

Reconstructing parton distribution functions from Ioffe time data: from Bayesian methods to Neural Networks

Joseph Karpie^{a,b}, Kostas Orginos^{a,b}, Alexander Rothkopf^c and Savvas Zafeiropoulos^d

^a*Department of Physics, The College of William & Mary, Williamsburg, VA 23187, USA*

^b*Thomas Jefferson National Accelerator Facility, Newport News, VA 23606, USA*

^c*Faculty of Science and Technology, University of Stavanger, 4021 Stavanger, Norway*

^d*Institute for Theoretical Physics, Heidelberg University, Philosophenweg 12, 69120 Heidelberg, Germany*

E-mail: jmkarpie@email.wm.edu, kostas@wm.edu,

alexander.rothkopf@uis.no, s.zafeiropoulos@thphys.uni-heidelberg.de

ABSTRACT: The computation of the parton distribution functions (PDF) or distribution amplitudes (DA) of hadrons from first principles lattice QCD constitutes a central open problem. In this study, we present and evaluate the efficiency of a selection of methods for inverse problems to reconstruct the full x -dependence of PDFs. Our starting point are the so called Ioffe time PDFs, which are accessible from Euclidean time calculations in conjunction with a matching procedure. Using realistic mock data tests, we find that the ill-posed incomplete Fourier transform underlying the reconstruction requires careful regularization, for which both the Bayesian approach as well as neural networks are efficient and flexible choices.

Contents

1	Introduction	1
2	Direct Inversion	3
3	Advanced PDF reconstructions	7
3.1	Backus-Gilbert Method	9
3.2	Neural Network Reconstruction	11
3.3	Bayesian PDF reconstruction	14
4	Mock Data Tests	18
4.1	Backus-Gilbert	18
4.2	Neural network reconstruction	20
4.3	Bayesian Analysis	27
4.4	Restricted χ^2 sampling	33
5	Summary and Conclusion	36

1 Introduction

Parton distributions form the core of our theoretical understanding of the inner workings of hadrons [1]. They encode how the momentum and angular momentum is distributed among quarks and gluons inside a hadron. As such there exists an intense high precision experimental program devoted to their determination and it is the concurrent goal of theoretical nuclear physics to compute these quantities from first principles QCD. As they constitute genuinely non-perturbative objects, lattice QCD calculations appear as a very promising route in achieving this goal.

All current methods for extracting parton distribution functions (PDF) or distribution amplitudes (DA) from lattice QCD require a Fourier transform or a quasi-Fourier transform of a certain class of hadronic position space matrix elements. These novel methods allow lattice QCD calculations to go beyond the traditional computation of the lowest moments [2–5]. Let us focus on two related methods, called the pseudo-PDF [6] and the quasi-PDF [7]. Both involve the hadron matrix element of space-like separated quark fields connected with a Wilson line and a γ matrix. For simplicity, we will use the example of the flavor iso-vector unpolarized quark PDFs. In this case the relevant matrix element has the following Lorentz decomposition

$$\langle h(p) | \bar{\psi}(z) \frac{\tau^3}{2} \gamma^\alpha W(z; 0) \psi(0) | h(p) \rangle = p^\alpha \mathcal{M}(\nu, z^2) + z^\alpha \mathcal{N}(\nu, z^2), \quad (1.1)$$

where p is an arbitrary hadron momentum, z is a space-like separation, τ^3 denotes a flavor Pauli matrix, γ^α refers to a gamma matrix acting in spin space, $W(z;0)$ to the $0 \rightarrow z$ Wilson line, and $\nu = p \cdot z$ represents a Lorentz invariant quantity known as the Ioffe time. Through a choice of p , z , and α one is able to isolate the term containing leading twist contributions $\mathcal{M}(\nu, z^2)$ [8].

Using pseudo-PDF formalism discussed in [9–11], the real component of this term, $\mathcal{M}_R(\nu, z^2)$, is matched to the \overline{MS} Ioffe time PDF, $\mathcal{Q}_R(\nu, \mu^2)$, via a perturbative kernel. The \overline{MS} Ioffe time PDF on the other hand is related to the valence quark PDF, $q_v(x, \mu^2)$, through the quasi-Fourier transformation

$$\mathcal{Q}_R(\nu, \mu^2) \equiv \int_0^1 dx \cos(\nu x) q_v(x, \mu^2), \quad (1.2)$$

as was shown in [8, 12].

Using the quasi-PDF formalism, originally proposed in [7], the lattice calculated $\mathcal{M}_R(\nu, z^2)$ is related to the quasi-PDF, $\tilde{q}_v(y, p_3)$, via the Fourier transformation

$$\mathcal{M}_R(\nu, \frac{\nu^2}{p_3^2}) \equiv \int_0^\infty dy \cos(\nu y) \tilde{q}_v(y, p_3). \quad (1.3)$$

The quasi-PDF can then be connected to the PDF using a perturbative kernel as discussed in [11, 13, 14]. Immediately after the presentation of the basic idea a plethora of works [15–21] explored the properties of the new methodologies as well as of other approaches [22–26]. Possible doubts against this approach that emanate from the need to invert the Fourier transform in Eq. (1.3) were raised in [27, 28]. These claims were refuted in [29–31], however as we discuss in this paper, apart from the theoretical issues raised, the numerical inversion of the Fourier transform is not straight forward. We refer the reader to [32, 33] for two detailed reviews of the topic.

This study will be concerned with the reconstruction of the PDF from the \overline{MS} Ioffe time PDF, but many of the conclusions are also applicable to the reconstruction of quasi-PDFs from the lattice calculated matrix element. In inverting Eq.(1.2) to obtain $q_v(x, \mu^2)$, there exist two challenges. First being that the required integral over ν does not extend over *the full Brillouin zone*. The second is that in practice only a *small number of points along ν* can be computed. Only the second challenge exists for the calculation of the quasi-PDF. As we will discuss in more detail below, taken together these issues render the extraction highly ill-posed and we explore several regularization strategies on how to nevertheless reliably estimate the PDF from the data at hand.

To assess the viability of different reconstruction approaches in practice and to elucidate their systematic uncertainties, we will carry out tests based on two different mock data sets. The first test scenario is based on experimentally determined PDF's for which it has been found that a simple ansatz is able to approximate their functional form quite well

$$p(x) = \frac{\Gamma(a+b+2)}{\Gamma(a+1)\Gamma(b+1)} x^a (1-x)^b. \quad (1.4)$$

Based on phenomenological fits the expectation is that, for scales $\mu > 2\text{GeV}$, the physical PDF shows a divergence close to $x = 0$, while vanishing at $x = 1$. This requires that $a < 0$

while $b > 0$. In mock scenario A, we insert into Eq. (1.2) such experimentally determined PDF's, which in turn tests the case $a < 0$.

On the other hand at very low scales μ^2 , lattice results in the quenched approximation and with heavy pions [8] suggest that a may become positive instead. Thus for mock scenario B we deform by hand the experimental PDF data so that it goes to zero at the origin. This scenario B with its different functional form serves as a gauge of the robustness of the methods we are testing. Also, scenario B is reminiscent of the quasi-PDF case which is known to converge in the limit of y going to zero.

This article starts out with the study of the direct inversion in section 2, followed by the recapitulation of the main ideas of advanced reconstructions in section 3. In section 4 we present extensive numerical experiments of all methods employing mock data and section 5 summarizes our conclusions.

2 Direct Inversion

While Eq. (1.2) is a Fourier transform and inverting may be considered a simple task, the facts that the range of ν extends only over a finite interval different from the full Brillouin zone and that our data for \mathcal{Q}_R are discrete, makes its inversion highly ill-posed. Let us explore this issue in more detail by naively discretizing the integral in the interval $x \in [0, 1]$, considering $N_x + 1$ points in a trapezoid integration rule. In that case

$$\Delta x = \frac{1}{N_x}, \quad x_k = k\Delta x = \frac{k}{N_x} \quad (2.1)$$

and

$$\mathcal{Q}_R(\nu) = \frac{\Delta x}{2} \cos(\nu x_0) q_v(x_0) + \sum_{k=1}^{N_x-1} \Delta x \cos(\nu x_k) q_v(x_k) + \frac{\Delta x}{2} \cos(\nu x_{N_x}) q_v(x_{N_x}). \quad (2.2)$$

Now let us further assume that we have an equal amount of $N_\nu + 1$ data points for $\mathcal{Q}_R(\nu)$. Then we can determine exactly the unknown values of the function $q_v(x_k)$ by solving a simple linear system of equations. We define the vector \mathfrak{Q} with components

$$\mathfrak{Q}_k = \mathcal{Q}_R(\nu_k) \quad (2.3)$$

where ν_k are the values of the Ioffe time for which data is available. Also let \mathfrak{q} be the vector with components the unknown values of $q_v(x_k)$ *i.e.*

$$\mathfrak{q}_k = q_v(x_k). \quad (2.4)$$

Then Eq. (2.2) can be written in matrix form as

$$\mathfrak{Q} = \mathfrak{C} \cdot \mathfrak{q}, \quad (2.5)$$

with \mathfrak{C} being the coefficient matrix with matrix elements,

$$\begin{aligned} \mathfrak{C}_{kl} &= \Delta x \cos(\nu_k x_l) = \frac{1}{N_x} \cos(\nu_k x_l) \quad \text{for } l \in [1, N_x - 1], \\ \mathfrak{C}_{kl} &= \frac{1}{2} \Delta x \cos(\nu_k x_l) = \frac{1}{2} \frac{1}{N_x} \cos(\nu_k x_l) \quad \text{for } l = 0, N_x. \end{aligned} \quad (2.6)$$

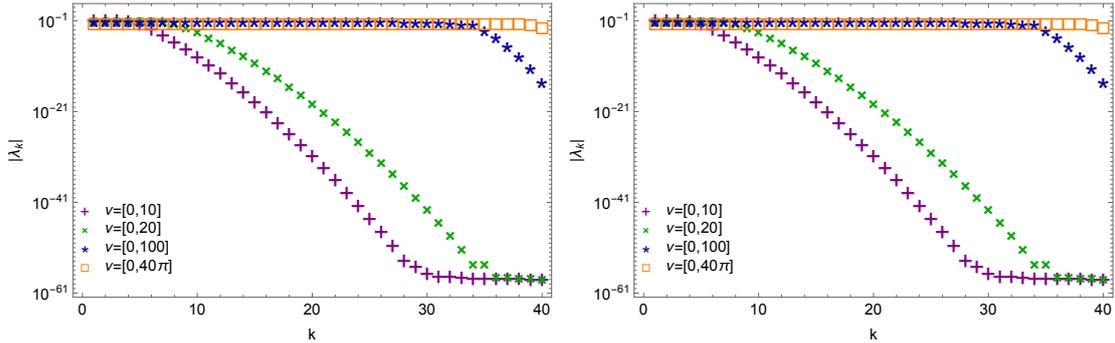


Figure 1. Eigenvalues λ_k of the Kernel matrix for the direct inversion method on the left and the derivative inversion method on the right calculated with different discretization intervals of ν . Note that only for the case corresponding to a genuine discrete Fourier transform, $\nu = [0, 40\pi]$, do all the eigenvalues remain of the same order. All other ranges show exponentially suppressed eigenvalues. The realistic case of $\nu = [0, 20]$ already shows a significant degradation of the eigenvalue spectrum.

However, Eq. (2.5) may be badly conditioned rendering the computation of \mathfrak{q} a difficult task. As a concrete example we take an idealized situation of having in our possession $N_\nu = 40$ data points for \mathcal{Q}_R over the four different intervals $I_0 = [0, 10]$, $I_1 = [0, 20]$, $I_2 = [0, 100]$, $I_3 = [0, 40\pi]$. Interval I_3 is special in the sense that with a redefinition of x one may rewrite Eq.(1.2) in the form of a genuine discrete Fourier transform, which we know is well conditioned. And indeed as shown in Fig. 1, for $\nu \in I_3$ the eigenvalues λ_k of the matrix \mathfrak{C} are all of the same order ($\mathcal{O}(0.1)$). Once the resolved ν region is shrunk below the full Brillouin zone, the cosine functions no longer constitute linearly independent basis functions and the columns of \mathfrak{C} become linearly dependent, which manifest itself in eigenvalues that are exactly zero and eigenvalues that are exponentially suppressed.

The exponential decay of the eigenvalues tells us that the inversion problem is ill-conditioned and that a direct inversion will become impractical once the ν range is significantly smaller than the full Brillouin zone. To make this explicit we carry out a direct inversion of mock data.

For this illustration we take Eq. (1.4) with parameters $a = -\frac{1}{4}$ and $b = 3$. Ideal data for \mathcal{Q}_R are obtained based on the three different discretization intervals I_1, I_2 and I_3 . The ideal data are then distorted by Gaussian noise corresponding to constant and uncorrelated relative errors on the averaged data of $\Delta\mathcal{Q}_R/\mathcal{Q}_R = \text{const}$. The matrix inverse is computed via a singular value decomposition, where only singular values which are larger than 10^{-4} are inverted.

In Fig. 2 the results of the direct inversion are shown, with the ideal data results depicted by red circles and the original $q(x)$ as gray dashed line. The leftmost panel corresponds to the well conditioned case of I_3 , where the reconstruction based on ideal data works flawlessly and even produces accurate results already for a relative large error on the input data of $\Delta\mathcal{Q}_R/\mathcal{Q}_R = 10^{-2}$. One exception is the point at $x = 0$, which formally would have to diverge.

The shorter the ν interval becomes the worse the reconstruction results, where even for

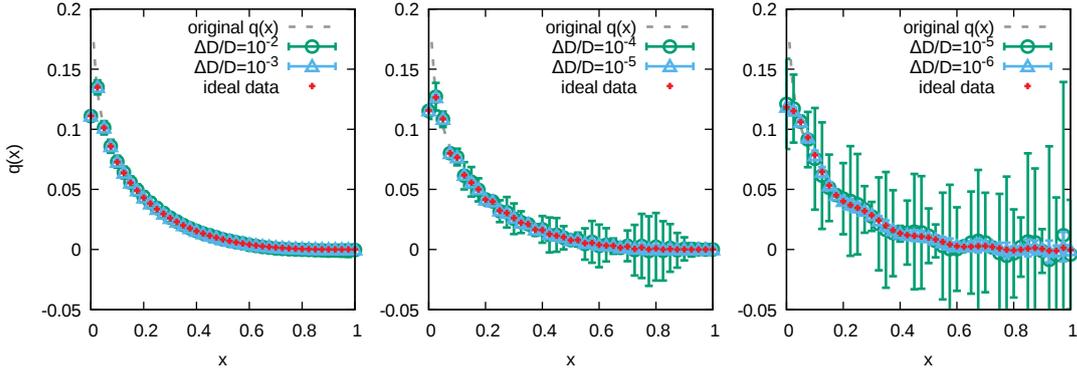


Figure 2. Direct inversion results for different discretization intervals of ν (left $\nu = [0, 40\pi]$, center $\nu = [0, 100]$, right $\nu = [0, 20]$). The matrix inverse is regularized by retaining only singular values of \mathfrak{C} that are larger than 10^{-4} . Note the different size of the relative errors needed, to obtain a well behaved result (left $\Delta\mathcal{Q}_R/\mathcal{Q}_R = 10^{-2}$, center $\Delta\mathcal{Q}_R/\mathcal{Q}_R = 10^{-5}$, right $\Delta\mathcal{Q}_R/\mathcal{Q}_R = 10^{-6}$).

ideal data having an error on the level of standard machine precision we obtain artificial fluctuations of the reconstructed PDF. At the same time, significantly smaller relative errors on the input data is required for a reasonably stable reconstruction to ensue. For $\nu = [0, 100]$ we already need $\Delta\mathcal{Q}_R/\mathcal{Q}_R = 10^{-5}$, while for $\nu = [0, 20]$ even the $\Delta\mathcal{Q}_R/\mathcal{Q}_R = 10^{-6}$ result only gives an approximate PDF with significant unphysical fluctuations.

One may attempt to improve the results of the direct inversion by considering a higher order integration scheme. At large values of ν , the integrand is a highly oscillatory function due to the presence of $\cos(\nu x)$. As a result, an integration algorithm that approximates the integrand by a low degree polynomial is bound to fail at large values of ν . We can improve on this by designing a better integration rule that performs similarly for all values of ν . We know that the oscillatory nature of the integrand is due to the $\cos(\nu x)$ term and that the unknown function $q_\nu(x)$ is slowly varying in the interval $[0, 1]$. Therefore we will approximate $q_\nu(x)$ with a linear interpolation and perform the integral exactly leaving the result to linearly depend on the unknown values of the function $q_\nu(x)$ on the grid points x_k . A linear interpolation $f(x)$ for the function $q_\nu(x)$ is given by

$$f(x) = \frac{x - x_k}{x_{k+1} - x_k} q_\nu(x_{k+1}) + \frac{x_{k+1} - x}{x_{k+1} - x_k} q_\nu(x_k) \text{ for } x \in [x_k, x_{k+1}]. \quad (2.7)$$

In order to compute the exact integral of the interpolating function we need to define

$$I_0(a, b, \nu) = \int_a^b dx \cos(\nu x) = \frac{1}{\nu} [\sin(\nu b) - \sin(\nu a)],$$

$$I_1(a, b, \nu) = \int_a^b dx x \cos(\nu x) = \frac{1}{\nu} [b \sin(\nu b) - a \sin(\nu a)] + \frac{1}{\nu^2} [\cos(\nu b) - \cos(\nu a)]. \quad (2.8)$$

Note that both integrals are finite for $\nu = 0$,

$$I_0(a, b, 0) = b - a,$$

$$I_1(a, b, 0) = \frac{1}{2} [b^2 - a^2]. \quad (2.9)$$

With these results at hand we can now write down the improved integration rule,

$$\mathcal{Q}_R(\nu) = \int_0^1 dx \cos(\nu x) q_v(x) \approx \int_0^1 dx \cos(\nu x) f(x). \quad (2.10)$$

The approximate integral is now a sum of integrals on the intervals $[x_k, x_{k+1}]$. To simplify expressions we introduce the notation

$$I_i(k, \nu) \equiv I_i(x_k, x_{k+1}, \nu) \text{ for } i = 0, 1. \quad (2.11)$$

With this notation we have

$$\mathcal{Q}_R(\nu) \approx \sum_{k=0}^{N_x-1} \left[\frac{I_1(k, \nu) - x_k I_0(k, \nu)}{x_{k+1} - x_k} q_v(x_{k+1}) + \frac{x_{k+1} I_0(k, \nu) - I_1(k, \nu)}{x_{k+1} - x_k} q_v(x_k) \right], \quad (2.12)$$

or

$$\begin{aligned} \mathcal{Q}_R(\nu) \approx & \frac{x_1 I_0(0, \nu) - I_1(0, \nu)}{x_1 - x_0} q_v(x_0) + \\ & + \sum_{k=1}^{N_x-1} \left[\frac{I_1(k-1, \nu) - x_{k-1} I_0(k-1, \nu)}{x_k - x_{k-1}} + \frac{x_{k+1} I_0(k, \nu) - I_1(k, \nu)}{x_{k+1} - x_k} \right] q_v(x_k) \\ & + \frac{I_1(N_x-1, \nu) - x_{N-1} I_0(N-1, \nu)}{x_N - x_{N-1}} q_v(x_N). \end{aligned} \quad (2.13)$$

Note that for $\nu = 0$ the above expression simplifies to the trapezoid rule for the parton density $q_v(x)$.

