## Deep-Inelastic Scattering and Collinear Physics

#### Second European School on the Physics of the EIC and Related Topics

#### Emanuele R. Nocera Università degli Studi di Torino and INFN, Torino

23-24 June 2025



# Summary of Lecture 1

- $\textcircled{0} Deep Inelastic Scattering has been, is, and will be a crucial laboratory of QCD \\ \longrightarrow hadronic structure encoded in unpolarised and polarised PDFs$ 
  - $\longrightarrow$  PDFs are related to physical observales via factorisation and evolution
  - $\longrightarrow$  (critically different) qualitative PDF features are driven by this theoretical framework
- PDFs are a limiting factor for precision and discovery
  - $\longrightarrow$  unpolarised PDFs: SM and BSM physics at the LHC
  - $\longrightarrow$  polarised PDFs: contribution of partons' spin to the proton spin
- PDFs are determined from experimental data by means of parametric regression → need to define data, theory, and methodology
- Oifferent physical observables constrain different PDF combinations
  - $\longrightarrow$  fixed-target NC DIS: u and d
  - $\longrightarrow$  fixed-target CC DIS: s and  $\bar{s}$
  - $\rightarrow$  HERA NC and CC DIS:  $u, \bar{u}, d, \bar{d}, g$  (scaling violations and tagged DIS)
  - $\longrightarrow$  fixed-target DY: u and d at large x
  - $\longrightarrow$  collider DY: u,  $\bar{u}$ , d,  $\bar{d}$ , s
  - $\longrightarrow$  collider DY+c: s (W) and c (Z)
  - $\longrightarrow Zp_T$ ,  $t\bar{t}$ , jets: g
  - $\longrightarrow$  only a small fraction of the above is available for polarised PDFs

Lecture 2: Data, theoretical, and methodological accuracy in PDF determination

## The ingredients of PDF determination



Each of these ingredients is a source of uncertainty in the PDF determination Each of these ingredients require to make choices which lead to different PDF sets

Emanuele R. Nocera (UNITO)

DIS and Collinear Physics

# Overview of current PDF determinations

	NNPDF4.0	MSHT20	CT18	HERAPDF2.0	CJ22	ABMP16
Fixed-target DIS	Ø	Ń	Ø	$\boxtimes$	Ø	Ø
JLAB	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\checkmark$	$\boxtimes$
HERA I+II	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
HERA jets	Ø	$\boxtimes$	$\boxtimes$	Ń	$\boxtimes$	$\boxtimes$
Fixed target DY	$\square$	$\square$	$\square$	$\boxtimes$	$\square$	$\square$
Tevatron $\boldsymbol{W}$ , $\boldsymbol{Z}$	Ø	Ø	Ø	$\boxtimes$	Ø	Ø
LHC vector boson	Ø	$\checkmark$	$\checkmark$	$\boxtimes$	$\square$	$\square$
LHC $W + c \ Z + c$	Ø	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$
Tevatron jets	Ø	Ø	Ø	$\boxtimes$	$\square$	$\boxtimes$
LHC jets	Ø	Ø	$\checkmark$	$\boxtimes$	$\boxtimes$	$\boxtimes$
LHC top	Ø	$\checkmark$	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\square$
LHC single $t$	Ø	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$
LHC prompt $\gamma$		$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$
statistical treatment	Monte Carlo	Hessian $\Delta\chi^2$ dynamical	Hessian $\Delta\chi^2$ dynamical	Hessian $\Delta \chi^2 = 1$	Hessian $\Delta \chi^2 = 1.645$	Hessian $\Delta \chi^2 = 1$
parametrisation	Neural Network	Chebyschev pol.	Bernstein pol.	polynomial	polynomial	polynomial
HQ scheme	FONLL	TR'	ACOT- $\chi$	TR'	ACOT- $\chi$	FFN
accuracy	aN <sup>3</sup> LO	aN <sup>3</sup> LO	NNLO	NNLO	NLO	NNLO
latest update	EPJ C82 (2022) 428	EPJ C81 (2021) 341	PRD 103 (2021) 014013	EPJ C82 (2022) 243	PRD 107 (2023) 113005	PRD 96 (2017) 014011
All PDF sets are available as $(x, Q^2)$ interpolation grids through the LHAPDF library						

Emanuele R. Nocera (UNITO)

DIS and Collinear Physics

4 / 75

Parton Distribution Functions Lecture 2: Data, Theoretical and Methodological Accuracy in PDF Determination

# Outline

- 2.1 Which data constrain which PDFs? DIS observables non-DIS observables
- 2.2 Can we improve the fit quality by improving the theory? heavy quarks missing higher order uncertainties
- 2.3 Why is the methodology important? parametrisation optimisation uncertainty representation validation of uncertainties PDF benchmarks

I will focus on a limited selection of recent results

I will not talk about some very interesting topics (*e.g.* aN<sup>3</sup>LO PDFs, interplay between fitting PDFs and New Physics, non parametric regression models, ...)

See also lectures by M. Ubiali and the ML hands-on session by P. Zurita

Emanuele R. Nocera (UNITO)

Parton Distribution Functions

#### 1.4 Data

# Overview of experimental data



 $N_{\rm dat} = 4618$ 

# Deep Inelastic Scattering

Kinematic coverage



#### Deep Inelastic Scattering

Re-write the cross section in terms of structure functions  $(F_2, F_3, F_L)$  $\frac{d^2\sigma^i}{dxdy} \propto Y_+ F_2^i \mp Y_- x F_3^i - y^2 F_L^i$   $Y_{\pm} = 1 \pm (1-y)^2$   $F_L^i = F_2^i - 2x F_1^i$   $i = \text{NC}(\gamma, Z, \gamma Z), \text{CC}(W^{\pm})$ 

NC DIS  $(\ell p \rightarrow \ell X)$  at LO (NMC, SLAC, BCDMS, HERA)

$$\begin{split} \left[ F_{2}^{\gamma}, F_{2}^{\gamma Z}, F_{2}^{Z} \right] &= x \sum_{q} \left[ e_{q}^{2}, 2e_{q}g_{V}^{q}, (g_{V}^{q})^{2} + (g_{A}^{q})^{2} \right] (q + \bar{q}) \\ \left[ F_{3}^{\gamma}, F_{3}^{\gamma Z}, F_{3}^{Z} \right] &= \sum_{q} \left[ 0, 2e_{q}g_{A}^{q}, 2g_{V}g_{A}^{q} \right] (q - \bar{q}) \\ \left[ F_{L}^{\gamma}, F_{L}^{\gamma Z}, F_{L}^{Z} \right] &= \left[ 0, 0, 0 \right] \end{split}$$

