

On the determination of uncertainties in parton densities

N. Hunt-Smith, **Alberto Accardi**, W. Melnitchouk,
N. Sato, A.W. Thomas, M.J. White

Phys. Rev. D 106 (2022) 036003

DIS 2023

March 29, 2023



This work is in part supported by the DOE Office of Science

Overview

- **Uncertainty Quantification: Parametric Methods**

- Monte Carlo Bayesian estimators
- Hessian approximation
- Data resampling

- **Description of Toy Model**

- Benchmark of Hessian and MC methods

- **Neural Network Comparison**

- Algorithmic modification of likelihood?
(see also N. Sato @ DIS 2018)

A whole session devoted to PDF uncertainties:

- P. Nadolsky – “Epistemic uncert. quant.”
- L. Kotz – “Bezier curve parametrizations”
- K. Mohan – A new statistical method”

Uncertainty quantification: parametric methods

Bayesian estimators

- **Bayes theorem** $p(\mathbf{a}|\mathbf{m}) = \frac{1}{\mathcal{Z}} p(\mathbf{m}|\mathbf{a}) p(\mathbf{a})$

with “evidence” $\mathcal{Z} = \int d\mathbf{a} p(\mathbf{m}|\mathbf{a}) p(\mathbf{a})$

and “likelihood” $p(\mathbf{m}|\mathbf{a}) = \mathcal{N} \exp \left[-\frac{1}{2} \chi^2(\mathbf{a}, \mathbf{m}) \right]$

Typical choice in PDF analyses

- Algorithms for sampling of likelihood $\rightarrow \{\mathbf{a}_k\}$
 - **HMC**: Hamiltonian Monte Carlo (an example of Markov-Chain MC methods)
 - **NS**: Nested Sampling, primarily aimed at estimating the evidence
 \rightarrow Samples the likelihood as a byproduct

- Expectation values $E_{\text{Bayes}}\{\mathcal{O}(\mathbf{a})\} = \frac{1}{n} \sum_{k=1}^n \mathcal{O}(\mathbf{a}_k)$,

and variance $V_{\text{Bayes}}\{\mathcal{O}(\mathbf{a})\} = \frac{1}{n} \sum_{k=1}^n [\mathcal{O}(\mathbf{a}_k) - E_{\text{Bayes}}\{\mathcal{O}(\mathbf{a})\}]^2$

Data resampling

- Data Resampling (**DR**) approximates Bayes' posterior using frequentist logic
 - Reshuffle data within data uncertainty (Gaussian distribution)
 - Maximize likelihood
 - Repeat n_{rep} times $\rightarrow \{\mathbf{a}_k\}$

- Estimate

$$E_{\text{freq}}\{\mathcal{O}(\mathbf{a})\} = \frac{1}{n_{\text{rep}}} \sum^{n_{\text{rep}}} \mathcal{O}(\mathbf{a}_{\text{rep}}),$$

$$V_{\text{freq}}\{\mathcal{O}(\mathbf{a})\} = \frac{1}{n_{\text{rep}}} \sum^{n_{\text{rep}}} [\mathcal{O}(\mathbf{a}_{\text{rep}}) - E_{\text{freq}}\{\mathcal{O}(\mathbf{a})\}]^2$$

- Good in parameter space region well constrained by data

Generalized Hessian Approximation

Hunt-Smith et al., PRD 106 (2022) 036003

- Start as usual:
 - Find minimum of likelihood
 - Diagonalize Hessian $\rightarrow e_k$ eigenvectors, w_k eigenvalues
- Change variables: $\mathbf{a}(t) = \mathbf{a}_0 + \sum_{k=1}^{n_{\text{par}}} t_k \frac{e_k}{\sqrt{w_k}}$, then $p(\mathbf{a}|\mathbf{m}) \rightarrow p(t|\mathbf{m})$
- Assume likelihood factorized along Hessian eigendirection, then

$$E_{\text{Hess}}\{\mathcal{O}(\mathbf{a})\} = \int d^n t p(t|\mathbf{m}) \mathcal{O}(\mathbf{a}(t)) \approx \mathcal{O}(\mathbf{a}_0)$$
$$V_{\text{Hess}}\{\mathcal{O}(\mathbf{a})\} \approx \sum_k T_k^2 \left(\left. \frac{\partial \mathcal{O}(\mathbf{a}(t))}{\partial t_k} \right|_{\mathbf{a}_0} \right)^2$$

