Statistical tools for global QCD analyses

GFI 1st miniworkshop

Emanuele R. Nocera

Università degli Studi di Torino and INFN — Torino

Oasi di Cavoretto - 29th June 2023







PDF fitting in Bayesian language

Inverse Problem

Given a data set F, determine p(f|F) in the sapce of functions $f:[0,1] \rightarrow \mathbb{R}$

$$E[\mathcal{O}] = \int \mathcal{D}f \, p(f|F)\mathcal{O}(f) \qquad V[\mathcal{O}] = \int \mathcal{D}f \, p(f|F)[\mathcal{O}(f) - E[\mathcal{O}(f)]]^2$$

Choose a basis in the space of functions

Choose a parametrisation for the functions and a fitting model where τ : kinematic variables; θ : parameters; h: hyperparameters

$$f = f(\tau, \theta, h) \qquad p(f|F) \to p(\theta|F, h)$$

Sample the statistical distribution of the data with Monte Carlo replicas $F^{(k)}$

The posterior distribution is the product of the likelihood and the prior (Bayes theorem)

$$p(\theta|F^{(k)},h) \propto p(F^{(k)}|\theta,h)p(\theta|h)$$

Each replica is the mode of the posterior distribution given one instance of the data

$$\theta^{*(k)} = \arg\max_{\theta} p(\theta|F^{(k)}, h)$$

Hyperoptimise h and determine $\theta^{*(k)}$, for each replica, by $\chi^{2(k)}$ minimisation

Importance sampling



 $N_{\rm rep} = 100$ replicas are sufficient to reproduce mean values and uncertainties to 1% Importance sampling is reproduced correctly (14 replicas out of 1000 vs 98.9% C.I.) Outliers in the distribution of fitted replicas are good fits to unlikely fluctuations of the data

Emanuele R. Nocera (U. Torino & INFN)

A tale of accuracy and precision

High accuracy

Low accuracy



High precision

Low precision

A tale of bias and variance

Low bias

High bias







High variance

The bias-variance trade-off



squared mean error = $bias^2 + variance + irreducible error$

Emanuele R. Nocera (U. Torino & INFN)

Error

Statistical tools

1. Goodness-of-fit

Defining the χ^2

$$\chi^{2} = \frac{1}{N_{\text{dat}}} \sum_{i,j}^{N_{\text{dat}}} \Delta_{i} \left(\cos v_{t_{0}}^{-1} \right)_{ij} \Delta_{j} \qquad \begin{cases} \chi^{2(k)} & \Delta_{i} = T_{i}^{(k)} - D_{i}^{(k)} & \text{t}_{0} \text{ used to avoid} \\ \chi^{2(k,c)} & \Delta_{i} = T_{i}^{(k)} - D_{i}^{(c)} & \text{JAgostini bias} \\ \chi^{2(0,c)} & \Delta_{i} = T_{i}^{0} - D_{i}^{(c)} & \text{[JHEP 05 (2010) 075]} \\ \chi^{2(0,c)} & \Delta_{i} = T_{i}^{0} - D_{i}^{(c)} & \text{[JHEP 04 (2013) 125]} \end{cases}$$

$$(\operatorname{cov}_{t_0})_{ij} = \delta_{ij} s_i^{(\mathsf{uncorr})} s_j^{(\mathsf{uncorr})} + \sum_{m=1}^{N_{\operatorname{norm}}} \sigma_{i,m}^{(\operatorname{norm})} \sigma_{j,m}^{(\operatorname{norm})} T_i^{(0)} T_j^{(0)} + \sum_{l=1}^{N_{\operatorname{corr}}} \sigma_{i,l}^{(\operatorname{corr})} \sigma_{j,l}^{(\operatorname{corr})} D_i D_j,$$



Normalisation uncertainties: D'Agostini bias [JHEP 1005 (2010) 075]

 $\textbf{O} \ \ \text{Consider one experiment with } N_{\mathrm{dat}} \ \text{data} \ d_i \ \text{of one theoretical quantity} \ t \\$

$$\chi^{2}(t) = \sum_{i,j}^{N_{\text{dat}}} (t - d_{i}) \left(\text{cov}^{-1} \right)_{ij} (t - d_{j})$$

2 The best-fit theoretical quantity t_0 and its variance v_t are given by

$$\frac{d\chi^2}{dt}\Big|_{t=t_0} = 0 \iff t_0 = \frac{\sum_{i,j}^{N_{\rm dat}} \left(\cos^{-1}\right)_{ij} d_j}{\sum_{i,j}^{N_{\rm dat}} \left(\cos^{-1}\right)_{ij}} \quad v_t = \left(\frac{1}{2} \frac{d^2 \chi^2}{dt^2}\right)^{-1} = \frac{1}{\sum_{i,j}^{N_{\rm dat}} \left(\cos^{-1}\right)_{ij}}$$

