

# A semi-intrusive Code Framework for Uncertainty Quantification

Jonas Kusch<sup>a</sup>, Jannick Wolters<sup>b</sup>, Martin Frank<sup>c</sup>

<sup>a</sup>Karlsruhe Institute of Technology, Karlsruhe, jonas.kusch@kit.edu

<sup>b</sup>Karlsruhe Institute of Technology, Karlsruhe, jannick.wolters@kit.edu

<sup>c</sup>Karlsruhe Institute of Technology, Karlsruhe, martin.frank@kit.edu

---

## Abstract

Methods for quantifying the effects of uncertainties in hyperbolic problems can be divided into intrusive and non-intrusive techniques. Intrusive methods require the implementation of a new code and yield

*Keywords:* conservation laws, hyperbolic, intrusive, UQ, IPM, SG, Collocation

---

## 1. Introduction

Hyperbolic equations play an important role in various research areas such as fluid dynamics or plasma physics. Efficient numerical methods combined with robust implementations are available for these problems, however they do not account for uncertainties which can arise in measurement data or modeling assumptions. Including the effects of uncertainties in differential equations has become an important topic in the last decades.

A general hyperbolic set of equations with random initial data can be written as

$$\partial_t \mathbf{u}(t, \mathbf{x}, \boldsymbol{\xi}) + \nabla \cdot \mathbf{F}(\mathbf{u}(t, \mathbf{x}, \boldsymbol{\xi})) = \mathbf{0}, \quad (1a)$$

$$\mathbf{u}(t = 0, \mathbf{x}, \boldsymbol{\xi}) = \mathbf{u}_{IC}(\mathbf{x}, \boldsymbol{\xi}), \quad (1b)$$

where the solution  $u \in \mathbb{R}^p$  depends on time  $t \in \mathbb{R}^+$ , spatial position  $\mathbf{x} \in D \subset \mathbb{R}^d$  as well as a vector of random variables  $\boldsymbol{\xi} \in \Theta \subset \mathbb{R}^s$  with given probability density functions  $f_{\Xi,i}(\xi_i)$  for  $i = 1, \dots, s$ . The physical flux is given by  $\mathbf{F} : \mathbb{R}^p \rightarrow \mathbb{R}^p$ . To simplify notation, we assume that the initial condition is random, i.e.  $\boldsymbol{\xi}$  enters through the definition of  $\mathbf{u}_{IC}$ . Equations (1) are usually supplemented with boundary conditions, which we will specify later for the individual problems.

Due to the randomness of the solution, one is interested in determining the expectation value or the variance of the solution, i.e.

$$\mathbb{E}[\mathbf{u}] = \langle \mathbf{u} \rangle, \quad \text{Var}[\mathbf{u}] = \langle (\mathbf{u} - \mathbb{E}[\mathbf{u}])^2 \rangle,$$

where we use the bracket operator  $\langle \cdot \rangle := \int_{\Theta} \cdot \prod_{i=1}^s f_{\Xi,i}(\xi_i) d\xi_1 \cdots d\xi_s$ . More generally, one is interested in determining the moments of the solution for a given set of orthonormal basis functions  $\varphi_i : \Theta \rightarrow \mathbb{R}$  such that for the multi-index  $i = (i_1, \dots, i_s)$  we have  $|i| \leq N$ . The moments are then given by  $\hat{\mathbf{u}}_i := \langle \mathbf{u} \varphi_i \rangle$ .

Numerical methods for approximating the moment  $\hat{\mathbf{u}}_i$  can be divided into intrusive and non-intrusive methods. A popular non-intrusive method is the stochastic-Collocation (SC) method, which computes the moments with the help of a numerical quadrature rule: For a given set of  $N_q$  quadrature weights  $w_k$  and

quadrature points  $\xi_k$ , the moments are approximated by

$$\hat{\mathbf{u}}_i = \langle \mathbf{u} \varphi_i \rangle \approx \sum_{k=1}^{N_q} w_k \mathbf{u}(t, \mathbf{x}, \xi_k) \varphi_i(\xi_k) f_{\Xi}(\xi_k).$$