- Here $T_k^2 = \int dt_k p_k(t_k|\mathbf{m}) t_k^2$ is the “tolerance” :
 - $T_k = 1$ where likelihood is Gaussian;
 - Approximates well the likelihood in non-Gaussian directions
 - Maintains a “68%” or “ 1σ ” kind of meaning also when $\neq 1$

CT, MSTW \rightarrow T=5-10

- Often T_k determined “ad hoc” to account for statistical inconsistency of data

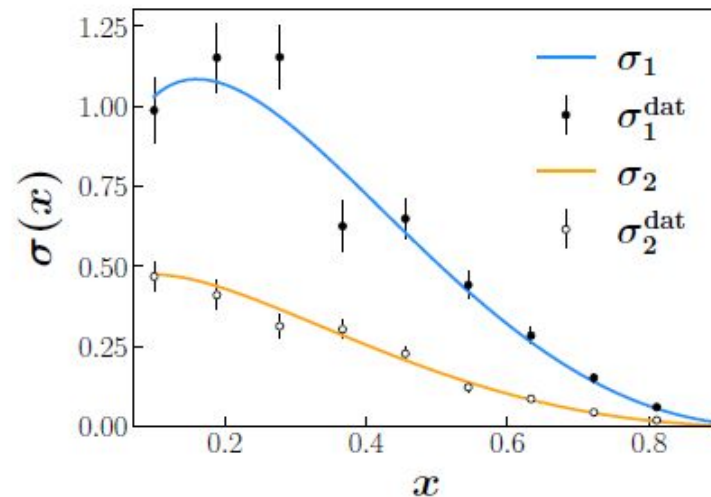
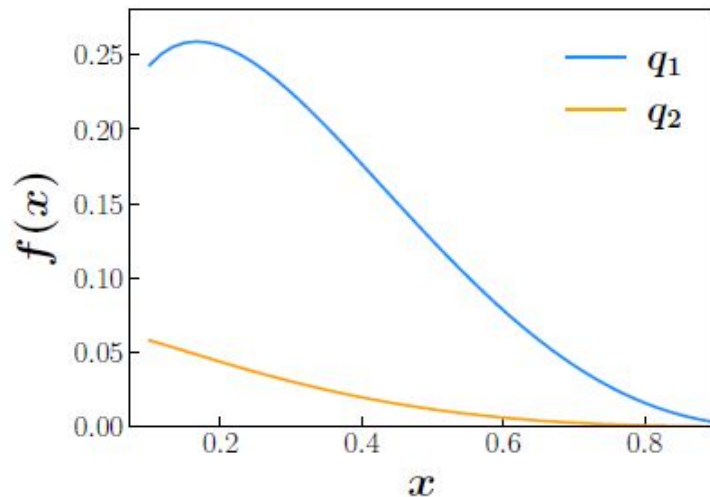
Toy Model

Toy model

- PDFs f : mimic up and down quarks
- Observables σ : mimic proton, neutron DIS cross section at fixed Q^2
 - Data randomly generated according to corresponding x distributions

