

No. 660

March 2023

**Finite Element approximation of
data-driven problems in conductivity**

A. Müller, C. Meyer

ISSN: 2190-1767

FINITE ELEMENT APPROXIMATION OF DATA-DRIVEN PROBLEMS IN CONDUCTIVITY

ANNIKA MÜLLER AND CHRISTIAN MEYER

ABSTRACT. This paper is concerned with the finite element discretization of the data driven approach according to [18] for the solution of PDEs with a material law arising from measurement data. To simplify the setting, we focus on a scalar diffusion problem instead of a problem in elasticity. It is proven that the data convergence analysis from [9] carries over to the finite element discretization as long as $H(\text{div})$ -conforming finite elements such as the Raviart-Thomas element are used. As a corollary, minimizers of the discretized problems converge in data in the sense of [9], as the mesh size tends to zero and the approximation of the local material data set gets more and more accurate. We moreover present several heuristics for the solution of the discretized data driven problems, which is equivalent to a quadratic semi-assignment problem and therefore NP-hard. We test these heuristics by means of two examples and it turns out that the “classical” alternating projection method according to [18] is superior w.r.t. the ratio of accuracy and computational time.

1. INTRODUCTION

In material science, empirically developed material models are commonly in use, i.e., material laws that describe the behavior of materials are derived from measured data. But, due to measuring errors and simplified models, this approach bears the risk of inaccuracies. For this reason, an alternative data-driven concept has been established in [18]. In a sense, this concept skips the modeling step and uses the measured data directly. The idea is to select that data point from the set of measurements that best fits axiomatic physical laws such as first principles.

Let us explain this data-driven approach in terms of a stationary diffusion process of the form

$$-\text{div } \kappa(\nabla u) = f \text{ in } \Omega, \quad u = 0 \text{ on } \Gamma := \partial\Omega. \quad (1.1)$$

Here and in the following, $\Omega \subset \mathbb{R}^d$, $d \in \mathbb{N}$, is a bounded domain, $f \in H^{-1}(\Omega) := H_0^1(\Omega)^*$, and $\text{div} : L^2(\Omega; \mathbb{R}^d) \rightarrow H^{-1}(\Omega)$ denotes the distributional divergence. Furthermore, $\kappa : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a given function which models the material law and is calibrated by $m \in \mathbb{N}$ measurements for the tuple $(\mathbf{q}, \nabla u)$ collected in the so-called *local material data set*

$$\mathcal{D}^{\text{loc}} := \{(\mathbf{r}_1, \mathbf{w}_1), \dots, (\mathbf{r}_m, \mathbf{w}_m)\} \subset \mathbb{R}^d \times \mathbb{R}^d. \quad (1.2)$$

Date: February 28, 2023.

2010 Mathematics Subject Classification. 65N30, 65K10, 49M41.

Key words and phrases. Data driven models, Raviart Thomas finite elements, data convergence, proximal gradient method.

This research was supported by the German Research Foundation (DFG) under grant number ME 3281/10-1.

As indicated above, the idea is now to use these measurements directly. For this purpose, we rewrite (1.1) equivalently as

$$(1.1) \quad \iff \quad (\mathbf{q}, \nabla u) \in \tilde{\mathcal{D}} \times \mathcal{E},$$

with the so-called *equilibrium set*

$$\mathcal{E} := \{(\mathbf{q}, \nabla u) \in L^2(\Omega; \mathbb{R}^d) \times L^2(\Omega; \mathbb{R}^d) : u \in H_0^1(\Omega), -\operatorname{div} \sigma = f\} \quad (1.3)$$

and the *material law set*

$$\tilde{\mathcal{D}} := \{(\mathbf{r}, \mathbf{w}) \in L^2(\Omega; \mathbb{R}^d) \times L^2(\Omega; \mathbb{R}^d) : \mathbf{r}(x) = \kappa(\mathbf{w}(x)) \text{ a.e. in } \Omega\}. \quad (1.4)$$

The data-driven approach now skips the modelling step and uses the measured data directly by replacing $\tilde{\mathcal{D}}$ with the *material data set*

$$\mathcal{D} := \{(\mathbf{r}, \mathbf{w}) \in L^2(\Omega; \mathbb{R}^d) \times L^2(\Omega; \mathbb{R}^d) : (\mathbf{r}, \mathbf{w}) \in \mathcal{D}^{\text{loc}} \text{ a.e. in } \Omega\}$$

where \mathcal{D}^{loc} is the collection of measurements from (1.2). Due to measurement errors and limited measuring capacities, the intersection $\mathcal{D} \cap \mathcal{E}$ is usually empty. One therefore resorts to a minimization problem of the form

$$\left. \begin{array}{l} \min_{(y,z) \in Z \times Z} \quad \|y - z\|_Z^2 \\ \text{s.t.} \quad y \in \mathcal{D}, z \in \mathcal{E}, \end{array} \right\} \quad (\text{DDP})$$

i.e., one searches for two elements of the sets \mathcal{E} and \mathcal{D} that have smallest distance to each other. Herein, we abbreviated $Z := L^2(\Omega; \mathbb{R}^d) \times L^2(\Omega; \mathbb{R}^d)$.

The optimization problem (DDP) is frequently called *data-driven problem* and gives rise to several questions and issues: First of all, while \mathcal{E} is easily seen to be weakly closed, the set \mathcal{D} is in general not. Hence, (DDP) does not necessarily admit a solution. Moreover, a natural question arising in context of (DDP) is its behavior for measurements getting more and more accurate. Does the arising data-driven limit recover the “true” material law and, if so, in which sense? Moreover, a numerical solution of (DDP) requires a discretization of (DDP) and one may ask how a discretization influences this data-driven limit. Finally, the material data set \mathcal{D} involves a discrete point set such that (DDP) is a mixed-integer optimization problem. Problems of this type are typically hard to handle such that the development of efficient optimization algorithms for (DDP) (and its discretized counterpart) is all but trivial.

With this work, we address the two latter questions, i.e., we discuss discretization schemes and optimization algorithms for the solution of (DDP). Let us put our work into perspective. In engineering science, the data-driven approach is meanwhile well accepted and has been applied to various problems, in particular in solid mechanics, we only refer to [18, 19, 21, 12, 7] for examples from elasticity, inelasticity, dynamics, and fracture mechanics. A rigorous mathematical analysis of this approach has only been initiated recently in [9], where the concept of *data convergence* has been introduced. This notion of convergence is especially tailored to the structure of (DDP) and allows to characterize the data-driven limit. In this way, it answers the above question what happens, if the measurement errors tend to zero. The precise characterization of data-driven limits strongly depends on the particular structure of \mathcal{D} and \mathcal{E} . This notion of convergence and the characterization of the associated limits have been investigated for several scenarios, we exemplarily refer to in [9, 10, 23]. To the best of our knowledge however, the discretization of \mathcal{E} and \mathcal{D} has not been incorporated into this convergence analysis so far and with this

work, we aim to fill this gap. This is an important issue, not only because (DDP) cannot be solved in infinite dimensional spaces, but also due to the general lack of existence of solutions to (DDP). If one turns to a discretized counterpart of (DDP), then the finite dimensional structure allows to establish the existence of optimal solutions under mild assumptions so that, not until then, it makes sense to look for efficient algorithms for their computation.