CC DIS ( $\ell^- p \to \nu X$  or  $\bar{\nu} p \to \ell^+ X$ ) at LO (CHORUS, NuTeV, HERA)

$$\begin{bmatrix} F_2^{W^-}, F_2^{W^-} \end{bmatrix} = 2x \left[ (u + \bar{d} + \bar{s} + c \dots), (d + \bar{u} + \bar{c} + s \dots) \right]$$
$$\begin{bmatrix} F_3^{W^-}, F_3^{W^+} \end{bmatrix} = 2 \left[ (u - \bar{d} - \bar{s} + c \dots), (d - \bar{u} - \bar{c} + s \dots) \right]$$
$$\begin{bmatrix} F_L^{W^+}, F_L^{W^-} \end{bmatrix} = [0, 0]$$

isospin  $(p \to n)$ :  $u^p = d^n$   $d^p = u^n$   $\bar{u}^p = \bar{d}^n$   $\bar{d}^p = \bar{u}^n$ 

deuteron target approximated as the average of one proton and one neutron  $Q^2 \ge Q_{\min}^2 \sim 1 \text{ GeV}^2$   $W^2 = Q^2(1-x)/x \ge W_{\min}^2 \sim 10 \text{ GeV}^2$ 

Emanuele R. Nocera (UNITO)

#### Deep Inelastic Scattering



NMC, NPB 487 (1997) 3; CHORUS, PLB 632 (2006) 65



Emanuele R. Nocera (UNITO)

Parton Distribution Functions

# Drell-Yan

#### Kinematic coverage



#### Drell-Yan

Work out the cross section differential in the rapidity of the lepton pair

$$\frac{d\sigma^{i}}{dy} \propto A^{i}\mathcal{L}^{i} \qquad i = \mathrm{NC}(\gamma, \mathbf{Z}, \gamma \mathbf{Z}), \mathrm{CC}(\mathbf{W}^{\pm})$$

$$\mathrm{NC} \ \mathrm{DY} \ (\gamma, Z) \ \mathrm{at} \ \mathrm{LO} \qquad \qquad \mathrm{CC} \ \mathrm{DY} \ (W^{\pm}) \ \mathrm{at} \ \mathrm{LO}$$

$$A^{\gamma}\mathcal{L}^{\gamma}, A^{Z}\mathcal{L}^{Z} \Big] = \left[\sum_{q} e_{q}^{2}q\bar{q}, \sum_{q,q'} |V_{qq'}^{\mathrm{CKM}}|q\bar{q'}\right] \qquad \left[A^{W^{\pm}}\mathcal{L}^{W^{\pm}}\right] = \left[\sum_{q} \left((g_{V}^{q})^{2} + (g_{A}^{q})^{2}\right) q\bar{q}\right]$$

isospin  $(p \to n)$ :  $u^p = d^n$   $d^p = u^n$   $\bar{u}^p = \bar{d}^n$   $\bar{d}^p = \bar{u}^n$ 

deuteron target approximated as the average of one proton and one neutron different experiments measure different cross section combinations

$$\begin{split} & \frac{\sigma_{pn}^Z}{\sigma_{pp}^P} \approx \frac{4/9u\bar{d}+1/9d\bar{u}}{4/9u\bar{u}+1/9d\bar{d}} \to \frac{\bar{d}}{\bar{u}} & \text{DY } p/d \text{ asymmetry (NuSea, SeaQuest)} \\ & \frac{\sigma_{p\bar{p}}^{W^+}}{\sigma_{p\bar{p}}^{W^+}} \approx \frac{ud+\bar{d}\bar{u}}{du+\bar{u}\bar{d}} \to \frac{ud}{du} & W^{\pm} \text{ asymmetry (CDF, D0)} \\ & \sigma_{pp}^{W^+} \approx u\bar{d}+c\bar{s} \quad \sigma_{pp}^Z \approx u\bar{u}+d\bar{d}+s\bar{s}\to s, \bar{s} & W^{\pm} \text{ and } Z \text{ production (ATLAS, CMS, LHCb)} \\ & \frac{\sigma_{pp}^{W^+}}{\sigma_{pp}^{W^-}} \approx \frac{u\bar{d}+\bar{d}u}{d\bar{u}+\bar{u}d} \to \bar{u}-\bar{d} & W^{\pm} \text{ muon asymmetry (ATLAS, CMS)} \end{split}$$

13 / 75

#### Drell-Yan



Emanuele R. Nocera (UNITO)

Parton Distribution Functions

28 August 2024 14 / 75

# $Zp_T$ , $t\bar{t}$ , Single-Inclusive Jet, and Dijet Production



Emanuele R. Nocera (UNITO)

# $Zp_T$ , $t\bar{t}$ , Single-Inclusive Jet, and Dijet Production

Various differential distributions, all proportional to the gluon PDF at LO

$$\underline{Zp_T \text{ production}}: \frac{d\sigma^2}{dp_T^2 dm_{\ell\ell}}, \frac{d\sigma^2}{dp_T^2 dy_Z} \text{ (ATLAS) } \frac{d\sigma}{dp_T^2} \text{ (CMS)}$$

need one final-state parton, then initial-state quark and gluon are on the same footing wide  $p_T$  range, constraints on a wide x (typically intermediate) and  $Q^2$  range

$$t\bar{t}$$
 production:  $\frac{d\sigma}{dp_T^t}$ ,  $\frac{d\sigma}{dy^t}$ ,  $\frac{d\sigma}{dy^{t\bar{t}}}$ ,  $\frac{d\sigma}{dm^{t\bar{t}}}$  (ATLAS and CMS)

process initiated by two gluons in the initial state

differential cross sections reconstructed at parton level (additional systematics) normalise by  $\sigma_{tot}^{t\bar{t}}$  (systematics largely cancel, but loose control on PDF shape) particle-level cross sections are theoretically more complicated

wide rapidity range, constraints on the large-x region

single-inclusive jet and dijet production:  $\frac{d\sigma^2}{dydp_T}$ ,  $\frac{d\sigma^2}{dy_{1,2}dm_{1,2}}$  (HERA, ATLAS and CMS)

process initiated by two gluons in the initial state

be careful with systematic uncertainties (mostly driven by jet energy reconstruction) can also be measured in DIS

wide rapidity range, constraints on the large-x region

Emanuele R. Nocera (UNITO)

Parton Distribution Functions

# $Zp_T$ , $t\bar{t}$ , Single-Inclusive Jet, and Dijet Production



Emanuele R. Nocera (UNITO)

28 August 2024 17 / 75

## Scaling Violations and Heavy Flavour Production

Scale dependence in DIS of singlet structure function

 $\frac{d}{d\ln(Q^2)}F_2^{\Sigma}\approx \frac{\alpha_s}{2\pi}\left[\gamma_{qq}f_{\Sigma}+2n_f\gamma_{gq}f_g\right]$ 

ANOMALOUS DIMENSIONS 20 15 10 5  $\gamma_{qg}$ 0 -5 ō **pp** 2

the gluon PDF can be determined at small x from DIS scaling violations (from HERA)

Heavy quark production in DIS initiated by gluons





the gluon PDF is determined by tagging a c or a b quark in the final state (HERA)