Since the solution at a fixed quadrature point can be computed by a standard deterministic solver, the SC method does not require a big implementation effort. Furthermore, it is embarrassingly parallel, since the required computations can be carried out in parallel on different cores. A downside of collocation methods are aliasing effects, which stem from the inexact approximation of the integral. Furthermore, collocation methods require more runs of the deterministic solver than intrusive methods [1, 2]. Despite their easy implementation, collocation methods can become cumbersome for unsteady problems, since in this case, the solution needs to be written out at all quadrature points in every time step to compute the time evolution of the moments. This problem does not occur in intrusive methods, since the computation is carried out on the moments, i.e. the time evolution can be recorded during the computation. Also for steady problems, the fact that the computation is carried out on the moments provided the ability to compute the stochastic residual, which indicates when to stop the iteration towards the steady state solution. However intrusive methods are in general more difficult to implement and come with higher numerical costs. The main idea of these methods is to derive a system of equations for the moments and implementing a numerical solver for this system: Testing the initial problem (1) with  $\varphi_i$  for  $|i| \leq M$  yields

$$\partial_t \hat{\mathbf{u}}_i(t, \mathbf{x}) + \nabla \cdot \langle \mathbf{F}(\mathbf{u}(t, \mathbf{x}, \cdot)) \varphi_i \rangle = 0, \quad (2a)$$

$$\hat{\mathbf{u}}_i(t=0, \mathbf{x}) = \langle \mathbf{u}_{IC}(\mathbf{x}, \cdot) \varphi_i \rangle. \quad (2b)$$

To obtain a closed set of equations, one needs to derive a closure  $\mathcal{U}$  such that

$$\mathbf{u}(t, \mathbf{x}, \xi) \approx \mathcal{U}(\hat{\mathbf{u}}_0, \dots, \hat{\mathbf{u}}_N; \xi).$$

A commonly used closure is given by stochastic-Galerkin, which represents the solution by a polynomial:

$$\mathcal{U}_{SG}(\hat{\mathbf{u}}_0, \dots, \hat{\mathbf{u}}_N; \xi) := \sum_{i=0}^N \hat{\mathbf{u}}_i \varphi_i(\xi) = \hat{\mathbf{u}}^T \boldsymbol{\varphi},$$

where we collect moments in  $\hat{\mathbf{u}} := (\hat{u}_0, \dots, \hat{u}_N)^T \in \mathbb{R}^{(N+1) \times p}$  and basis functions in  $\boldsymbol{\varphi} := (\varphi_0, \dots, \varphi_N)^T \in \mathbb{R}^{N+1}$ . When using the stochastic-Galerkin method to close (2), the resulting moment system is not necessarily hyperbolic, meaning that it cannot be solved with standard methods. A generalization of stochastic-Galerkin, which ensures hyperbolicity is the Intrusive Polynomial Moment (IPM) closure [3]: The closure is given by a constraint optimization problem. For a given convex entropy  $s : \mathbb{R}^p \rightarrow \mathbb{R}$  for the original problem (1), this optimization problem is

$$\mathcal{U}_{ME}(\hat{\mathbf{u}}) = \arg \min_{\mathbf{u}} \langle s(\mathbf{u}) \rangle \quad \text{subject to } \hat{\mathbf{u}}_i = \langle \mathbf{u} \varphi_i \rangle \text{ for } i = 0, \dots, N. \quad (3)$$

This problem can be rewritten as

$$\hat{\boldsymbol{\lambda}}(\hat{\mathbf{u}}) := \arg \min_{\boldsymbol{\lambda} \in \mathbb{R}^{(N+1) \times p}} \left\{ \langle s_*(\boldsymbol{\lambda}^T \boldsymbol{\varphi}) \rangle - \sum_{i=0}^N \boldsymbol{\lambda}_i^T \hat{\mathbf{u}}_i \right\}, \quad (4)$$

where  $s_* : \mathbb{R}^p \rightarrow \mathbb{R}$  is the Legendre transformation of  $s$ , and  $\hat{\boldsymbol{\lambda}} \in \mathbb{R}^{(N+1) \times p}$  are called the dual variables. The solution to (3) is then given by

$$\mathcal{U}_{ME}(\hat{\mathbf{u}}) = (s_{\hat{\mathbf{u}}})^{-1}(\hat{\boldsymbol{\lambda}}(\hat{\mathbf{u}})^T \boldsymbol{\varphi}). \quad (5)$$

Plugging the derived closure into the moment system (2), one obtains the IPM moment system

$$\partial_t \hat{\mathbf{u}}_i(t, \mathbf{x}) + \nabla \cdot \langle \mathbf{F}(\mathcal{U}_{ME}(\hat{\mathbf{u}})) \varphi_i \rangle = \mathbf{0}, \quad (6a)$$

$$\hat{\mathbf{u}}_i(t = 0, \mathbf{x}) = \langle \mathbf{u}_{IC}(\mathbf{x}, \cdot) \varphi_i \rangle. \quad (6b)$$