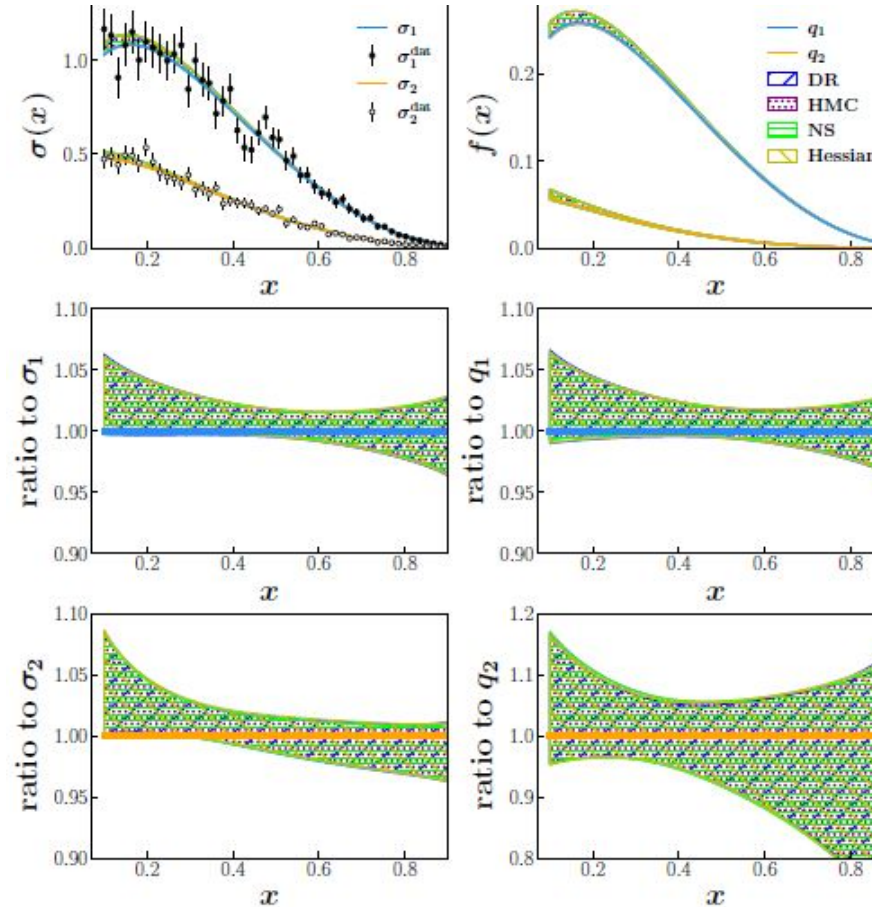
$$q_i(x) = x^{\alpha_i}(1-x)^{\beta_i},$$
$$i = 1, 2.$$

$$\sigma_j = \sum_{i=1,2} c_{ji} q_i,$$
$$c_{11} = 4c_{12} = 4c_{21} = c_{22}.$$



Equivalency of parametric methods

- **Bayesian MC estimators** used as benchmark
- **Hessian approximation is good!**
 - Generalized tolerance marginally needed even in this simplified example
- Crucially, **data resampling provides same likelihood estimation** as Bayesian MC methods



Neural Network Fits

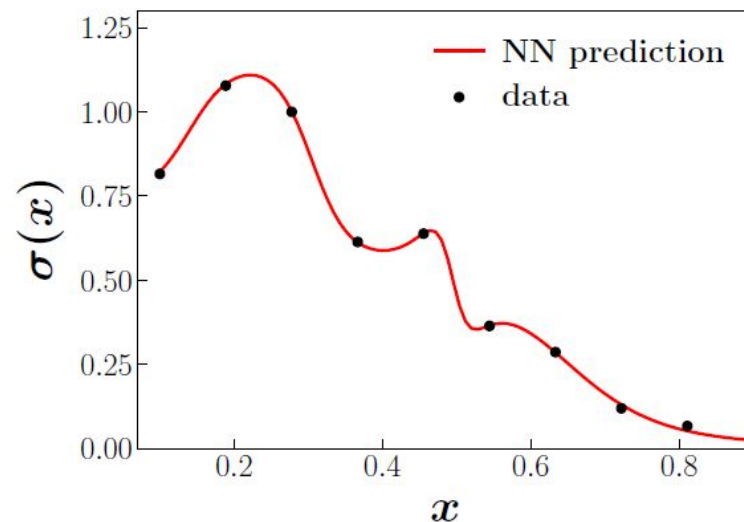
Neural Networks and overfitting

- Neural networks provide:
 - Efficient, very flexible parametrizations
 - Hundreds of parameters
 - Essentially a parameter free functional form

- Aim at maximizing the same likelihood

$$p(\mathbf{m}|\mathbf{a}) = \mathcal{N} \exp \left[-\frac{1}{2} \chi^2(\mathbf{a}, \mathbf{m}) \right]$$

