



PDFs WITH 1% ACCURACY

STEFANO FORTE UNIVERSITÀ DI MILANO & INFN



UNIVERSITÀ DEGLI STUDI DI MILANO

DIPARTIMENTO DI FISICA



JULY 6, 2021

CMS SMP-COM MEETING

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





PDFs WITH 1% ACCURACY

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

AMSTERDAM-CAMBRIDGE-EDINBURGH-INFN-MILAN-NIKHEF-SINGAPORE







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

- DELIVERY
- PHENOMENOLOGY
- STABILITY
- VALIDATION
- UNCERTAINTIES
- DATASET SELECTION
- METHODOLOGY
- THEORY
- DATA

NNPDF4.0





DATA



- ABOUT 50 NEW DATASETS & 400 EXTRA DATAPOINTS
- FULL DIS AND FT DY DATASET
 - AS IN NNPDF3.1: FINAL HERA, NMC, BCDMS, CHORUS, NUTEV
 - NOW ALSO NOMAD NEUTRINO
 - SEAQUEST DY
- FULL 7 TEV AND 8 TEV DATASET & EXTENSIVE USE OF 13 TEV DATA:
 - W, Z production: rapidity distributions, asymmetries, $Z p_T$ distributions
 - TOP PAIR PRODUCTION: ALL AVAILABLE DISTRIBUTIONS
 - SINGLE-INCLUSIVE JETS
- SEVERAL NEW PROCESSES:
 - PROMPT PHOTON
 - SINGLE TOP
 - DIJETS
 - HERA JETS

LHC DATA

LHCB

Data set	NNPDF4.0	NNPDF3.1	ABMP16	CT18	MSHT20
LHCb Z 940 pb	1	1	×	×	1
LHC b $Z \to ee$ 2 fb	1	1	1	1	1
LHC b $W,Z \to \mu$ 7 TeV	1	1	1	1	1
LHC b $W\!,Z\to\mu$ 8 TeV	1	 Image: A second s	1	1	 Image: A second s
LHC b $Z \to \mu \mu, ee$ 13 TeV	1	×	×	×	×

ATLAS

Data set	NNPDF4.0	NNPDF3.1	ABMP16	CT18	MSHT20
ATLAS W, Z 7 TeV (2010)	1	1	1	1	1
ATLAS W, Z 7 TeV (2011)	1	1	×	1	1
ATLAS low-mass DY 7 TeV	1	1	×	×	×
ATLAS high-mass DY 7 TeV	1	1	×	×	1
ATLAS W 8 TeV	 Image: A second s	×	×	×	1
ATLAS DY 2D 8 TeV	1	×	×	×	1
ATLAS high-mass DY 2D 8 TeV	1	×	×	×	1
ATLAS $\sigma_{W,Z}$ 13 TeV	1	×	1	×	×
ATLAS W^+ +jet 8 TeV	1	×	×	×	1
ATLAS $Z p_T$ 8 TeV	1	1	×	1	1
ATLAS σ_{tt}^{tot} 7, 8 TeV	1	1	1	×	×
ATLAS σ_{tt}^{tot} 13 TeV	1	1	1	×	×
ATLAS $t\bar{t}$ lepton+jets 8 TeV	1	1	×	1	1
ATLAS $t\bar{t}$ dilepton 8 TeV	1	×	×	×	1
ATLAS single-inclusive jets 7 TeV, $R=0.6$	×	1	×	1	1
ATLAS single-inclusive jets 8 TeV, $R=0.6$	1	×	×	×	×
ATLAS dijets 7 TeV, $R=0.6$	1	×	×	×	×
ATLAS direct photon production 13 TeV	1	×	×	×	×
ATLAS single top R_t 7, 8, 13 TeV	1	×	1	×	×
ATLAS single top diff. 7, 8 TeV	1	×	×	×	×
ATLAS single top diff. 8 TeV	1	×	×	×	×

• CUTOFF DATE AROUND 06/2020

• DIJETS NOW INCLUDED ALONG WITH JETS CANNOT INCLUDE SIMULTANEOUSLY FROM SAME UNDERLYING DATASET CMS

Data set	NNPDF4.0	NNPDF3.1	ABMP16	CT18	MSHT20
CMS W electron asymmetry 7 TeV	1	1	×	1	1
CMS W muon asymmetry 7 TeV	1	1	1	1	×
CMS Drell-Yan 2D 7 TeV	1	1	×	×	1
CMS W rapidity 8 TeV	1	1	1	1	1
CMS $Z p_T$ 8 TeV	1	1	×	1	×
CMS $W + c$ 7 TeV	1	1	×	×	1
CMS $W + c$ 13 TeV	1	×	×	×	×
CMS single-inclusive jets $2.76~{\rm TeV}$	×	✓	×	×	 Image: A second s
CMS single-inclusive jets 7 ${\rm TeV}$	×	1	×	1	1
CMS dijets 7 TeV	1	×	×	×	×
CMS single-inclusive jets 8 ${\rm TeV}$	1	×	×	1	 Image: A second s
CMS 3D dijets 8 TeV	×	×	×	×	×
CMS σ_{tt}^{tot} 5 TeV	1	×	1	×	×
CMS σ_{tt}^{tot} 7, 8 TeV	1	✓	1	×	1
CMS σ_{tt}^{tot} 13 TeV	1	1	1	×	×
CMS $t\bar{t}$ lepton+jets 8 TeV	1	1	×	×	1
CMS $t\bar{t}$ 2D dilepton 8 TeV	1	×	×	1	1
CMS $t\bar{t}$ lepton+jet 13 TeV	1	×	×	×	×
CMS $t\bar{t}$ dilepton 13 TeV	1	×	×	×	×
CMS single top $\sigma_t + \sigma_{\bar{t}}$ 7 TeV	1	×	1	×	×
CMS single top R_t 8, 13 TeV	1	×	1	×	×