The IPM method has several advantages: Choosing the entropy  $s(\mathbf{u}) = \frac{1}{2} \mathbf{u}^T \mathbf{u}$  yields the stochastic-Galerkin method, i.e. IPM generalizes different intrusive method. Furthermore, at least for scalar problems, IPM is significantly less oscillatory compared to SG [4]. Also, as discussed in [3], when choosing  $s(\mathbf{u})$  to be a physically correct entropy of the deterministic problem, the IPM solution dissipates the entropy

$$S(\hat{\mathbf{u}}) := \langle s(\mathcal{U}_{ME}(\hat{\mathbf{u}})) \rangle,$$

i.e. the IPM method yields a physically correct entropy solution. This again underlines a weakness of stochastic-Galerkin: If  $s(\mathbf{u}) = \frac{1}{2} \mathbf{u}^T \mathbf{u}$  is not a correct entropy of the original problem, the SG method can lead to non-physical solution, which can then cause a failure of the method. The main weakness of the IPM method is its run time, since it requires the repeated evaluation of (5), which involves the computation of the optimization problem (4). Hence, the desirable costs of IPM (and intrusive methods in general) require a significantly increased run time. A further advantages which we will not use in this work is the possibility of using adaptivity [5], i.e. refining the stochastic space at certain spatial position and time steps in which complex structures such as discontinuities occur. The main task here is to find an adequate refinement indicator [5, 6].

In this paper, we present a semi-intrusive code framework, which facilitates the task of implementing general intrusive methods. The framework only requires the numerical flux of the deterministic problem. If IPM is used, the entropy of the deterministic problem needs to be provided. We thereby provide the ability to recycle existing implementations of deterministic solvers. Furthermore, we investigate intrusive methods for steady problems and compare them to collocation methods. The steady setting provides different opportunities to take advantage of features of intrusive methods:

- Apply intrusive methods as a post-processing step for collocation methods: We converge the moments of the solution to a steady state with an inaccurate, but cheap collocation method and then use the resulting collocation moments as starting values for an expensive but accurate intrusive method, which we then again converge to steady state.
- Since the moments during the iteration process are inaccurate, i.e. they are not the steady state solution, we propose to not fully converge the dual iteration, which computes (4).

The paper is structured as follows: After the introduction, we discuss the numerical discretization as well as the implementation and structure of the semi-intrusive framework in section 2. Section 3 discusses the IPM acceleration with a non-intrusive method and in section 4. A comparison of results computed with the presented methods is given in 5, followed by a summary and outlook in section 6.

## 2. Discretization and code framework

### 2.1. Discretization

In the following, we discretize the moment system in space and time following according to [4]. Due to the fact, that stochastic-Galerkin can be interpreted as IPM with a quadratic entropy, it suffices to only derive a discretization of the IPM moment system (6). Omitting initial conditions and assuming a one-dimensional spatial domain, we can write this system as

$$\partial_t \hat{\mathbf{u}} + \partial_x \mathbf{F}(\hat{\mathbf{u}}) = \mathbf{0}$$

with the flux  $\mathbf{F} : \mathbb{R}^{(N+1) \times p} \rightarrow \mathbb{R}^{(N+1) \times p}$ ,  $\mathbf{F}(\mathbf{u}) = \langle f(\mathcal{U}_{ME}(\hat{\mathbf{u}}))\boldsymbol{\varphi} \rangle$ , where we omit the index ME in the following for efficiency of exposition. Due to hyperbolicity of the IPM moment system, one can use a finite-volume method to approximate the time evolution of the IPM moments. First we perform the discretization of the spatial domain: We choose the discrete unknowns to be the spatial averages over each cell at time  $t_n$ , given by

$$\hat{\mathbf{u}}_{ij}^n \simeq \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_{j+1/2}} \hat{\mathbf{u}}_i(t_n, x) dx.$$

If a moment vector in cell  $j$  at time  $t_n$  is denoted as  $\hat{\mathbf{u}}_j^n = (\hat{\mathbf{u}}_{0j}^n, \dots, \hat{\mathbf{u}}_{N_j}^n)^T \in \mathbb{R}^{N+1}$ , the finite-volume scheme can be written in conservative form with the numerical flux  $\mathbf{G} : \mathbb{R}^{(N+1) \times p} \times \mathbb{R}^{(N+1) \times p} \rightarrow \mathbb{R}^{(N+1) \times p}$  as

$$\hat{\mathbf{u}}_j^{n+1} = \hat{\mathbf{u}}_j^n - \frac{\Delta t}{\Delta x} (\mathbf{G}(\hat{\mathbf{u}}_j^n, \hat{\mathbf{u}}_{j+1}^n) - \mathbf{G}(\hat{\mathbf{u}}_{j-1}^n, \hat{\mathbf{u}}_j^n)) \quad (7)$$

for  $j = 1, \dots, N_x$  and  $n = 0, \dots, N_t$ , where  $N_x$  is the number of spatial cells and  $N_t$  is the number of time steps. The numerical flux is assumed to be consistent, i.e., that  $\mathbf{G}(\mathbf{u}, \mathbf{u}) = \mathbf{F}(\mathbf{u})$ . To ensure stability, a CFL condition has to be derived by investigating the eigenvalues of  $\nabla \mathbf{F}$ .