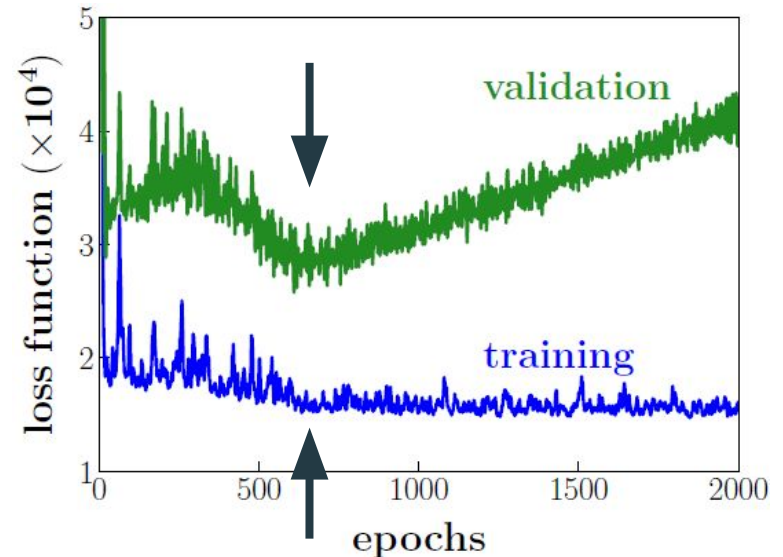
- Without intervention, will overfit the data
 - The plot shows an extreme example



Cross-validation (CV) and stopping

- Needs a “stopping criterion”
 - to avoid fitting statistical noise instead of physics
- Randomly separate the data into 2 groups, say
 - 70% → training (T)
 - 30% → validation (V)
- Fit the training, calculate $\chi^2(T)$ and $\chi^2(V)$
- Resample data, repeat
- “Stop” training when $\chi^2(V)$ is minimum:

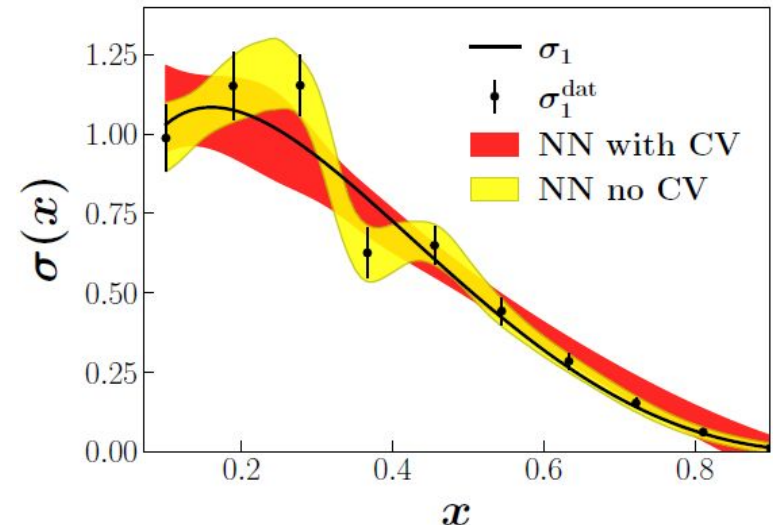
$$\sigma = E[\sigma_{\text{fit}}]$$
$$\delta\sigma = V[\sigma_{\text{fit}}]$$