Using the same notation as before, we can now compute the matrix elements of the coefficient matrix \mathfrak{C} as

$$\begin{aligned} \mathfrak{C}_{lk} = & \frac{I_1(k-1, \nu_l) - x_{k-1} I_0(k-1, \nu_l)}{x_k - x_{k-1}} + \frac{x_{k+1} I_0(k, \nu_l) - I_1(k, \nu_l)}{x_{k+1} - x_k} \text{ for } k \in [1, N-1], \\ & \mathfrak{C}_{lk} = \frac{x_1 I_0(0, \nu_l) - I_1(0, \nu_l)}{x_1 - x_0} \text{ for } k = 0, \\ & \mathfrak{C}_{lk} = \frac{I_1(N-1, \nu_l) - x_{N-1} I_0(N-1, \nu_l)}{x_N - x_{N-1}} \text{ for } k = N, \end{aligned} \quad (2.14)$$

We have tested this improved integration method and compared it to the trapezoid rule. It turns out that the relative integration error for typical functions $q_v(x)$ with 10 interpolating points continuously grows with ν for the trapezoid rule and reaches 100% at $\nu = 15$. On the other hand the improved integration, which performs similarly with the trapezoid rule at small ν , at $\nu = 15$ has 65% relative error which remains almost unchanged up to $\nu = 100$. The improved integration achieves a constant relative error as a function of ν at large ν for any number of interpolating points. On the other hand the trapezoid rule has a relative error that always increases with ν for any number of interpolating points. Therefore, this improved integration scheme achieves its design goal of having a constant integration error for all values of ν . We could further improve the integration by using a second order interpolation formula therefore producing a generalized

Simpson's rule. Perhaps this is needed because the number of points we have is very limited. However, it may be that the largest systematic in our problem is the truncation at relatively small values of ν and further improvement of the integration scheme will not affect this systematic. The proposed integration scheme can prove valuable because it can significantly reduce the number of points required to discretize the integral resulting a smaller maximum value of ν for which the problem is no longer ill-defined.

It has been proposed in [34], that the unphysical oscillations in the related quasi-PDF inverse problem can be controlled by considering the derivative of the integral equation with respect to ν or z_3 . Even if the derivative method results in a smoother curve, it does not alleviate the ill-posed nature of the problem. Assuming that the derivative can be calculated explicitly or accurately, Eq. (1.2) becomes

$$\partial\mathcal{Q}(\nu) = - \int_0^1 dx x \sin(\nu x) q(x), \quad (2.15)$$

which can be discretized as before to

$$\partial\mathcal{Q} = (\partial\mathfrak{C}) \cdot \mathfrak{q}, \quad (2.16)$$

with $\partial\mathfrak{C}$ being the coefficient matrix with matrix elements,

$$\begin{aligned} [\partial\mathfrak{C}]_{kl} &= -\Delta x x_l \sin(\nu_k x_l) = -\frac{1}{N_x} x_l \sin(\nu_k x_l) \quad \text{for } l \in [1, N_x - 1], \\ [\partial\mathfrak{C}]_{kl} &= -\frac{1}{2}\Delta x x_l \sin(\nu_k x_l) = -\frac{1}{2} \frac{1}{N_x} x_l \sin(\nu_k x_l) \quad \text{for } l = 0, N_x. \end{aligned} \quad (2.17)$$

Using the same intervals of ν and x mentioned above, one can find that the eigenvalues of this matrix, shown in Fig. 1, have the same pattern of exponential decay which characterizes an ill-conditioned inverse problem.

Our numerical experiments indicate that for realistic scenarios encountered in lattice studies, for which we foresee $\nu \in [0, 20]$ and $\Delta\mathcal{Q}_R/\mathcal{Q}_R = 10^{-2} - 10^{-3}$, a direct inversion is not feasible and more sophisticated regularization schemes of the ill-defined inversion problem are needed. In particular we have observed that the truncation in ν , results in a reconstructed PDF with unphysical fluctuations. This is an effect similar to that observed in [34]. In addition we also see that any non-analytic behavior of $q(x)$, e.g. a divergence at $x = 0$ will pose a difficulty for (direct) inversion methods and needs to be considered with care. Therefore in the main part of this paper we will explore several modern methods for treating inverse problems and compare their efficiency in dealing with the uncertainties introduced by the truncation of the integration regime.

3 Advanced PDF reconstructions

The fundamental difficulty of solving an ill-posed inverse problem lies in the fact that the input data by themselves do not single out a unique answer. To give meaning to the PDF extraction the inversion needs to be regularized, i.e. we need to provide criteria on how to choose a single PDF from an infinite number of possible solutions reproducing the discrete and noisy input.

Naïve methods, such as the direct inversion and the derivative method, introduce a regularization by removing small singular values of the kernel or by setting data to zero at large distances. However, these approaches introduce uncontrolled systematic errors. A straightforward example of a stable solution to the inverse problem, which was used in [8, 12], is to parametrize the solution with a functional form, such as (1.4), containing a small number of parameters. One can then use this functional form to fit the data using a χ^2 minimization. This approach typically results in smaller statistical errors, however, it introduces a model systematic that one may control by varying the functional form used in the model. As the flexibility of the parametrization increases the model systematic is reduced and the ill-posed nature of the inverse problem becomes more apparent. More sophisticated methods often use some additional prior information to constrain the results in order to improve the robustness of the result and provide controlled systematic error estimates.

In the following, we discuss three extraction methods that do not presuppose a functional form of the encoded PDFs or have a very flexible functional form parametrizing the solution. Each of these methods regulates the inverse problem in different ways. The first, the Backus-Gilbert method, regulates the problem by minimizing the statistical variance of the solution. The second, a neural network parametrization, that provides a flexible parametrization given by specific choices of network geometry and activation functions. The third class of methods, is based on Bayesian inference, that rely on Bayes theorem to systematically incorporate prior information, such as e.g. positivity or smoothness of the solution. In addition, the Bayesian inference methods introduce a default model, which represent the prior information in that this model is the correct solution in the absence of any data. Non Bayesian approaches, while not explicitly mentioning prior information often incorporate additional knowledge about the problem at hand in an indirect fashion. One such way is to apply a preconditioning to the inverse problem, in case that a relatively good guess of the form of the solution is available. In a Bayesian method such information would instead be supplied in the form of a default model.

Let us discuss the preconditioning, which we will employ in order to improve the results of the Backus-Gilbert and neural network methods. The main ingredient is to define a rescaled kernel and rescaled PDF $h(x)$

$$K_j(x) \equiv \cos(\nu_j x) p(x) \quad \text{and} \quad h(x) \equiv \frac{q_v(x)}{p(x)}, \quad (3.1)$$

where $p(x)$ corresponds to an appropriately chosen function that makes the problem easier to solve. In particular, we wish to incorporate into $p(x)$ most of the non-trivial structure of $q(x)$, such that $h(x)$ is a slowly varying function of x and contains only the deviation of $q(x)$ from $p(x)$.

As discussed in the introduction, it has been found that the Ansatz Eq. (1.4) describes phenomenological PDFs rather well. Integrated over the cosine, it yields the function

$$\mathcal{Q}_p(\nu, a, b) = \pi 2^{-a-b-1} \Gamma(a+1) \Gamma(b+1) {}_2\tilde{F}_3 \left(\frac{a+1}{2}, \frac{a+2}{2}; \frac{1}{2}, \frac{1}{2}(a+b+2), \frac{1}{2}(a+b+3); -\frac{\nu^2}{4} \right) \quad (3.2)$$

where ${}_2\tilde{F}_3$ is the generalized hypergeometric function. One can find a good choice for the parameters a and b by first fitting the Ioffe-time PDF data with Eq. (3.2). The choice of preconditioning function or default model will influence the result of the reconstruction and it is therefore necessary to explore the stability of the end result on that choice.

3.1 Backus-Gilbert Method

One approach to inverse problems, which has been used in a number of engineering and physics applications, is the Backus-Gilbert method [35–38]. Like many approaches to the inverse problem, the Backus-Gilbert method provides a unique solution to the ill posed problem given some condition. This method differs from other regularizations, which impose a smoothness criterion on the resulting function $q(x)$, by imposing a minimization condition of the variance of the resulting function. The Backus-Gilbert method has been studied previously in [39] as a solution to an inverse problem for extracting PDFs from other types of lattice calculated hadronic matrix elements.

Let us start from the preconditioned expression (3.1) with a rescaled PDF $h(x)$ that is only a slowly varying function of x . Hence our inverse problem becomes

$$d_j \equiv \mathcal{Q}_R(\nu_j) = \int_0^1 dx K_j(x) h(x). \quad (3.3)$$

The Backus-Gilbert method introduces a function $\Delta(x - \bar{x})$ that is written as

$$\Delta(x - \bar{x}) = \sum_j a_j(\bar{x}) K_j(x), \quad (3.4)$$

where $a_j(\bar{x})$ are unknown functions to be determined. It then estimates the unknown function as a linear combination of the data,

$$\hat{h}(\bar{x}) = \sum_j a_j(\bar{x}) d_j, \quad \text{or} \quad \hat{q}_v(\bar{x}) = \sum_j a_j(\bar{x}) d_j p(\bar{x}). \quad (3.5)$$

Given the above definitions, if $\Delta(x - \bar{x})$ were to be a Dirac δ function then $\hat{h}(\bar{x})$ would be equal to $h(\bar{x})$. If $\Delta(x - \bar{x})$ approximates a δ -function with a width σ , then the smaller σ is, the better the approximation of $\hat{h}(x)$ to $h(x)$. In other words, if $\hat{h}_\sigma(x)$ is the approximation resulting from $\Delta(x)$ with a width σ then

$$\lim_{\sigma \rightarrow 0} \hat{h}_\sigma(x) = h(x). \quad (3.6)$$

With this in mind the Backus-Gilbert method minimizes the width σ given by

$$\sigma = \int_0^1 dx (x - \bar{x})^2 \Delta(x - \bar{x})^2. \quad (3.7)$$

Note that other choices for the definition of the width can be used. This choice makes the resulting integrals easy to compute and the minimization problem becomes quadratic in the unknown values a_j . Furthermore, if $\Delta(x)$ is to approximate a δ -function then one has to impose the constraint

$$\int_0^1 dx \Delta(x - \bar{x}) = 1. \quad (3.8)$$

Using a Lagrange multiplier λ one can minimize the width and impose the constraint by minimizing

$$\chi[a] = \int_0^1 dx (x - \bar{x})^2 \sum_{j,k} a_j(\bar{x}) K_j(x) K_k(x) a_k(\bar{x}) + \lambda \int_0^1 dx \sum_j K_j(x) a_j(\bar{x}). \quad (3.9)$$

By setting the derivative to zero with respect to $a_j(\bar{x})$ we get

$$\frac{\partial \chi[a]}{\partial a_j(\bar{x})} = 2 \int_0^1 dx (x - \bar{x})^2 \sum_k K_j(x) K_k(x) a_k(\bar{x}) + \lambda \int_0^1 dx K_j(x) = 0, \quad (3.10)$$

which results to

$$\int_0^1 dx (x - \bar{x})^2 \sum_k K_j(x) K_k(x) a_k(\bar{x}) = -\frac{1}{2} \lambda \int_0^1 dx K_j(x). \quad (3.11)$$

Let's now define the matrix \mathbf{M} with matrix elements

$$M_{jk} = \int_0^1 dx (x - \bar{x})^2 K_j(x) K_k(x), \quad (3.12)$$

the vector \mathbf{u} with components

$$u_j = \int_0^1 dx K_j(x), \quad (3.13)$$

and promote $a_j(\bar{x})$ to a vector \mathbf{a} . With these definitions the minimization condition takes the matrix equation form

$$\mathbf{M}\mathbf{a} = -\frac{1}{2} \lambda \mathbf{u} \quad (3.14)$$

or

$$\mathbf{a} = -\frac{1}{2} \lambda \mathbf{M}^{-1} \mathbf{u}. \quad (3.15)$$

Imposing the normalization condition for $\Delta(x)$ we get

$$\lambda = -2 [\mathbf{u}^T \mathbf{M}^{-1} \mathbf{u}]^{-1}, \quad (3.16)$$

therefore,

$$\mathbf{a} = \frac{1}{\mathbf{u}^T \mathbf{M}^{-1} \mathbf{u}} \mathbf{M}^{-1} \mathbf{u}. \quad (3.17)$$

We can now obtain an estimate of the unknown function using Eq. (3.5). We can also compute the width of the function $\Delta(x)$ which is otherwise known as the resolution function. The width of the distribution $\Delta(x - \bar{x})$ is given by

$$\sigma(\bar{x}) = \mathbf{a}^T \mathbf{M} \mathbf{a}. \quad (3.18)$$

Note that the width is dependent on \bar{x} because \mathbf{a} depends on \bar{x} and represents the resolution of the method at $x = \bar{x}$. Features of the unknown function at scales shorter than $\sigma(\bar{x})$ will not be resolved and thus the method works better when the unknown function is smooth. This approach is known to fail when the matrix \mathbf{M} becomes singular, and thus various techniques for regularizing the inversion of \mathbf{M} have been proposed. As mentioned before,

minimization of the variance of the resulting solution leads to a regularization of the matrix \mathbf{M} with the covariance matrix of the data. This is implemented by the substitution

$$\mathbf{M} \rightarrow \mathbf{M} + \rho \mathbf{C}, \quad (3.19)$$

where \mathbf{C} is the covariance matrix of the data and ρ is a small free parameter. The larger ρ the better conditioned the matrix \mathbf{M} and the lower the resolution of the method is. A similar method is the Tikhonov regularization [40, 41], where one makes the substitution

$$\mathbf{M} \rightarrow \mathbf{M} + \rho \mathbf{1}, \quad (3.20)$$

where $\mathbf{1}$ is the identity matrix and ρ a small adjustable parameter. In this case all the singular values of \mathbf{M} are lifted resulting in a well-defined matrix \mathbf{M} . However, the Tikhonov regularization does not minimize the variance of the resulting solution of the inverse problem as the covariance matrix regularization does. Again larger values of ρ result in better regulated matrices at the cost of reduced resolution.

Another regularization method similar to Tikhonov is the method we adopted for our experiments. We just project out the singular vectors with singular values smaller than a given cut-off ρ which we choose. For our problem we have noticed that the Tikhonov regularization and the approach we chose result in nearly identical solutions of the inverse problem. Furthermore, the covariance matrix regularization for appropriate values of the parameter ρ produces very similar results as the other two methods for our problem. Therefore, we decided to use our simple approach for regulating the matrix \mathbf{M} and present results of our experiments using this SVD cutoff method.

The Backus-Gilbert method provides a unique solution to the inverse problem demanding that the solution maximizes a stability criterion. The method has a tunable free parameter, ρ , which provides a trade off between stability and resolution. This freedom can lead to a bias from the user, but the Backus-Gilbert method provides the variance and the resolution function as quantitative measures to assure a proper analysis is performed. It should be noted that the forms of preconditioning functions, $p(x)$, are restricted in the Backus-Gilbert method. The preconditioning function must be chosen such that integrals of the preconditioned kernels which define \mathbf{M} and \mathbf{u} remain finite. Furthermore, it should be noted that a good preconditioning function makes the remainder smooth, as a result the demand for a small resolution is smaller the better the preconditioning is.

3.2 Neural Network Reconstruction

Neural networks can be used as a rather flexible parametrization of functions and thus have been used in the literature to address various inverse problems. Neural networks do not provide any explicit constraints such as a smoothness condition or minimization of a variance, as the Backus-Gilbert method does but instead it is the depth and structure of their layers that limits what functional forms may be encoded. Hence, neural networks provide a very general parametrization of the unknown function, without forcing a model dependent form such as Eq. (3.2), for a statistical regression. Neural networks have e.g. been used in the literature to extract PDFs from experimental data. The pioneering work

of [42] is now one of the established approaches in obtaining PDFs from cross section data by the NNPDF collaboration [43–46]. In this work, we explore the use of neural networks in solving the inverse problem at hand. The specific neural network implementation used in this work is known as a multilayer feedforward neural network.

A neural network is composed of a system of interconnected nodes, called the neurons. The connection between neuron pairs is called the synapsis. The output of each neuron, called the activation, is typically a non-linear transformation of its input, which is called the activation function of that neuron. The input for each neuron is a weighted sum of the activations of the connected neurons shifted by some real number, which is called the threshold or bias of the neuron. The weights and the thresholds are the free parameters which will be determined by some process which is called the training procedure.