When a consistent numerical flux  $\mathbf{g} : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}^p$ ,  $\mathbf{g} = \mathbf{g}(\mathbf{u}_\ell, \mathbf{u}_r)$  is available for the deterministic problem (??), then for the IPM system we can simply take

$$\mathbf{G}(\hat{\mathbf{u}}_j^n, \hat{\mathbf{u}}_{j+1}^n) = \langle g(u_{ME}(\hat{\Lambda}(\hat{\mathbf{u}}_j^n)), u_{ME}(\hat{\Lambda}(\hat{\mathbf{u}}_{j+1}^n)))\boldsymbol{\varphi} \rangle.$$

This choice of the numerical flux is a common choice in kinetic theory and is called kinetic flux. The time update of the moment vector now becomes

$$\hat{\mathbf{u}}_j^{n+1} = \hat{\mathbf{u}}_j^n - \frac{\Delta t}{\Delta x} \left( \langle g(u_{ME}(\hat{\Lambda}_j^n), u_{ME}(\hat{\Lambda}_{j+1}^n)))\boldsymbol{\varphi} \rangle - \langle g(u_{ME}(\hat{\Lambda}_{j-1}^n), u_{ME}(\hat{\Lambda}_j^n)))\boldsymbol{\varphi} \rangle \right), \quad (8)$$

where  $\hat{\Lambda}_j^n := \hat{\Lambda}(\hat{\mathbf{u}}_j^n)$  for all  $j$ . Note that the computation of  $\hat{\Lambda}_j^n$  requires solving the dual problem (4) for the moment vector  $\hat{\mathbf{u}}_j^n$ .

Unfortunately (8) cannot be implemented because the dual problem cannot be solved exactly.<sup>1</sup> Instead, it must be solved numerically, for example with Newton's method. The stopping criterion for the numerical optimizer ensures that the approximate multiplier vector it returns, which we denote  $\bar{\boldsymbol{\lambda}}_j^n \in \mathbb{R}^{N+1}$  for the moment vector  $\hat{\mathbf{u}}_j^n$ , satisfies the stopping criterion

$$\left\| \hat{\mathbf{u}}_j^n - \left\langle u_{ME} \left( (\bar{\boldsymbol{\lambda}}_j^n)^T \boldsymbol{\varphi} \right) \right\rangle \right\| < \tau. \quad (9)$$

This is derived from the first-order necessary conditions for the dual problem. Once the numerical optimizer finds such a  $\bar{\boldsymbol{\lambda}}_j^n$ , the corresponding dual state  $\bar{\Lambda}_j^n := (\bar{\boldsymbol{\lambda}}_j^n)^T \boldsymbol{\varphi} \in \mathbb{P}(\Theta)$  can be used in (8) for the unknown  $\hat{\Lambda}_j^n$ . This gives Algorithm 1.

---

<sup>1</sup>Equation (8) also includes integral evaluations which cannot be computed in closed form. Their approximation by numerical quadrature, however, does not play a role in the realizability problems we discuss below.

---

**Algorithm 1:** IPM for Uncertainty Quantification

---

```

1: for  $j = 0$  to  $N_x + 1$  do
2:    $\mathbf{u}_j^0 = \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_{j+1/2}} \langle u_0(x, \cdot) \varphi \rangle dx$ 
3: end for
4: for  $n = 0$  to  $N_t$  do
5:   for  $j = 0$  to  $N_x + 1$  do
6:      $\bar{\lambda}_j^n \approx \arg \min_{\lambda} (\langle s_*(\lambda^T \varphi) \rangle - \lambda^T \hat{\mathbf{u}}_j^n)$  such that (9) holds
7:      $\bar{\Lambda}_j^n = (\bar{\lambda}_j^n)^T \varphi$ 
8:   end for
9:   for  $j = 1$  to  $N_x$  do
10:     $\hat{\mathbf{u}}_j^{n+1} = \hat{\mathbf{u}}_j^n - \frac{\Delta t}{\Delta x} (\langle g(u_{ME}(\bar{\Lambda}_j^n), u_{ME}(\bar{\Lambda}_{j+1}^n)) \varphi \rangle - \langle g(u_{ME}(\bar{\Lambda}_{j-1}^n), u_{ME}(\bar{\Lambda}_j^n)) \varphi \rangle)$ 
11:   end for
12: end for

