



PDF DETERMINATION AS MACHINE LEARNING

STEFANO FORTE UNIVERSITÀ DI MILANO & INFN



UNIVERSITÀ DEGLI STUDI DI MILANO DIPARTIMENTO DI FISICA



MAINZ, JULY 5, 2022

DEEP-LEARNING ERA OF PARTICLE THEORY

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

SUMMARY

- THE STATE OF THE ART
 - PDFs and data
 - UNCERTAINTIES: DATA AND EXTRAPOLATION
- PDF UNCERTAINTIES
 - BACKWARD AND FORWARD COMPATIBILITY
 - METHODOLOGY AND THEORY BIAS
 - CLOSURE TESTS
- THE NNPDF METHDOLOGY
 - THE FUNCTIONAL MONTECARLO
 - NEURAL NETS AND GENETIC MINIMIZATION
- MACHINE LEARNINING PDFs
 - CODE STRUCTURE
 - HYPEROPTIMIZATION
 - K-folding
- VALIDATION
 - CLOSURE TESTS
 - FUTURE TESTS
- DELIVERY
 - GENETIC COMPRESSION
 - GAN OPTIMIZATION
- FUTURE DEVELOPMENTS
 - FEATURE SCALING
 - OVERLEARNING METRICS

PDFs: THE STATE OF THE ART

PDFs AND DATA

THE NNPDF4.0 DATASET

Kinematic coverage

- LHC CROSS SECTION: $- \sigma = \sum_{ij} \hat{\sigma}_{ij} \otimes f_i^{(1)} \otimes f_j^{(2)}$ - $\hat{\sigma}_{ij}$ partonic cross section, INCOMING PARTONS i, j- $f_i^{(j)}(x, Q^2)$ PDF for parton of species iIN j-TH INCOMING PROTON - \otimes CONVOLUTION OVER x- PDF DEPENDS ON Q^2 AND x, OTHER KINE-MATIC VARIABLES IN $\hat{\sigma}$ PARTONIC CROSS SECTION COMPUTED PERTUR-BATIVELY • PDFs determined comparing σ to data - About 4600 datapoints
 - LEPTOPRODUCTION & HADROPRODUCTION, COLLIDER & FIXED-TARGET





UNCERTAINTIES: RECENT (NNPDF3.1, 2017)

- TYPICAL UNCERTAINTIES IN DATA REGION: SINGLET $\sim 3\%$, NONSINGLET $\sim 5\%$
- DATA REGION: $10^2 \lesssim M_X \lesssim 10^3$ TeV, $-2 \lesssim y \lesssim 2$



UNCERTAINTIES: STATE OF THE ART (NNPDF4.0, 2021)

TYPICAL UNCERTAINTIES IN DATA REGION: SINGLET $\sim 1\%$, NONSINGLET $\sim 2-3\%$

• DATA REGION: $10 \lesssim M_X \lesssim 3 \cdot 10^3$ TeV, $-4 \lesssim y \lesssim 4$





UNCERTAINTY ESTIMATION

THE PDF UNCERTAINTY PROBLEM: THE HERA-LHC BENCHMARK (2005)

- RESTRICTED AND VERY CONSISTENT DATASET USED
- RESULTS COMPARED TO THEN-BEST RESULT FROM FULL DATASET



BENCHMARK VS DEFAULT GLUON

"...the partons extracted using a very limited data set are completely incompatible, even allowing for the uncertainties, with those obtained from a global fit with an identical treatment of errors...The comparison illustrates the problems in determining the true uncertainty on parton distributions." (R.Thorne, HERALHC, 2005)



- CTEQ5 2002: $xg(x, Q_0^2) = A_0 x^{A_1} (1-x)^{A_2} (1+A_3 x^{A_4})$
- MRST-HERALHC 2005: $xg(x, Q_0^2) = A_g x^{\delta g} (1-x)^{\eta g} (1+\epsilon_g x^{0.5} + \gamma_g x) + A_{g'} x^{\delta g'} (1-x)^{\eta g'}$
- CT18: $g(x, Q = Q_0) = x^{a_1 1} (1 x)^{a_2} \left[a_3 (1 y)^3 + a_4 3y (1 y)^2 + a_5 3y^2 (1 y) + y^3 \right];$ $y = \sqrt{x}; a_5 = (3 + 2a_1)/3.$

BIAS?