In a multilayer neural network, the neurons are organized into distinct layers whose neurons are not directly connected to each other. In a multilayer feedforward network, the layers are ordered and each neuron on a given layer is only connected to the neurons on the previous and subsequent layers. The geometry of these networks, shown in Fig. 3, is described by a list of numbers, N_1, \dots, N_L , giving the number of neurons on each of the L layers. The first layer, $\xi^{(1)}$, is a given input vector of length N_1 and the final layer, $\xi^{(L)}$ is the output vector of length N_L , also called the response of the neural network. All other layers are referred to as hidden layers. The activation of the i -th neuron on the l -th layer of the network is given by the recursive relationship

$$\xi_i^{(l)} = g_i^{(l)} \left(\sum_j^{N_{l-1}} w_{ij}^{(l)} \xi_j^{(l-1)} - \theta_i^{(l)} \right). \quad (3.21)$$

Specified by the geometry, activation functions, thresholds, and weights, a multilayer feedforward neural network can be considered a parametrization of a function from $\mathbb{R}^{N_1} \rightarrow \mathbb{R}^{N_L}$. For a given geometry and a set of activation functions, the thresholds and weights can be chosen for the network to approximate a given continuous function. For the case of reconstructing the rescaled PDF $h(x)$, which is a single valued function with a single argument, the geometry is restricted to have $N_1 = 1$ and $N_L = 1$. It should be noted that NNPDF has also used $\log \frac{1}{x}$ as a second argument in the input layer [47], which is important to reproduce the small x behavior of the PDF. However, in our case our data are not sensitive to small x and thus using $\log \frac{1}{x}$ at the input layer is not essential.

A neural network can be used to perform a regression by choosing the thresholds and weights with a training procedure. During the training procedure, the weights and thresholds are modified to minimize some error function, which describes the difference between the response of the neural network and some desired output. When using a neural network to perform a statistical regression, a common choice of error function is the χ^2 function, e.g.

$$\chi^2(\{w\}, \{\theta\}) = \sum_{k=1}^N \left(Q_k - \int_0^1 dx K_k(x) h(x; \{w\}, \{\theta\}) \right)^2 / \sigma_k^2, \quad (3.22)$$

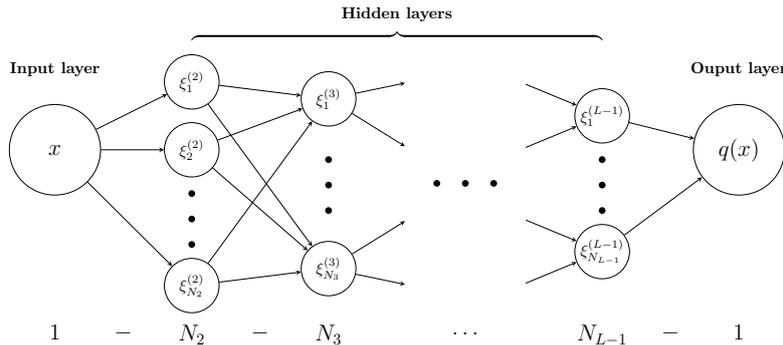


Figure 3. A neural network can be used as a general parametrization of an unknown function from $\mathbb{R}^{N_1} \rightarrow \mathbb{R}^{N_L}$. For the case of a PDF, a single valued function of a single argument, the input and output layers have only one neuron.

where \mathcal{Q}_k are N data points with standard deviations σ_k and h is the output layer of the neural network given an input layer x , weights $\{w\}$, and thresholds $\{\theta\}$. Using a neural network to parametrize the unknown function may result in a $\chi^2(\{w\}, \{\theta\})$ with a large number of local minima. Some these local minima are trivial multiplicities due to symmetries a neural network has under permutations of the weights and thresholds. However, the possibility of multiple non-trivial local minima exists resulting in many realizations of the network that reproduce the data equally well. In these cases special care has to be taken to avoid “over-fitting” and several methods to do so have been developed in the literature [47].

These roughly equivalent minima can be found by a training procedure such as a genetic algorithm. A genetic algorithm is an iterative process based upon the idea of natural selection. Each iteration, also called a generation, begins with a sample of possible networks, called a population. A fitness function is evaluated for each of the networks, which in this case is the error function $\chi^2(\{w\}, \{\theta\})$. Those networks which are the “least fit”, i.e. largest χ^2 , are removed from the population. The surviving population is then “mutated” by randomly changing their parameters, i.e. weights and thresholds, to create the starting population for a new generation. This procedure is iterated for enough generations that a final population covers a sufficient number of minima with sufficiently small values of the error function.

The genetic algorithm used in this study is based upon simulated annealing. The initial population is created from N_{rep}^0 sets of initial weights and thresholds which are generated from a random normal distribution with a wide initial search width, σ_0 .

$$\{w\}_i^{(0)}, \{\theta\}_i^{(0)} \sim \mathcal{N}(0, \sigma_0) \quad (3.23)$$

where $\{w\}_i^{(0)}$ and $\{\theta\}_i^{(0)}$ are the parameters of the i -th neural network in the initial population. The weights and thresholds of this initial population of neural networks is the starting point for minimizing with respect to χ^2 to some new values $\{w'\}_i^{(0)}$ and $\{\theta'\}_i^{(0)}$. The resulting neural networks which are the least fit, i.e. that have the largest χ^2 , are removed from the population and not used in the next generation. The surviving neural networks are then mutated by adjusting their weights and thresholds by some small,

Gaussian distributed amount with a mutation width σ ,

$$\{w\}_i^{(g)} \sim \{w'\}_i^{(g-1)}(1 + \eta\mathcal{N}(0, \sigma)) \quad \{\theta\}_i^{(g)} \sim \{\theta'\}_i^{(g-1)}(1 + \eta\mathcal{N}(0, \sigma)), \quad (3.24)$$

where η is the relative size of the mutation, $\{w\}_i^{(g)}$ and $\{\theta\}_i^{(g)}$ are the initial parameters of the i -th neural network in the g -th generation and $\{w'\}_i^{(g)}$ and $\{\theta'\}_i^{(g)}$ are the parameters of the i -th surviving neural network in the g -generation after minimization. This procedure of minimization and removal is repeated for N_{gen} generations. After enough generations have passed, the population of neural networks will all exist in various minima of χ^2 which can be used to estimate the PDF.

Many other training procedures exist which all have different methods of identifying minima, both stochastic and deterministic. An example of such a procedure is the approach the NNPDF collaboration is taking in obtaining PDFs from experimental data [47]. Though different training procedures vary in their efficiency for finding minima, they should eventually all converge on the same set of minima of the error function.

Neural networks provide a rather general solution to the inverse problem. The large number of highly interconnected parameters allows for model independent result, but the complexity of the hidden layers do not give any insight to what effects underlie that result.

3.3 Bayesian PDF reconstruction

The third approach to inverse problems we discuss, which has proven to be versatile in practice, is Bayesian inference. It acknowledges the fact that the inverse problem is ill-defined and a unique answer may only be provided once further information about the system has been made available. This method does not require any explicit constraints, though it may benefit from them, nor does it require any functional form, even one as flexible as a neural network. This method finds the most probable value of $q(x)$ given the data and whatever prior information is provided. Other regularization methods may often be rewritten in terms of Bayesian reconstruction with the constraints treated as prior information.

Formulated in terms of probabilities, one finds in the form of Bayes theorem that

$$P[q|\mathcal{Q}, I] = \frac{P[\mathcal{Q}|q, I]P[q|I]}{P[\mathcal{Q}|I]}. \quad (3.25)$$

It states that the so called posterior probability, $P[q|\mathcal{Q}, I]$, for a test function q to be the correct x -space PDF, given lattice calculated Ioffe time PDF \mathcal{Q} and additional prior information I may be expressed in terms of three quantities. The likelihood probability $P[\mathcal{Q}|q, I]$ denotes how probable it is to find the data \mathcal{Q} if q were the correct PDF. Obtaining the most probable q by maximizing the likelihood is nothing but a χ^2 fit to the \mathcal{Q} data, which as we saw from the direct inversion is by itself ill-defined. The second term of importance is the prior probability $P[q|I]$, which quantifies, how compatible our test function q is with respect to any prior information we have. In particular such information can be related to the appearance of non-analytic behavior of $q(x)$ at the boundaries of the x interval or the positivity of the PDF. $P[\mathcal{Q}|I]$, the so called evidence represents a q independent normalization.

For sampled (lattice) data, due to the central limit theorem, the likelihood probability may be written as the quadratic distance functional $P[\mathcal{Q}|q, I] = \exp[-L]$,

$$L = \frac{1}{2} \sum_{k,l} (\mathcal{Q}_k - \mathcal{Q}_k^q) C_{kl}^{-1} (\mathcal{Q}_l - \mathcal{Q}_l^q), \quad (3.26)$$

where \mathcal{Q}_k^q are the Ioffe-time data arising from inserting the test function q into Eq.(1.2) and

$$C_{kl} = \frac{1}{N_m(N_m - 1)} \sum_h (\mathcal{Q}_k^h - \mathcal{Q}_k) (\mathcal{Q}_l^h - \mathcal{Q}_l), \quad (3.27)$$

denotes the covariance matrix of the N_m measurements of simulation data \mathcal{Q}_k^h .

For the regularization of the inversion task, we now have to specify an appropriate prior probability $P[q|I] = \exp[\alpha S[q, m]]$. Prior information enters in two ways here: on the one hand the (unique) extremum of the S functional is given by the so called default model m , which, by means of Bayes theorem, represents the most probable answer in the absence of data. On the other hand the functional form of S encodes which solutions are admissible to the inverse problem, e.g. enforcing positivity.

In a Bayesian approach the preconditioning discussed in the beginning of this section is naturally incorporated not via a modification of the Kernel but via the specification of an appropriate default model. Actually dividing out a fitted Ansatz $p(x)$ from the Kernel and using a constant default model is equivalent to working with the original Kernel and simply assigning $m(x) = p(x)$. In other words, in the Bayesian approach we select the most appropriate default model from a best fit of Eq. (3.2) to the Ioffe time data and estimate the dependence of the end result $q_{\text{Bayes}}(x)$ on the choice of m by repeating the reconstruction with default models arising from a slight variation of the best fit parameters in Eq. (3.2).

Based on different sets of prior information, different regularization functionals have been derived in the literature. Note that due to Bayes theorem, in the "Bayesian continuum limit" i.e. the combined limit of the number of supplied datapoints becoming infinite and the uncertainty on the data going to zero, all reconstructions will lead to the same solution. A good choice of regulator functional and the availability of valid prior information can help us approach this solution more closely even for data that are coarse and noisy. Indeed if an accurate approximate solution is already available we can use Bayesian methods that imprint this prior information strongly onto the end result (steep $S[q]$), while in cases where only unreliable prior information is available one wishes to use a method that let's the data "speak" as freely as possible without distorting them through the default model (weaker $S[q]$).

The reconstruction of PDF's from Ioffe-time PDF data benefits from the availability of good prior information in the form of the empirically obtained Eq.(1.4). This situation is quite different from other inverse problems, e.g. the reconstruction of hadronic spectral functions, where little to no relevant information about the spectral structures of interest is known. We therefore foresee that methods with a steep prior probability will fare better than those that are explicitly designed to minimize the impact of the default model.

The first example is the well known Maximum Entropy Method [48, 49], which features the Shannon-Jaynes entropy as regulator

$$S_{SJ}[q, m] = \sum_n \Delta x_n \left(q_n - m_n - q_n \log\left(\frac{q_n}{m_n}\right) \right). \quad (3.28)$$

Its functional form has been derived using arguments from two-dimensional image reconstruction and it is designed to introduce as least as possible additional correlations into the end result, beyond what is encoded in the input data. The standard MEM implementation restricts the space of solutions to those close to the default model, meaning that for a small number of available input points its result will lie close to the function m . While it has been shown [50] that this implementation of the MEM in general fails to recover the global extremum of the posterior, for a default model that is close to the correct result, the outcome is often satisfactorily accurate.

We contrast the MEM to another Bayesian approach simply named Bayesian reconstruction (BR) [51], which has been derived by requiring positive definiteness of the resulting q , smoothness of q , where the data \mathcal{Q} do not provide constraints on its form, as well as independence of the resulting q on the units used to express \mathcal{Q} . The resulting regulator reads

$$S_{BR}[q, m] = \sum_n \Delta x_n \left(1 - \frac{q_n}{m_n} + \log\left(\frac{q_n}{m_n}\right) \right). \quad (3.29)$$

It has been shown that this regulator also leads to a unique extremum of the posterior and since in contrast to the MEM no flat direction appears in S_{BR} , the search space does not need to be restricted. The BR method had been developed in particular for inverse problems, where only scarce or unreliable prior information is available, imprinting the form of m as weakly as possible onto the end result. We expect thus that in the present case with good prior information available, this method will produce less accurate results than the MEM, as it does not leverage the information in m as strongly.

Note that in the definition of $P[q|I]$ we introduced a further parameter α , a so called hyperparameter, which weighs the influence of simulation data and prior information. It has to be taken care of self-consistently. In the MEM it is selected, such that the evidence has an extremum. In the BR method, we marginalize the parameter α a priori, i.e. we integrate the posterior with respect to the hyperparameter, assuming complete ignorance of its values $P[\alpha] = 1$.

Up to this point we have only considered regulators for the reconstruction of positive definite functions. It is fathomable though that one has to deal with PDF reconstructions where positive definiteness does not hold, such as for spin dependent PDFs. In that case the regulator (3.29) is not applicable. A choice of regulator often employed in the literature in such a situation is the quadratic one, which corresponds to a modified form of Tikhonov regularization

$$S_{QDR}[q, m] = \sum_n \Delta x_n \left(q_n - m_n \right)^2. \quad (3.30)$$

It is a comparatively strong regulator and imprints the form of the default model significantly onto the end result. Marginalizing the hyperparameter α with a flat prior probability $P[\alpha] = 1$ is hence not possible in this case. Indeed integrating over the normalized prior probability with respect to α leads to a delta function that fixes the end result to $\rho = m$ everywhere. Instead one uses the "historic MEM" approach [48, 49], where α is chosen such that for the corresponding reconstruction q^α the likelihood takes on the value $L[q^\alpha] = N_\nu$. As in the following we are able to utilize default models, which are already close to the correct result, the use of this regulator is justified.

Again we can compare this approach to one specifically designed to keep the influence of the default model to a minimum, i.e. to imprint its functional form as weakly as possible on the end result. To this end we resort to an extension of the BR method (3.29) to non-positive functions. The corresponding

$$S_{BRg}[q, m] = \sum_n \Delta x_n \left(-\frac{|q_n - m_n|}{h_n} + \log\left(\frac{|q_n - m_n|}{h_n} + 1\right) \right), \quad (3.31)$$

keeps the advantageous properties of the original BR prior, e.g. smoothness and scale invariance at the price of having to introduce another default model related function h . It may be thought of as the confidence we have in our default model. Note that for $h = m$ one obtains a regulator, which for $q > 0$ and $q > m$ takes on the form of the standard BR prior and mirrors it for $q < m$. In order to obtain the full uncertainty budget of the end result obtained with the generalized BR prior, one has to vary not only the default model but also the function h . Eq. (3.31) also admits for an a priori integrating out of α .

Once L , S and m have been provided, the most probable PDF q , given simulation data and prior information is obtained by numerically finding the extremum of the posterior

$$\left. \frac{\delta P[q|\mathcal{Q}, I]}{\delta q} \right|_{q=q_{\text{Bayes}}} = 0. \quad (3.32)$$

It has been proven that if the regulator is strictly concave, as is the case for all the regulators discussed above, there only exists a single unique extremum in the space of functions q on a discrete ν interval.

Note that as long as positive definiteness is imposed on the end result, the space of admissible solutions is significantly reduced. On the other hand, regulators admitting also q functions with negative contributions have to distinguish between a multitude of oscillatory functions, which if integrated over, resemble a monotonous function to high precision. We will observe the emergence of ringing artifacts for the quadratic and generalized BR prior below.

Other regularizations may be reformulated in the spirit of the Bayesian strategy. To give meaning to the inversion task they introduce in addition to the likelihood another functional which encodes prior information or prior expectations on the end result. In the standard Tikhonov approach e.g. the condition is that the values of q shall have the smallest possible magnitude. While they do not rely on an explicitly formulated default model, they contain the prior information implicitly in the form of their regularization functional. Otherwise they could not provide a unique regularized answer to the task. The

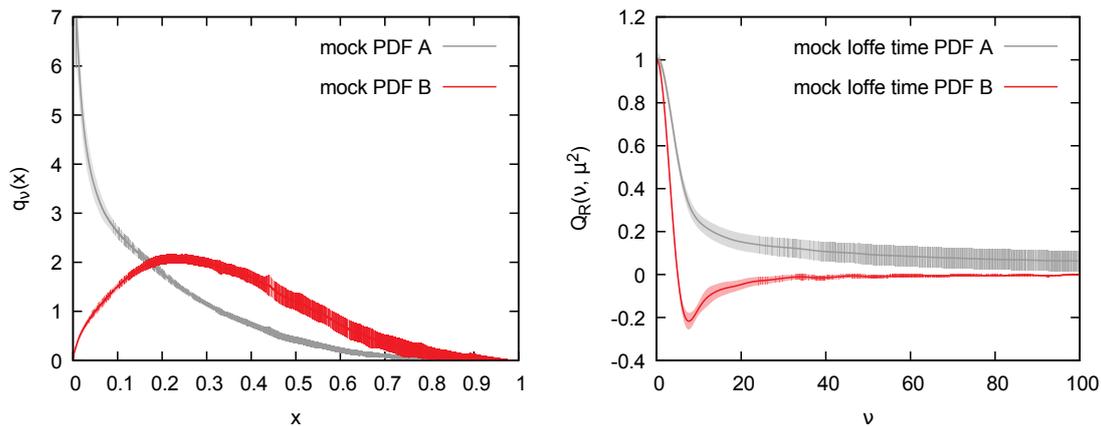


Figure 4. The two mock datasets used for testing the methods. The PDF on the left. The Ioffe time PDF on the right. Scenario A (red band) will test the case of a PDF which diverges at small x and scenario B (blue band) will test the case of a PDF which converges at small x .

benefit of the Bayesian method is that all prior information is made explicit, so that the dependence of the end result on the regularization may be thoroughly and fully tested.

