
ATLAS $\sigma_{tt}^{tot}$ 13 TeV ( $\mathcal{L} = 3.2 \text{ fb}^{-1}$ )	[66]	1	×	1	×	×
ATLAS $\sigma_{tt}^{\text{tot}}$ 13 TeV ( $\mathcal{L} = 139 \text{ fb}^{-1}$ )	[134]	×	1	×	×	×
ATLAS $\sigma_{tt}^{tot}$ and Z ratios	[267]	×	×	×	×	(🗸)
ATLAS tr lepton+jets 8 TeV	[67]	1	1	x	1	1
ATLAS $t\bar{t}$ dilepton 8 TeV	(89)	×	1	x	×	1
ATLAS single-inclusive jets 7 TeV, R=0.6	[73]	1	(✔)	×	1	1
ATLAS single-inclusive jets 8 TeV, R=0.6	(86)	×	1	×	×	×
ATLAS dijets 7 TeV, R=0.6	[148]	×	1	×	×	×
ATLAS direct photon production 8 TeV	[100]	×	(✔)	×	×	×
ATLAS direct photon production 13 TeV	[101]	×	1	×	×	×
ATLAS single top $R_t$ 7, 8, 13 TeV	[94, 96, 98]	×	1	1	×	×
ATLAS single top diff. 7 TeV	[94]	×	1	×	×	×
ATLAS single top diff. 8 TeV	[96]	×	1	×	×	×

Data set	Ref.	NNPDF3.1	NNPDF4.0	ABMP16	CT18	MSHT20
CMS W asym. 7 TeV ( $\mathcal{L} = 36 \text{ pb}^{-1}$ )	[268]	×	×	×	×	1
CMS Z 7 TeV ( $\mathcal{L} = 36 \text{ pb}^{-1}$ )	[269]	×	×	×	×	1
CMS $W$ electron asymmetry 7 TeV	[55]	1	1	×	1	1
CMS W muon asymmetry 7 TeV	[56]	1	1	1	1	×
CMS Drell-Yan 2D 7 TeV	[57]	100	1	×	(✔)	1
CMS Drell-Yan 2D 8 TeV	[270]	(✔)	×	×	×	×
CMS W rapidity 8 TeV	[58]	1	1	1	1	1
CMS W, Z p <sub>T</sub> 8 TeV ( $\mathcal{L} = 18.4 \text{ fb}^{-1}$ )	[271]	×	×	×	(✔)	×
CMS Z p <sub>T</sub> 8 TeV	[64]	1	1	×	(✔)	×
CMS $W + c$ 7 TeV	[76]	1	1	×	(🖌)	1
CMS $W + c$ 13 TeV	[84]	×	1	×	×	(✔)
CMS single-inclusive jets 2.76 TeV	[75]	1	×	×	×	1
CMS single-inclusive jets 7 TeV	[147]	1	(✔)	×	1	1
CMS dijets 7 TeV	[74]	×	1	×	×	×
CMS single-inclusive jets 8 TeV	[87]	×	1	×	1	1
CMS 3D dijets 8 TeV	[149]	×	(✔)	×	×	×
CMS $\sigma_{tt}^{tot}$ 5 TeV	[88]	×	1	×	×	×
CMS $\sigma_{tt}^{tot}$ 7, 8 TeV	[146]	1	1	×	×	×
CMS $\sigma_{tt}^{1ot}$ 8 TeV	[272]	×	×	×	×	100
CMS $\sigma_{t\bar{t}}^{tot}$ 5, 7, 8, 13 TeV	[68, 273 - 281]	×	×	1	×	×
CMS $\sigma_{tt}^{\text{tot}}$ 13 TeV	[69]	1	1	1	×	×
CMS $t\bar{t}$ lepton+jets 8 TeV	[70]	1	1	×	×	1
CMS $t\bar{t}$ 2D dilepton 8 TeV	[90]	×	1	×	1	1
CMS $t\bar{t}$ lepton+jet 13 TeV	[91]	×	1	×	×	×
CMS $t\bar{t}$ dilepton 13 TeV	[92]	×	1	×	×	×
CMS single top $\sigma_t + \sigma_{\bar{t}}$ 7 TeV	[95]	×	1	1	×	×
CMS single top $R_t$ 8, 13 TeV	[97, 99]	×	1	1	×	×
CMS single top 13 TeV	[282, 283]	×	×	×	× .	(✔)
Data set	Ref	NNPDE3.1	NNPDE4.0	ABMP16	CT18	MSHT20
	(22)			-	erno	
LHCb Z 7 TeV ( $\mathcal{L} = 940 \text{ pb}^{-1}$ )	[59]	1.1	- 1	- <u>*</u>	<u>.</u>	1
LHCD $Z \rightarrow cc$ 8 TeV ( $L = 2$ fb <sup>-1</sup> )	[61]		- (j	- ( )	÷.	1
Encow riev (L = 3r pb <sup>-1</sup> )	[284]		- 1	- 1	<u></u>	
LHCD $W, Z \rightarrow \mu \tau$ lev	[60]	1	- 1	- C	1	- 1
LINCO $W, z \rightarrow \mu 8$ TeV	[62]	1	10	1	1	1
LINCO $W \rightarrow e 8$ 1ev	[80]	1 A	(2)	<u> </u>	<u> </u>	1
LHCb $Z \rightarrow \mu\mu, ee$ 13 TeV	[82]	×		*	×	*