PDF UNCERTAINTIES AND NEW PHYSICS

- DISCREPANCY BETWEEN QCD CALCULATION AND CDF JET DATA (1995)
- EVIDENCE FOR QUARK COMPOSITENESS?
- RESULT STRONGLY DEPENDS ON GLUON AT $x \gtrsim 0.1$
- PDF MUST VANISH AT x = 0, BUT (THEN) NO DATA FOR $x \ge 0.05!$



DISCREPANCY REMOVED IF JET DATA USED FOR GLUON DETERMINATION



NOW: NO DATA FOR $x \gtrsim 0.5 \Rightarrow$ **DISCOVERY** (THRESHOLD) REGION!





х

.4 .5 .6 .7 .8

(Scale is linear in $x^{1/3}$)

.3

A. de Roeck, Cracow epiphany conf. 1996

• RISE OF F_2 AT HERA CAME \Rightarrow SURPRIZE

0.5 0

1.5

1

0.5

HINTED BY PRE-HERA DATA; VETOED BY THEORETICAL BIAS •

THE NNPDF METHODOLOGY

PROTON STRUCTURE AS AN AI PROBLEM: NNPDF



THE FUNCTIONAL MONTE CARLO

REPLICA SAMPLE OF FUNCTIONS ⇔ PROBABILITY DENSITY IN FUNCTION SPACE KNOWLEDGE OF LIKELIHHOD SHAPE (FUNCTIONAL FORM) NOT NECESSARY



FINAL PDF SET: $f_i^{(a)}(x,\mu)$; i =up, antiup, down, antidown, strange, antistrange, charm, gluon; $j = 1, 2, ... N_{\text{rep}}$



- 37 parameters \times 8 PDFs
- PREPROCESSING EXPONENTS RANDOMIZED REPLICA PER REPLICA, RANGE DETERMINED SELF-CONSISTENTLY

GENETIC MINIMIZATION

- NODAL MUTATION REPLACED POINT MUTATION
- SINGLE EPOCH WITH NO REWEIGHTING UNLIKE PREVIOUS

OPTIMAL FIT

- **CROSS-VALIDATION** WITH 50-50 TRAINING & VALIDATION FRACTIONS
- LOOKBACK STOPPING REPLACED THRESHOLD ON DERIVATIVE

MACHINE LEARNING PDFs

LEARNING THE METHODOLOGY

THE N3FIT PROJECT



HOW DO WE KNOW THAT THE METHODOLOGY IS THE BEST? "ACCUMULATED WISDOM" INEFFICIENT AND SLOW

CHANGE OF PHILOSOPHY \Rightarrow DETERMINISTIC MINIMIZATION (GRADIENT DESCENT) GO FOR THE ABSOLUTE MINIMUM, AND (HYPER)OPTIMIZE



- PYTHON-BASED KERAS + TENSORFLOW FRAMEWORK
- EACH BLOCK INDEPENDENT LAYER
- CAN VARY ALL ASPECTS OF METHODOLOGY

THE NNPDF CODE STRUCTURE

- MODULAR PYTHON-BASED CODE
- HIGH DEGREE PARALLELIZATION & HARDWARE ACCELERATION

Average fitting time per replica and use of resources
SAME DATASET FOR OLD AND NEW METHODOLOGIES IN CPU AND GPU
CPU: INTEL(R) CORE(TM) 17-4770 AT 3.40GHZ; GPU: NVIDIA TITAN V

	NNPDF31 CODEBASE	NNPDF40 CODEBASE IN CPU	NNPDF40 CODEBASE IN GPU
Тіме	15.2 н.	38 ± 5 Min.	6.6 MIN.
RAM USE	1.5 GB	6.1 GB	NA



MINIMIZATION AND CROSS-VALIDATION

- DATA REPLICAS \Rightarrow PDF REPLICAS
- EACH PDF REPLICA: PREPROCESSED NEURAL NET
- NEURAL NET \Rightarrow OBSERVABLES
- RANDOM TRAINING-VALIDATION SPLIT, χ^2 to training data replicas minimized
- TRAINING STOPS IF VALIDATION χ^2 GROWS FOR A WHILE (PATIENCE)
- LOWEST VALIDATION $\chi^2 \Rightarrow {\rm OPTIMAL}\; {\rm FIT}$





- SCAN PARAMETER SPACE
- OPTIMIZE FIGURE OF MERIT: VALIDATION χ^2
- BAYESIAN UPDATING





- NOT HYPEROPTIMIZED: WIGGLES: FINITE SIZE \Rightarrow WILL GO AWAY AS N_{rep} GROWS
- N3FIT: WIGGLY PDFS \Leftrightarrow OVERFITTING \Rightarrow WILL NOT GO AWAY ($\chi^2_{train} \ll \chi^2_{valid}$!!)

WHAT HAPPENED?



CROSS-VALIDATION SELECTS THE OPTIMAL MINIMUM

WHAT HAPPENED?