```

---

### 3. Collocation accelerated IPM for steady state problems

#### 4. One-shot IPM

For classical IPM, the iteration scheme for the moment vectors is

$$\mathbf{u}_j^{n+1} = \mathbf{c}(\lambda(\mathbf{u}_{j-1}^n), \lambda(\mathbf{u}_j^n), \lambda(\mathbf{u}_{j+1}^n)), \quad (10)$$

where  $\mathbf{c} : \mathbb{R}^{N+1} \times \mathbb{R}^{N+1} \times \mathbb{R}^{N+1} \rightarrow \mathbb{R}^{N+1}$  is given by

$$\mathbf{c}(\lambda_\ell, \lambda_c, \lambda_r) := \langle u(\lambda_c^T \varphi) \varphi \rangle - \frac{\Delta t}{\Delta x} (\langle g(u(\lambda_c^T \varphi), u(\lambda_r^T \varphi)) \varphi \rangle - \langle g(u(\lambda_\ell^T \varphi), u(\lambda_r^T \varphi)) \varphi \rangle).$$

The map from the moment vector to the dual variables is given by the exact fix point of  $\mathbf{d} : \mathbb{R}^{N+1} \times \mathbb{R}^{N+1} \rightarrow \mathbb{R}^{N+1}$ ,

$$\mathbf{d}(\lambda, \mathbf{u}) := \lambda - \mathbf{B} \cdot (\langle u(\lambda^T \varphi) \varphi \rangle - \mathbf{u}),$$

where  $\mathbf{B}$  is a preconditioner, which is used in the numerical scheme to ensure convergence of the fix point iteration (which we call dual iteration in the following)

$$\lambda_j^{m+1} = \mathbf{d}(\lambda_j^m, \mathbf{u}_j^n). \quad (11)$$

Note that here,  $\mathbf{u}_j^n$  is a fixed parameter. By letting the dual iteration converge, i.e.  $m \rightarrow \infty$ , we obtain  $\lim_{m \rightarrow \infty} \mathbf{d}(\lambda_j^m, \mathbf{u}_j^n) =: \lambda_j^n \equiv \lambda(\mathbf{u}_j^n)$ , which are called the dual variables of  $\mathbf{u}_j^n$ . Hence, the numerical scheme for IPM converges the dual iteration (11) to compute the mapping from the moment vectors to their dual variables and then performs one iteration of (10) to obtain the time update of the moment vector. Note that this can become extremely expensive, since in each iteration of the moment vectors, we fully converge the dual iteration for all spatial cells.

To simplify notation, we leave out the dependency on  $\lambda$  when writing the moment scheme in the following: Hence with

$$\tilde{\mathbf{c}}(\mathbf{u}_\ell^n, \mathbf{u}_c^n, \mathbf{u}_r^n) := \mathbf{c}(\lambda(\mathbf{u}_\ell^n), \lambda(\mathbf{u}_c^n), \lambda(\mathbf{u}_r^n))$$

we can rewrite the moment iteration as

$$\mathbf{u}_j^{n+1} = \tilde{\mathbf{c}}(\mathbf{u}_{j-1}^n, \mathbf{u}_j^n, \mathbf{u}_{j+1}^n). \quad (12)$$

For steady problems, we assume that the IPM scheme converges to a fix point, i.e. we must have that  $\rho(\tilde{\mathbf{c}}_{\mathbf{u}}) < 1$ . The main idea of *One-Shot IPM* is to not fully converge the dual iteration, since the moment vectors are not yet converged to the exact steady solution. So if we successively perform one update of the moment iteration and one update of the dual iteration, we obtain

$$\boldsymbol{\lambda}_j^{n+1} = \mathbf{d}(\boldsymbol{\lambda}_j^n, \mathbf{u}_j^n) \quad \text{for all } j \quad (13a)$$

$$\mathbf{u}_j^{n+1} = \mathbf{c}(\boldsymbol{\lambda}_{j-1}^{n+1}, \boldsymbol{\lambda}_j^{n+1}, \boldsymbol{\lambda}_{j+1}^{n+1}). \quad (13b)$$

In the following, we study convergence of this scheme. For this we take an approach commonly chosen to prove the local convergence properties of Newton's method: In Theorem 1, we show that the iteration function is contractive at its fix point and conclude in Theorem 2 that this yields local convergence.