As already indicated, the design of reliable and efficient solvers for (DDP) (and its discretization, respectively) is a delicate issue due to the discrete structure of the local material data set. In [18], a fixed-point type heuristic based on projections has been introduced, which is able to handle extensive measurement data, but may fail to converge or converges to spurious fixed-points that are not optimal, as shown in [17]. In the latter reference, a standard mixed-integer programming solver has been employed to solve a data-driven problem of (unrealistically) small size. Due to the vast amount of measurement data, it is in principle impossible to use exact mixed-integer programming solvers for the solution of (DDP). Therefore, various heuristics have been developed and applied such as kernel regression [15], local regression [16], tensor voting [13], and neural networks [22]. In the second part of the paper, we present some new heuristics and compare them with existing methods. Some of our algorithms are based on projection heuristic from [18], but we also tested a method, which employs an exact mixed-integer solver in combination with a local search algorithm.

We point out that we restrict ourselves to the conductivity example from (1.1) in order to keep the discussion as concise as possible. An extension of the finite element convergence analysis as well as the algorithmic approaches to problems in elasticity should be possible and is subject to future research.

The plan of the paper reads as follows: After introducing our standing assumptions and some well known results from saddle point theory in Section 2, we focus on the discretization of the equilibrium set by means of Raviart-Thomas type finite elements in Section 3. Afterwards, in Section 4, we recall the notion of data convergence from [9] and adapt it to our setting. Section 5 is then devoted to our main results, incorporating the finite element discretization of \mathcal{E} and \mathcal{D} into the data convergence analysis. In Section 6, we discuss the need for $H(\text{div})$ -conforming finite elements like the Raviart-Thomas element for the discretization of (DDP). Section 7 is dedicated to the algorithms and their implementation and finally, in Section 8, we present some numerical results.

2. PRELIMINARIES AND STANDING ASSUMPTIONS

As usual we define

$$H(\text{div}) := \{\mathbf{w} \in L^2(\Omega; \mathbb{R}^d) : \text{div } \mathbf{w} \in L^2(\Omega)\},$$

where $\text{div} : L^2(\Omega; \mathbb{R}^d) \rightarrow H^{-1}(\Omega)$ denotes the distributional divergence. The following two lemmas concerning the space $H(\text{div})$ will be useful in the rest of the paper. For their proofs, we refer to [25, ...].

Lemma 2.1. *If the complement of $\bar{\Omega}$ satisfies the cone condition according to [?], then $C^\infty(\bar{\Omega}; \mathbb{R}^d)$ is dense in $H(\text{div})$.*

Lemma 2.2. *Let $F \in H(\operatorname{div})^*$ and $f \in L^2(\Omega)$ be given. Then there exists a unique solution $(\boldsymbol{\tau}, \lambda) \in H(\operatorname{div}) \times L^2(\Omega)$ to the saddle point problem*

$$\int_{\Omega} \boldsymbol{\tau} \cdot \boldsymbol{w} \, dx + \int_{\Omega} \lambda \operatorname{div} \boldsymbol{w} \, dx = \langle F, \boldsymbol{w} \rangle \quad \forall \boldsymbol{w} \in H(\operatorname{div}), \quad (2.1a)$$

$$- \int_{\Omega} v \operatorname{div} \boldsymbol{\tau} \, dx = \int_{\Omega} f v \, dx \quad \forall v \in L^2(\Omega). \quad (2.1b)$$

If Ω satisfies the regularity assumptions from Lemma 2.1 and $F \in L^2(\Omega; \mathbb{R}^d)$, then $\lambda \in H_0^1(\Omega)$. If moreover Ω is H^2 -regular and $F \in H^1(\Omega; \mathbb{R}^d)$, then $\boldsymbol{\tau} \in H^1(\Omega; \mathbb{R}^d)$ and

$$\|\boldsymbol{\tau}\|_{H^1(\Omega; \mathbb{R}^d)} \leq c(\|f\|_{L^2(\Omega)} + \|F\|_{H^1(\Omega; \mathbb{R}^d)})$$

with a constant $c > 0$, which only depends on Ω .

Assumption 2.3. *Throughout this paper, we assume the following regularity assumptions on the domain and the inhomogeneity f :*

- (i) *The domain Ω is supposed to be bounded and convex with a polygonal resp. polyhedral boundary.*
- (ii) *The right-hand side in the divergence constraint satisfies $f \in L^2(\Omega)$.*

Remark 2.4. The regularity of the domain can be relaxed in the sense that we can drop the convexity, see Remark 3.9 below. In contrast to this, we need the higher regularity of f and cannot work with inhomogeneities in $H^{-1}(\Omega)$, since our discretization requires $H(\operatorname{div})$ -conforming finite elements as demonstrated in Section 6 below.

We underline that Assumption 2.3 is tacitly supposed to hold throughout the rest of the paper without mentioning it every time.

3. DISCRETIZATION

The equilibrium constraint set \mathcal{E} from (1.3) is discretized by

$$\mathcal{E}_h := \{(\boldsymbol{q}_h, \nabla u_h) : \boldsymbol{q}_h \in Q_h, u_h \in U_h, -\operatorname{div}_h \boldsymbol{q}_h = f\}, \quad (3.1)$$

where U_h and Q_h are finite dimensional function spaces satisfying the following

Assumption 3.1. *The discretization of \mathcal{E} in (3.1) is supposed to fulfill the following conditions:*

- (i) *For all $h > 0$, the discrete spaces U_h and Q_h are conforming, i.e., they are finite dimensional linear subspaces of $H_0^1(\Omega)$ and $H(\operatorname{div})$. Moreover, $\bigcup_{h>0} U_h$ is dense in $H_0^1(\Omega)$ w.r.t. the $H^1(\Omega)$ -norm.*
- (ii) *We set*

$$V_h := \operatorname{div}(Q_h) \subset L^2(\Omega). \quad (3.2)$$

Then the discrete divergence $\operatorname{div}_h : H(\operatorname{div}) \rightarrow V_h^$ from (3.1) is defined by*

$$\langle \operatorname{div}_h \boldsymbol{\tau}, v_h \rangle := \int_{\Omega} v_h \operatorname{div} \boldsymbol{\tau} \, dx, \quad \boldsymbol{\tau} \in H(\operatorname{div}), v_h \in V_h.$$