HYPEROPTIMIZATION



WE ARE MISSING A SELECTION CRITERION



- HANDPICKED: WIGGLES: FINITE SIZE \Rightarrow WILL GO AWAY AS $N_{\rm rep}$ GROWS
- N3FIT: WIGGLY PDFS \Leftrightarrow OVERFITTING \Rightarrow WILL NOT GO AWAY ($\chi^2_{train} \ll \chi^2_{valid}$!!)
- CORRELATIONS BETWEEN TRAINING AND VALIDATION DATA

THE SOLUTION

TUNED HYPEROPTIMIZATION



TESTS GENERALIZATION POWER

THE TEST SET METHOD

- COMPLETELY UNCORRELATED TEST SET
- OPTIMIZE ON WEIGHTED AVERAGE OF VALIDATION AND TEST \Rightarrow NO OVERLEARNING



- NO OVERFITTING
- COMPARED TO HANDPICKED
 - MUCH GREATER STABILITY \Rightarrow FEWER REPLICAS FOR EQUAL ACCURACY
 - UNCERTAINTIES SOMEWHAT REDUCED

*K***-FOLDING** THE BASIC IDEA:

- DIVIDE THE DATA INTO n REPRESENTATIVE SUBSETS EACH CONTAINING PROCESS TYPES, KINEMATIC RANGE OF FULL SET
- FIT n 1 SETS AND USE n-TH SET AS TEST $\Rightarrow n$ VALUES OF $\chi^2_{\text{test, i}}$
- HYPEROPTIMIZE ON NON FITTED $\chi^2_{\rm test,\ i}$ \rightarrow GOOD & STABLE GENERALIZATION



K-FOLDING IMPLEMENTATION



	Fold 1	
CHORUS σ_{CC}^{ν}	HERA I+II inc NC e^+p 920 GeV	BCDMS p
LHCb Z 940 pb	ATLAS W, Z 7 TeV 2010	CMS Z p_T 8 TeV (p_T^{ll}, y_{ll})
DY E605 σ_{DY}^{p}	CMS Drell-Yan 2D 7 TeV 2011	CMS 3D dijets 8 TeV
ATLAS single- $\bar{t} y$ (normalised)	ATLAS single top R_t 7 TeV	CMS $t\bar{t}$ rapidity $y_{t\bar{t}}$
CMS single top R_t 8 TeV		
	Fold 2	
HERA I+II inc CC e^-p	HERA I+II inc NC e^+p 460 GeV	HERA comb. $\sigma_{b\bar{b}}^{red}$
NMC p	NuTeV σ_c^p	LHCb $Z \rightarrow ee \ 2 \ fb$
CMS W asymmetry 840 pb	ATLAS Z p_T 8 TeV (p_T^{ll}, M_{ll})	D0 $W \rightarrow \mu\nu$ asymmetry
DY E886 σ_{DY}^{p}	ATLAS direct photon 13 TeV	ATLAS dijets 7 TeV, R=0.6
ATLAS single antitop y (normalised)	CMS σ_{tt}^{tot}	CMS single top $\sigma_t + \sigma_{\bar{t}}$ 7 TeV
	Fold 3	
HERA I+II inc CC e^+p	HERA I+II inc NC e^+p 575 GeV	NMC d/p
NuTeV σ_c^{ν}	LHCb $W, Z \rightarrow \mu$ 7 TeV	LHCb $Z \rightarrow ee$
ATLAS W, Z 7 TeV 2011 Central selection	ATLAS W^+ +jet 8 TeV	ATLAS HM DY 7 TeV
CMS W asymmetry 4.7 fb	DYE 866 $\sigma_{DY}^d / \sigma_{DY}^p$	CDF Z rapidity (new)
ATLAS σ_{tt}^{tot}	ATLAS single top y_t (normalised)	CMS σ_{tt}^{tot} 5 TeV
CMS $t\bar{t}$ double diff. $(m_{t\bar{t}},y_t)$		
	Fold 4	
CHORUS σ^{p}_{CC}	HERA I+II inc NC e^+p 820 GeV	LHC b $W,Z \to \mu$ 8 TeV
	ATLAS W Z 7 TeV 2011 Feed	ATLAS W ⁻ +iet 8 TeV
LHCb $Z \rightarrow \mu\mu$	ALLAS W, Z 7 IEV 2011 FWG	111110 11 1 100 101
LHCb $Z \rightarrow \mu\mu$ ATLAS low-mass DY 2011	ATLAS $W, Z T$ fev 2011 Fwd ATLAS $Z p_T$ 8 TeV (p_T^{ll}, y_{ll})	CMS W rapidity 8 TeV
LHCb $Z \rightarrow \mu\mu$ ATLAS low-mass DY 2011 D0 Z rapidity	ATLAS $W, Z \neq 10V$ 2011 Pwd ATLAS $Z p_T $ 8 TeV (p_T^{ll}, y_{ll}) CMS dijets 7 TeV	$\begin{array}{c} \text{CMS } W \text{ rapidity 8 TeV} \\ \text{ATLAS single top } y_t \text{ (normalised)} \end{array}$