Cross-validation (CV) and stopping

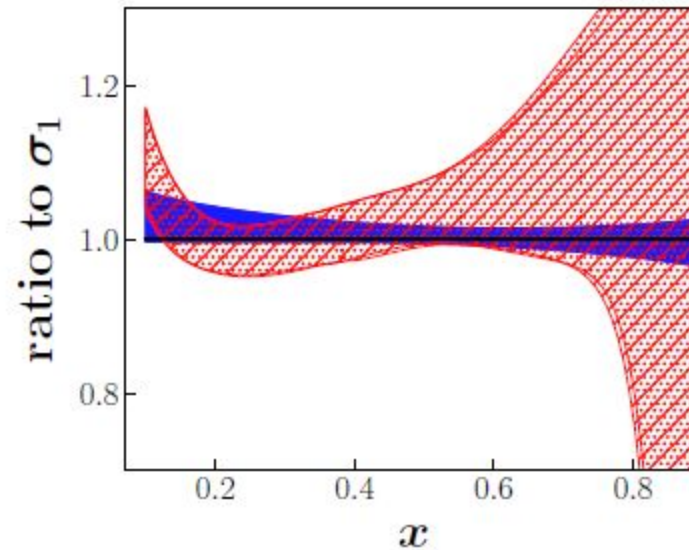
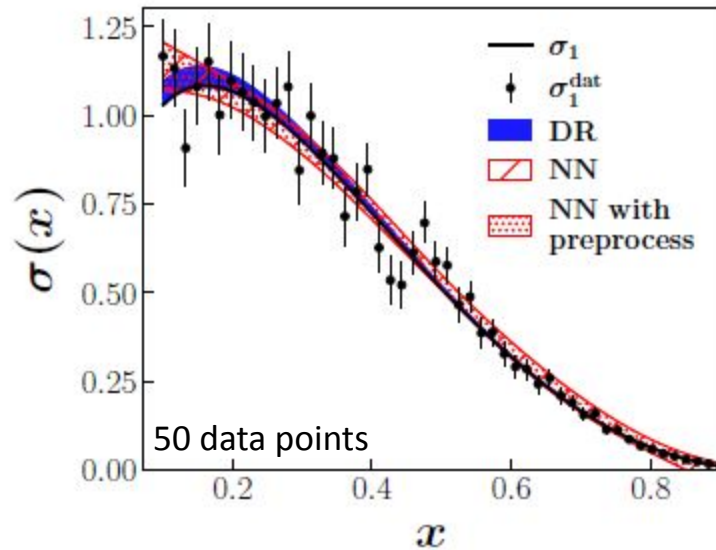
- Needs a “stopping criterion”
 - to avoid fitting statistical noise instead of physics
- Randomly separate the data into 2 groups, say
 - 70% → training (T)
 - 30% → validation (V)
- Fit the training, calculate $\chi^2(T)$ and $\chi^2(V)$
- Resample data, repeat
- “Stop” training when $\chi^2(V)$ is minimum:

$$\sigma = E[\sigma_{\text{fit}}]$$
$$\delta\sigma = V[\sigma_{\text{fit}}]$$



Comparison of NN to parametric methods

- DR as representative of parametric methods
- **Neural Network fits:**
 - Comparable in shape
 - **Quite larger uncertainty!**

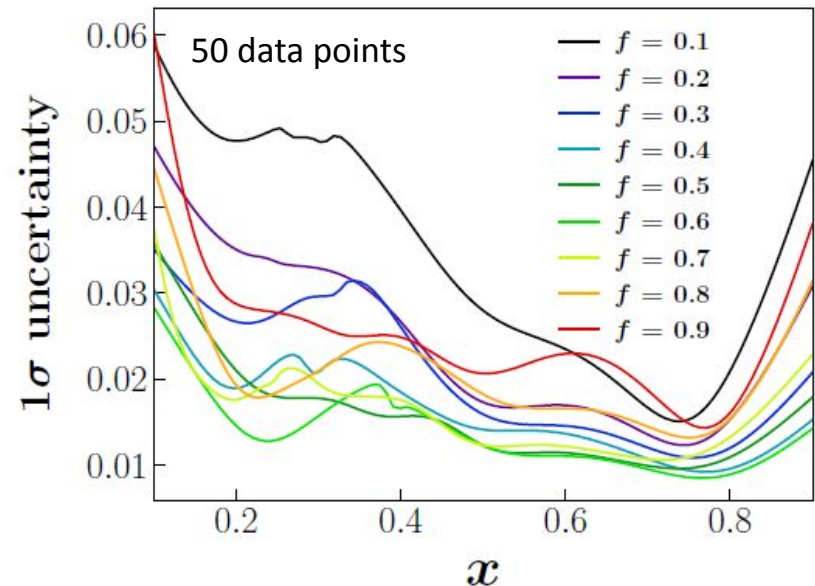


Dependence on training fraction f

- **The fit is quite independent of the T/V partitioning**
 - b/c in each replica the training data is randomly chosen
 - So it spans the whole x range

BUT

- **The uncertainty strongly depends on the fraction of training data**
 - with $f \approx 0.6$ providing the smallest uncertainty