ELECTROWEAK CORRECTIONS

- PineAPPL FAST INTERFACE TO Madgraph5_aMC@NLO AVAILABLE (Schwan, Carrazza, Nocera, Zaro 2020)
 ⇒ FULL NLO EW+QCD POSSIBLE
- DATA W/O FSR & PHOTON-INITIATED SUBTRACTION OFTEN NOT AVAILABLE
- CURRENTLY USED FOR DATASET SELECTION: \Rightarrow DISCARDED IF EW CORRNS EXCEED THRESHOLD



NUCLEAR CORRECTIONS

- INCLUDED AS CONTRIBUTION TO COVARIANCE MATRIX (FULLY CORRELATED) (Ball, Nocera, Pearson, 2019)
- COMPUTED AS SHIFT BETWEEN NUCLEAR & STANDARD PDF
- DEUTERIUM PDF DETERMINED FROM SELF-CONSISTENT NNPDF FIT (Ball, Nocera, Pearson, 2019)
- NUCLEAR PDFS FROM NNPDF2.0 (Abdul Khalek, Ethier, Rojo, van Weelden, 2020)



PDF POSITIVITY & INTEGRABILITY

- MS PDFs ARE NON-NEGATIVE!(Candido, Hekhorn, Forte, 2020)
- PDF POSITIVITY IMPOSED (PREVIOUSLY: OBSERVABLE POSITIVITY) \Rightarrow SMALLER LARGE x UNCERTAINTIES
- SEA NONSINGLET COMBINATIONS INTEGRABLE: GOTTFRIED $u + \bar{u} - (d + \bar{d})$ STRANGENESS $u + \bar{u} + (d + \bar{d}) - 2(s + \bar{s})$ \Rightarrow SMALLER SMALL x UNCERTAINTIES



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 $\chi^2 \Rightarrow {\rm OPTIMAL}\; {\rm FIT}$



HYPEROPTIMIZATION

- PARAMETRIZATION AND MINIMIZATION PARAMETERS VARIED
- SCAN OF PARAMETER SPACE
- BAYESIAN UPDATING LEADS TO BEST METHODOLOGY



K-FOLDING



	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 $R_t \ 8 \ {\rm TeV}$		
	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 $\sigma_{tt}^{\rm tot}$	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_{CC}^{p}	HERA I+II inc NC e^+p 820 GeV	LHC b $W,Z \to \mu$ 8 TeV
LHCb $Z \rightarrow \mu \mu$	ATLAS W, Z 7 TeV 2011 Fwd	ATLAS W ⁻ +jet 8 TeV
ATLAS low-mass DY 2011	ATLAS Z p_T 8 TeV (p_T^{ll}, y_{ll})	CMS W rapidity 8 TeV
D0 Z rapidity	CMS dijets 7 TeV	ATLAS single top y_t (normalised
ATLAS single top R_t 13 TeV	CMS single top R_t 13 TeV	





- HYPEROPTIMIZATION \Rightarrow OVERFITTING (χ^2 TOO GOOD)
- CHECK GENERALIZATION POWER: K-FOLDING
 - DIVIDE DATA IN FOLDS
 - EXCLUDE ONE FOLD IN TURN FROM FIT
 - Optimize on the χ^2 of the excluded folds
 - BEST AVERAGE OR BEST WORST

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 $\Lambda^{(int)}$	10	10	10

HYPEROPTIMIZED PARAMETERS



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

DATASET SELECTION



- MISSING HIGHER-ORDER CORRECTIONS
- NO RESUMMATION WHERE NEEDED
- ILL-CONDITIONED COVARIANCE MATRIX
- EXPERIMENTAL ISSUES

THE WEIGHTED FIT METHOD

- FLAG PROBLEMATIC DATASETS:
 - LARGE χ^2
 - LARGE FROBENIUS NUMBER OF COVMAT (EIGENVALUES TOO SMALL)
- REPEAT GLOBAL FIT WITH LARGE WEIGHT GIVEN TO EACH PROBLEMATIC DATASET IN TURN
- χ^2 OF DATASET
 - − UNCHANGED \Rightarrow INTERNAL INCONSISTENCY
 - **DECREASES** \Rightarrow TENSION
- GLOBAL χ^2
 - UNCHANGED \Rightarrow CONSISTENT, KEEP
 - − INCREASES \Rightarrow INCONSISTENT, DISCARD



INCONSISTENT!