- EACH FOLD REPRODUCES FEATURES OF FULL DATASET
- DIFFERENT CHOICES POSSIBLE FOR LOSS (NON-FITTED)
 - BEST WORST
 - BEST AVERAGE
- RESULTS **STABLE**

THE ML METHODOLOGY

Parameter	NNPDF4.0	L as in Eq. (3.21)	Flavour basis Eq. (3.2)	
Architecture	25-20-8	70-50-8	7-26-27-8	
Activation function	hyperbolic tangent	hyperbolic tangent	sigmoid	
Initializer	glorot_normal	glorot_uniform	glorot_normal	
Optimizer	Nadam	Adadelta	Nadam	
Clipnorm	6.0×10^{-6}	5.2×10^{-2}	2.3×10^{-5}	
Learning rate	2.6×10^{-3}	2.5×10^{-1}	2.6×10^{-3}	
Maximum # epochs	17×10^3	45×10^{3}	45×10^{3}	
Stopping patience	10% of max epochs	12% of max epochs	16% of max epochs	
Initial positivity $\Lambda^{(pos)}$	185	106	2	
Initial integrability $\Lambda^{(\rm int)}$	10	10	10	

HYPEROPTIMIZED PARAMETERS



- HYPEROPT ADAPTS TO EXTERNAL CHOICES (E.G. PARAMETRIZATION BASIS)
- SIMILAR RESULTS CAN BE OBTAINED WITH RATHER DIFFERENT SETTINGS
- ~ 800 free parameters

VALIDATION

CLOSURE TESTS FAITHFUL UNCERTAINTIES IN DATA REGION?

- Assume "true" underlying PDF \Rightarrow E.G. some random PDF replica
- GENERATE DATA DISTRIBUTED ACCORDING TO EXPERIMENTAL COVARIANCE MATRIX
- RUN WHOLE METHDOLOGY ON THESE DATA
- DO STATISTICS ON "RUNS OF THE UNIVERSE", POSSIBLE THANKS TO EFFICIENT METHDOLOGY: COMPARE TO TRUE PDFS, OR TO TRUE VALUES OF OBSERVABLES (NOT FITTED)
 - BIAS/VARIANCE: MEAN SQUARE DEVIATION WR TO TRUTH VS UNCERTAINTY
 - is truth within one sigma 68% of times?

0.40	deviation from truth	Dataset	$\sqrt{\text{bias/variance}}$	$\xi_{1\sigma}^{(m data)}$
0.30		DY	0.99 ± 0.08	0.69 ± 0.02
0.25		Top-pair	0.75 ± 0.06	0.75 ± 0.03
0.20		Jets	1.14 ± 0.05	0.63 ± 0.03
0.15		Dijets	0.99 ± 0.07	0.70 ± 0.03
0.10		Direct photon	0.71 ± 0.06	0.81 ± 0.03
		Single top	0.87 ± 0.07	0.69 ± 0.04
0.00	-4 -2 0 2 4 Difference to underlying prediction	Total	1.03 ± 0.05	0.68 ± 0.02

RESULTS

FUTURE TESTS FAITHFUL UNCERTAINTIES IN EXTRAPOLATION?

- DETERMINE PDFs FROM A SUBSET OF CURRENT DATA: "PRE-HERA", "PRE-LHC",...
- COMPUTE χ^2 to the full current dataset:
 - WITHOUT PDF UNCERTAINTIES \Rightarrow IF \gg 1, MISSING INFORMATION
 - with PDF uncertainty \Rightarrow if \sim 1, missing info reproduced by uncertainty



VALENCE: PRE-HERA VS

PROCESS	PRE-HERA	PRE-LHC	3.1-like	4.0-glob
FT DIS (NC)	1.04	1.17	1.17	1.26
FT DIS (CC)	0.80	0.86	0.88	0.90
FT DY	0.93	1.27	1.43	1.59
HERA	24.01/1.12	1.22	1.21	1.21
Coll. DY (Tev.)	5.31/ 1.08	0.96	0.95	1.13
Coll. DY (LHC)	15.50/ 1.37	2.64/1.54	1.39	1.54
TOP QUARK	23.35/1.08	1.29/ 0.86	0.82	0.98
JETS	6.18/ 1.21	3.66/ 1.29	2.07/1.37	1.26
TOTAL	9.70	1.44	1.22	1.17

 χ^2 WITHOUT/WITH PDF UNC.