4 Mock Data Tests

In this section, the methods constructed above are tested using data for the Ioffe time PDFs constructed from the phenomenological PDFs. Using the software package LHAPDF [52] and the dataset NNPDF31_nn1o_as_0118 from the NNPDF collaboration [43], we Fourier transform numerically the iso-vector quark PDFs and obtain data for the $Q_R(\nu)$ at $Q^2 = 2 \text{ GeV}^2$. This data set will be called scenario A. A second dataset, scenario B, is created by multiplying the NNPDF data by Nx , where N is a normalization to fix the valence quark sum rule. This modification forces the PDF to vanish as x approaches zero. This scenario is reminiscent of the quasi-PDF case which is finite at low y before the matching procedure is applied. These data sets contain a set of replicas that can be used to obtain errors on the PDFs. The standard deviation over the ensemble of replicas represent the 68% confidence interval. The PDFs and Ioffe time PDFs for these two cases are plotted in Fig. 4. For the various tests, the Ioffe time PDF will be sampled in the ranges mentioned above to study the dependence on the maximum available Ioffe time.

4.1 Backus-Gilbert

The first advanced reconstruction technique to be tested is the Backus-Gilbert method. All integrals are computed numerically in double precision to tolerance 1.0×10^{-16} . In order to regularize the matrix \mathbf{M} an SVD cut-off of $\rho = 1 \times 10^{-8}$ is used as discussed in previous sections ¹. With this set up, a resolution function width, $\sigma(\bar{x})$, of $\mathcal{O}(10^{-2})$ is obtained for most of $x < 0.75$ gradually increasing after that reaching $\mathcal{O}(10^{-1})$ for $x = 1.0$. With the ensemble of replicas, one can also compute the covariance matrix of the data and therefore

¹We have also tested using the Moore-Penrose pseudo inverse instead and got similar results.

use the covariance matrix regularization for the Backus-Gilbert method. We have done that as well as employ the Tikhonov regularization and obtained results that are similar with the results we discuss here where the SVD cutoff method was used to regularize the matrix inverse. In order to obtain the statistical error-band and the mean for the reconstructed curve we take the standard deviation and the mean of the ensemble of reconstructed curves. The resulting statistical error-band represents a 68% confidence level.

After performing several numerical experiments on the Backus-Gilbert method, we have concluded that the use of a preconditioning function $p(x)$ is essential in obtaining an accurate determination of the PDF. The reconstructions of the PDF from data with $\nu_{max} = 10$ and no preconditioning function are shown in Fig. 5. The reconstructed PDF is consistent with the mock PDF for $x > 0.3$, but deviates from it in the low x region. This result is expected, because the low x region is dominated by the Ioffe time distribution at large Ioffe times which are not present in the input data. Also in Fig. 5, the fidelity of the reconstruction is tested by taking the Fourier transform of the reconstructed PDF and comparing to the data used to generate it which shows agreement across the range of Ioffe time.

As argued earlier, preconditioning is essential in improving the Backus-Gilbert extraction of the PDF. In order to test its effectiveness, the preconditioning function defined in Eq. (1.4) is used with a range of exponents a and b that allows for all the integrals that define \mathbf{M} and \mathbf{u} to be convergent. For scenario A we achieved the best reproduction with $a = -0.35$ and $b = 2$ while for scenario B with $a = 0.3$ and $b = 2$. In order to get a handle of possible systematics due to preconditioning we varied the exponents a and b in the intervals $[-0.25, -0.4]$ and $[1, 4]$ respectively for scenario A. Correspondingly for scenario B we varied the exponents in the intervals $[0.2, 0.35]$ and $[1.5, 3]$.

As one can see, for both scenarios, a wide range in x is well reproduced and deviations from the original data begins to appear for $x < 0.2$. Perhaps, these deviations are expected due to the rather small maximum value of ν used in this example, $\nu_{max} = 10$. However, this aggressive cutoff for ν is realistic for lattice calculations. In addition, the number of selected points, which is 12, is also a number that is plausible for present lattice QCD calculations. From these results we can conclude that the fidelity of the reconstruction can be improved significantly by choosing appropriate preconditioning functions. The best reconstruction is obtained using a preconditioning function that is chosen to roughly fit the data. In this case the resulting reconstruction is indistinguishable from the original data for nearly all the range of x with the largest deviations occurring for $x < 0.1$. Therefore, we conclude that the Backus-Gilbert reconstruction with an appropriate preconditioning function is well suited for the reconstruction of the PDFs from the limited data for the Ioffe time PDFs that are provided by present day lattice QCD calculations.

Finally, we explore the dependence of the reconstruction on the maximum value of ν . Experiments are performed with $\nu_{max} = 20$ and $\nu_{max} = 100$. In all cases the data points were selected to be uniformly spread in the interval with a separation $\delta\nu = 10/11$. The results are shown in Fig. 7. The best result that is in agreement with the input PDF within errors across the whole range of x is obtained from the experiment with $\nu_{max} = 100$. However, this is an unrealistically large range in ν unlikely to be accessible in contemporary

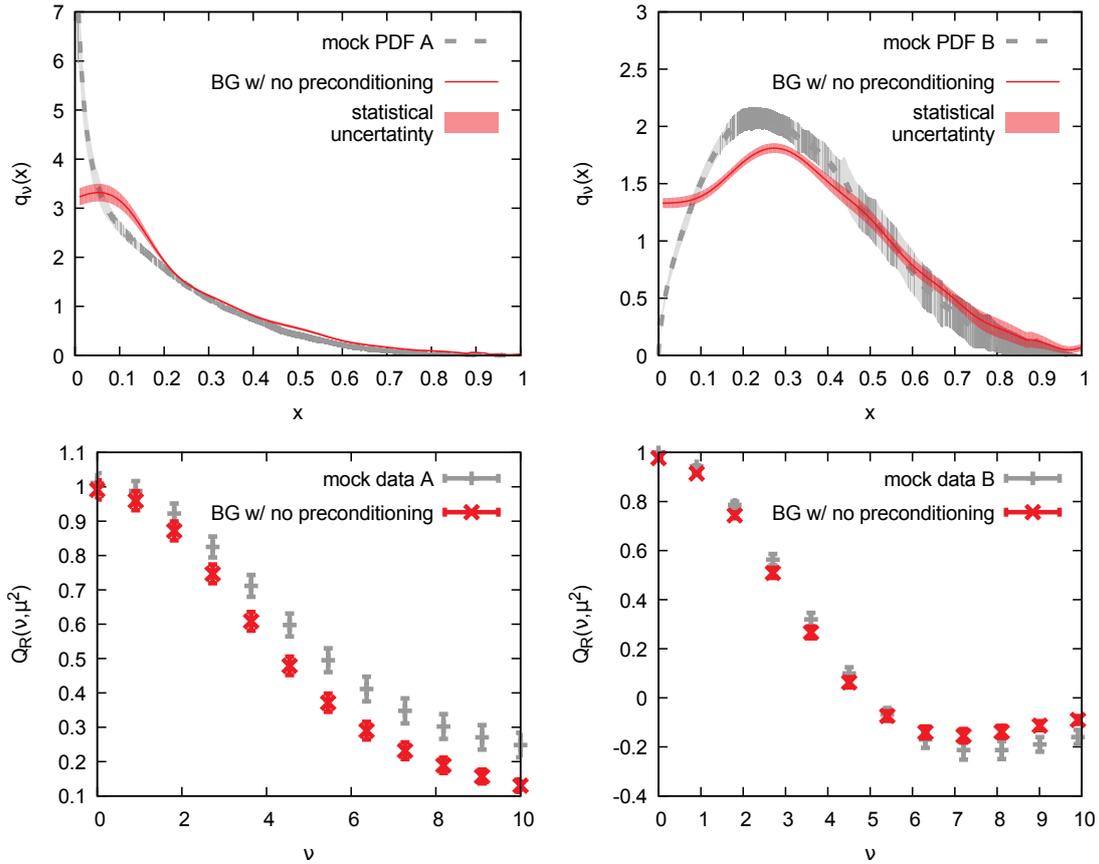


Figure 5. x -space PDF's reconstructed using the Backus-Gilbert (BG) method from $N_\nu = 12$ Ioffe-time data points on the interval $\nu = [0, 10]$ (top) as well as the input data (gray crosses) compared to the data arising from the reconstructed PDF (red crosses) in the bottom panels. The plots in the left column show the results for mock data based on a phenomenological PDF, while the right column from the modified scenario where the PDF vanishes at the origin. In both scenarios no preconditioning has been employed.

lattice QCD calculations. As ν_{max} is lowered then the reconstruction deteriorates. However, even for $\nu_{max} = 20$ the results seem to agree with the input with the exception of one small region $x = [0.05, 0.1]$ where we observe a one σ deviation. Overall, the $\nu_{max} = 20$ result is only very slightly better than the case of $\nu_{max} = 10$ which was presented in Fig. 6. Therefore, we conclude that extending the range of ν slightly beyond $\nu = 10$ will have only small improvement in the determination of the PDF. Our conclusion from this analysis is that 12 points for $\nu \in [0, 10]$ seems to be adequate to obtain a good reconstruction of the underlying PDF. This finding is rather encouraging for present day calculations as it is possible to cover the range of $\nu \in [0, 10]$ with about 12 points in lattice QCD calculations.

4.2 Neural network reconstruction

In this section, the neural network method described above is tested to reconstruct the mock PDFs. The data used for this reconstruction is the mock scenario A and B in the

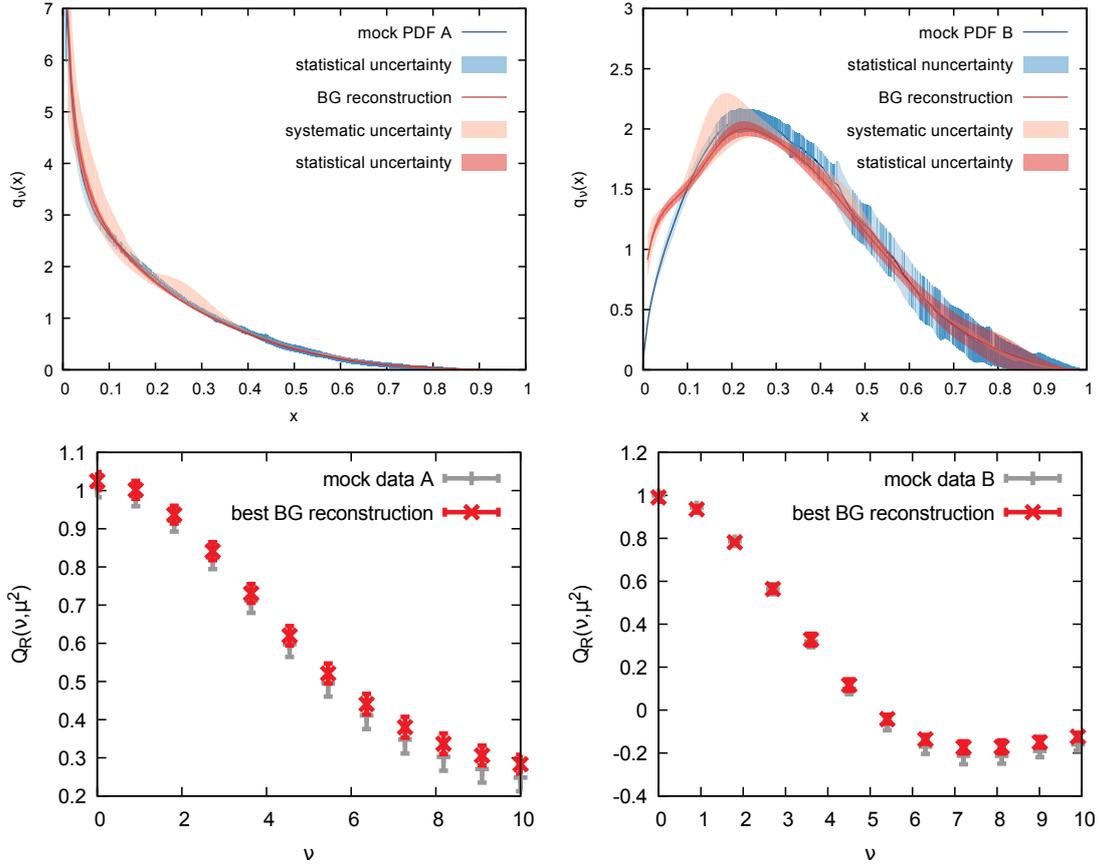


Figure 6. x -space PDF's reconstructed using the Backus-Gilbert (BG) method from $N_\nu = 12$ Ioffe-time data points on the interval $\nu = [0, 10]$ (top) as well as the input data (gray crosses) compared to the data arising from the reconstructed PDF (red crosses) in the bottom panels. The plots in the left column show the results for mock data based on a phenomenological PDF, while the right column from the modified scenario where the PDF vanishes at the origin. The reconstruction was performed with preconditioning exponents $a = -0.35$ and $b = 2$ for scenario A and $a = 0.3$ and $b = 2$ for scenario B.

smallest range of ν , $\nu \in [0, 10]$, discretized into $N_\nu = 12$ points. For this study, we chose the hyperbolic tangent as activation function for all the nodes in the hidden layers. The activation function for the final layer is linear with the threshold value fixed to zero. For a single hidden layer of size N_2 , the neural network parametrizes the PDF as

$$q(x) = \sum_{i=1}^{N_2} w_{1,i}^{(3)} \tanh(w_{i,1}^{(2)} x + \theta_i^{(2)}), \quad (4.1)$$

while for two hidden layers of sizes N_2 and N_3

$$q(x) = \sum_{j=1}^{N_3} w_{1,j}^{(4)} \tanh \left(\sum_{i=1}^{N_2} w_{j,i}^{(3)} \tanh(w_{i,1}^{(2)} x + \theta_i^{(2)}) + \theta_j^{(3)} \right). \quad (4.2)$$

The neural network behavior is governed by the collective interactions of all the weights, thresholds, and connections, not on any individual neurons. Slight modifications of the

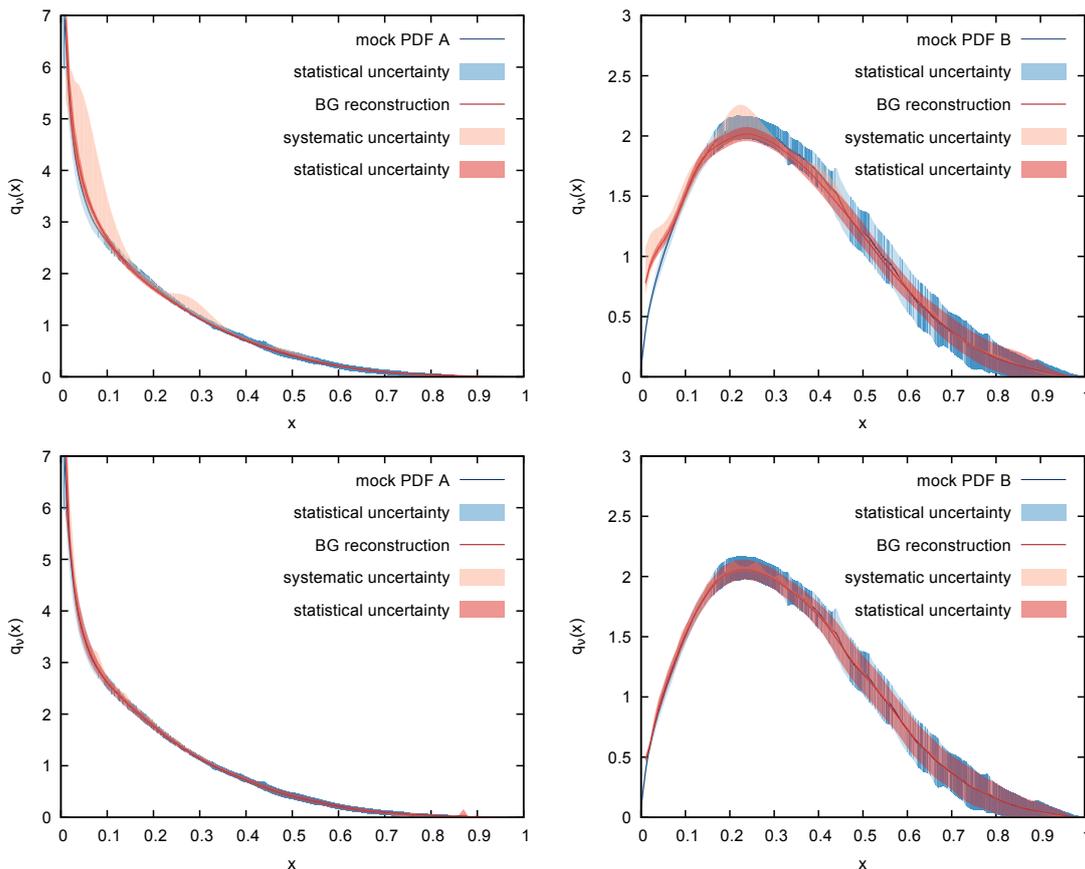


Figure 7. x -space PDF's reconstructed using the Backus-Gilbert (BG) method from $N_\nu = 23$ Ioffe-time data points on the interval $\nu = [0, 20]$ (top) and from $N_\nu = 112$ Ioffe-time data points on the interval $\nu = [0, 100]$ (bottom). The plots in the left column show the results for mock data based on a phenomenological PDF, while the right column from the modified scenario where the PDF vanishes at the origin. The reconstruction was performed with preconditioning exponents $a = -0.35$ and $b = 2$ for scenario A and $a = 0.3$ and $b = 2$ for scenario B. It is noteworthy that by having lattice data up to $\nu_{\max} = 20$ one can achieve a very good reconstruction with the BG method.

geometry, such as adding or removing a neuron, will have little change in the final result. To test this point, 1-3-1, 1-4-1, and 1-2-2-1 geometries are tried, all resulting in similar PDFs. No preconditioning function will be used in this test. A preconditioning function of the form in Eq. (1.4) can be set to prefer certain features, such as the PDF vanishing in the limit $x \rightarrow 1$ by using $b > 0$ or a divergence as $x \rightarrow 0$ by using $a < 0$. Remarkably the neural network was able to reproduce both these high and low x features without the need of any constraint or prior information.