 $\textbf{S} \quad \text{Consider completely uncorrelated additive errors: } (\operatorname{cov})_{ij} = s_i^2 \delta_{ij}$

$$t_0 = w = \Sigma^2 \sum_i^{N_{ ext{dat}}} rac{d_i}{s_i^2}$$
 $v_t = \Sigma^2$ with $rac{1}{\Sigma^2} = \sum_i^{N_{ ext{dat}}} rac{1}{s_i^2}$

) Consider an additional common normalisation error: $({
m cov})_{ij}=(s_i^2+\sigma^2 d_i^2)\delta_{ij}$

$$t_0 = \frac{w}{1 + r^2 \sigma^2 w^2 / \Sigma^2} \qquad v_t = \frac{\Sigma^2 + \sigma^2 w^2 (1 + r^2)}{1 + r^2 \sigma^2 w^2 / \Sigma^2} \qquad \text{with } r^2 = \frac{\Sigma^2}{w^2} \sum_{i}^{N_{\text{dat}}} \frac{(d_i - w)^2}{s_i^2}$$

(3) Both t_0 and v_t are affected by a downward bias smaller values of d_i have a smaller normalization uncertainties σd_i and are thus preferred

Normalisation uncertainties: D'Agostini bias [JHEP 1005 (2010) 075]

The penalty trick: redefine the fit quality

$$\chi^2(t) \rightarrow \chi^2(t, \mathcal{N}) = \sum_i^{N_{\text{dat}}} \frac{(t/\mathcal{N} - d_i)^2}{s_i^2} + \frac{(\mathcal{N} - 1)^2}{\sigma^2}$$

$$\frac{\partial \chi^2}{\partial t}\Big|_{t=t_0} = \frac{\partial \chi^2}{\partial \mathcal{N}} = 0 \iff t_0 = w \qquad v_t = \left(\frac{1}{2}\frac{d^2\chi^2}{dt^2}\right)^{-1} = \Sigma^2 + \sigma^2 w^2$$

 \longrightarrow recover the unbiased estimators for t_0 and v_t

2 The t_0 method: redefine the covariance matrix

$$(\operatorname{cov})_{ij} \to (\operatorname{cov}_{t_0})_{ij} \iff (s_i^2 + \sigma^2 d_i^2) \delta_{ij} \to s_i^2 \delta_{ij} + \sigma^2 t_0^2$$

$$\left(\operatorname{cov}_{t_0}^{-1}\right)_{ij} = \frac{\delta_{ij}}{s_i^2} - \frac{\sigma^2 t_0^2}{s_i^2 s_j^2} \frac{\Sigma^2}{\Sigma^2 + \sigma^2 t_0^2} \Longleftrightarrow t_0 = w \qquad v_t = \Sigma^2 + \sigma^2 w^2$$

 \longrightarrow recover the unbiased estimators for t_0 and $v_t,$ provided that w is tuned to t_0 \longrightarrow w can be tuned to t_0 via an iterative procedure

The d'Agostini bias and its solution can be generalised to more than one experiment with different normalisation errors (per experiment/per data point)

2. Consistency of data

Data inconsistency: tensions between data sets

Give more weight to a data set p $\chi^2 \rightarrow \chi^2 + w \chi_p^2$

Refit: the total χ^2 will increase Which data sets get worse? How much?

Refit: the data set χ_p^2 will decrese Self-consistency? Inconsistency?

Examples: ATLAS W, Z and $t\bar{t}$

Inconsistency clearly spotted unnatural PDF shapes appear error in other data sets increases

Otherwise global fit quality and PDFs remain unaltered

Data set	baseline	$rw\ W, Z$	rw $t\bar{t}$
ATLAS W, Z 7 TeV ATLAS $t\bar{t}$ 8 TeV	1.86 4.11	1.23	 1.21
Total	1.20	1.21	1.73



Data inconsistency: experimental correlations

Single inclusive jet data from ATLAS 7 TeV default correlations: terrible χ^2

(correlations across rapidity bins)

decorrelation models: improve the fit a lot

$n_{\rm dat}$	default	part. decorr.	full decorr.
140	1.89	1.28	0.83

no significant effect on the extracted gluon similar gluon irrespective of the rapidity bin



[EPJ C78 (2018) 248; EPJ C80 (2020) 797]

Top pair production from ATLAS 8 TeV

default correlations: terrible χ^2

(correlations across different spectra)

decorrelation models: improve the fit a lot

$n_{\rm dat}$	default	stat. uncorr.	p.s. uncorr
25	7.00	3.28	1.80

appreciable effect on the extracted gluon different gluon depending on the top spectrum



[EPJ C80 (2020) 1; Les Houches proceedings, 2019]

Good knowledge of experimental correlations is important

Consider the COMPASS π^{\pm} multiplicities [PLB 764 (2017) 1]

Only 80% of the systematic uncertainty is bin-by-bin correlated

What if you incorporate a different piece of information in a FF fit?