**Theorem 1.** *Assume that the classical IPM iteration (12) is contractive at its fix point  $\mathbf{u}^*$ . Then the Jacobi matrix  $\mathbf{J}$  of the One-Shot IPM iteration (13) has a spectral radius  $\rho(\mathbf{J}) < 1$  at the fix point  $(\boldsymbol{\lambda}^*, \mathbf{u}^*)$  if the preconditioner  $\mathbf{B} = \langle u'(\boldsymbol{\varphi}^T \boldsymbol{\lambda}_j^n) \boldsymbol{\varphi} \boldsymbol{\varphi}^T \rangle^{-1}$ , i.e. the Hessian of the dual problem is used.*

*Proof.* We first point out that since we assume that the classical IPM scheme is contractive at its fix point, we have  $\rho(\tilde{\mathbf{c}}_{\mathbf{u}}(\mathbf{u}^*)) < 1$  with  $\tilde{\mathbf{c}}_{\mathbf{u}} \in \mathbb{R}^{(N+1) \cdot N_x \times (N+1) \cdot N_x}$  defined by

$$\tilde{\mathbf{c}}_{\mathbf{u}} = \begin{pmatrix} \partial_{\mathbf{u}_c} \tilde{\mathbf{c}}_1 & \partial_{\mathbf{u}_r} \tilde{\mathbf{c}}_1 & 0 & 0 & \dots \\ \partial_{\mathbf{u}_\ell} \tilde{\mathbf{c}}_2 & \partial_{\mathbf{u}_c} \tilde{\mathbf{c}}_2 & \partial_{\mathbf{u}_r} \tilde{\mathbf{c}}_2 & 0 & \dots \\ 0 & \partial_{\mathbf{u}_\ell} \tilde{\mathbf{c}}_3 & \partial_{\mathbf{u}_c} \tilde{\mathbf{c}}_3 & \partial_{\mathbf{u}_r} \tilde{\mathbf{c}}_3 & \\ \vdots & & & \ddots & \\ 0 & \dots & 0 & \partial_{\mathbf{u}_\ell} \tilde{\mathbf{c}}_{N_x} & \partial_{\mathbf{u}_c} \tilde{\mathbf{c}}_{N_x} \end{pmatrix},$$

where we define  $\tilde{\mathbf{c}}_j := \tilde{\mathbf{c}}(\mathbf{u}_{j-1}^*, \mathbf{u}_j^*, \mathbf{u}_{j+1}^*)$  for all  $j$ . Now we have for each term inside the matrix  $\tilde{\mathbf{c}}_{\mathbf{u}}$

$$\partial_{\mathbf{u}_\ell} \tilde{\mathbf{c}}_j = \frac{\partial \mathbf{c}_j}{\partial \boldsymbol{\lambda}_\ell} \frac{\partial \boldsymbol{\lambda}(\mathbf{u}_{j-1}^*)}{\partial \mathbf{u}}, \quad \partial_{\mathbf{u}_c} \tilde{\mathbf{c}}_j = \frac{\partial \mathbf{c}_j}{\partial \boldsymbol{\lambda}_c} \frac{\partial \boldsymbol{\lambda}(\mathbf{u}_j^*)}{\partial \mathbf{u}}, \quad \partial_{\mathbf{u}_r} \tilde{\mathbf{c}}_j = \frac{\partial \mathbf{c}_j}{\partial \boldsymbol{\lambda}_r} \frac{\partial \boldsymbol{\lambda}(\mathbf{u}_{j+1}^*)}{\partial \mathbf{u}}. \quad (14)$$

We first wish to understand the structure of the terms  $\partial_{\mathbf{u}} \boldsymbol{\lambda}(\mathbf{u})$ . For this, we note that the exact dual state fulfills

$$\mathbf{u} = \langle u(\boldsymbol{\lambda}^T \boldsymbol{\varphi}) \boldsymbol{\varphi} \rangle =: \mathbf{h}(\boldsymbol{\lambda}), \quad (15)$$

which is why we have the mapping  $\mathbf{u} : \mathbb{R}^{N+1} \rightarrow \mathbb{R}^{N+1}$ ,  $\mathbf{u}(\boldsymbol{\lambda}) = \mathbf{h}(\boldsymbol{\lambda})$ . Since the solution of the dual problem for a given moment vector is unique, this mapping is bijective and therefore we have an inverse function

$$\boldsymbol{\lambda} = \mathbf{h}^{-1}(\mathbf{u}(\boldsymbol{\lambda})) \quad (16)$$

Now we differentiate both sides w.r.t.  $\boldsymbol{\lambda}$  to get

$$\mathbf{I}_d = \frac{\partial \mathbf{h}^{-1}(\mathbf{u}(\boldsymbol{\lambda}))}{\partial \mathbf{u}} \frac{\partial \mathbf{u}(\boldsymbol{\lambda})}{\partial \boldsymbol{\lambda}}.$$