The training procedure described above has a number of tunable parameters, but the end result does not depend very strongly on any particular parameter. For this study, the initial search width, σ_0 , is set to 5, the mutation width, σ , is set to 1, and the mutation size, η , is set to 0.05. From most starting points in any given generation, the minimization

procedure sends the χ^2 to very small values. In case this does not occur, the minimization is restarted. Three checks on χ^2 are performed and if the value is too large then the mutation and minimization steps are performed again. First if the χ^2 is greater than 10^{-2} , then the mutation and minimization steps are repeated with the same mutation width and size. Second if the χ^2 is greater than 10^{-1} , then the mutation and minimization steps are repeated with the same mutation width and a size of 0.025. Third if the χ^2 is still greater than 10^{-1} , then the mutation and minimization steps are repeated with the same mutation width and a size of 0.0125. The above minimization strategy was used in two stages. First we used the global average of all available data and began with 45 initial neural networks and removed 2 every generation for 10 generations until 25 neural network replicas remained, $\{q_i(x)\}$. The response of the surviving population is shown in Fig. 8 for each of the geometries studied. The standard deviation of the response of these neural networks is used to estimate the systematic error in this method of regression.

$$\sigma_{sys} = \text{StdDev}[q_i(x)] \quad ; \quad i = 1 \dots N_{rep} \quad (4.3)$$

In the second stage, we perform the minimization on each of the $N = 100$ mock Ioffe time PDF replicas² in order to estimate the statistical fluctuations of the resulting function. Starting with each of the networks in the final population, each PDF replica is used as the training data for another minimization to create the PDF replica dependent networks of the final population, $\{q_i^{(b)}(x)\}$. The variance of the response averaged over the members of the final population across the different PDF replica dependent networks is used to estimate the statistical error,

$$\sigma_{stat} = \text{StdDev}\left[\frac{1}{N_{rep}} \sum_{i=1}^{N_{rep}} q_i^{(b)}(x)\right] \quad ; \quad b = 1 \dots N. \quad (4.4)$$

The average response of the PDF replica dependent networks and the combined systematic and statistical errors are shown in Fig. 9. As anticipated, the errors tend to grow at small- x . Since the input data is over a truncated region of ν , the information of the low- x behavior is lost. The neural networks are able to precisely reconstruct the large- x region while in the low- x region the flexible parametrization allows for a wide range of functions that can reproduce the data. One can also see that the various network geometries all give comparable results for the reconstructed PDF.

Next we consider how a preconditioning function affects the reconstructions. Using the function defined in Eq. (1.4), the preconditioning function for scenario A with $a = -0.25$ and $b = 2$ and for scenario B with $a = 0.3$ and $b = 2$. Both of these cases were tested on the 1-3-1 geometry. The reconstructed PDFs, shown in Fig. 11, are not significantly different for the majority of large x region. For scenario A, the low x behavior is slight improved compared to the case without preconditioning. This behavior is perhaps to be expected, because the preconditioned neural network is parameterizing a much slower varying function in that region compared to the divergent PDF.

²These are the replicas the NNPDF collaboration provides in their data set. In a realistic lattice calculations these replicas would be a set of jackknife or bootstrap samples of the matrix elements.

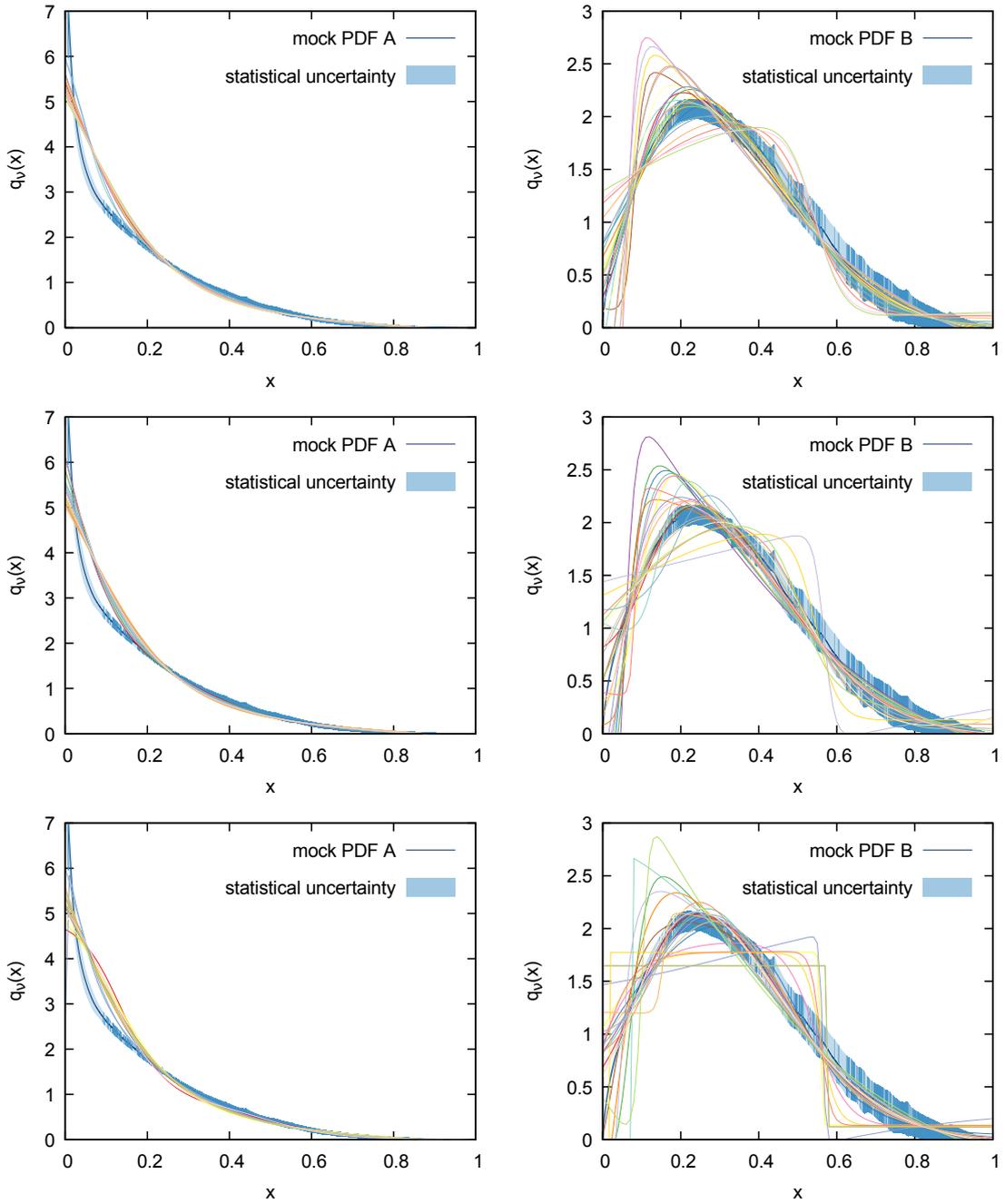


Figure 8. The genetically trained neural nets. The blue band is the original data. The curves are the responses for the final population of genetically trained neural networks. The left column is with NNPDF data. The right column is with modified data. The first row has a network geometry of 1-3-1. The second row has a geometry of 1-4-1. The third row has a geometry of 1-2-2-1.

Our conclusion is that for ranges of ν which are realistic in modern lattice calculations, a neural network is capable of reconstructing the PDF for a wide range of x without the need of preconditioning, though the neural networks can benefit from it. This study of

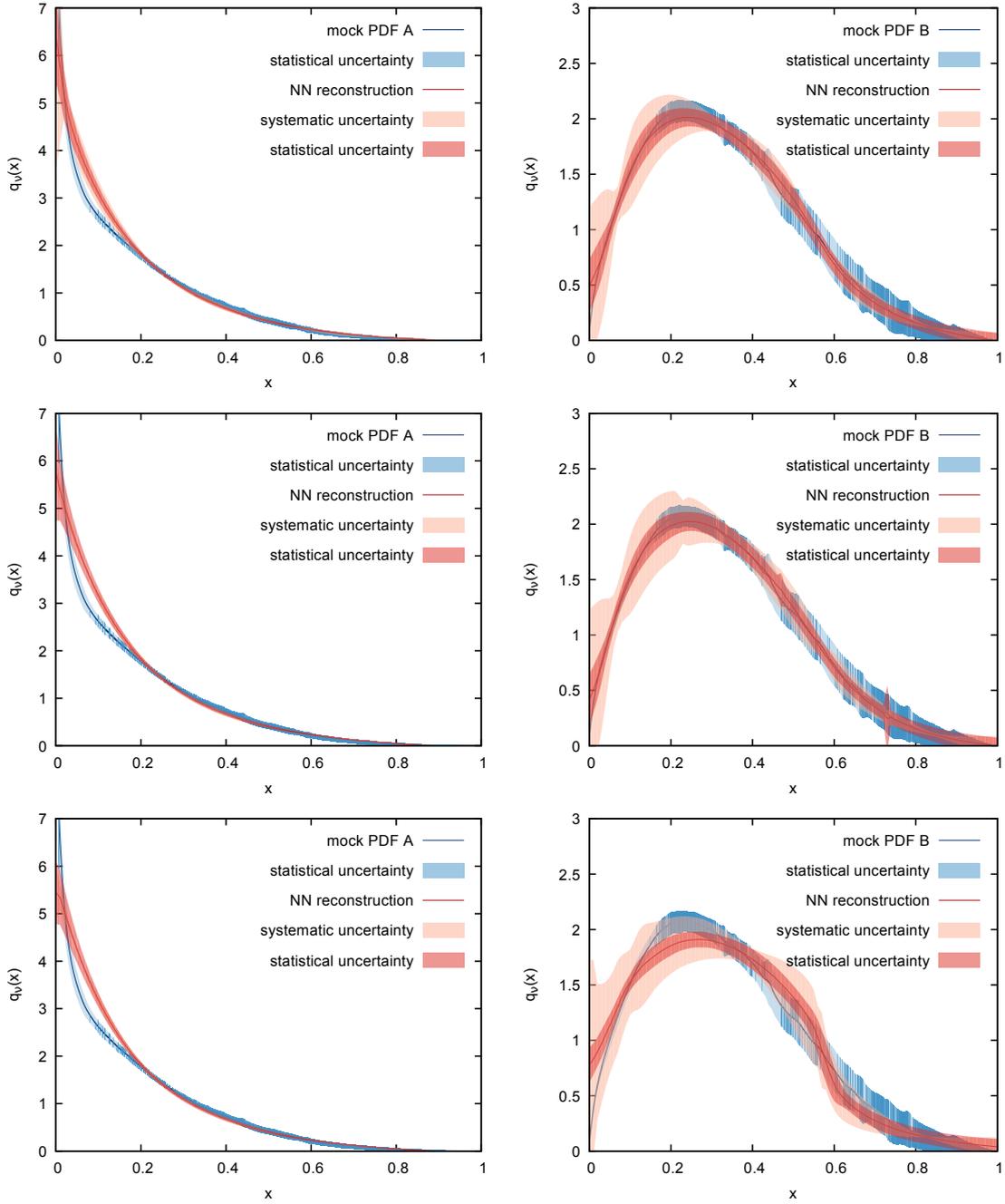


Figure 9. The genetically trained neural nets. The blue band is the original data. The red band is the reconstructed PDF with statistical and systematic errors. The left column is with NNPDF data. The right column is with modified data. The first row has a network geometry of 1-3-1. The second row has a geometry of 1-4-1. The third row has a geometry of 1-2-2-1.

neural network parametrization of the PDF is by no means complete. There are a number of choices that were made which may not be optimal. It is clear that a more dedicated study of this approach is required in order to understand its full potential. However, we

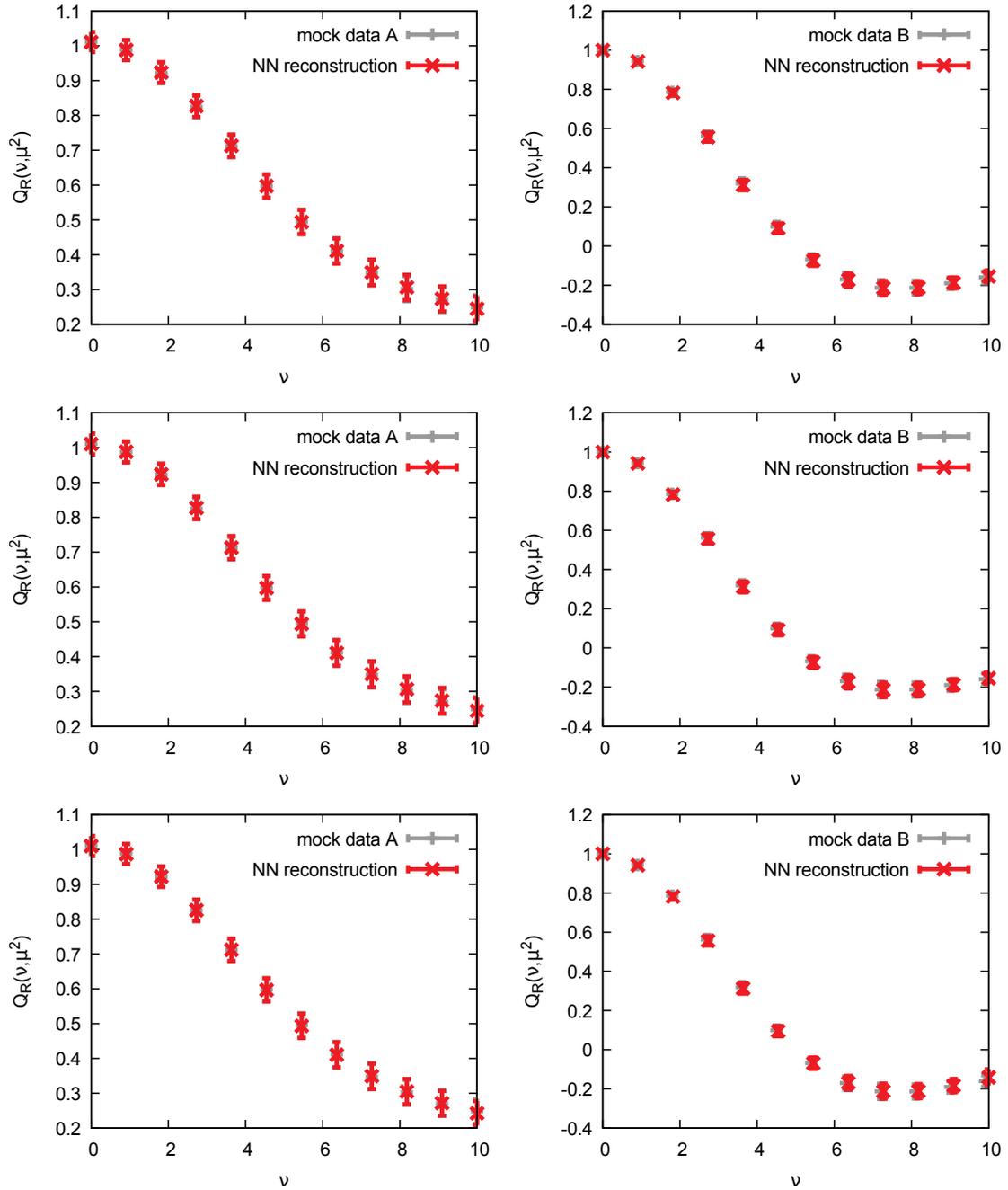


Figure 10. The genetically trained neural nets. The grey points are the original data while the red points are the ones reconstructed by the NN. The left column is with NNPDF data. The right column is with modified data. The first row has a network geometry of 1-3-1. The second row has a geometry of 1-4-1. The third row has a geometry of 1-2-2-1.

find the results obtained here very encouraging and thus we plan to further investigate this approach in a future publication.