- (iii) *There exists an interpolation operator $\Pi_h : H^1(\Omega; \mathbb{R}^d) \rightarrow Q_h$ such that, for all $\boldsymbol{\tau} \in H^1(\Omega; \mathbb{R}^d)$ there holds*

$$\operatorname{div}_h \boldsymbol{\tau} = \operatorname{div}_h \Pi_h \boldsymbol{\tau} \quad (3.3)$$

and

$$\Pi_h \boldsymbol{\tau} \rightarrow \boldsymbol{\tau} \quad \text{in } L^2(\Omega; \mathbb{R}^d) \text{ as } h \searrow 0. \quad (3.4)$$

Note that, in view of Assumption 3.1(ii), the constraint $-\operatorname{div}_h \mathbf{q}_h = f$ in the definition of \mathcal{E}_h is short for

$$-\int_{\Omega} v_h \operatorname{div} \mathbf{q}_h \, dx = \int_{\Omega} f v_h \, dx \quad \forall v_h \in V_h.$$

Assumption 3.1 implies the well known LBB-condition, as we shortly sketch in the following (by arguments analogous to the discussion after [4, Lemma 5.4]). To this end, let $\mathbf{w} \in H(\operatorname{div})$ be arbitrary and solve (2.1) with $(F, f) = (0, -\operatorname{div} \mathbf{w}) \in H^1(\Omega; \mathbb{R}^d) \times L^2(\Omega)$. Now, since Ω is supposed to be convex and thus H^2 -regular, we obtain a solution $(\boldsymbol{\tau}, \lambda)$ with $\boldsymbol{\tau} \in H^1(\Omega; \mathbb{R}^d)$. Thus we can apply the interpolation operator Π_h from Assumption 3.1(iii) and obtain

$$\operatorname{div}_h \mathbf{w} = \operatorname{div}_h \boldsymbol{\tau} = \operatorname{div}_h \Pi_h \boldsymbol{\tau}. \quad (3.5)$$

By construction, the mapping $H(\operatorname{div}) \ni \mathbf{w} \mapsto \Pi_h \boldsymbol{\tau} \in Q_h$ is linear and continuous and, in view of (3.5), it is a Fortin interpolation operator. Therefore, by [4, 4.8 Fortin's Criterion], the tuple (Q_h, V_h) satisfies the LBB-condition, i.e., we have shown the following:

Corollary 3.2. *Under Assumption 3.1, (Q_h, V_h) satisfies the LBB-condition, i.e., there exists a constant $\beta > 0$, independent of $h > 0$, such that*

$$\inf_{v_h \in V_h} \sup_{\mathbf{w}_h \in Q_h} \frac{\int_{\Omega} v_h \operatorname{div} \mathbf{w}_h \, dx}{\|\mathbf{w}_h\|_{H(\operatorname{div})} \|v_h\|_{L^2(\Omega)}} \geq \beta. \quad (3.6)$$

Based on Assumption 3.1(i)–(ii), the standard theory for mixed finite elements yields the following lemma. For the corresponding proof, we refer to [4, Section 4].

Lemma 3.3. *Let Assumption 3.1 be fulfilled. Then, the following is valid:*

- (i) *For all $\mathbf{q}_h, \boldsymbol{\tau}_h \in Q_h$ satisfying $\operatorname{div}_h \mathbf{q}_h = \operatorname{div}_h \boldsymbol{\tau}_h$, it holds that $\operatorname{div} \mathbf{q}_h = \operatorname{div} \boldsymbol{\tau}_h$.*
- (ii) *For every $F \in H(\operatorname{div})^*$ and $f \in L^2(\Omega)$ there exists a unique solution $(\boldsymbol{\tau}_h, \lambda_h) \in Q_h \times V_h$ to the saddle point problem*

$$\int_{\Omega} \boldsymbol{\tau}_h \cdot \mathbf{w}_h \, dx + \int_{\Omega} \lambda_h \operatorname{div} \mathbf{w}_h \, dx = \langle F, \mathbf{w}_h \rangle \quad \forall \mathbf{w}_h \in Q_h \quad (3.7a)$$

$$-\int_{\Omega} v_h \operatorname{div} \boldsymbol{\tau}_h \, dx = \int_{\Omega} f v_h \, dx \quad \forall v_h \in V_h. \quad (3.7b)$$

- (iii) *This solution satisfies the following best approximation result*

$$\|\boldsymbol{\tau} - \boldsymbol{\tau}_h\|_X \leq 2 \inf\{\|\boldsymbol{\tau} - \mathbf{w}_h\|_X : \mathbf{w}_h \in Q_h, -\operatorname{div}_h \mathbf{w}_h = f\}, \quad (3.8)$$

where $\boldsymbol{\tau}$ is the solution of (2.1) and $X = H(\operatorname{div})$ or $X = L^2(\Omega; \mathbb{R}^d)$.

- (iv) *There exists a constant $C > 0$ depending only on the LBB-constant such that the solution of (3.7) satisfies.*

$$\|\boldsymbol{\tau}_h\|_{H(\operatorname{div})} + \|\lambda_h\|_{L^2(\Omega)} \leq C(\|F\|_{H(\operatorname{div})^*} + \|f\|_{L^2(\Omega)}). \quad (3.9)$$

Moreover, the following best approximation result holds true

$$\|\boldsymbol{\tau} - \boldsymbol{\tau}_h\|_{H(\operatorname{div})} \leq 2(1 + C) \inf_{\mathbf{w}_h \in Q_h} \|\boldsymbol{\tau} - \mathbf{w}_h\|_{H(\operatorname{div})}, \quad (3.10)$$

where $\boldsymbol{\tau}$ again denotes the solution of (2.1).

Remark 3.4. Note that the assertions of Lemma 3.3(i)–(iii) also hold without the LBB-condition, i.e., without the existence of the interpolation operator from Assumption 3.1(iii).