- MISSING HIGHER-ORDER CORRECTIONS
- NO RESUMMATION WHERE NEEDED
- ILL-CONDITIONED COVARIANCE MATRIX
- EXPERIMENTAL ISSUES



THE WEIGHTED FIT METHOD

INCONSISTENT!

- FLAG PROBLEMATIC DATASETS:
 - LARGE χ^2
 - LARGE FROBENIUS NUMBER OF COVMAT (EIGEN-VALUES TOO SMALL)
- REPEAT GLOBAL FIT WITH LARGE WEIGHT GIVEN TO EACH PROBLEMATIC DATASET IN TURN



- χ^2 OF DATASET
 - UNCHANGED \Rightarrow INTERNAL INCONSISTENCY
 - **DECREASES** \Rightarrow TENSION
- GLOBAL χ^2
 - UNCHANGED \Rightarrow CONSISTENT, KEEP
 - − INCREASES \Rightarrow INSONSISTENT, DISCARD

UNCERTAINTIES

UNCERTAINTIES: FROM NNPDF3.1...



- 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: ...TO NNPDF4.0



- TYPICAL UNCERTAINTIES IN DATA REGION: SINGLET $\sim 1\%$, NONSINGLET $\sim 2-3\%$
- DATA REGION: $10 \lesssim M_X \lesssim 3 \cdot 10^3$ TeV, $-4 \lesssim y \lesssim 4$



- MOST χ^2 VALUES OF ORDER ONE PER DATAPOINT
- OUTLIERS \Rightarrow FLAGGED DATASETS
- LARGE DATASETS (DIS) WELL FITTED

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

BIAS/VARIANCE

- 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?

experiment	\mathcal{R}_{bv}	bootstrap error	experiment	$\xi_{1\sigma}^{(\rm data)}$	bootstrap error	$\operatorname{erf}(\mathcal{R}_{bv}/\sqrt{2})$	bootstrap error
SeaQuest ATLAS CMS LHCb Total	$\begin{array}{c} 0.99 \\ 1.08 \\ 1.04 \\ 0.92 \\ 1.03 \end{array}$	0.10 0.04 0.06 0.07 0.05	SeaQuest ATLAS CMS LHCb Total	$0.67 \\ 0.69 \\ 0.68 \\ 0.68 \\ 0.69$	$\begin{array}{c} 0.05 \\ 0.02 \\ 0.02 \\ 0.04 \\ 0.02 \end{array}$	$0.69 \\ 0.64 \\ 0.67 \\ 0.72 \\ 0.67$	$\begin{array}{c} 0.05 \\ 0.02 \\ 0.03 \\ 0.03 \\ 0.03 \end{array}$

RESULTS

ONE SIGMA

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

2				
χ^2	WITHOUT/	WITH	PDF	UNC.

STABILITY

PARAMETRIZATION BASIS

- PDFS BY DEFAULT PARAMETRIZED IN "EVOLUTION BASIS": SINGLET $\Sigma = \sum_i q_i + \bar{q}_i$, VALENCE $V = \sum_i q_i - \bar{q}_i$, TRIPLET $T_3 = u + \bar{u} - (d + \bar{d})$, ...
- WHAT IF ONE CHOOSES THE "FLAVOR BASIS": u, \bar{u}, d, \bar{d} ?
- COMPLETE STABILITY OF RESULTS!



NNPDF4.0 vs. NNPDF3.1

- FULL BACKWARD COMPATIBILITY
- SUBSTANTIAL REDUCTION IN UNCERTAINTY





NNPDF4.0 vs DIS-only

- DIS-ONLY FIT NO LONGER COMPETITIVE
- HADRONIC DATA NEEDED FOR PRECISION





NNPDF4.0 VS COLLIDER ONLY

- COLLIDER ONLY COMPETITIVE!
- ONLY DEUTERIUM FIXED-TARGET DATA STILL RELEVANT





PHENOMENOLOGY

REPRESENTATIVE PROCESSES

- EW CORRECTIONS UNDER CONTROL
- **SMALL** UNCERTAINTIES





AN OPEN SOURCE CODE!

- THE FULL NNPDF CODE WILL BE MADE PUBLIC!
- INCLUDING HYPEROPTIMIZATION, EVOLUTION, THEORY, FITTING, VISUALIZATION
- FULLY DOCUMENTED CODE

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

OUTLOOK

THE ACCURACY CHALLENGE

- AT 1%, NOT ALL REDUCED DATASETS AGREE
- MUST INCLUDE MISSING HIGHER-ORDER CORRECTION UNCERTAINTIES
- INCLUDE **EW** CORRECTIONS
- GO BEYOND K-Factors
- WEIGHT SMALL DATASETS

THE ACCURACY CHALLENGE

- AT 1%, NOT ALL REDUCED DATASETS AGREE
- MUST INCLUDE MISSING HIGHER-ORDER CORRECTION UNCERTAINTIES
- INCLUDE **EW** CORRECTIONS
- GO BEYOND K-Factors
- WEIGHT SMALL DATASETS

STAY TUNED FOR NNPDF4.1!!