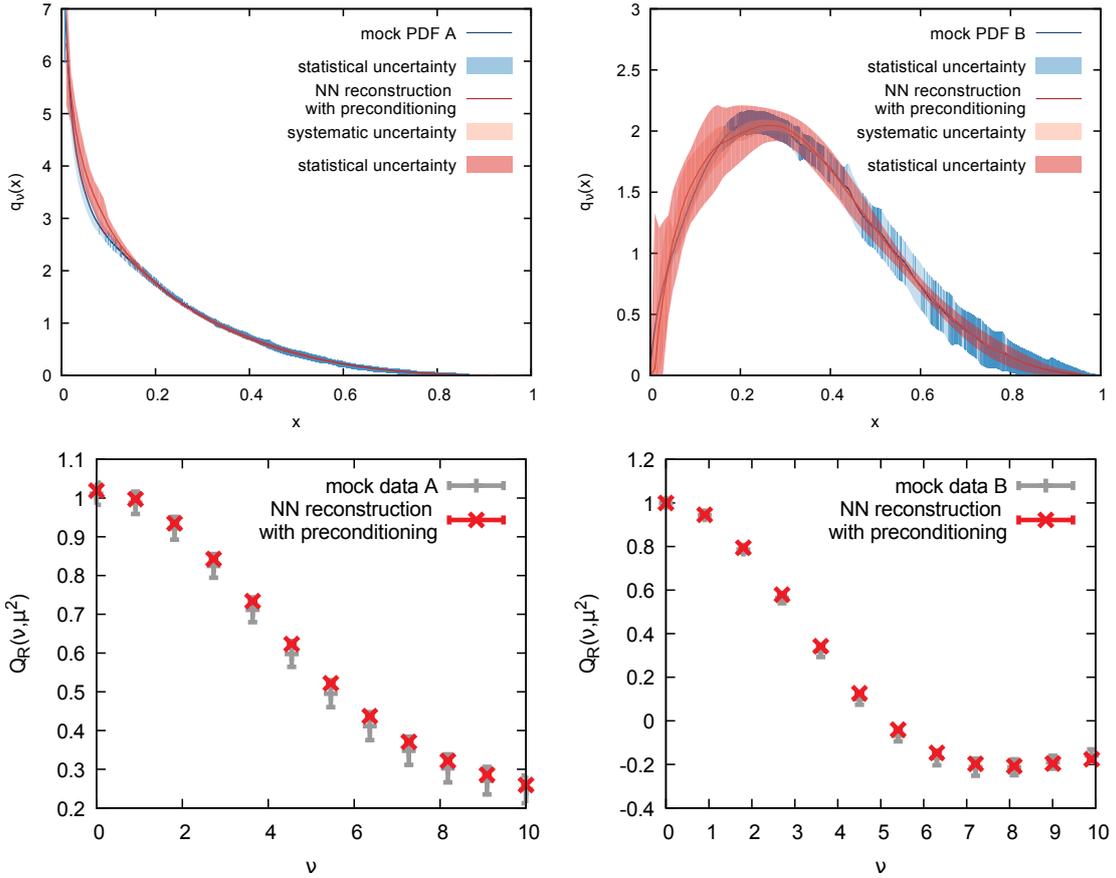


Figure 11. The PDF reconstructions from preconditioned Neural Networks (top). The grey points are the original data while the red points are the ones reconstructed by the NN (bottom). The left is scenario A and the right is scenario B.

4.3 Bayesian Analysis

Finally, mock data tests of the Bayesian strategy outlined above are carried out to determine the feasibility of extracting the x -space PDF $q(x)$ from the Ioffe-time data $Q(v)$ in realistic settings for applications to lattice generated data. The x interval is discretized in $N_x = 2000$ steps to have fine enough resolution of the possibly highly oscillating cosine. To remain as close as possible to the situation encountered with real data, we first carry out a simple fit on the noisy mock Ioffe-time data, based on the simple Ansatz Eq. (1.4) via the expression in Eq. (3.2).

In both scenario A and B we find that the simple functional form (3.2) allows us to capture the overall features of the Ioffe-time data very well, as shown in Fig. 12. The best fit result for $q(x)$ is shown in Fig. 13 as a red curve and compared to the actual mock PDF, given as gray dashed line. At intermediate- x values, the best fit in both scenarios deviates from the correct $q(x)$ but close to $x = 0$ provides a rather good description. This is exactly what we have in mind: the fit will provide us with prior information about the non-analytic behavior of the function, which we can use as default model in the Bayesian reconstruction.

The Bayesian approach will then imprint the information encoded in the simulation data as deviations from the default model onto the end result.

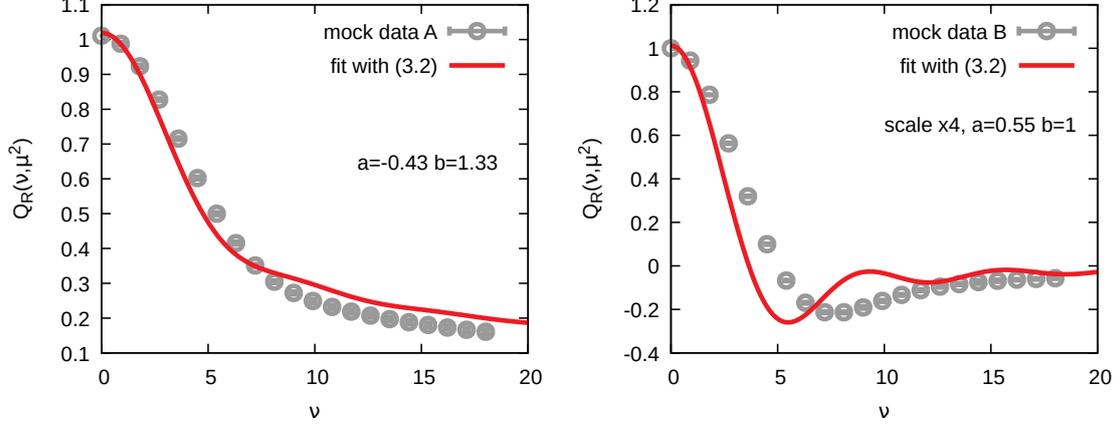


Figure 12. Mock data $Q(\nu)$ in the interval I_3 from (left) realistic PDF data [mock scenario A] and (right) a modified scenario with the PDF vanishing at the origin [mock scenario B]. In red we show the corresponding best fit of the noisy data with Eq.(3.2).

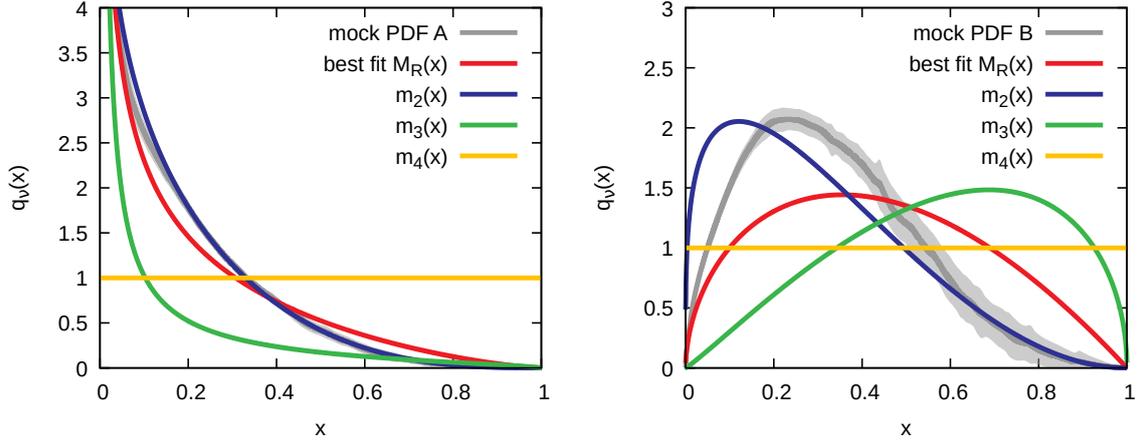


Figure 13. Best fit PDF (red solid line) from (left) realistic PDF data [mock scenario A] and (right) a modified scenario with the PDF vanishing at the origin [mock scenario B]. The actual mock PDF in the former cases is given as gray solid line. To determine the dependence of our results on the choice of default model, three further choices for m are plotted, two arising from varying the best fit parameters by factors of 2, one being the constant default model $m=1$.

Note that in order to use the best fit $q(x)$ as default model we need to modify its functional form at the boundaries of the x -interval in both cases. If there is e.g. a pole at the origin, then the value of m is not well defined there. On the other hand if the function vanishes exactly at the origin, then the prior assumption of positive definiteness in some of the Bayesian methods is not fulfilled. Therefore we set $m(0) = q_{\text{best fit}}(\Delta x)/2$. To further understand how strongly our end result depends on the choice of default model,

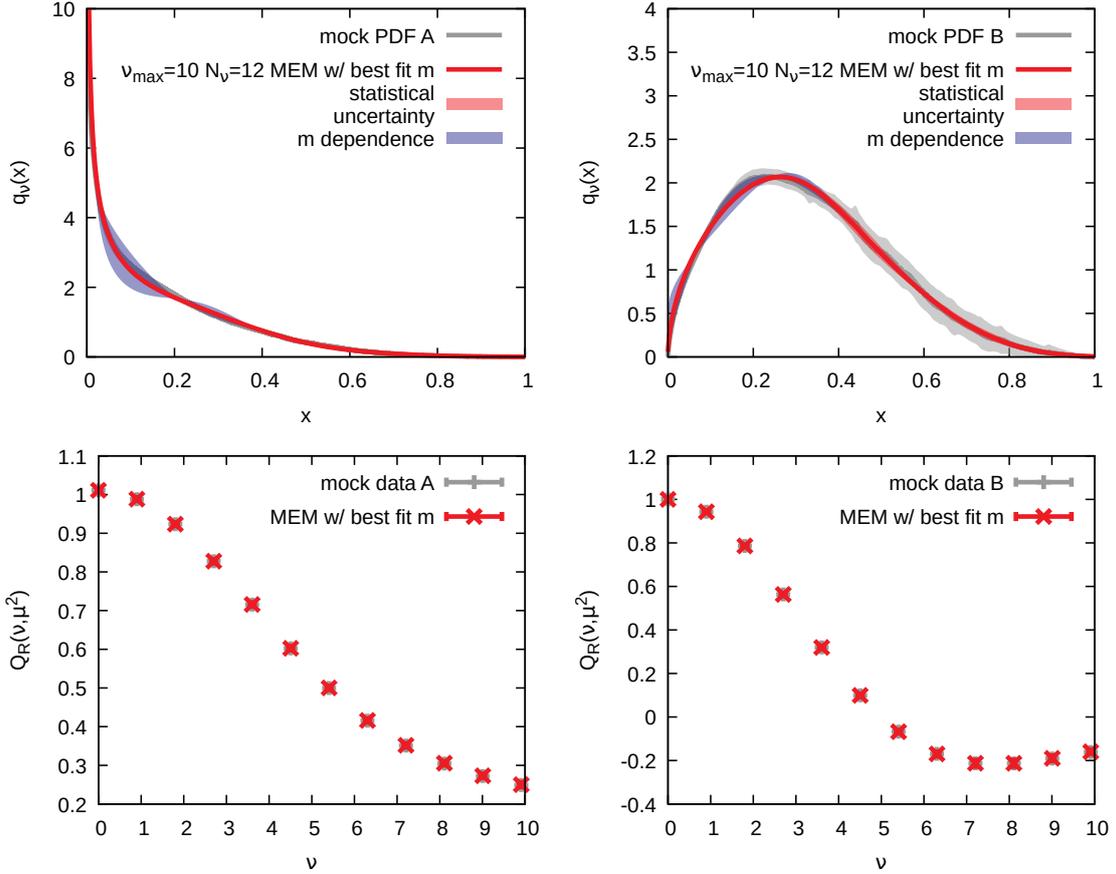


Figure 14. x -space PDF's reconstructed using the Maximum Entropy method (MEM) from $N_\nu = 12$ Ioffe-time data points on the interval $\nu = [0, 10]$ (top) as well as the input data (gray crosses) compared to the data arising from the reconstructed PDF (red crosses) in the bottom panels. The plots in the left column denote the results for mock data based on a phenomenological PDF, while the right column from the modified scenario where the PDF vanishes at the origin. (top) The original mock PDF is given as gray dashed line and the result based on the BR method with the best fit default model as solid red line. Statistical uncertainty is visualized by a red, systematic uncertainty from the choice of default model by a blue errorband respectively. (bottom) The red solid line corresponds to the relative deviation of the reconstruction from the input mock PDF, the red errorband encodes the combined uncertainty of the BR reconstruction, the gray errorband the uncertainty of the original mock PDF.

we also carry out reconstructions with other functions than the best fit one, two of which are obtained from varying the best fit parameters by factors of two and one function simply being the constant $m(x) = 1$.

We start with the most realistic and most challenging setting, where the input data ν is in the range I_1 , i.e. $\nu = [0, 10]$ discretized with $N_\nu = 12$ points, very similar to what is currently available in lattice QCD simulations. Due to the availability of a good approximate solution obtained from the fit with eq.(1.4), provided as default model, we first consider the Maximum Entropy Method, whose results are presented in Fig.14

The original mock $q(x)$ is shown on the top of Fig. 14 as a gray dashed line, and the MEM reconstruction based on the best fit default model in red. The red errorband denotes the statistical uncertainty, while the blue band arises from the default model dependence. In the bottom row panels the input Ioffe-time data (gray) as well as the data corresponding to the Bayesian reconstruction (red) are shown. We find that the reproduction for both scenario A and scenario B is excellent, i.e. the MEM is able to utilize the approximate solution provided by the default model to regularize the inverse problem and reproduce the input PDF within uncertainties. If one inspects the reconstruction in scenario B carefully around the maximum one finds that the MEM solution slightly underestimates the mean below $x = 0.25$ and overestimates it for $x > 0.25$. The deviation however is minute.

Let us compare how the BR method fares. In this method the information of the default model is more weakly incorporated. The corresponding results are shown in Fig.15 and yield, as expected, a less accurate reproduction of the input data. While for scenario A it reproduces the PDF within uncertainties, the reconstruction of scenario B shows undulations around the correct result. Close to the maximum at intermediate x the mean is found to be underestimated.

This result is not surprising but the appearance of differences in the results from different Bayesian regulators urges us to figure out how far we have to improve the input data for both results to provide a reproduction of the correct result within uncertainties. To this end we have also carried out reconstructions based on the larger interval $\nu = [0, 20]$ with the same number of $N_\nu = 12$ Ioffe-time data-points. The results of the MEM and BR method are shown in Fig.16 and reflect that the improvement in input data has significantly reduced the errors of the reconstruction so that now both methods lie within the uncertainty of the original input data. It is still the case that the MEM utilizes the provided prior information to a stronger extent than the BR method and produces slightly more accurate results in this setting. A further extension of the interval to $\nu = [0, 10]$ and increasing the number of available points to $N_\nu = 100$ confirms this trend, improving the accuracy of the reconstructions in particular around the maximum of the PDF in scenario B.

Let us turn our attention now to the Bayesian reconstruction methods, which do not require the positivity of the PDF. The first uses the quadratic prior combined with the "historic MEM" procedure to choose the hyperparameter α . Results for both the realistic PDF mock scenario A (left column) and the mock scenario B (right column) are given in Fig. 17. Since this method also imprints the default model information strongly on the end result, it is not surprising that the reconstruction of scenario A also works excellently here. For scenario B we find that while the overall shape has been reproduced in an acceptable manner, the small $x < 0.1$ region is overestimated, while close to the maximum of the PDF the solution underestimates the correct value. It is reassuring to see that with only $N_\nu = 12$, the presence of accurate prior information is able to efficiently suppress the many possible oscillatory solutions, which could also reproduce the provided input data.

Let us consider as a last item the reconstruction based on the generalized BR method with $h = m$, whose results are given in Fig. 18. To understand the performance of this method, we note that it has been designed for cases, where no reliable prior model is available and where thus the influence of m on the end result needs to be held as small as

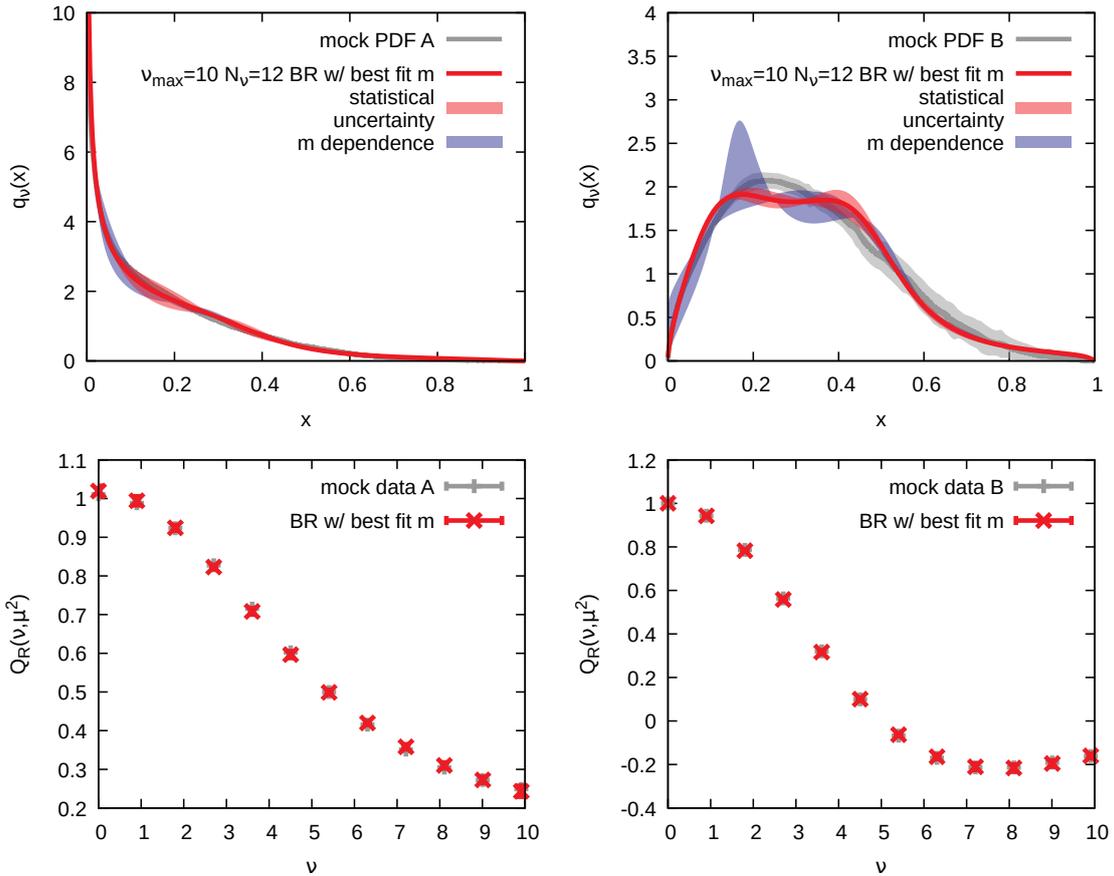


Figure 15. x -space PDF's reconstructed using the Bayesian reconstruction (BR) method from $N_\nu = 12$ Ioffe-time data points on the interval $\nu = [0, 10]$ (top) as well as the input data (gray crosses) compared to the data arising from the reconstructed PDF (red crosses) in the bottom panels.

possible. We find that for both scenarios A and B the reconstructions are, as expected, less accurate than for the quadratic prior. In particular we see that excursions into negative values occur at some x , which in turn lead to overestimation of values at different x . Note that the obtained reconstructions however lead to Ioffe-time data that reproduce the input data very accurately.