Lemma 3.5. *Let Assumption 3.1 hold and let a sequence $\{(\mathbf{q}_h, \nabla u_h)\}_{h>0}$ with $(\mathbf{q}_h, \nabla u_h) \in \mathcal{E}_h$ be given. Then, there exists a sequence $\{(\hat{\mathbf{q}}_h, \nabla \hat{u}_h)\}_{h>0} \subset \mathcal{E}$ such that*

$$\|(\hat{\mathbf{q}}_h - \mathbf{q}_h, \nabla(\hat{u}_h - \nabla u_h))\|_{L^2(\Omega; \mathbb{R}^d)^2} \rightarrow 0 \quad \text{as } h \searrow 0. \quad (3.11)$$

Proof. We consider the system (2.1) with $(F, f) = (0, f)$ and denote the corresponding solution by $(\boldsymbol{\tau}, \lambda) \in H(\text{div}) \times L^2(\Omega)$. Moreover, we solve (3.7) with the same right hand side and denote this solution by $(\boldsymbol{\tau}_h, \lambda_h) \in Q_h \times V_h$. Furthermore, we set

$$\hat{\mathbf{q}}_h := \mathbf{q}_h + \boldsymbol{\tau} - \boldsymbol{\tau}_h \in H(\text{div}).$$

Due to Lemma 3.3(i), there holds $\text{div}(\mathbf{q}_h - \boldsymbol{\tau}_h) = 0$, and, consequently, $\text{div} \hat{\mathbf{q}}_h = \text{div} \boldsymbol{\tau} = -f$ and therefore, $(\hat{\mathbf{q}}_h, \nabla u_h) \in \mathcal{E}$ by the conformity of U_h by Assumption 3.1(i).

Now, since Ω is convex and $f \in L^2(\Omega)$, Lemma 2.2 implies $\boldsymbol{\tau} \in H^1(\Omega; \mathbb{R}^d)$ with $\|\boldsymbol{\tau}\|_{H^1(\Omega; \mathbb{R}^d)} \leq c\|f\|_{L^2(\Omega)}$. Thus, Lemma 3.3(iii) gives

$$\|\hat{\mathbf{q}}_h - \mathbf{q}_h\|_{L^2(\Omega; \mathbb{R}^d)} = \|\boldsymbol{\tau} - \boldsymbol{\tau}_h\|_{L^2(\Omega; \mathbb{R}^d)} \leq 2\|\boldsymbol{\tau} - \Pi_h \boldsymbol{\tau}\|_{L^2(\Omega; \mathbb{R}^d)} \rightarrow 0 \quad \text{as } h \searrow 0.$$

Therefore, if we set $\hat{u}_h := u_h$, we obtain $(\hat{\mathbf{q}}_h, \nabla \hat{u}_h) \in \mathcal{E}$ and (3.11). \square

Lemma 3.6. *Let Assumption 3.1 be fulfilled and let $(\mathbf{q}, \nabla u) \in \mathcal{E}$ be given. Then there is a sequence $\{(\mathbf{q}_h, \nabla u_h)\}_{h>0} \subset H(\text{div}) \times L^2(\Omega; \mathbb{R}^d)$ such that $(\mathbf{q}_h, \nabla u_h) \in \mathcal{E}_h$ and*

$$(\mathbf{q}_h, \nabla u_h) \rightarrow (\mathbf{q}, \nabla u) \quad \text{in } L^2(\Omega; \mathbb{R}^d) \times L^2(\Omega; \mathbb{R}^d) \quad \text{as } h \searrow 0.$$

Proof. Let $\varepsilon > 0$ be arbitrary. By Lemma 2.1, there is a function $\mathbf{q}_\varepsilon \in C^\infty(\bar{\Omega}; \mathbb{R}^d)$ such that

$$\|\mathbf{q} - \mathbf{q}_\varepsilon\|_{H(\text{div})} \leq \min\{1, C^{-1}\} \frac{\varepsilon}{3}, \quad (3.12)$$

where $C > 0$ is the constant from Lemma 3.3(iv). Since \mathbf{q}_ε is smooth, we are allowed to apply Π_h , which yields

$$\|\mathbf{q}_\varepsilon - \Pi_h \mathbf{q}_\varepsilon\|_{L^2(\Omega; \mathbb{R}^d)} \leq \frac{\varepsilon}{3}$$

provided that $h > 0$ is chosen sufficiently small. Define now $f_\varepsilon := -\text{div} \mathbf{q}_\varepsilon$ and denote the solution of (3.7) with right hand side $(0, f - f_\varepsilon)$ by $(\boldsymbol{\tau}_h^\varepsilon, \lambda_h^\varepsilon)$. Then we set

$$\mathbf{q}_h^\varepsilon := \Pi_h \mathbf{q}_\varepsilon + \boldsymbol{\tau}_h^\varepsilon \in Q_h.$$

Then, (3.3) implies for every $v_h \in V_h$ that

$$\begin{aligned} \int_{\Omega} \text{div} \mathbf{q}_h^\varepsilon v_h \, dx &= \int_{\Omega} \text{div}(\Pi_h \mathbf{q}_\varepsilon) v_h \, dx + \int_{\Omega} \text{div} \boldsymbol{\tau}_h^\varepsilon v_h \, dx \\ &= \int_{\Omega} \text{div} \mathbf{q}_\varepsilon v_h \, dx - \int_{\Omega} (f - f_\varepsilon) v_h \, dx = - \int_{\Omega} f v_h \, dx, \end{aligned}$$

i. e., $-\text{div}_h \mathbf{q}_h^\varepsilon = f$. According to Lemma 3.3(iv), we deduce from (3.12) that

$$\|\boldsymbol{\tau}_h^\varepsilon\|_{L^2(\Omega; \mathbb{R}^d)} \leq C\|f - f_\varepsilon\|_{L^2(\Omega)} = C\|\text{div} \mathbf{q} - \text{div} \mathbf{q}_\varepsilon\|_{L^2(\Omega)} \leq \frac{\varepsilon}{3}.$$

Altogether, we obtain

$$\|\mathbf{q} - \mathbf{q}_h^\varepsilon\|_{L^2(\Omega; \mathbb{R}^d)} \leq \|\mathbf{q} - \mathbf{q}_\varepsilon\|_{L^2(\Omega; \mathbb{R}^d)} + \|\mathbf{q}_\varepsilon - \Pi_h \mathbf{q}_\varepsilon\|_{L^2(\Omega; \mathbb{R}^d)} + \|\boldsymbol{\tau}_h^\varepsilon\|_{L^2(\Omega; \mathbb{R}^d)} \leq \varepsilon.$$

Finally, since $\bigcup_{h>0} U_h$ is dense in $H_0^1(\Omega)$ by assumption, there exist $h > 0$ and $u_h \in U_h$ such that $\|\nabla u - \nabla u_h\|_{L^2(\Omega; \mathbb{R}^d)} \leq \varepsilon$. As $\varepsilon > 0$ was arbitrary, this proves the claim. \square

Proposition 3.7. *Let $\{\mathcal{T}_h\}_{h>0}$ be a family of shape regular triangulations of Ω according to [4, Definition 5.1]. Then the Raviart-Thomas space of order $k \in \mathbb{N} \cup \{0\}$ given by*

$$\mathcal{RT}_k(\mathcal{T}_h) := \{\mathbf{w} \in H(\operatorname{div}) : \mathbf{w}|_T \in \mathcal{RT}_k(T) \forall T \in \mathcal{T}_h\}$$

with $\mathcal{RT}_k(T) := \mathcal{P}_k(T)^d + x \mathcal{P}_k(T)$, where $\mathcal{P}_k(T)$ denotes the space of polynomials of order k on T , is a feasible choice for Q_h fulfilling Assumption 3.1.