# NNPDF4.0 model

For more information see EPJC79 (2019) 676



#### Main changes:

- Python codebase
  - Easier and faster development
- Freedom to use external libraries (default: TensorFlow)
- $\bullet\,$  Modularity  $\Rightarrow$  ability to vary all aspects of the methodology

## Performance benefit - time per replica

	NNPDF3.1	NNPDF4.0 (CPU)	NNPDF4.0 (GPU)
Fit timing per replica	15.2 h	38 min	6.6 min
Speed up factor	1	24	140
RAM use	1.5 GB	6.1 GB	NA

## Hyperoptimization: the reward function

Choosing as the hyperoptimization target the  $\chi^2$  of fitted data results in overfitting.



#### Hyperoptimization: the reward function

Choosing as the hyperoptimization target the  $\chi^2$  of fitted data results in overfitting.

#### We solve this using k-fold cross-validation:

 $\Rightarrow$  The hyperoptimization target is not based on data that entered the fit.



- No overfitting
- Compared to NNPDF3.1:
  - Increased stability
  - Reduced uncertainties

#### The (negligible) impact of datasets with tension

Excluding datasets with large  $(\chi^2 - 1)/\sigma_{\chi^2}$  one at a time and combining the resulting PDFs following the conservative PDF4LHC15 prescription shows stability at the level of statistical fluctuations.



#### Envelope of fits with different arametrization bases

#### Different strategies to parametrize the PDF flavour combinations lead to the same result



# Understanding the $\chi^2$ distribution



# Impact of positivity on the PDFs



#### Trusting uncertainties outside the data region

- The improved methodology and extended dataset result in a reduction of the PDF uncertainties
- 'Closure test' to validate uncertainty in the data region: arxiv:1410.8849
- Can we trust the uncertainties in the extrapolation region?

#### Idea:

- Take a historic dataset e.g. pre-HERA or pre-LHC
- Perform fit
- Ompare predictions to "future" data



#### Future tests

For more information see arxiv:2103.08606

 $\chi^2/N$  (only exp. covmat)

.

(dataset)	I	NNPDF4.0	pre-LHC	pre-Hera
pre-HERA		1.09	1.01	0.90
pre-LHC		1.21	1.20	23.1
NNPDF4.0		1.29	<b>3.30</b>	23.1





#### Future tests

For more information see arxiv:2103.08606

 $\chi^2/N$  (exp. and PDF covmat)



The total uncertainty increases, and accommodates for difference between predictions and new data.

#### Closure test See Eur.Phys.J.C 82 (2022); arxiv:2111.05787

Closure test of a known input assumption

- Assume a "true" underlying PDF (e.g. a single PDF replica)
- Produce data distributed according to the experimental covariance matrices
- Perform a fit to this data

Examples of statistical estimators:

- Bias: squared difference between central value and true observable Variance: variance of the model predictions Faithful uncertainties require *E*[bias] = variance
- Is truth within one sigma in 68% of cases?

$\sqrt{\text{bias}/\text{variance}}$	$\xi_{1\sigma}^{( m data)}$
$1.03\pm0.05$	$0.68\pm0.02$