It is possible to improve the accuracy of the generalized BR method by changing the value of the default model confidence function h , which essentially makes the regulator steeper and steeper, i.e. imprinting the default model more and more strongly. Our intention of showing the generalized BR method with a weak regulator is to emphasize the role accurate prior information has in approaching the correct reconstruction result.

In summary, we find that due to the availability of accurate prior information on the true shape of the PDF via fits with eq.(1.4), Bayesian methods that are designed to imprint prior information strongly on the end result outperform those that minimize the influence of the default model. I.e. in case of a rather small interval $\nu = [0, 10]$ and number of

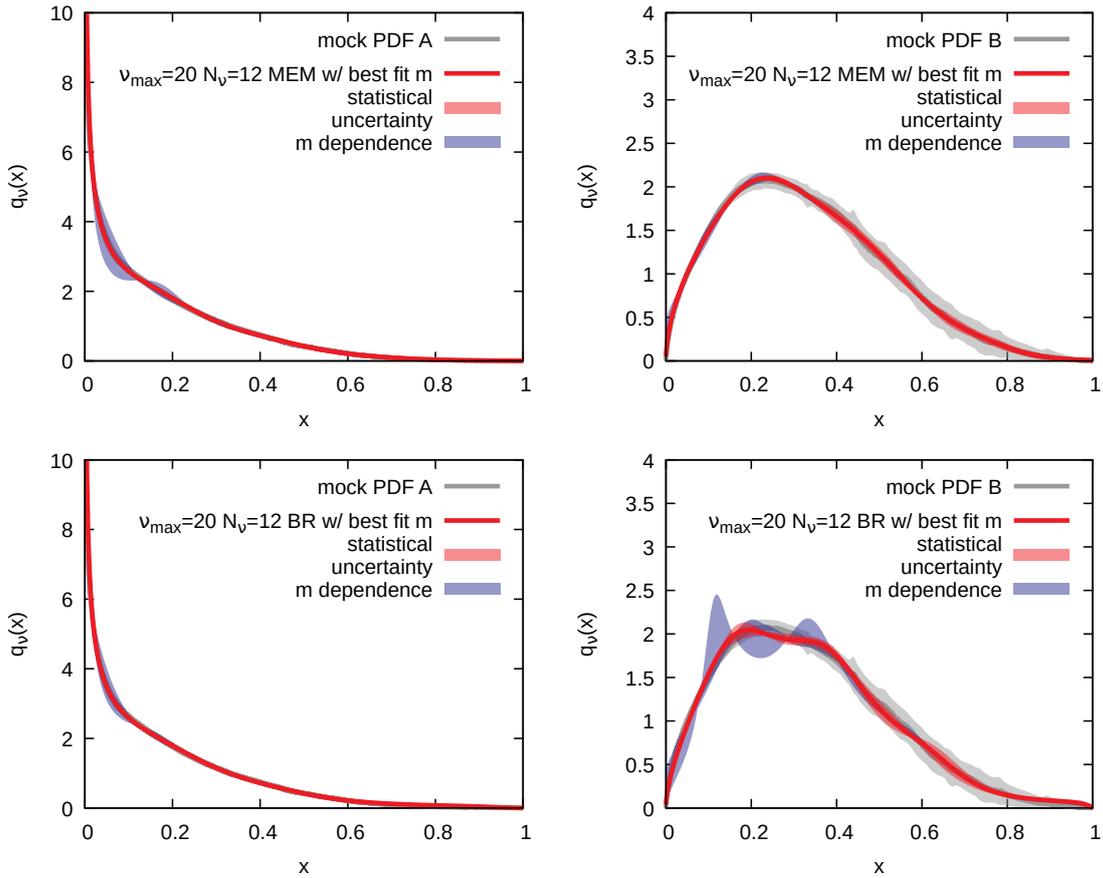


Figure 16. x -space PDF's reconstructed using the Maximum Entropy Method (top) as well as the Bayesian reconstruction (BR) method (bottom) from $N_\nu = 12$ Ioffe-time data points on the interval $\nu = [0, 20]$.

data points $N_\nu = 12$, the Maximum Entropy Method provides the most accurate results for both scenario A and B (Fig. 14). The BR method is competitive in scenario A but shows a larger deviation from the true result for scenario B (Fig.15). Increasing the extent of the ν interval and the number of data points consistently shows that both methods approach the correct solution in the "Bayesian continuum limit".

When comparing methods without the requirement of positivity we find that the quadratic prior still fares excellently with scenario A but shows some deviations, in particular at small x values in scenario B. As expected and by construction, the generalized BR method with the choice $h = m$ shows the weakest dependence on the used default model but at the same time also provides the least accurate reconstructions.

The success of the MEM relies on the availability of very good prior information. On the other hand the BR method is constructed to imprint the default model on the end result in a weaker fashion than the MEM (lower curvature of the prior functional). Hence we plan to deploy it together with the MEM when investigating actual lattice QCD data in the future. The reason is that as long as the asymptotics of the PDF at $x = 0$ and $x = 1$ are

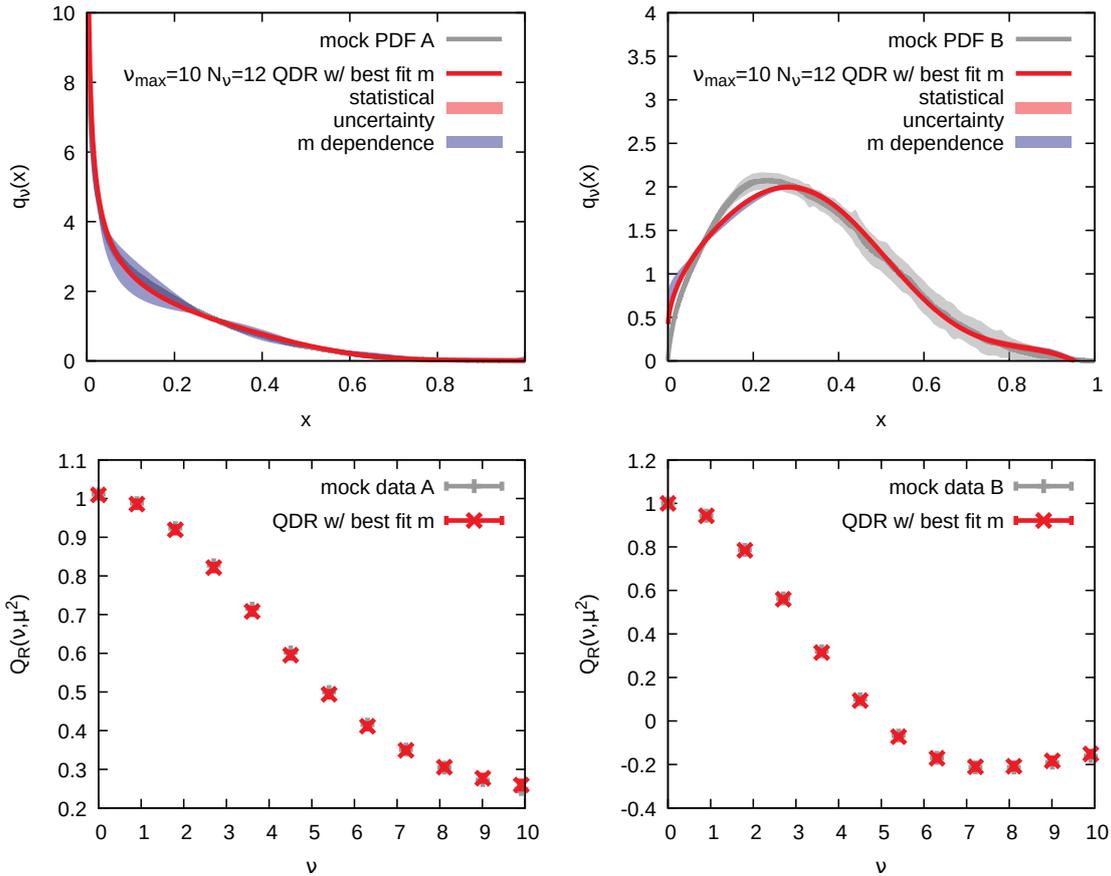


Figure 17. x -space PDF's reconstructed using a quadratic prior Bayesian (QDR) method from $N_\nu = 12$ Ioffe-time data points on the interval $\nu = [0, 10]$ (top) as well as the input data (gray crosses) compared to the data arising from the reconstructed PDF (red crosses) in the bottom panels.

adequately provided, the BR method is expected to be able to maintain its reconstruction quality in case that in a genuine lattice data based reconstruction the quality of the default model at intermediate x values will be worse than in the above tests.

If on the other hand a non-positive q function is concerned and accurate prior information is available, then we can utilize the quadratic prior even in cases that only small number of data points are available. The generalized BR method to become competitive requires a more careful treatment of the default model confidence function h in this context.

4.4 Restricted χ^2 sampling

In order to clarify the role of different prior information on the success of the reconstruction, we perform a further numerical experiment on our mock data sets. Here we consider positivity as the only property of the PDF known a priori. The ill-posedness of the inversion is related to the many local extrema in the likelihood. If we work with a constant prior probability, the resulting posterior probability will just be a flat distribution, unable to provide meaningful insight. This fact has been explicitly tested and confirmed. It is

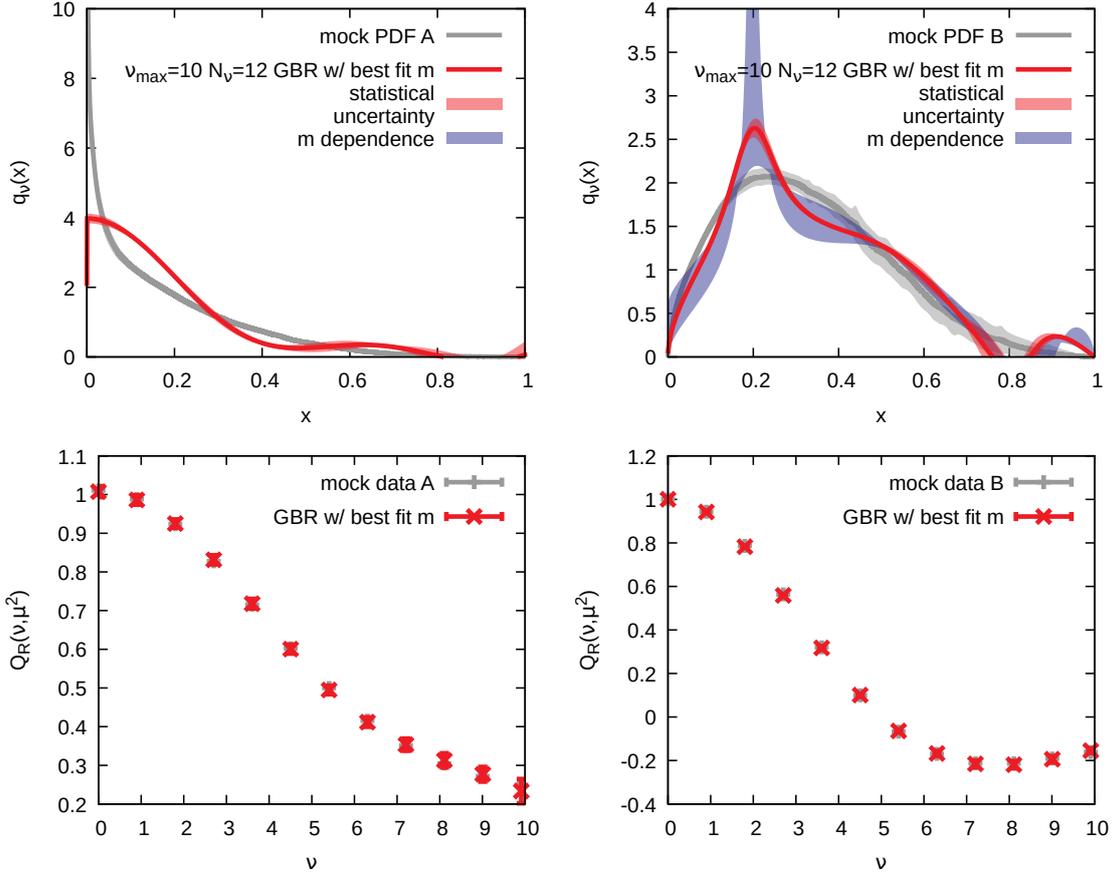


Figure 18. x -space PDF's reconstructed using the generalized Bayesian reconstruction (BRg) method from $N_\nu = 12$ Ioffe-time data points on the interval $\nu = [0, 10]$ (top) as well as the input data (gray crosses) compared to the data arising from the reconstructed PDF (red crosses) in the bottom panels.

interesting to then ask, whether the restriction to positive PDF's limits the number of local extrema in L , such that all of their contributions taken together in a statistical fashion will lead to a meaningful posterior?

To explore this idea, the MC-Stan library, a modern tool for statistical inference, can be used. It implements an efficient hybrid Monte-Carlo algorithm, specifically a no-U-Turn sampler, which allows the sampling of a wide variety of joint probability distributions of random variables. The applicability of the machinery of HMC arises from identifying its Hamiltonian with the logarithm of the joint probability distribution (for technical details see Ref. [53]). Using its high level programming language, the test can be formulated by identifying each of the (appropriately decorrelated) Ioffe-time mock data points with an individual Gaussian distribution. Its mean is expressed in term of the PDF $q(x)$ and the cosine integral Kernel, the spread by the corresponding data uncertainty. Positivity of the PDF is enforced by restricting its sampling to values larger than zero a priori. The complete STAN model reads

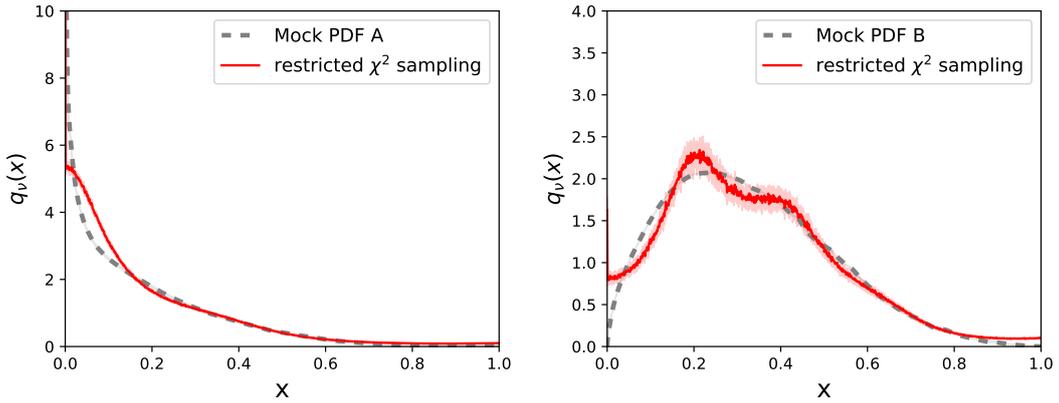


Figure 19. Reconstructed PDF (red dashed) from restricted sampling of the χ^2 functional based on $N_\nu = 12$ Ioffe-time data points on the interval $\nu = [0, 10]$. Mock scenario A is shown on the left, scenario B on the right. While positivity alone already provides a powerful regularization, for this realistic scenario we do observe deviations from the correct result, especially for scenario B.

```

data {
  int NNu; int Nx; matrix[NNu, Nx] Kernel;
  vector[NNu] Q;
  vector[NNu] Uncertainty;
}
parameters {
  vector<lower=0>[Nx] q;
}
model {
  Q ~ normal(Kernel * q, Uncertainty);
}

```

For the Kernel we discretize the x -range with $N_x = 1000$ points and consider two sets of input data for each scenario A and B. As realistic test case we use again the a $\nu \in [0, 10]$ with $N_\nu = 12$ points (see Fig. 19), exactly as has been used in the Bayesian reconstruction, as well as the ideal case of $N_\nu = 100$ points with $\nu_{\max} = 100$ (see Fig. 20). The sampling is performed over 20 chains of HMC trajectories, each with 1500 steps in MC time, 500 of which are discarded as warmup phase. The mean value of $q(x)$ over the 20 chains is shown as the red data points, while its erroband is obtained from the usual variance among different chains. As we do not imprint any further prior information, except positivity, and we assume that positive holds exactly, there is no further systematic error to be considered for this reconstruction.