# Single t, Direct $\gamma$ , W, Z+c-jets



# Single t, Direct $\gamma$ , W, Z+c-jets

Other processes, currently limited by experimental uncertainties Single t production (t-channel):  $\frac{d\sigma}{dp_{\pi}^{t}}$ ,  $\frac{d\sigma}{dy^{t}}$ ,  $\frac{d\sigma}{dy^{t}}$ ,  $\frac{d\sigma}{dy^{t}}$  (ATLAS, CMS) partonic cross sections similar to CC DIS; t reconstructed from Wb decay potential sensitivity to  $\bar{u}$  and  $\bar{d}$ , also through ratios of t and  $\bar{t}$  production potential currently limited by large experimental uncertainties <u>Prompt  $\gamma$  production</u>:  $\frac{d\sigma^2}{dE_{\gamma}^2 u^{\gamma}}$  (ATLAS and CMS) gluon-quark-initiated Compton scattering potential sensitivity to the gluon PDF potential currently limited by large experimental uncertainites W, Z + charm-tagged jets:  $\frac{d\sigma}{du}$  (ATLAS and CMS) W + c: sensitivity to strange PDF (and  $s - \bar{s}$  asymmetry)

W + Z: sensitivity to charm PDF (including intrinsic charm and  $c - \bar{c}$  asymmetry) be careful with systematic uncertainties (due to jet tagging algorithm)

More exclusive processes: double gauge boson production, multijet production, ...

generally less precise and potentially contaminated by BSM physics

Emanuele R. Nocera (UNITO)

Parton Distribution Functions

# Single t, Prompt $\gamma$ , W, Z+c-jets





	Hadronic Process	Partonic Process	PDFs probed	x coverage
Lepton-nucleon	$\ell^{\pm}\{p,n\} \to \ell^{\pm} + X$	$\gamma^* q \to q$	q, ar q, g	$x \gtrsim 0.01$
	$\ell^{\pm} n/p \to \ell^{\pm} + X$	$\gamma^* d/u  o d/u$	d/u	$x\gtrsim 0.01$
	$\nu(\bar{\nu})N \to \mu^-(\mu^+) + X$	$W^*q \rightarrow q'$	$q,ar{q}$	$0.01 \lesssim x \lesssim 0.5$
	$\nu N \to \mu^- \mu^+ + X$	$W^*s \rightarrow c$	s	$0.01 \lesssim x \lesssim 0.2$
	$\bar{\nu}N \to \mu^+\mu^- + X$	$W^* \bar{s} \rightarrow \bar{c}$	$\overline{s}$	$0.01 \lesssim x \lesssim 0.2$
	$e^{\pm}p \rightarrow e^{\pm} + X$	$\gamma^* q \to q$	$g,q,ar{q}$	$0.0001 \lesssim x \lesssim 0.1$
	$e^+p \to \bar{\nu} + X$	$W^+\{d,s\} \to \{u,c\}$	d, s	$x\gtrsim 0.01$
	$e^{\pm}p \rightarrow e^{\pm}c\bar{c} + X$	$\gamma^* c \to c, \gamma^* g \to c \bar{c}$	c, g	$0.0001 \lesssim x \lesssim 0.1$
	$e^{\pm}p \rightarrow jet(s) + X$	$\gamma^*g  o q \bar{q}$	g	$0.01 \lesssim x \lesssim 0.1$
Proton-(anti)proton	$pp \to \mu^+ \mu^- + X$	$u\bar{u}, d\bar{d} \to \gamma^*$	$\bar{q}$	$0.015 \lesssim x \lesssim 0.35$
	$pn/pp \to \mu^+\mu^- + X$	$(u\bar{d})/(u\bar{u}) \rightarrow \gamma^*$	$ar{d}/ar{u}$	$0.015 \lesssim x \lesssim 0.35$
	$p\bar{p}(pp) \rightarrow jet(s) + X$	gg, qg, qq  ightarrow jets	g, q	$0.005 \lesssim x \lesssim 0.5$
	$p\bar{p} \to (W^{\pm} \to \ell^{\pm}\nu) + X$	$ud  ightarrow W^+$ , $ar{u}ar{d}  ightarrow W^-$	$u, d, ar{u}, ar{d}$	$x\gtrsim 0.05$
	$pp \to (W^{\pm} \to \ell^{\pm} \nu) + X$	$u \bar{d}  ightarrow W^+$ , $d \bar{u}  ightarrow W^-$	$u,d,\bar{u},\bar{d},(g)$	$x\gtrsim 0.001$
	$p\bar{p}(pp) \to (Z \to \ell^+ \ell^-) + X$	$uu, dd(u\bar{u}, d\bar{d}) \to Z$	u, d(g)	$x\gtrsim 0.001$
	$pp \rightarrow (W+c) + X$	$gs \to W^- c, g\bar{s} \to W^+ \bar{c}$	$s, \bar{s}$	$x \sim 0.01$
	$pp \rightarrow (Z+c) + X$	$gc \rightarrow Zc, g\bar{c} \rightarrow Z\bar{c}$	$c,ar{c}$	$x \sim 0.01$
	$pp \rightarrow t\bar{t} + X$	$gg  ightarrow t \bar{t}$	g	$x \sim 0.1$
	$pp \rightarrow t, \bar{t} + X$	$gq  ightarrow t, ar{t}q$	u, d	$x \sim 0.01$
	$pp \rightarrow \gamma + X$	$gq \rightarrow \gamma q$	g	$x \sim 0.01$

#### All the data together

Emanuele R. Nocera (UNITO)

	NNPDF4.0	MSHT20	CT18	HERAPDF2.0	CJ22	ABMP16
Fixed-target DIS	Ø	Ø	Ń	$\boxtimes$	Ø	Ø
JLAB	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\checkmark$	$\boxtimes$
HERA I+II	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
HERA jets	$\checkmark$	$\boxtimes$	$\boxtimes$	$\checkmark$	$\boxtimes$	$\boxtimes$
Fixed target DY	$\checkmark$	$\checkmark$	$\checkmark$	$\boxtimes$	$\checkmark$	$\checkmark$
Tevatron $W$ , $Z$	$\checkmark$	$\checkmark$	$\checkmark$	$\boxtimes$	$\checkmark$	$\checkmark$
LHC vector boson	$\checkmark$	$\checkmark$	$\checkmark$	$\boxtimes$	$\checkmark$	$\checkmark$
LHC $W + c \ Z + c$	$\checkmark$	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$
Tevatron jets	$\checkmark$	$\checkmark$	$\checkmark$	$\boxtimes$	$\checkmark$	$\boxtimes$
LHC jets	$\checkmark$	$\checkmark$	$\checkmark$	$\boxtimes$	$\boxtimes$	$\boxtimes$
LHC top	$\square$	$\checkmark$	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\checkmark$
LHC single $t$	$\checkmark$	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$
LHC prompt $\gamma$	$\checkmark$	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$
statistical treatment	Monte Carlo	Hessian $\Delta\chi^2$ dynamical	Hessian $\Delta\chi^2$ dynamical	Hessian $\Delta \chi^2 = 1$	Hessian $\Delta \chi^2 = 1.645$	Hessian $\Delta \chi^2 = 1$
parametrisation	Neural Network	Chebyschev pol.	Bernstein pol.	polynomial	polynomial	polynomial
HQ scheme	FONLL	TR'	ACOT- $\chi$	TR'	ACOT- $\chi$	FFN
accuracy	$aN^3LO$	$aN^3LO$	NNLO	NNLO	NLO	NNLO
latest update	EPJ C82 (2022) 428	EPJ C81 (2021) 341	PRD 103 (2021) 014013	EPJ C82 (2022) 243	PRD 107 (2023) 113005	PRD 96 (2017) 014011
All PDF sets are available as $(x, Q^2)$ interpolation grids through the LHAPDF library						