We multiply with the matrix inverse of  $\frac{\partial \mathbf{u}(\boldsymbol{\lambda})}{\partial \boldsymbol{\lambda}}$  to get

$$\left( \frac{\partial \mathbf{u}(\boldsymbol{\lambda})}{\partial \boldsymbol{\lambda}} \right)^{-1} = \frac{\partial \mathbf{h}^{-1}(\mathbf{u}(\boldsymbol{\lambda}))}{\partial \mathbf{u}}.$$

Note that on the left-hand-side we have the inverse of a matrix and on the right-hand-side, we have the inverse of a multi-dimensional function. By rewriting  $\mathbf{h}^{-1}(\mathbf{u}(\boldsymbol{\lambda}))$  as  $\boldsymbol{\lambda}(\mathbf{u})$  and simply computing the term  $\frac{\partial \boldsymbol{\lambda}(\mathbf{u})}{\partial \mathbf{u}}$  by differentiating (15) w.r.t.  $\boldsymbol{\lambda}$ , one obtains

$$\partial_{\mathbf{u}} \boldsymbol{\lambda}(\mathbf{u}) = \langle u'(\boldsymbol{\lambda}^T \boldsymbol{\varphi}) \boldsymbol{\varphi} \boldsymbol{\varphi}^T \rangle^{-1}. \quad (17)$$

Now we begin to derive the spectrum of the *One-Shot IPM* iteration (13). Note that in its current form this iteration is not really a fix point iteration, since it uses the time updated dual variables in (13b). To obtain a fix point iteration, we plug the dual iteration step (13a) into the moment iteration (13b) to obtain

$$\begin{aligned}\lambda_j^{n+1} &= d(\lambda_j^n, \mathbf{u}_j^n) \quad \text{for all } j \\ \mathbf{u}_j^{n+1} &= \mathbf{c}(d(\lambda_{j-1}^n, \mathbf{u}_{j-1}^n), d(\lambda_j^n, \mathbf{u}_j^n), d(\lambda_{j+1}^n, \mathbf{u}_{j+1}^n))\end{aligned}$$

The Jacobian  $\mathbf{J} \in \mathbb{R}^{2(N+1) \cdot N_x \times 2(N+1) \cdot N_x}$  has the form

$$\mathbf{J} = \begin{pmatrix} \partial_\lambda \mathbf{d} & \partial_{\mathbf{u}} \mathbf{d} \\ \partial_\lambda \mathbf{c} & \partial_{\mathbf{u}} \mathbf{c} \end{pmatrix}, \quad (18)$$

where each block has entries for all spatial cells. We start by looking at  $\partial_\lambda \mathbf{d}$ . For the columns belonging to cell  $j$ , we have

$$\begin{aligned}\partial_\lambda d(\lambda_j^n, \mathbf{u}_j^n) &= \mathbf{I}_d - \mathbf{B} \cdot \langle u'(\varphi^T \lambda_j^n) \varphi \varphi^T \rangle - \partial_\lambda \mathbf{B} \cdot (\langle u(\varphi^T \lambda_j^n) \varphi \rangle - \mathbf{u}) \\ &= -\partial_\lambda \mathbf{B} \cdot (\langle u(\varphi^T \lambda_j^n) \varphi \rangle - \mathbf{u}),\end{aligned}$$

where we used  $\mathbf{B} = \langle u'(\varphi^T \lambda_j^n) \varphi \varphi^T \rangle^{-1}$ . Recall that at the fix point  $(\lambda^*, \mathbf{u}^*)$ , we have  $\langle u(\varphi^T \lambda_j^n) \varphi \rangle = \mathbf{u}$ , hence one obtains  $\partial_\lambda \mathbf{d} = \mathbf{0}$ . For the block  $\partial_{\mathbf{u}} \mathbf{d}$ , we get

$$\partial_{\mathbf{u}} d(\lambda_j^n, \mathbf{u}_j^n) = \mathbf{B},$$

hence  $\partial_{\mathbf{u}} \mathbf{d}$  is a block diagonal matrix. Let us now look at  $\partial_\lambda \mathbf{c}$  at a fixed spatial cell  $j$ :

$$\frac{\partial \mathbf{c}}{\partial \lambda_\ell} \frac{\partial d(\lambda_{j-1}^n, \mathbf{u}_{j-1}^n)}{\partial \lambda} = \mathbf{0},$$

since we already showed that by the choice of  $\mathbf{B}$  the term  $\partial_\lambda \mathbf{d}$  is zero. We can show the same result for all spatial cells and all inputs of  $\mathbf{c}$  analogously, hence  $\partial_\lambda \mathbf{c} = \mathbf{0}$ . For the last block, we have that