For U_h one can choose the classical finite element space

$$U_h := \{u \in C(\overline{\Omega}) \cap H_0^1(\Omega) : u|_T \in \mathcal{P}_k(T) \forall T \in \mathcal{T}_h\}$$

in order to fulfill Assumption 3.1.

Proof. The conformity of $\mathcal{RT}_k(\mathcal{T}_h)$ and U_h is already part of their definition. The density of $\bigcup_{h>0} U_h$ in $H_0^1(\Omega)$ follows by smooth approximation and standard interpolation error estimates. In case of Raviart-Thomas finite elements, the space $V_h = \operatorname{div}(Q_h)$ equals $\mathcal{P}_k(\mathcal{T}_h) := \{v \in L^2(\Omega) : v|_T \in \mathcal{P}_k(T) \forall T \in \mathcal{T}_h\}$, see e.g. [11, Lemma 3.5]. The existence of an interpolation operator Π_h fulfilling Assumption 3.1(iii) is established in [11, Theorem 3.1, Lemma 3.5]. \square

Remark 3.8. There are several other elements satisfying Assumption 3.1, for instance the \mathcal{BDM} -element or the Raviart-Thomas element on quadrilateral meshes. We refer to [11] and the references therein.

Remark 3.9. The regularity assumptions on Ω in Assumption 2.3 can be relaxed. In fact, it is sufficient to require that Ω is polygonally resp. polyhedrally bounded, i.e., we can drop the convexity of Ω . This is due to the fact that convexity is only needed for the regularity of the solution of the saddle point problem (2.1) for the construction of Fortin's interpolation operator for Corollary 3.2 and for the solution $(\boldsymbol{\tau}, \lambda)$ in the proof of 3.5. In both cases however, one can resort to a larger convex domain B containing Ω and solve the continuous saddle point problem there such that the regularity result from Lemma 2.2 applies. The function $\boldsymbol{\tau}_h$ in the proof of Lemma 3.5 is then defined on B and for this reason, one needs to assume that the meshes can be extended in a shape regular way from Ω to B so that the results of Lemma 3.3 also hold on B instead of Ω . In order to avoid these technical issues, we restrict ourselves to the case of a convex domain Ω .

4. DATA TOPOLOGY

Let us recall the concept of data convergence, which was first introduced in [9]. It represents an intermediate convergence between weak and strong convergence and is especially tailored to the structure of the data driven problem (DDP).

Definition 4.1 (Data convergence). Let Z be a reflexive, separable Banach space. A sequence $\{(y_k, z_k)\}_{k \in \mathbb{N}}$ in $Z \times Z$ is said to converge to $(y, z) \in Z \times Z$ in the data topology, denoted $(y, z) = \Delta\text{-}\lim_{k \rightarrow \infty} (y_k, z_k)$, if

$$y_k \rightharpoonup y, \quad z_k \rightharpoonup z \quad \text{and} \quad y_k - z_k \rightarrow y - z \quad \text{in } Z.$$

The concept of data convergence can be transferred to sets.

Definition 4.2 (Data convergence of sets). Let Z be a reflexive, separable Banach space and $\mathcal{D}, \mathcal{D}_k, \mathcal{E}, \mathcal{E}_k \subset Z$, $k \in \mathbb{N}$. We write $\mathcal{D} \times \mathcal{E} = \Delta\text{-}\lim_{k \rightarrow \infty} (\mathcal{D}_k \times \mathcal{E}_k)$, if

(DC1) for each $(y, z) \in \mathcal{D} \times \mathcal{E}$ there is a sequence $\{(y_k, z_k)\}_{k \in \mathbb{N}}$ with $(y_k, z_k) \in \mathcal{D}_k \times \mathcal{E}_k$ for each $k \in \mathbb{N}$ such that $(y, z) = \Delta\text{-}\lim_{k \rightarrow \infty} (y_k, z_k)$,

(DC2) for each sequence $\{(y_j, z_j)\}_{j \in \mathbb{N}}$ with $(y_j, z_j) \in \mathcal{D}_{k_j} \times \mathcal{E}_{k_j}$ for each $j \in \mathbb{N}$, $\{k_j\}_{j \in \mathbb{N}}$ strictly monotonically increasing, and $(y, z) = \Delta\text{-}\lim_{j \rightarrow \infty} (y_j, z_j)$ it holds that $(y, z) \in \mathcal{D} \times \mathcal{E}$.

Note that the above definition of data convergence of sets corresponds to Kuratowski convergence of sets with respect to data convergence. The notion of data convergence of sets is especially well suited to the approximation of data-driven problems of the form (DDP), as the following proposition shows. Its proof is along the lines of [9, Theorem 3.2], where the equilibrium set \mathcal{E} is fixed. Here we additionally consider the approximation of \mathcal{E} is with a sequence of sets \mathcal{E}_k . Though the proof is a straightforward adaptation of the one in [9], we present it for convenience of the reader.

Proposition 4.3. *Let Z be a reflexive and separable Banach space and suppose that subsets $\mathcal{D}, \mathcal{E} \subset Z$ and sequences of subsets $\{\mathcal{D}_k\}_{k \in \mathbb{N}}, \{\mathcal{E}_k\}_{k \in \mathbb{N}}, \mathcal{D}_k, \mathcal{E}_k \subset Z$ for all $k \in \mathbb{N}$, are given such that*

$$\mathcal{D} \times \mathcal{E} = \Delta\text{-}\lim_{j \rightarrow \infty} (\mathcal{D}_k \times \mathcal{E}_k). \quad (4.1)$$

Assume moreover that there are constants $c > 0$ and $b \geq 0$, independent of $k \in \mathbb{N}$, such that, for all $k \in \mathbb{N}$,

$$\|y - z\|_Z \geq c(\|y\|_Z + \|z\|_Z) - b \quad \forall (y, z) \in \mathcal{D}_k \times \mathcal{E}_k. \quad (4.2)$$