#### Overview of current unpolarised PDF determinations

Emanuele R. Nocera (UNITO)

Parton Distribution Functions

28 August 2024

23 / 75

	Hadronic Process	Partonic Process		PDFs probed	x coverage	
Lepton-nucleon	$ \ell^{\pm}\{p,n\} \to \ell^{\pm} + \\ \ell^{\pm}\{p,n\} \to \ell^{\pm} + h $	$ \begin{array}{ccc} X & \gamma^* q \to q \\ + X & \gamma^* q \to q \end{array} $		$egin{array}{lll} \Delta q, \Delta ar q, \Delta g \ \Delta q, \Delta ar q, \Delta ar q, \Delta g \end{array}$	$\begin{aligned} x \gtrsim 0.003 \\ 0.005 \lesssim x \lesssim 0.5 \end{aligned}$	
Proton-proton	$pp \to jet(s) + X$ $pp \to (W^{\pm} \to \ell^{\pm} \nu)$ $pp \to h + X$	$\begin{array}{ccc} gg, qg, \\ +X & u\bar{d} \rightarrow W \\ gg, qg \end{array}$	$gg, qg, qq \rightarrow jets$ $\cdot X  u\bar{d} \rightarrow W^+, d\bar{u} \rightarrow W^-  \Delta u,$ $gg, qg, qq \rightarrow h$		$\begin{array}{l} 0.05 \lesssim x \lesssim 0.2 \\ 0.05 \lesssim x \lesssim 0.4 \\ 0.05 \lesssim x \lesssim 0.4 \end{array}$	
		NNPDFpol2.0	MAPPDFpol1.0	BDSSV25	JAM25	
Fixed-target DIS Fixed-target SIDI RHIC vector boso RHIC jets RHIC hadrons	5 n				1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	
statistical treatme	ent	Monte Carlo	Monte Carlo	Monte Carlo	Monte Carlo	
parametrisation		Neural Network	polynomial	polynomial	polynomial	
HQ scheme		FONLL	ZM	ZM	ZM	
accuracy		NNLO	NNLO	NNLO	NLO	
latest update		arXiv:2503.11814	PLB 865 (2025) 1394	197 PRL 133 (2024) 15	arXiv:2506.13616	

And polarised PDFs?

#### Emanuele R. Nocera (UNITO)

Parton Distribution Functions

24 / 75

#### Unpolarised and polarised PDFs: data



Thousands of data for unpolarised PDFs; hundreds of data for unpolarised PDFs

Emanuele R. Nocera (UNITO)

Parton Distribution Functions

#### Impact of future data: HL-LHC



26 / 75

#### Impact of future data: unpolarised EIC



PRD 103 (2021) 096005; PRD 109 (2024) L091501

#### Impact of future data: polarised EIC



#### Impact of future data: FPF



EPJ C84 (2024) 369

## 2.1 Theory

# Can we improve the fit quality by improving the theory?



31 / 75

2.1.1 Heavy Quarks

## Heavy Quarks in DIS

Two possible factorisation schemes for DIS structure functions  $\overline{\rm MS} \ \ \ scheme$ 

Heavy quarks are treated as massless (zero-mass scheme) corrections proportional to  $\ln(Q^2/m_h^2)$  are resummed to all orders by DGLAP corrections that are  $\mathcal{O}(m_h^2/Q^2)$  are neglected This scheme is appropriate when  $Q^2 \gg m_h^2$ <u>Decoupling scheme</u> Heavy quarks are treated as massive (massive scheme) corrections proportional to  $\ln(Q^2/m_h^2)$  are treated at fixed order corrections that are  $\mathcal{O}(m_h^2/Q^2)$  are included

This scheme is appropriate when  $Q^2 \sim m_h^2$ 

The third way: match the two schemes

General-mass variable-flavour number schemes (ACOT, S-ACOT, TR, FONLL, ...)

use  $\overline{\mathrm{MS}}$  for  $Q^2 \gg m_h^2$  with full mass dependence retained

keep all flavour sin running DGLAP

subtract double counting terms

Emanuele R. Nocera (UNITO)

Parton Distribution Functions

# Example: the $g_1$ and $g_1^c$ structure functions at the EIC



Emanuele R. Nocera (UNITO)

# 2.1.2 Missing Higher Order Uncertainties

#### Perturbative Accuracy in PDF Determination

NNLO is the precision frontier for PDF determination

N3LO is the precision frontier for partonic cross sections

Mismatch between perturbative order of partonic cross sections and accuracy of PDFs may become a significant source of uncertainty

$$\hat{\sigma} = \alpha_s^p \hat{\sigma}_0 + \alpha_s^{p+1} \hat{\sigma}_1 + \alpha_s^{p+2} \hat{\sigma}_2 + \mathcal{O}(\alpha_s^{p+3}) \qquad \delta(\text{PDF} - \text{TH}) = \frac{1}{2} \left| \frac{\sigma_{\text{NNLO-PDFs}}^{(2)} - \sigma_{\text{NLO-PDFs}}^{(2)}}{\sigma_{\text{NNLO-PDFs}}^{(2)}} \right|$$



Emanuele R. Nocera (UNITO)

28 August 2024 36 / 75
## MHOUs and Scale Variations

As an example, let us consider the NS DIS structure function

$$F_2^{\rm NS}(N,Q^2) = xC_{\rm NS}(N,\alpha_s(Q^2)) \exp\left[\int_{Q_0^2}^{Q^2} \frac{d\lambda^2}{\lambda} \gamma_{\rm NS}\left(N,\alpha_s(\mu^2)\right)\right] f_{\rm NS}(Q_0^2)$$

Sources of MHOUs  

$$\gamma_{\rm NS}^{\rm N^kLO}(N,\alpha_s) = \alpha_s \gamma_{\rm NS}^{(0)} + \alpha_s^2 \gamma_{\rm NS}^{(1)} + \dots \alpha_s^{k+1} \gamma_{\rm NS}^{(k)}$$

$$C_{\rm NS}^{\rm N^kLO}(N,\alpha_s) = 1 + \alpha_s C_{\rm NS}^{(1)} + \dots \alpha_s^k C_{\rm NS}^{(k)}$$