NNPDF4.0 vs. NNPDF3.1

- FULL BACKWARD COMPATIBILITY
- SUBSTANTIAL REDUCTION IN UNCERTAINTY







MONTECARLO COMPRESSION CAN WE REDUCE THE NUMBER OF REPLICAS?

- START WITH LARGE REPLICA SAMPLE
- SELECT BY GENETIC ALGORITHM SUBSET OF REPLICAS \Rightarrow STATISTICAL FEATURES OPTIMIZED TO PRIOR
- MINIMIZE LOSS: DIFFERENCE OF MOMENTS, KL DIVERGENCE, ...
- 50 COMPRESSED REPLICA REPRODUCE 1000 REPLICA SET TO PRECENT ACCURACY





- TRAIN A NETWORK TO SIMULATE THE TRUE DISTRIBUTION (GENERATOR)
- TRAIN A NETWORK TO **DISCRIMINATE** TRUTH FROM SIMULATION (**DISCRIMINATOR**)
- TRAIN THE GENERATOR TO TRICK THE DISCRIMINATOR

GAN ENHANCEMENT

- ENHANCE THE STARTING PDF SET BY ADDING GAN-PDFS TO IT
- PERFORM COMPRESSION OF THE ENHANCED SET



ENHANCED: NUMBER OF REPLICAS CUT IN HALF FOR SAME TARGET ACCURACY

NEW IDEAS

PREPROCESSING

- PDF EQUAL TO PREPROCESSED NN: $f_i = x^{\alpha_i} (1-x)_i^{\beta} NN(i, x, \ln x)$
- PREPROCESSING EXPONENTS α_i , β_i varied randomly replica by replica
- RANGE OF VARIATION DETERMINED SELF-CONSISTENTLY
- PDF TAKES x, $\ln x$ as inputs



QUARK SINGLET EFFECTIVE α

- NEED TO ITERATE FITS
- POSSIBLE SOURCE OF BIAS?



- RESCALE \Rightarrow UNIFORMLY DISTRIBUTED INPUT (ECDF+INTERPOLATION)
- RERUN]BLUE HYPEROPT
- ONLY ONE INPUT NEEDED



OVERFITTING IN HYPEROPT

- HYPEROPT \Rightarrow OVERFITTING?
- WHAT IS AN "UNNATURAL" PDF?



OVERFITTING METRIC

- RECOMPUTE VALIDATION χ^2
 - SAME TRAINING-VALIDATION SPLIT
 - DIFFERENT FLUCTUATED VALIDATION DATA
- COMPUTE AVERAGE χ^2 & DETERMINE DIFFERENCE TO VALIDATION $\mathcal{R}_O = \langle \chi^2_{val} \chi^2_{val'} \rangle$ (OVERFITNESS)
- **NEGATIVE** OVERFITNESS $\mathcal{R}_O \Rightarrow$ OVERFIT



FOOD FOR THOUGHT

- OPTIMIZED FOLDS
- OVERFITTING AND UNDERFITTING METRICS IN HYPEROPT
- HYPEROPT ON FUTURE TESTS
- BEYOND NEURAL NETS

AN OPEN SOURCE CODE!

- THE FULL NNPDF CODE IS PUBLIC!
- INCLUDING HYPEROPTIMIZATION, EVOLUTION, THEORY, FITTING, VISUALIZATION
- FULLY DOCUMENTED CODE
- LINKS TO CODE (GITHUB), DOCUMENENTATION, INSTALLATION BINARY PACKAGES AVAILABLE FROM http://nnpdf.mi.infn.it/nnpdf-open-source-code/

https://arxiv.org/abs/2109.02671

An open-source machine learning framework for global analyses of parton distributions

The NNPDF Collaboration: Richard D. Ball · Stefano Carrazza · Juan Cruz-Martinez · Luigi Del Debbio · Stefano Forte · Tommaso Giani · Shayan Iranipour · Zahari Kassabov · Jose I. Latorre · Emanuele R. Nocera · Rosalyn L. Pearson · Juan Rojo · Roy Stegeman · Christopher Schwan · Maria Ubiali · Cameron Voisey · Michael Wilson



Fig. 2.1. Workflow for an NNPDF fit

THE WORK OF MANY PEOPLE



NNPDF collaboration and N³PDF team meeting, Gargnano, Italy, September 2021





N. Laurenti



- T. Rabemananjara