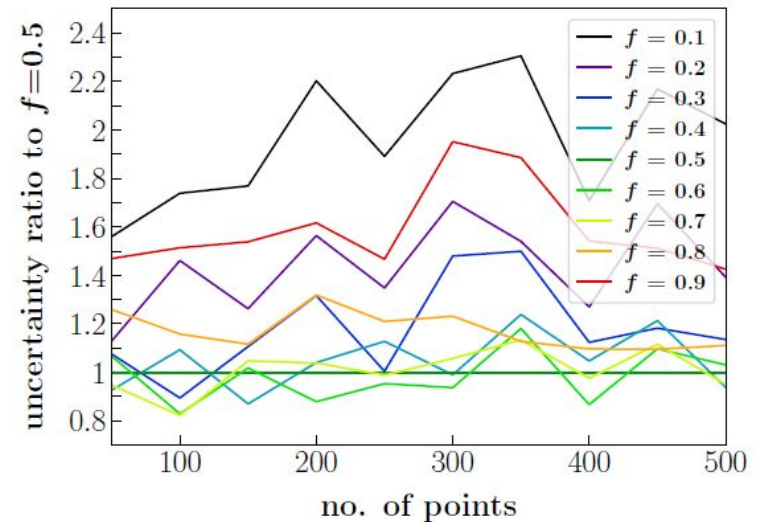


Dependence on training fraction f

- **The fit is quite independent of the T/V partitioning**
 - b/c in each replica the training data is randomly chosen
 - So it spans the whole x range

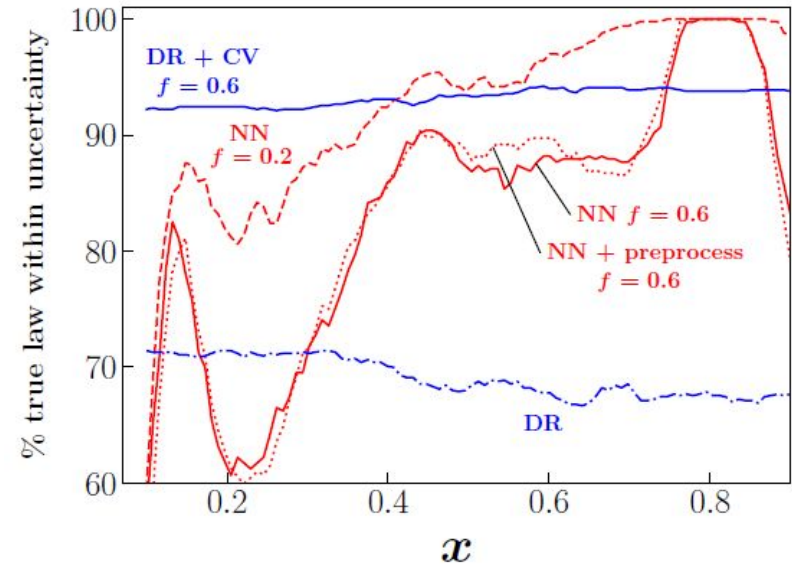
BUT

- **The uncertainty strongly depends on the fraction of training data**
 - with $f \approx 0.6$ providing the smallest uncertainty
 - Independently of how dense the data is in x



Comparison of NN to parametric methods

- **NN fits inflate the uncertainty estimate!**
 - Partly due to cross-validation
 - Structure in x difficult to understand
 - Uncertainty explodes at large x
- Data resampling + Cross Validation also inflates the uncertainty
 - Validation set “pulls” against training set
 - But in the same way across x



The algorithms have effectively modified the nominal $\exp(-\chi^2)$ likelihood!

In conclusion...

Food for thought

- Reliable quantification of PDF uncertainties needed for QCD and HEP applications
- **Parametric methods produce the same likelihood estimates**
 - Bayesian MC methods
 - Hessian approximation
 - Data resampling
- **Neural Network fits**
 - **Algorithmically modify the nominal likelihood**
 - The resulting uncertainties are not directly comparable to parametric estimates
 - *Enlarged uncertainties do not look like a natural replacement for tolerance criterion to account for tension in the data sets*
- **In what sense can NNPDF be combined with others in, say, PDF4LHC fits?**