Scale variations  
Idea: 
$$\alpha_s(\kappa^2\mu^2) = \alpha_s(\mu^2)[1 + \mathcal{O}(\alpha_s)]$$

at N^kLO differences due to higher orders are related to the QCD  $\beta$  function up to  $\beta_k$ 

$$\begin{split} \bar{C}_{\rm NS}(\alpha_s(\kappa_r^2\mu^2,\kappa_r^2)) &= C_{\rm NS}(\alpha_s(\mu^2))[1+\mathcal{O}(\alpha_s)] \text{ fixes } \bar{C}^{(k)} \text{ in terms of } C^{(k)} \\ \bar{\gamma}_{\rm NS}(\alpha_s(\kappa_f^2\mu^2,\kappa_f^2)) &= \gamma_{\rm NS}(\alpha_s(\mu^2))[1+\mathcal{O}(\alpha_s)] \text{ fixes } \bar{\gamma}^{(k)} \text{ in terms of } \gamma^{(k)} \\ \Delta C_{\rm NS} &= \bar{C}_{\rm NS}(\alpha_s(\kappa_r^2\mu^2,\kappa_r^2)) - C_{\rm NS}(\alpha_s(\mu^2)) \end{split}$$

renormalisation scale (at which UV divergences are subtracted)  $\mu_r = \kappa_r \mu$ 

$$\Delta \gamma_{\rm NS} = \bar{\gamma}_{\rm NS} (\alpha_s (\kappa_f^2 \mu^2, \kappa_f^2)) - \gamma_{\rm NS} (\alpha_s (\mu^2))$$
  
factorisation scale (at which collinear divergences are factorised)  $\mu_f = \kappa_f \mu$ 

Propagate  $\Delta C$  and  $\Delta \gamma$  into  $\Delta f$ 

#### Scale Variations: Prescriptions



Vary  $\mu_r$  and  $\mu_f$  about  $\mu_0$ 

Pick a set of possible variations

3-points:  $\mu_r = \mu_f$ ,  $\kappa_{r,f} = 2, 1/2$ 

7-points:  $\mu_r, \mu_f$  varied independently,  $\kappa_{r,f} = 2, 1/2$ , remove  $\mu_r/\mu_f = 4$ 

9-points:  $\mu_r, \mu_f$  varied independently,  $\kappa_{r,f} = 2, 1/2$ 

To estimate MHOUs, take the envelope, *i.e.* the difference between the largest and smallest predictions

# A Theory Covariance Matrix

Assuming that theory uncertainties are (a) Gaussian and (b) independent from experimental uncertainties, modify the figure of merit to account for theory errors

$$\chi^{2} = \sum_{i,j}^{N_{\text{dat}}} (D_{i} - T_{i}) (\operatorname{cov}_{\exp} + \operatorname{cov}_{\operatorname{th}})_{ij}^{-1} (D_{j} - T_{j}); \ (\operatorname{cov}_{\operatorname{th}})_{ij} = \frac{1}{N} \sum_{k}^{N} \Delta_{i}^{(k)} \Delta_{j}^{(k)}; \ \Delta_{i}^{(k)} \equiv T_{i}^{(k)} - T_{i}$$

Problem reduced to estimate the th. cov. matrix, e.g. in terms of nuisance parameters

$$\Delta_i^{(k)} = T_i(\mu_R, \mu_F) - T_i(\mu_{R,0}, \mu_{F,0});$$
 vary scales in  $\frac{1}{2} \le \frac{\mu_F}{\mu_{F,0}}, \frac{\mu_R}{\mu_{R,0}} \le 2$ 



# A Theory Covariance Matrix

Assuming that theory uncertainties are (a) Gaussian and (b) independent from experimental uncertainties, modify the figure of merit to account for theory errors

$$\chi^{2} = \sum_{i,j}^{N_{\text{dat}}} (D_{i} - T_{i}) (\operatorname{cov}_{\exp} + \operatorname{cov}_{\operatorname{th}})_{ij}^{-1} (D_{j} - T_{j}); \ (\operatorname{cov}_{\operatorname{th}})_{ij} = \frac{1}{N} \sum_{k}^{N} \Delta_{i}^{(k)} \Delta_{j}^{(k)}; \ \Delta_{i}^{(k)} \equiv T_{i}^{(k)} - T_{i}$$

Problem reduced to estimate the th. cov. matrix, e.g. in terms of nuisance parameters

$$\Delta_i^{(k)} = T_i(\mu_R, \mu_F) - T_i(\mu_{R,0}, \mu_{F,0}); \text{ vary scales in } \frac{1}{2} \le \frac{\mu_F}{\mu_{F,0}}, \frac{\mu_R}{\mu_{R,0}} \le 2$$



# A Theory Covariance Matrix

Assuming that theory uncertainties are (a) Gaussian and (b) independent from experimental uncertainties, modify the figure of merit to account for theory errors

$$\chi^{2} = \sum_{i,j}^{N_{\text{dat}}} (D_{i} - T_{i}) (\operatorname{cov}_{\exp} + \operatorname{cov}_{\operatorname{th}})_{ij}^{-1} (D_{j} - T_{j}); \ (\operatorname{cov}_{\operatorname{th}})_{ij} = \frac{1}{N} \sum_{k}^{N} \Delta_{i}^{(k)} \Delta_{j}^{(k)}; \ \Delta_{i}^{(k)} \equiv T_{i}^{(k)} - T_{i}$$

Problem reduced to estimate the th. cov. matrix, e.g. in terms of nuisance parameters



#### Impact on Parton Distributions



Faster perturbative convergence when MHOU are incorporated into PDFs

EPJ C79 (2019) 838; ibid. 931; EPJ C84 (2024) 517

## Impact on Uncertainties and Fit Quality



Dataset	N	NEO		ININEO	
	1'dat	no MHOU	MHOU	no MHOU	MHOU
DIS NC	2100	1.30	1.22	1.23	1.20
DIS CC	989	0.92	0.87	0.90	0.90
DY NC	736	2.01	1.71	1.20	1.15
DY CC	157	1.48	1.42	1.48	1.37
Top pairs	64	2.08	1.24	1.21	1.43
Single-inclusive jets	356	0.84	0.82	0.96	0.81
Dijets	144	1.52	1.84	2.04	1.71
Prompt photons	53	0.59	0.49	0.75	0.67
Single top	17	0.36	0.35	0.36	0.38
Total	4616	1.34	1.23	1.17	1.13

Overall (rather small) variation of uncertainties. Tensions relieved: improvement in  $\chi^2$ [EPJC79 (2019) 838; ibid. 931; EPJC84 (2024) 517]

Emanuele R. Nocera (UNITO)