Consider two cases: [arXiv:2204.10331] full correlation; full decorrelation





3. Consistency of methodology

Closure tests [EPJ C77 (2017) 663; EPJ C82 (2022) 330]

A test to validate PDF uncertainties in the data region

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
$$\label{eq:Ebias} \begin{split} \text{E[bias]} = \text{variance} \\ \text{25 fits, 40 replicas each} \end{split}$$



Future tests [Acta Phys.Polon. B52 (2021) 243]

A test to validate PDF uncertainties in the extrapolation regions

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



Future tests [Acta Phys.Polon. B52 (2021) 243]

A test to validate PDF uncertainties in the extrapolation regions

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



How to test for overfitting?

Metric 1: Kinetic Energy

$$\mathrm{KE} = \sqrt{1 + \left(\frac{d}{d\ln x}xf(x,Q^2)\right)^2}$$

A measure of PDF wiggliness:

the higher the KE, the longer the curve that joins two fixed points



How to test for overfitting?

Metric 2: Overfit metric

$$\mathcal{R}_O = \chi_{\mathrm{val}}^2 \left[T^{(k)}, D^{(k)} \right] - \overline{\chi_{\mathrm{val}}^2} \left[T^{(k)}, D^{(k)} \right]$$

$$\overline{\chi^2_{\text{val}}}\left[T^{(k)}, D^{(k)}\right] \equiv \frac{1}{N} \sum_{k'=1}^N \chi^2_{\text{val}}\left[T^{(k)}, D^{(k')}\right]\Big|_{\text{fixed mask}}$$

A measure of PDF overfitting:

for each replica, there is a different fluctuation and training/validation split (mask) if $\mathcal{R}_O < 0$ the replica is overfitted, *i.e.* it contains information specific to D^k



4. Consistency of theory

Theory uncertainties in PDF determination

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

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



Experimental + Theory Correlation Matrix (3 pt)

Emanuele R. Nocera (U. Torino & INFN)

Theory uncertainties in PDF determination



PDF uncertainty increase encapsulates NLO-NNLO shift Overall (rather small) increase in uncertainties Increase in PDF uncertainties due to replica generation is counteracted by extra correlations in fitting minimisation

 $\begin{array}{ll} \mbox{Tensions relieved: improvement in } \chi^2 \\ \mbox{exp only: } \chi^2/N_{\rm dat} = 1.139 & \mbox{exp+th: } \chi^2/N_{\rm dat} = 1.110 \end{array}$

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

EPJ C79 (2019) 838; ibid. 931

Theory uncertainty in PDF determination

Experimental+Nuclear correlation matrix CHORUS 1.00 0.75 0.50 0.25 0.00 -0.25 -0.50NUTEV -0.75 DYE605 -1 00 1.00 SLAC -075 - 0 50 BCDMS - 0.25 - 0.00 - -0.25 -0.50 NMC -0.75 NuSea -1.00

Effect of nuclear uncertainties relevant at large xto reconcile FT DIS with LHC DY data $\chi^2_{tot} = 1.17 \rightarrow \chi^2_{tot} = 1.26$ (no nucl. uncs.) $\chi^2_{LHCb} = 1.54 \rightarrow \chi^2_{tot} = 1.76$ (no nucl. uncs.) The bulk of the effect is due to nuclear uncertainties for heavy nuclei deuteron uncertainties have a comparatively

smaller effect at inermediate values of \boldsymbol{x}



Emanuele R. Nocera (U. Torino & INFN)

29th June 2023 22 / 24

5. To conclude

Summary and outlook

Performing a global QCD analysis is a daunting task (in psychological terms)

Performing a global QCD analysis is an inverse problem (in Bayesian terms)

There is an inherent bias-variance trade-off

In order to maximise precision and accuracy one has to minimise bias and variance

One has to take into account the experimental, methodological and theoretical inputs

Summary and outlook

Performing a global QCD analysis is a daunting task (in psychological terms)

Performing a global QCD analysis is an inverse problem (in Bayesian terms)

There is an inherent bias-variance trade-off

In order to maximise precision and accuracy one has to minimise bias and variance

One has to take into account the experimental, methodological and theoretical inputs

Thank you