Furthermore, define $F_k : Z \times Z \rightarrow [0, \infty]$ by

$$F_k(y, z) := I_{\mathcal{D}_k}(y) + I_{\mathcal{E}_k}(z) + \|y - z\|_Z^2,$$

where $I_{\mathcal{D}_k} : Z \rightarrow \{0, \infty\}$ is the indicator functional of \mathcal{D}_k , i.e.,

$$I_{\mathcal{D}_k}(y) := \begin{cases} 0, & y \in \mathcal{D}_k, \\ \infty, & y \notin \mathcal{D}_k \end{cases}$$

and $I_{\mathcal{E}_k}$ is defined analogously. Then, the following is valid:

- (a) If $F_k(y_k, z_k) \rightarrow 0$, there exists $z \in \mathcal{D} \cap \mathcal{E}$ such that, up to subsequences, $(z, z) = \Delta\text{-}\lim_{k \rightarrow \infty} (y_k, z_k)$;
- (b) If $z \in \mathcal{D} \cap \mathcal{E}$, there exists a sequence $\{(y_k, z_k)\}_{k \in \mathbb{N}}$ in $Z \times Z$ such that $(z, z) = \Delta\text{-}\lim_{k \rightarrow \infty} (y_k, z_k)$ and $F_k(y_k, z_k) \rightarrow 0$.

Proof. ad (a): Let $F_k(y_k, z_k) \rightarrow 0$. Then, it follows that $y_k \in \mathcal{D}_k, z_k \in \mathcal{E}_k$ for k sufficiently large and $\|y_k - z_k\| \rightarrow 0$ as $k \rightarrow \infty$. By (4.2), $\{y_k\}_{k \in \mathbb{N}}$ and $\{z_k\}_{k \in \mathbb{N}}$ are bounded. Therefore, there are subsequences $\{y_{k_j}\}_{j \in \mathbb{N}}$ and $\{z_{k_j}\}_{j \in \mathbb{N}}$ and $y \in Z$ and $z \in Z$ such that $y_{k_j} \rightarrow y$ and $z_{k_j} \rightarrow z$. By weak lower-semicontinuity of the norm, we have that

$$0 \leq \|y - z\|_Z \leq \liminf_{j \rightarrow \infty} \|y_{k_j} - z_{k_j}\|_Z = 0.$$

Hence $y = z$ and $(z, z) = \Delta\text{-}\lim_{j \rightarrow \infty} (y_{k_j}, z_{k_j})$ and therefore, (4.1) yields $z \in \mathcal{D} \cap \mathcal{E}$ as claimed.

ad (b): Let $z \in \mathcal{D} \cap \mathcal{E}$ be given. Then, thanks to (4.1), there exists a sequence $\{(y_k, z_k)\}_{k \in \mathbb{N}}$ with

$$(y_k, z_k) \in \mathcal{D}_k \times \mathcal{E}_k \quad \text{and} \quad (z, z) = \Delta\text{-}\lim_{k \rightarrow \infty} (y_k, z_k).$$

This in particular implies $y_k - z_k \rightarrow z - z = 0$ and hence, by continuity of the norm,

$$\lim_{k \rightarrow \infty} F_k(y_k, z_k) = \lim_{k \rightarrow \infty} (I_{\mathcal{D}_k}(y_k) + I_{\mathcal{E}_k}(z_k) + \|y_k - z_k\|_Z^2) = 0,$$

as required. \square

Proposition 4.3 shows that, if a sequence of sets $\{(\mathcal{D}_k, \mathcal{E}_k)\}_{k \in \mathbb{N}}$ satisfies (4.2) and more importantly (4.1), then the data-driven problem with limit sets \mathcal{D} and \mathcal{E} admits a solution, which can be approximated (w.r.t. data convergence) with solutions of the respective data-driven problems subject to the sets \mathcal{D}_k and \mathcal{E}_k . The crucial question is of course now, which (sequences of sets) satisfy (4.1). This will be answered for our conductivity example in the following section.

5. CONVERGENCE RESULTS

As in [9, Theorem 3.3], we aim at giving sufficient conditions under which the assumptions of Proposition 4.3 are fulfilled. Recall again the setting in our conductivity example, where

$$Z = L^2(\Omega; \mathbb{R}^d) \times L^2(\Omega; \mathbb{R}^d) \quad (5.1)$$

and

$$\mathcal{E} = \{(\mathbf{q}, \nabla u) \in L^2(\Omega; \mathbb{R}^d) \times L^2(\Omega; \mathbb{R}^d) : u \in H_0^1(\Omega), -\operatorname{div} \mathbf{q} = f\}. \quad (5.2)$$

For the approximation of \mathcal{E} , we choose the discretized equilibrium constraint sets \mathcal{E}_{h_k} from (3.1). The following theorem shows that such a discretization can be included in the convergence analysis of [9, Theorem 3.3].

Theorem 5.1. *Let Z and \mathcal{E} be given as in (5.1) and (5.2), respectively, and assume that a global material data set $\mathcal{D} \subset Z$ and approximations thereof, denoted by $\mathcal{D}_k \subset Z$, $k \in \mathbb{N}$, are given. Suppose moreover the following to hold:*

- (i) (Data closure) $\overline{\mathcal{D}} \times \mathcal{E} = \overline{\mathcal{D} \times \mathcal{E}}^\Delta$, i. e., $\overline{\mathcal{D}} \times \mathcal{E}$ is the closure of $\mathcal{D} \times \mathcal{E}$ w.r.t. data convergence;
- (ii) (Fine approximation) For each $\xi \in \mathcal{D}$, there is a sequence $\{\xi_k\}_{k \in \mathbb{N}}$ with $\xi_k \in \mathcal{D}_k$ for all $k \in \mathbb{N}$ such that $\xi_k \rightarrow \xi$ as $k \rightarrow \infty$;
- (iii) (Uniform approximation) There is a sequence $\{t_k\}_{k \in \mathbb{N}} \subset \mathbb{R}_{>0}$ with $t_k \searrow 0$ such that

$$d(\xi, \mathcal{D}) := \inf_{y \in \mathcal{D}} \|y - \xi\|_Z \leq t_k \quad \forall \xi \in \mathcal{D}_k;$$

- (iv) (Transversality) There are constants $c > 0$ and $b \geq 0$ such that, for all $y \in \mathcal{D}$ and $z \in \mathcal{E}$,

$$\|y - z\|_Z \geq c(\|y\|_Z + \|z\|_Z) - b;$$

- (v) (Conforming discretization) There is a monotonically decreasing sequence of mesh sizes $\{h_k\}_{k \in \mathbb{N}} \subset \mathbb{R}_{>0}$ with $h_k \searrow 0$ as $k \rightarrow \infty$ such that the discrete spaces $Q_k := Q_{h_k}$ and $U_k := U_{h_k}$ from the discrete equilibrium set $\mathcal{E}_k := \mathcal{E}_{h_k}$ in (3.1) satisfy Assumption 3.1.