# What Happens at aN<sup>3</sup>LO?

Dataset	$N_{\mathrm{dat}}$	NLO no MHOU	мнои	$N_{\mathrm{dat}}$	NNLO no MHOU	мнои	$N_{\mathrm{dat}}$	aN <sup>3</sup> LO no MHOU	мнои
DIS NC	1980	1.30	1.22	2100	1.22	1.20	2100	1.22	1.20
DIS CC	988	0.92	0.87	989	0.90	0.90	989	0.91	0.92
DY NC	667	1.49	1.32	736	1.20	1.15	736	1.17	1.16
DY CC	193	1.31	1.27	157	1.45	1.37	157	1.37	1.36
Top pairs	64	1.90	1.24	64	1.27	1.43	64	1.23	1.41
Single-inclusive jets	356	0.86	0.82	356	0.94	0.81	356	0.84	0.83
Dijets	144	1.55	1.81	144	2.01	1.71	144	1.78	1.67
Prompt photons	53	0.58	0.47	53	0.76	0.67	53	0.72	0.68
Single top	17	0.35	0.34	17	0.36	0.38	17	0.35	0.36
Total	4462	1.24	1.16	4616	1.17	1.13	4616	1.15	1.14

Fit quality improves with perturbative order

Fit quality almost independent from perturbative order when MHOU are included

Data whose theoretical description is affected by large scale uncertainties are deweighted in favour of more perturbatively stable data



## Impact on Inclusive Cross Sections



Effect of using aN<sup>3</sup>LO PDFs instead of NNLO PDFs in N<sup>3</sup>LO predictions is small Good consistency between NNPDF4.0 [EPJ C84 (2024) 659] and MSHT20 [EPJ C83 (2023) 185]

Emanuele R. Nocera (UNITO)

#### 2.2 Methodology

## Why is the methodology important?



The methodology is crucial if we aim at percent-level accurate PDFs

Emanuele R. Nocera (UNITO)

#### Accuracy vs precision or bias vs variance



# What are the ingredients of a fitting methodology?

parametrisation

polynomials/neural network(s) is there a bias due to the parametrisation?

optimisation

(adaptive) gradient descent is the parameter space explored efficiently?

uncertainty representation

Hessian/bootstrap of experimental uncertainties what is the statistical meaning of uncertainties?

validation

closure tests (what happens if I know in advance the underlying law that I am fitting?) are interpolation and extrapolation uncertainties statistically faithful?

<u>benchmark</u>

PDF4LHC working group

are PDFs obtained independently by various groups equivalent?

Emanuele R. Nocera (UNITO)

# 2.2.1 Parametrisation

#### Parametrisation: general features

Problem projected onto the finite-dimensional space of parameters

Choose a parametrisation at an initial scale  $Q_0^2$  for each independent parton *i* (or a combination of them)

$$xf_i(z, Q_0^2) = A_i x^{a_i} (1-x)^{b_i} \mathscr{F}_i(x, \{c_{f_i}\})$$



The problem is reduced to the determination of the finite set of parameters  $\{c_{f_i}\}$ 

The interpolating function  $\mathscr{F}_i(x, \{c_{f_i}\})$  should be sufficiently GENERAL (the range of PDF behaviours in the space of functions should not be limited) SMOOTH (PDFs are implicitly assumed ot be smooth functions) FLEXIBLE (it should be able to adapt to a variety of data and processes) to describe the data with minimal bias

Emanuele R. Nocera (UNITO)

#### Parametrisation: two alternative choices

Delynomial (Bernstein, Chebyschev) parametrisation, e.g.

$$\mathscr{F}_i = 1 + \sum_{i=1}^n a_i T_i^{\mathrm{Ch}} \left( y(x) \right) \qquad y = 1 - 2\sqrt{x}$$

in terms of a (relatively) small set of parameters ( $\mathcal{O}(30)$  per PDF set)

$$\{\mathbf{a}\} = \{a_i, b_i, \gamma_i, \delta_i\}$$

 $\Rightarrow$  smooth behavior (a desirable feature for a PDF)

 $\Rightarrow$  potential source of bias if the parametrisation is too rigid

Bedundant parametrisation, e.g.

a neural network

in terms of a huge set of parameters ( $\mathcal{O}(200)$  per PDF set)

$$\{\mathbf{a}\} = \{\omega_{ij}^{(L-1),f_i}, \theta_i^{(L),f_i}\}$$

- $\Rightarrow$  potentially non-smooth
- $\Rightarrow$  bias due to the parametrisation reduced as much as possible

# Parametrisation: what a neural network exactly is?

A convenient functional form providing a flexible parametrization used as a generator of random functions in the PDF space

#### EXAMPLE: MULTY-LAYER FEED-FORWARD PERCEPTRON



$$\begin{split} \xi_i^{(l)} &= g\left(\sum_{j}^{n_l-1} \omega_{ij}^{(l-1)} \xi_j^{(l-1)} - \theta_i^{(l)}\right) \\ g(y) &= \frac{1}{1+e^{-y}} \end{split}$$

- made of neurons grouped into layers (define the architecture)
- each neuron receives input from neurons in the preceding layer (feed-forward NN)
- activation  $\xi_i^{(l)}$  determined by a set of parameters (weights and thresholds)
- activation determined according to a non-linear function (except the last layer)

#### Parametrisation: what a neural network exactly is?

EXAMPLE: THE SIMPLEST 1-2-1 MULTI-LAYER FEED-FORWARD PERCEPTRON



$$f(z) \equiv \xi_1^{(3)} = \left\{ 1 + \exp\left[ \theta_1^{(3)} - \frac{\omega_{11}^{(2)}}{1 + e^{\theta_1^{(2)} - x\omega_{11}^{(1)}}} - \frac{\omega_{12}^{(2)}}{1 + e^{\theta_2^{(2)} - x\omega_{21}^{(1)}}} \right] \right\}^{-1}$$

$$\text{Recall:} \qquad \xi_i^{(l)} = g\left(\sum_{j}^{n_l-1} \omega_{ij}^{(l-1)} \xi_j^{(l-1)} - \theta_i^{(l)}\right) \ ; \qquad g(z) = \frac{1}{1+e^{-z}}$$

Emanuele R. Nocera (UNITO)

## Parametrisation: standard vs redundant

HERA-LHC 2009 PDF benchmark



#### 2.2.2 Optimisation

# Fit quality

**1** Define the fit quality (the  $\chi^2$  function)

$$\chi^{2} = \sum_{i,j}^{N_{dat}} \left( T_{i}[\{\mathbf{a}\}] - D_{i} \right) \left( \operatorname{cov}^{-1} \right)_{ij} \left( T_{j}[\{\mathbf{a}\}] - D_{j} \right)$$

with the experimental covariance matrix

$$(\text{cov})_{ij} = \delta_{ij} s_i^2 + \left(\sum_{\alpha}^{N_c} \sigma_{i,\alpha}^{(c)} \sigma_{j,\alpha}^{(c)} + \sum_{\alpha}^{N_{\mathcal{L}}} \sigma_{i,\alpha}^{(\mathcal{L})} \sigma_{j,\alpha}^{(\mathcal{L})}\right) D_i D_j$$