$$\frac{\partial \mathbf{c}}{\partial \lambda_\ell} \frac{\partial d(\lambda_{j-1}^n, \mathbf{u}_{j-1}^n)}{\partial \mathbf{u}} = \frac{\partial \mathbf{c}}{\partial \lambda_\ell} \mathbf{B} = \frac{\partial \mathbf{c}}{\partial \lambda_\ell} \langle u'(\varphi^T \lambda_{j-1}^n) \varphi \varphi^T \rangle^{-1} = \partial_{\mathbf{u}_\ell} \tilde{\mathbf{c}}_j$$

by the choice of  $\mathbf{B}$  as well as (14) and (17). Hence, we have that  $\partial_{\mathbf{u}} \mathbf{c} = \tilde{\mathbf{c}}_{\mathbf{u}}$ , which only has eigenvalues between  $-1$  and  $1$  by the assumption that the classical IPM iteration is contractive. Since  $\mathbf{J}$  is an upper triangular block matrix, the eigenvalues are given by  $\lambda(\partial_\lambda \mathbf{d}) = 0$  and  $\lambda(\partial_{\mathbf{u}} \mathbf{c}) \in (-1, 1)$ , hence the One-Shot IPM is contractive around its fix point.  $\square$

**Theorem 2.** *With the assumptions from Theorem 1, the One-Shot IPM converges locally, i.e. there exists a  $\delta > 0$  s.t. for all starting points  $(\lambda^0, \mathbf{u}^0) \in B_\delta(\lambda^*, \mathbf{u}^*)$  we have*

$$\|(\lambda^n, \mathbf{u}^n) - (\lambda^*, \mathbf{u}^*)\| \rightarrow 0 \quad \text{for } n \rightarrow \infty.$$

*Proof.* By Theorem 1, the One-Shot scheme is contractive at its fix point. Since we assumed convergence of the classical IPM scheme, we can conclude that all entries in the Jacobian  $\mathbf{J}$  are continuous functions. Furthermore, the determinant of  $\tilde{\mathbf{J}} := \mathbf{J} - \lambda \mathbf{I}_d$  is a polynomial of continuous functions, since

$$\det(\tilde{\mathbf{J}}) = \sum_{\sigma} \text{sgn}(\sigma) \prod_{i=1}^{2N_x(N+1)} \tilde{J}_{\sigma(i), i}.$$

Since the roots of a polynomial vary continuously with its coefficients, the eigenvalues of  $\mathbf{J}$  are continuous w.r.t  $(\lambda, \mathbf{u})$ . Hence there exists an open ball with radius  $\delta$  around the fix point in which the eigenvalues remain in the interval  $(-1, 1)$ .  $\square$

**Remark 3.** *Theorem 2 gives some insights in how to set up the One-Shot IPM iteration: If we look at the Jacobian (18) at an arbitrary point, we obtain*

*Let us assume that the classical IPM scheme converges globally.*

- *Choose a moment vector  $\mathbf{u}^0$  and compute the corresponding exact dual variables  $\boldsymbol{\lambda}^0 := \boldsymbol{\lambda}(\mathbf{u}^0)$  as starting vector. In this case, the term  $\partial_{\boldsymbol{\lambda}} \mathbf{d}$  is equal to zero for the starting conditions.*

## 5. Results

## 6. Summary and outlook

## References

- [1] Dongbin Xiu. Fast numerical methods for stochastic computations: a review. *Communications in computational physics*, 5(2-4):242–272, 2009.
- [2] AK Alekseev, IM Navon, and ME Zelentsov. The estimation of functional uncertainty using polynomial chaos and adjoint equations. *International Journal for numerical methods in fluids*, 67(3):328–341, 2011.
- [3] Gaël Poëtte, Bruno Després, and Didier Lucor. Uncertainty quantification for systems of conservation laws. *Journal of Computational Physics*, 228(7):2443–2467, 2009.
- [4] Jonas Kusch, Graham W Alldredge, and Martin Frank. Maximum-principle-satisfying second-order intrusive polynomial moment scheme. *arXiv preprint arXiv:1712.06966*, 2017.
- [5] Ilja Kröker and Christian Rohde. Finite volume schemes for hyperbolic balance laws with multiplicative noise. *Applied Numerical Mathematics*, 62(4):441–456, 2012.
- [6] Jan Giesselmann, Fabian Meyer, and Christian Rohde. A posteriori error analysis for random scalar conservation laws using the stochastic galerkin method. *arXiv preprint arXiv:1709.04351*, 2017.