Then the assumptions of Proposition 4.3 are fulfilled, i. e.,

- (a) (Data convergence) $\overline{\mathcal{D}} \times \mathcal{E} = \Delta\text{-}\lim_{k \rightarrow \infty} (\mathcal{D}_k \times \mathcal{E}_k)$;
- (b) (Equi-transversality) There are constants $c > 0$ and $b \geq 0$ such that, for all $k \in \mathbb{N}$ and all $(y, z) \in \mathcal{D}_k \times \mathcal{E}_k$, there holds

$$\|y - z\|_Z \geq c(\|y\|_Z + \|z\|_Z) - b.$$

Proof. ad (a), condition (DC1): Let $(y, z) \in \overline{\mathcal{D}} \times \mathcal{E}$ be fixed but arbitrary. Our goal is to find a sequence $\{(y_k^*, z_k^*)\}_{k \in \mathbb{N}}$ with $(y_k^*, z_k^*) \in \mathcal{D}_k \times \mathcal{E}_k$ such that $(y, z) = \Delta\text{-}\lim_{k \rightarrow \infty} (y_k^*, z_k^*)$. By (i), there is a sequence $\{(\hat{y}_n, \hat{z}_n)\}_{n \in \mathbb{N}} \subset \mathcal{D} \times \mathcal{E}$ such that $(y, z) = \Delta\text{-}\lim_{n \rightarrow \infty} (\hat{y}_n, \hat{z}_n)$. Due to (ii) and Lemma 3.6, for each $n \in \mathbb{N}$, there are sequences $\{y_{n,k}\}_{k \in \mathbb{N}}$ with $y_{n,k} \in \mathcal{D}_k$ and $\{z_{n,k}\}_{k \in \mathbb{N}}$ with $z_{n,k} \in \mathcal{E}_k$ and a finite number $m_n \in \mathbb{N}$ with $m_n \geq m_{n-1} + 1$ such that

$$\|y_{n,k} - \hat{y}_n\|_Z < \frac{1}{n} \quad \text{and} \quad \|z_{n,k} - \hat{z}_n\|_Z < \frac{1}{n} \quad \forall k \geq m_n.$$

This of course gives rise to a diagonal sequence $\{y_{n,m_n}, z_{n,m_n}\}$ with the desired properties, but, for each $(y, z) \in \overline{\mathcal{D}} \times \mathcal{E}$, one obtains a different sequence $\{m_n\}_{n \in \mathbb{N}}$ with different approximations \mathcal{D}_{m_n} and discretizations \mathcal{E}_{m_n} . To overcome this issue, let us define

$$\{(y_k^*, z_k^*)\}_{k \in \mathbb{N}} := \underbrace{\{(\hat{y}_1, \hat{z}_1), \dots, (\hat{y}_1, \hat{z}_1)\}}_{(m_1 - 1)\text{-times}}, \underbrace{\{(\hat{y}_1, \hat{z}_1), \dots, (\hat{y}_1, \hat{z}_1), (\hat{y}_2, \hat{z}_2), \dots, (\hat{y}_2, \hat{z}_2)\}}_{(m_2 - m_1)\text{-times}}, \dots$$

as well as

$$\begin{aligned} \{(y_k^*, z_k^*)\}_{k \in \mathbb{N}} := & \{(y_{1,1}, z_{1,1}), \dots, (y_{1,m_1-1}, z_{1,m_1-1}), \\ & (y_{1,m_1}, z_{1,m_1}), \dots, (y_{1,m_2-1}, z_{1,m_2-1}), \\ & (y_{2,m_2}, z_{2,m_2}), \dots, (y_{2,m_3-1}, z_{2,m_3-1}), \dots\}. \end{aligned}$$

Then, by construction, $(y_k^*, z_k^*) \in \mathcal{D}_k \times \mathcal{E}_k$ for all $k \in \mathbb{N}$. Moreover, we have

$$(y, z) = \Delta\text{-}\lim_{k \rightarrow \infty} (y_k^*, z_k^*) \tag{5.3}$$

and, since, for each $n \in \mathbb{N}$ and all $k \geq m_n$, it holds

$$\|\hat{y}_k^* - y_k^*\|_Z \leq \frac{1}{n} \quad \text{and} \quad \|\hat{z}_k^* - z_k^*\|_Z \leq \frac{1}{n},$$

we obtain

$$\|\hat{y}_k^* - y_k^*\|_Z \rightarrow 0 \quad \text{and} \quad \|\hat{z}_k^* - z_k^*\|_Z \rightarrow 0 \quad \text{as } k \rightarrow \infty. \tag{5.4}$$

By the definition of data convergence, (5.3) and (5.4) yield $y_k^* \rightarrow y$, $z_k^* \rightarrow z$, and

$$\|y_k^* - z_k^* - (y - z)\|_Z \leq \|y_k^* - \hat{y}_k^*\|_Z + \|\hat{y}_k^* - \hat{z}_k^* - (y - z)\|_Z + \|\hat{z}_k^* - z_k^*\|_Z \rightarrow 0,$$

which is nothing else than

$$(y, z) = \Delta\text{-}\lim_{k \rightarrow \infty} (y_k^*, z_k^*)$$

with $(y_k^*, z_k^*) \in \mathcal{D}_k \times \mathcal{E}_k$ for all $k \in \mathbb{N}$. Since $(y, z) \in \overline{\mathcal{D}} \times \mathcal{E}$ was arbitrary, this implies (DC1).

ad (a), condition (DC2): Suppose that $(y, z) = \Delta\text{-}\lim_{j \rightarrow \infty} (y_j, z_j)$ in $Z \times Z$ with $(y_j, z_j) \in \mathcal{D}_{k_j} \times \mathcal{E}_{k_j}$ for all $j \in \mathbb{N}$ and a strictly monotonically increasing sequence $\{k_j\}_{j \in \mathbb{N}}$. We need to prove $(y, z) \in \overline{\mathcal{D}} \times \mathcal{E}$. By (iii) and Lemma 3.5, there exist $\hat{y}_j \in \mathcal{D}$ and $\hat{z}_j \in \mathcal{E}$ such that

$$\|\hat{y}_j - y_j\|_Z \leq t_{k_j} \quad \text{and} \quad \|\hat{z}_j - z_j\|_Z \rightarrow 0 \quad \text{as } j \rightarrow \infty.$$