- $s_i$  are  $N_{\rm dat}$  uncorrelated uncertainties (statistic + uncorrealted systematic ucnertainties)  $\sigma_{i,\alpha}^{(c)}$  are  $N_{\rm dat} \times N_c$  additive correlated uncertainties  $\sigma_{i,\alpha}^{(\mathcal{L})}$  are  $N_{\rm dat} \times N_{\mathcal{L}}$  multiplicative uncertainties
- 2) Find the best-fit configuration of parameters  $\{{f a}_{f 0}\}$  which minimise the  $\chi^2$
- 3 Treat conveniently
  - uncorrelated/correlated uncertainties need not to overestimate uncertainties and to let the  $\chi^2$  be faithful
  - additive/multiplicative uncertainties need to avoid the D'Agostini bias

# Parameter optimisation: general framework

Optimisation usually performed by means of simple gradient descent: compute and minimise the gradient of the fit quality with respect to the fit parameters

$$\frac{\partial \chi^2}{\partial a_i}$$
, for  $i = 1, \dots, N_{\text{par}}$ 

Optimisation should minimise the noise in the  $\chi^2$  driven by noisy experimental data

Additional complications in case of a redundant parametrisation (huge parameter space)

- need to explore the parameter space as uniformly as possible (in order to avoid stopping the fit in a local minimum)
- need for a computationally efficient minimisation (non-trivial relationship between FFs and observables via convolution)
- Ineed to define a criterion for minimisation stopping (avoid learning statistical fluctuations of the data)

Alternative algorithms: genetic algorithms, adaptive algorithms, ...



$$\chi^{2} = \sum_{i,j}^{N_{dat}} (T_{i}[\{\mathbf{a}\}] - D_{i}) (\operatorname{cov}^{-1})_{ij} (T_{j}[\{\mathbf{a}\}] - D_{j})$$
$$\operatorname{cov})_{ij} = \delta_{ij} s_{i}^{2} + \left(\sum_{\alpha}^{N_{c}} \sigma_{i,\alpha}^{(c)} \sigma_{j,\alpha}^{(c)} + \sum_{\alpha}^{N_{\mathcal{L}}} \sigma_{i,\alpha}^{(\mathcal{L})} \sigma_{j,\alpha}^{(\mathcal{L})}\right) D_{i} D_{j}$$

Emanuele R. Nocera (UNITO)

#### Optimisation: stopping criterion

Divide the data into two subsets (training & validation) Train the NN on the training subset and compute  $\chi^2$  for each subset Stop when the training loss reaches the absolute minimum



The best fit does not coincide with the absolute minimum of the  $\chi^2$ 

Emanuele R. Nocera (UNITO)

# Hypertoprimisation: fitting the methodology



Compare to a Test Set (new set of data previously not used at all) Who picks the Test Set? Automatic generalisation based on K foldings Divide the data into n representative sets, fit n-1 sets and use the n-th set as test set Hyperoptimise on mean and standard deviation of  $\chi^2_{{\rm test},i}$ ,  $i=1\ldots n$ 

# Hyperoptimisation: K-folding



Compare to a Test Set (new set of data previously not used at all) Who picks the Test Set? Automatic generalisation based on K foldings Divide the data into n representative sets, fit n-1 sets and use the n-th set as test set Hyperoptimise on mean and standard deviation of  $\chi^2_{{\rm test},i}$ ,  $i=1\ldots n$ 

#### Hyperparameters



# 2.2.3 Uncertainty representation

#### The Hessian method: general strategy

**(**) Expand the  $\chi^2$  about its global minimum at first (nontrivial) order

$$\chi^{2}\{\mathbf{a}\} \approx \chi^{2}\{\mathbf{a}_{0}\} + \delta a^{i} H_{ij} \delta a^{j}, \qquad H_{ij} = \frac{1}{2} \left. \frac{\partial^{2} \chi^{2}\{\mathbf{a}\}}{\partial a_{i} \partial a_{j}} \right|_{\{\mathbf{a}\} = \{\mathbf{a}_{0}\}}$$

2) Assume linear error propagation for any observable  ${\mathcal O}$  depending on  $\{{\mathbf a}\}$ 

$$\left|\mathcal{O}\{\mathbf{a}\}\right\rangle \approx \mathcal{O}\{\mathbf{a}_{\mathbf{0}}\} + a_{i} \left. \frac{\partial \mathcal{O}\{\mathbf{a}\}}{\partial a_{i}} \right|_{\{\mathbf{a}\} = \{\mathbf{a}_{\mathbf{0}}\}} \qquad \sigma_{\mathcal{O}\{\mathbf{a}\}} \approx \sigma_{ij} \frac{\partial \mathcal{O}\{\mathbf{a}\}}{\partial a_{i}} \left. \frac{\partial \mathcal{O}\{\mathbf{a}\}}{\partial a_{j}} \right|_{\{\mathbf{a}\} = \{\mathbf{a}_{\mathbf{0}}\}}$$

**③** Determine  $\sigma_{ij}$  from  $H_{ij}$  from maximum likelihood (under Gaussian hypothesis)

$$\sigma_{ij}^{-1} = \left. \frac{\partial^2 \chi^2 \{ \mathbf{a} \}}{\partial a_i \partial a_j} \right|_{\{ \mathbf{a} \} = \{ \mathbf{a}_0 \}} = H_{ij}$$

• A C.L. about the best fit is obtained as the volume (in parameter space) about  $\chi^2$ {a<sub>0</sub>} that corresponds to a fixed increase of the  $\chi^2$ ; for Gaussian uncertainties:

68% C.L. 
$$\iff \Delta \chi^2 = \chi^2 \{ \mathbf{a} \} - \chi^2 \{ \mathbf{a_0} \} = 1$$

#### The Hessian method: some remarks

Compact representation and computation of observables and their uncertainties

 $\langle \mathcal{O}[f(x,Q^2)] \rangle = \mathcal{O}[f_0(x,Q^2)]$ 

$$\sigma_{\mathcal{O}}[f(x,Q^2)] = \frac{1}{2} \left[ \sum_{i=1}^{N_{\text{par}}} \left( \mathcal{O}[f_i(x,Q^2)] - \mathcal{O}[f_0(x,Q^2)] \right)^2 \right]^{1/2}$$

Parameters can always be adjusted so that all eigenvalues of  $H_{ij}$  are equal to one (diagonalise  $H_{ij}$  and rescale the eigenvectors by their eigenvalues)

$$\delta a_i H_{ij} \delta a_j = \sum_{i=1}^{N_{\text{par}}} \left[ a'_i(a_i) \right]^2 \Longleftrightarrow \sigma_{\mathcal{O}\{\mathbf{a}'\}} = \left| \nabla' \mathcal{O}\{\mathbf{a}'\} \right|$$

The total contribution to the uncertainty due to two different sources (possibly correlated) is obtained by simply adding them in quadrature