We find that positivity alone is a powerful regulator, as sampling of the restricted χ^2 functional reveals. Even though we did not specify any further constraints, the posterior probability for each $q(x_i)$ leads to a well defined mean, its variance of course depends on the number and quality of the provided input data. In the ideal case of $\nu_{\max} = 100$ we manage to reproduce the PDF for scenario A excellently, for scenario B however ringing artifacts

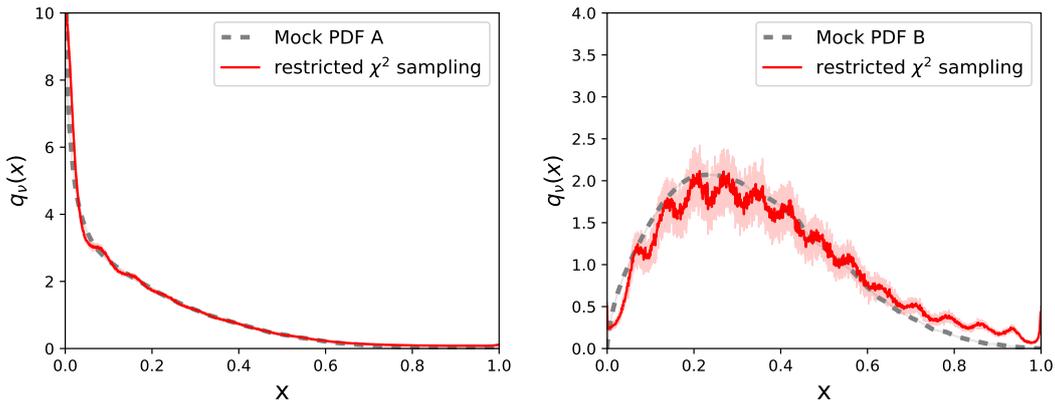


Figure 20. Reconstructed PDF (red dashed) from restricted sampling of the χ^2 functional based on $N_\nu = 100$ Ioffe-time data points on the interval $\nu = [0, 100]$. Mock scenario A is shown on the left, scenario B on the right. In this almost ideal setting, positivity alone provides a regularization that is sufficient for a reasonable determination of the PDF of scenario A, while it is unable to suppress ringing artifacts in scenario B.

remain clearly visible. On the other hand $\nu_{\max} = 20$, as expected leads to results that are less accurate than those obtained by the advanced Bayesian reconstruction methods or the neural network.

We thus learn that the additional prior information encoded in the Bayesian reconstruction manages to provide meaningful regularization to the problem beyond positivity. At the same time it also tells us that the neural network approach is able to extract more relevant prior information for a regularization than just positivity, i.e. it goes beyond a simple χ^2 fit.

5 Summary and Conclusion

In this work, we studied the general problem of extracting the full x -dependence of PDFs or DAs which are defined by the Fourier transform of a position space hadronic matrix element computed by means of lattice QCD. We have identified the two main challenges that one encounters during the computation of the aforementioned Fourier transform. The first one is related to the range of integration being restricted compared to a genuine Fourier transform and the second issue is that one has to perform such a task with only a few data points at hand. We explicitly have showed direct inversion methods fail at performing their task and the same also holds for trivial modifications thereof. However, as we demonstrated, advanced reconstruction methods open new paths towards obtaining PDFs and DAs from lattice QCD data. The methods that we tested here are the Backus-Gilbert method, a reconstruction which is based on neural networks, and a set of Bayesian reconstructions, including the MEM and the BR method.

We explicitly tested all these methods using mock input data computed using PDFs extracted from experiment that diverge at $x = 0$ as well as an artificial variation of these PDFs where the behavior at small x was modified to make the PDF vanish at $x = 0$. It is

clear that one can not probe the $x = 0$ region employing calculations at finite lattice spacing, which imposes a natural cutoff to the highest available momentum, in plain analogy to a scattering experiment. However, one can use advanced reconstruction methods that can determine the unknown PDFs with well defined estimates of the uncertainties.

In order to obtain a realistic impression of the method efficiencies we carried out an analysis using $\mathcal{O}(10)$ points for $\nu \in [0, 10]$. We find that such a limited set of data, which can be reached using current lattice QCD calculations, already leads to satisfactory reconstruction results for the underlying PDF for $x > 0.1$. The goal to double the available amount of Ioffe time input points to $\mathcal{O}(20)$ is realistic and achievable in the near future, given reasonable computational resources. However, tripling or quadrupling the maximum Ioffe time accessible would require a concerted effort by the lattice community working on PDFs. We have explicitly shown that an increase in the available range of ν values to $\nu_{\max} = 20$ significantly further reduces the uncertainty of the reconstruction for all three tested approaches, the Backus-Gilbert method, the neural network reconstruction, as well as the Bayesian methods.

For the case of the Backus-Gilbert reconstruction we saw that preconditioning the ill-defined problem was mandatory for obtaining the correct result in the small- x region. With a preconditioning that incorporated many of the relevant features of the final result, already a naive application of the algorithm was sufficient for obtaining a good reconstruction of the intermediate and large- x region. An advantage of the Backus-Gilbert method is its simplicity (linear problem) and the fact that it is numerically the cheapest of all to implement. On the other hand the relative lack of freedom in terms of preconditioning due to restrictions being imposed by the requirement of convergence of the involved integrals represents a minor downside of this approach.

The method of neural network parametrization also provided faithful reconstructions and has thus been shown to be a competitive candidate for solving the inverse problem for PDFs on the lattice. In this study, we have only explored a limited number of its facets, noticing that all three different geometries that we tested lead to fully equivalent results for the reconstructions of scenarios A and B. Furthermore, a small improvement of the reconstruction in the small x region can be achieved when preconditioning is employed as in the Backus-Gilbert method.

The Bayesian methods included the traditional Maximum Entropy Method (MEM) and the Bayesian reconstruction (BR) method. A fit of the input data with eq.(1.4) already provides a very good estimate of the shape of the underlying PDF and this information can be incorporated into the Bayesian methods as a default model. Hence, as expected, those Bayesian methods that are designed to imprint the information of the default model more strongly onto the end result (e.g. MEM) outperform those designed to keep the influence of the default model to a minimum (e.g. BR). We found that already with $\mathcal{O}(10)$ points for $\nu \in [0, 10]$ the MEM excellently reproduced the PDF of both scenario A and scenario B. The BR method is competitive in scenario A but showed larger deviations from the correct result for scenario B. Extending the range of ν to larger values and providing more input points consistently improves the reconstructions, approaching the correct result in the "Bayesian continuum limit".

When utilizing Bayesian methods that do not presuppose positive PDFs, the inverse problem is even more ill-posed. We have found that the quadratic prior Bayesian method performs remarkably well, as it also imprints the information provided by the default model strongly on the end result. It thus outperformed the generalized BR method for the case of a small number of input points $\mathcal{O}(10)$. Also in for general PDFs, increasing the range in available ν significantly reduced the uncertainties and made the results of different methods approach each other, as well as the correct result.

In order to further explore the role of prior information we performed an HMC sampling of the associated χ^2 functional employing the MC-Stan library. These studies indicated that the positivity constraint alone is a powerful regulator. At the same time we conclude that the additional prior information included in Bayesian methods, the Backus-Gilbert method and also the neural network approach are essential in producing a faithful reconstruction on the underlying PDF. Further studies using the HMC sampling approach are underway both in BR-MEM and the neural network methods.

We plan to apply the reconstruction approaches discussed in this paper in future publications that will employ realistic results from lattice calculations. We encourage other lattice practitioners to implement these methods, because in our opinion the systematic artifacts related to the Fourier transform among the numerous systematics that one has to face in the lattice studies of PDFs, are the ones that can not be dealt by brute force methods. This point has been discussed detail in [27, 28, 30, 31, 54], however our methods presented here completely avoid the difficulties that arise from the inverse Fourier transform that is required in other approaches.

Acknowledgments

We thank Anatoly Radyushkin for especially stimulating and enlightening discussions throughout this work. Fruitful discussions with Luigi del Debbio are also acknowledged. This work has been supported by the U.S. Department of Energy through Grant Number DE-FG02-04ER41302, and through contract Number DE-AC05-06OR23177, under which JSA operates the Thomas Jefferson National Accelerator Facility. SZ acknowledges support by the DFG Collaborative Research Centre SFB 1225 (ISOQUANT). K.O. acknowledges support in part by STFC consolidated grant ST/P000681/1, and the hospitality from DAMTP and Clare Hall at Cambridge University, where this work was performed. JK acknowledges support from the U.S. Department of Energy, Office of Science, Office of Workforce Development for Teachers and Scientists, Office of Science Graduate Student Research (SCGSR) program. The SCGSR program is administered by the Oak Ridge Institute for Science and Education for the DOE under contract number DE-SC0014664. This work was performed in part using computing facilities at the College of William & Mary which were provided by contributions from the National Science Foundation (MRI grant PHY-1626177), the Commonwealth of Virginia Equipment Trust Fund and the Office of Naval Research. In addition, this work used resources at NERSC, a DOE Office of Science User Facility supported by the Office of Science of the U.S. Department of Energy under Contract # DE-AC02-05CH11231 as well as computing resources provided by UNINETT Sigma2 - the National

References

- [1] R. Feynman, *Photon-hadron Interactions*, Advanced Books Classics (Avalon Publishing, 1998), ISBN 9780201360745.
- [2] C. Best, M. Gockeler, R. Horsley, E.-M. Ilgenfritz, H. Perlt, P. E. L. Rakow, A. Schafer, G. Schierholz, A. Schiller, and S. Schramm, Phys. Rev. **D56**, 2743 (1997), [hep-lat/9703014](#).
- [3] M. Guagnelli, K. Jansen, F. Palombi, R. Petronzio, A. Shindler, and I. Wetzorke (Zeuthen-Rome (ZeRo)), Eur. Phys. J. **C40**, 69 (2005), [hep-lat/0405027](#).
- [4] C. Alexandrou, M. Constantinou, K. Hadjiyiannakou, K. Jansen, C. Kallidonis, G. Koutsou, A. Vaquero Avilés-Casco, and C. Wiese, Phys. Rev. Lett. **119**, 142002 (2017), [1706.02973](#).
- [5] M. Oehm, C. Alexandrou, M. Constantinou, K. Jansen, G. Koutsou, B. Kostrzewa, F. Steffens, C. Urbach, and S. Zafeiropoulos (2018), [1810.09743](#).
- [6] A. V. Radyushkin, Phys. Rev. **D96**, 034025 (2017), [1705.01488](#).
- [7] X. Ji, Phys. Rev. Lett. **110**, 262002 (2013), [1305.1539](#).
- [8] K. Orginos, A. Radyushkin, J. Karpie, and S. Zafeiropoulos, Phys. Rev. **D96**, 094503 (2017), [1706.05373](#).
- [9] A. Radyushkin (2018), [1801.02427](#).
- [10] J.-H. Zhang, J.-W. Chen, and C. Monahan (2018), [1801.03023](#).
- [11] T. Izubuchi, X. Ji, L. Jin, I. W. Stewart, and Y. Zhao, Phys. Rev. **D98**, 056004 (2018), [1801.03917](#).
- [12] J. Karpie, K. Orginos, A. Radyushkin, and S. Zafeiropoulos (2017), [1710.08288](#).
- [13] X. Xiong, X. Ji, J.-H. Zhang, and Y. Zhao, Phys. Rev. **D90**, 014051 (2014), [1310.7471](#).
- [14] I. W. Stewart and Y. Zhao, Phys. Rev. **D97**, 054512 (2018), [1709.04933](#).
- [15] H.-W. Lin, J.-W. Chen, S. D. Cohen, and X. Ji, Phys. Rev. **D91**, 054510 (2015), [1402.1462](#).
- [16] C. Alexandrou, K. Cichy, V. Drach, E. Garcia-Ramos, K. Hadjiyiannakou, K. Jansen, F. Steffens, and C. Wiese, Phys. Rev. **D92**, 014502 (2015), [1504.07455](#).
- [17] J.-W. Chen, S. D. Cohen, X. Ji, H.-W. Lin, and J.-H. Zhang, Nucl. Phys. **B911**, 246 (2016), [1603.06664](#).
- [18] C. Alexandrou, K. Cichy, M. Constantinou, K. Jansen, A. Scapellato, and F. Steffens, Phys. Rev. Lett. **121**, 112001 (2018), [1803.02685](#).
- [19] J.-H. Zhang, J.-W. Chen, X. Ji, L. Jin, and H.-W. Lin, Phys. Rev. **D95**, 094514 (2017), [1702.00008](#).
- [20] J.-W. Chen, T. Ishikawa, L. Jin, H.-W. Lin, Y.-B. Yang, J.-H. Zhang, and Y. Zhao, Phys. Rev. **D97**, 014505 (2018), [1706.01295](#).
- [21] W. Broniowski and E. Ruiz Arriola, Phys. Rev. **D97**, 034031 (2018), [1711.03377](#).
- [22] Y.-Q. Ma and J.-W. Qiu, Phys. Rev. Lett. **120**, 022003 (2018), [1709.03018](#).

- [23] A. J. Chambers, R. Horsley, Y. Nakamura, H. Perlt, P. E. L. Rakow, G. Schierholz, A. Schiller, K. Somfleth, R. D. Young, and J. M. Zanotti, *Phys. Rev. Lett.* **118**, 242001 (2017), [1703.01153](#).
- [24] K.-F. Liu and S.-J. Dong, *Phys. Rev. Lett.* **72**, 1790 (1994), [hep-ph/9306299](#).
- [25] G. S. Bali, V. M. Braun, B. Gläsel, M. Göckeler, M. Gruber, F. Hutzler, P. Korcyl, A. Schäfer, P. Wein, and J.-H. Zhang, *Phys. Rev.* **D98**, 094507 (2018), [1807.06671](#).
- [26] R. S. Sufian, J. Karpie, C. Egerer, K. Orginos, J.-W. Qiu, and D. G. Richards (2019), [1901.03921](#).
- [27] G. C. Rossi and M. Testa, *Phys. Rev.* **D96**, 014507 (2017), [1706.04428](#).
- [28] G. Rossi and M. Testa, *Phys. Rev.* **D98**, 054028 (2018), [1806.00808](#).
- [29] X. Ji, J.-H. Zhang, and Y. Zhao, *Nucl. Phys.* **B924**, 366 (2017), [1706.07416](#).
- [30] A. V. Radyushkin, *Phys. Lett.* **B788**, 380 (2019), [1807.07509](#).
- [31] J. Karpie, K. Orginos, and S. Zafeiropoulos, *JHEP* **11**, 178 (2018), [1807.10933](#).
- [32] H.-W. Lin et al., *Prog. Part. Nucl. Phys.* **100**, 107 (2018), [1711.07916](#).
- [33] K. Cichy and M. Constantinou (2018), [1811.07248](#).
- [34] H.-W. Lin, J.-W. Chen, T. Ishikawa, and J.-H. Zhang (LP3), *Phys. Rev.* **D98**, 054504 (2018), [1708.05301](#).
- [35] G. Backus and F. Gilbert, *Geophysical Journal International* **16**, 169 (1968).
- [36] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery (1992).
- [37] B. B. Brandt, A. Francis, H. B. Meyer, and D. Robaina, *Phys. Rev.* **D92**, 094510 (2015), [1506.05732](#).
- [38] R.-A. Tripolt, P. Gubler, M. Ulybyshev, and L. Von Smekal (2018), [1801.10348](#).
- [39] J. Liang, K.-F. Liu, and Y.-B. Yang, *EPJ Web Conf.* **175**, 14014 (2018), [1710.11145](#).
- [40] M. V. Ulybyshev, C. Winterowd, and S. Zafeiropoulos, *EPJ Web Conf.* **175**, 03008 (2018), [1710.06675](#).
- [41] M. Ulybyshev, C. Winterowd, and S. Zafeiropoulos, *Phys. Rev.* **B96**, 205115 (2017), [1707.04212](#).
- [42] S. Forte, L. Garrido, J. I. Latorre, and A. Piccione, *JHEP* **05**, 062 (2002), [hep-ph/0204232](#).
- [43] R. D. Ball et al. (NNPDF), *Eur. Phys. J.* **C77**, 663 (2017), [1706.00428](#).
- [44] R. D. Ball et al. (NNPDF), *JHEP* **04**, 040 (2015), [1410.8849](#).
- [45] R. D. Ball, L. Del Debbio, S. Forte, A. Guffanti, J. I. Latorre, J. Rojo, and M. Ubiali, *Nucl. Phys.* **B838**, 136 (2010), [1002.4407](#).
- [46] R. D. Ball, V. Bertone, F. Cerutti, L. Del Debbio, S. Forte, A. Guffanti, J. I. Latorre, J. Rojo, and M. Ubiali (NNPDF), *Nucl. Phys.* **B849**, 112 (2011), [Erratum: *Nucl. Phys.* **B855**, 927 (2012)], [1012.0836](#).
- [47] J. Rojo, in *13th Conference on Quark Confinement and the Hadron Spectrum (Confinement XIII) Maynooth, Ireland, July 31-August 6, 2018* (2018), [1809.04392](#).

- [48] J. Skilling and S. F. Gull, *Bayesian maximum entropy image reconstruction* (Institute of Mathematical Statistics, Hayward, CA, 1991), vol. Volume 20 of *Lecture Notes–Monograph Series*, pp. 341–367, URL <https://doi.org/10.1214/lnms/1215460511>.
- [49] M. Asakawa, T. Hatsuda, and Y. Nakahara, *Prog. Part. Nucl. Phys.* **46**, 459 (2001), [hep-lat/0011040](https://arxiv.org/abs/hep-lat/0011040).
- [50] A. Rothkopf, *J. Comput. Phys.* **238**, 106 (2013), [1110.6285](https://arxiv.org/abs/1110.6285).
- [51] Y. Burnier and A. Rothkopf, *Phys. Rev. Lett.* **111**, 182003 (2013), [1307.6106](https://arxiv.org/abs/1307.6106).
- [52] A. Buckley, J. Ferrando, S. Lloyd, K. Nordström, B. Page, M. Rüfenacht, M. Schönherr, and G. Watt, *Eur. Phys. J.* **C75**, 132 (2015), [1412.7420](https://arxiv.org/abs/1412.7420).
- [53] B. Carpenter, A. Gelman, M. D. Hoffman, D. Lee, B. Goodrich, M. Betancourt, M. Brubaker, J. Guo, P. Li, and A. Riddell, *Journal of Statistical Software* **76** (2017).
- [54] G. Rossi and M. Testa (2018), [1811.10267](https://arxiv.org/abs/1811.10267).