- Any rotation in the space of parameters preserves the gradient (one can diagonalise a chosen observable without spoiling the result)
- Unmanageable Hessian matrix if the numer of parameters is huge

# The Monte Carlo method: general strategy

I Generate (art) replicas of (exp) data according to the distribution

$$\mathcal{O}_i^{(art)(k)} = \mathcal{O}_i^{(exp)} + r_i^{(k)} \sigma_{\mathcal{O}_i}, \qquad i = 1, \dots N_{\text{dat}}, \qquad k = 1, \dots, N_{\text{rep}}$$

where  $r_i^{(k)}$  are (Gaussianly distributed) random numbers for each k-th replica  $(r_i^{(k)}$  can be generated with any distribution, not neccesarily Gaussian)

- 2 Perform a fit for each replica  $k = 1, \ldots, N_{rep}$
- Compact computation of observables and their uncertainties (PDF replicas are equally probable members of a statistical ensemble)

$$\langle \mathcal{O}[f(x,Q^2)]\rangle = \frac{1}{N_{\rm rep}}\sum_{k=1}^{N_{\rm rep}}\mathcal{O}[f^{(k)}(x,Q^2)]$$

$$\sigma_{\mathcal{O}}[f(x,Q^2)] = \left[\frac{1}{N_{\mathsf{rep}} - 1} \sum_{k=1}^{N_{\mathsf{rep}}} \left(\mathcal{O}[f^{(k)}(x,Q^2)] - \langle \mathcal{O}[f(x,Q^2)] \rangle \right)^2\right]^{1/2}$$

 $\Rightarrow$  no need to rely on linear approximation

 $\Rightarrow$  computational expensive: need to perform  $N_{\rm rep}$  fits instead of one

## The Monte Carlo method: determining the sample size

Require that the average over the replicas reproduces the central value of the original experimental data to a desired accuracy (the standard deviation reproduces the error and so on)



Accuracy of few % requires  $\sim 100$  replicas

Emanuele R. Nocera (UNITO)

## 2.2.4 Validation

#### Accuracy vs precision or bias vs variance



#### Closure tests: general idea [JHEP 1504 (2015) 040]

Validation and optimisation of the fitting strategy with known underlying physical law



#### Full control of procedural uncertainties

Emanuele R. Nocera (UNITO)

# Closure Tests: Levels

Level 0 no fluctuations



interpolation uncertainty



Level 1 Gaussian fluctuation



#### fuctional uncertainty



Level 2 Monte Carlo replicas



data uncertainty


## Closure tests at work

Data region: closure tests

Fit PDFs to pseudodata generated assuming a known underlying law

Define bias and variance bias difference of central prediction and truth variance uncertainty of replica predictions

> If PDF uncertainty faithful, then E[bias] = variance25 fits, 40 replicas each





Emanuele R. Nocera (UNITO)

## Future tests



Extrapolation regions: future test

Test PDF uncertainties on data sets not included in a given PDF fit that cover unseen kinematic regions

Data set	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	3.30	23.1

Only exp. cov. matrix



Acta Phys.Polon. B52 (2021) 243

## Future tests



Extrapolation regions: future test

Test PDF uncertainties on data sets not included in a given PDF fit that cover unseen kinematic regions

Data set	NNPDF4.0	pre-LHC	pre-HERA
pre-HERA pre-LHC NNPDF4.0	1.12	1.17 1.30	0.86 1.22 1.38

Exp+PDF cov. matrix

0.7





Acta Phys.Polon. B52 (2021) 243

# 2.3 Summary of Lecture 2

# Overview of current PDF determinations

	NNPDF4.0	MSHT20	CT18	HERAPDF2.0	CJ22	ABMP16		
Fixed-target DIS	Ø	Ø	Ø	$\boxtimes$	Ø	Ø		
JLAB	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\checkmark$	$\boxtimes$		
HERA I+II	$\checkmark$	$\checkmark$	$\checkmark$	$\square$	$\checkmark$	$\checkmark$		
HERA jets	$\square$	$\boxtimes$	$\boxtimes$	$\square$	$\boxtimes$	$\boxtimes$		
Fixed target DY	$\square$		$\square$	$\boxtimes$	$\square$			
Tevatron $W$ , $Z$	Ø	Ø	Ø	$\boxtimes$	Ø	Ø		
LHC vector boson	Ø		$\square$	$\boxtimes$	$\square$	$\square$		
$LHC\ W + c\ Z + c$	Ø	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$		
Tevatron jets	Ø	Ø	Ø	$\boxtimes$	$\square$	$\boxtimes$		
LHC jets	Ø	Ø	$\square$	$\boxtimes$	$\boxtimes$	$\boxtimes$		
LHC top	Ø		$\boxtimes$	$\boxtimes$	$\boxtimes$	$\square$		
LHC single $t$	Ø	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$		
LHC prompt $\gamma$		$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$	$\boxtimes$		
statistical treatment	Monte Carlo	Hessian $\Delta\chi^2$ dynamical	Hessian $\Delta\chi^2$ dynamical	Hessian $\Delta \chi^2 = 1$	Hessian $\Delta\chi^2 = 1.645$	Hessian $\Delta \chi^2 = 1$		
parametrisation	Neural Network	Chebyschev pol.	Bernstein pol.	polynomial	polynomial	polynomial		
HQ scheme	FONLL	TR'	ACOT- $\chi$	TR'	ACOT- $\chi$	FFN		
accuracy	$aN^3LO$	$aN^3LO$	NNLO	NNLO	NLO	NNLO		
latest update	EPJ C82 (2022) 428	EPJ C81 (2021) 341	PRD 103 (2021) 014013	EPJ C82 (2022) 243	PRD 107 (2023) 113005	PRD 96 (2017) 014011		
All PDF sets are available as $(x, Q^2)$ interpolation grids through the LHAPDF library								

Parton Distribution Function

74 / 75

# Summary of Lecture 2

**1** PDF accuracy can be improved by improving the theory

- $\longrightarrow$  proper treatment of heavy quarks is mandatory to describe DIS data
- $\longrightarrow$  evidence for intrinsic charm in the proton
- $\longrightarrow$  MHOUs can be estimated by scale variations
- $\longrightarrow$  inclusion of MHOUs stabilises fit quality
- $\longrightarrow$  electroweak corrections modify DGLAP equations
- $\longrightarrow$  the photon PDF is determined very precisely
- $\longrightarrow$  inclusion of photon PDFs impacts the gluon PDF
- 2 Devising the methodology is essential to minimise bias and variance
  - $\longrightarrow$  bias is a measure of accuracy, variance is a measure of precision
  - $\rightarrow$  choices of parametrisation (polynomial vs neural network)
  - $\rightarrow$  choices of uncertainty representation (Hessian vs Monte Carlo)
  - $\longrightarrow$  are all sources of bias and variance
  - $\longrightarrow$  closure tests are a way to validate PDF uncertainites in the  $\mathit{seen}$  region
  - $\longrightarrow$  future tests validate the generalisation power in the unseen region

#### Thank you