Consequently, $\hat{y}_j \rightarrow y$, $\hat{z}_j \rightarrow z$ and $\hat{y}_j - \hat{z}_j \rightarrow y - z$ so that $(y, z) = \Delta\text{-}\lim_{j \rightarrow \infty} (\hat{y}_j, \hat{z}_j)$. Thus (i) implies $(y, z) \in \overline{\mathcal{D}} \times \mathcal{E}$.

ad (b): Let $k \in \mathbb{N}$ and $(y, z) \in \mathcal{D}_k \times \mathcal{E}_{h_k}$ be arbitrary. By the uniform approximation property (iii), there is $\hat{y} \in \mathcal{D}$ with $\|y - \hat{y}\|_Z < t_k$ and by Lemma 3.5 there exists $\hat{z} \in \mathcal{E}$ with $\|z - \hat{z}\|_Z \leq c h_k \|f\|_{L^2(\Omega)} =: r_k$. Therefore, (iv) implies

$$\|y - z\|_Z \geq c(\|y\|_Z + \|z\|_Z) - b - (c+1)r_k - (c+1)t_k.$$

Since the sequences t_k and r_k are bounded, equi-transversality holds with $b' := b + (1+c)(\max_{k \in \mathbb{N}} t_k + \max_{k \in \mathbb{N}} r_k)$. \square

Remark 5.2. Since \mathcal{E} as defined in (5.2) is closed and convex and thus weakly closed and data convergence implies weak convergence, the set \mathcal{E} itself arises in the data closure in (i). The situation changes, if one turns to the material data set \mathcal{D} . Of course, if the constitutive law coupling \mathbf{q} and ∇u is linear such as in case of Fourier's law for instance, \mathcal{D} is weakly closed, too, such that $\mathcal{D} \times \mathcal{E} = \overline{\mathcal{D}} \times \mathcal{E}^\Delta$. By contrast, if the constitutive law is nonlinear, then $\overline{\mathcal{D}}$ will in general differ from the Z -closure of \mathcal{D} , but also from the closure of its convex hull. The latter is due to the fact that data convergence provides more information than just weak convergence. The computation of data closures is a field of active research, we only refer to [23] and the references therein.

So far we have focused on the discretization of the set \mathcal{E} . A possible discretization of the set \mathcal{D} is given by piecewise constant functions. To fulfill condition (ii) of Theorem 5.1, we need to bound the distance between the values of those piecewise constant functions and the values of the functions in \mathcal{D} by a monotonically decreasing sequence that converges to zero, which is done in the following

Proposition 5.3. *Let a monotonically decreasing sequence $\{h_k\}_{k \in \mathbb{N}} \subset \mathbb{R}_{>0}$ with $h_k \rightarrow 0$ and a corresponding sequence of shape regular triangulations \mathcal{T}_{h_k} of Ω be given. Let*

$$\mathcal{D} := \{y \in Z : y(x) \in \mathcal{D}^{\text{loc}} \text{ a.e. in } \Omega\} \quad (5.5)$$

with $\mathcal{D}^{\text{loc}} \subset \mathbb{R}^d \times \mathbb{R}^d$ and

$$\mathcal{D}_k := \{y \in Z : y(x) \in \mathcal{D}_k^{\text{loc}} \text{ a.e. in } \Omega, y|_T \in \mathcal{P}_0(T) \forall T \in \mathcal{T}_{h_k}\} \quad (5.6)$$

with $\mathcal{D}_k^{\text{loc}} \subset \mathbb{R}^d \times \mathbb{R}^d$ be given. Moreover, assume that there is a sequence $\{\rho_k\}_{k \in \mathbb{N}}$ with $\rho_k \searrow 0$ such that the Hausdorff distance between \mathcal{D}^{loc} and $\mathcal{D}_k^{\text{loc}}$ satisfies

$$d_{\text{H}}(\mathcal{D}^{\text{loc}}, \mathcal{D}_k^{\text{loc}}) = \max \left\{ \sup_{\xi \in \mathcal{D}^{\text{loc}}} d(\xi, \mathcal{D}_k^{\text{loc}}), \sup_{\eta \in \mathcal{D}_k^{\text{loc}}} d(\eta, \mathcal{D}^{\text{loc}}) \right\} \leq \rho_k. \quad (5.7)$$

Then \mathcal{D} and \mathcal{D}_k as defined in (5.5) and (5.6), respectively, satisfy the fine and uniform approximation assumption (ii) and (iii) in Theorem 5.1.

Proof. Define $V_{h_k} := \{v : \Omega \rightarrow \mathbb{R}^d \times \mathbb{R}^d : v|_T \equiv \text{const.} \forall T \in \mathcal{T}_{h_k}\}$. Then, by standard interpolation error analysis, $\bigcup_{k \in \mathbb{N}} V_{h_k}$ is dense in Z . Therefore, for $y \in \mathcal{D}$ and $\varepsilon > 0$ fixed, but arbitrary, there exist $k_1^\varepsilon \in \mathbb{N}$ such that for all $k \geq k_1^\varepsilon$ there is a $v_k \in V_{h_k}$ with

$$\|y - v_k\|_Z \leq \varepsilon. \quad (5.8)$$

Define $\bar{y}_k \in V_{h_k}$, $k \in \mathbb{N}$, by

$$\begin{aligned} \bar{y}_k(x) &\in \mathcal{D}^{\text{loc}} \quad \text{a.e. in } \Omega, \\ \|\bar{y}_k(x) - v_k(x)\|_{\mathbb{R}^d \times \mathbb{R}^d} &\leq \inf_{\xi \in \mathcal{D}^{\text{loc}}} \|\xi - v_k(x)\|_{\mathbb{R}^d \times \mathbb{R}^d} + \varepsilon \quad \text{a.e. in } x \in \Omega. \end{aligned}$$

Note that \bar{y}_k is well defined, since v_k is constant on each $T \in \mathcal{T}_{h_k}$. Then, $\bar{y}_k \in \mathcal{D}$ and

$$\|v_k - \bar{y}_k\|_Z \leq \|v_k - y\|_Z + \sqrt{|\Omega|} \varepsilon \quad (5.9)$$

for all $k \geq k_1^\varepsilon$. Moreover, there is $k_2^\varepsilon \in \mathbb{N}$ such that $\rho_k \leq \varepsilon$ for all $k \geq k_2^\varepsilon$. Hence, according to 5.7, for each $k \geq k_2^\varepsilon$, there is an $y_k \in \mathcal{D}_k$ such that

$$\|\bar{y}_k - y_k\|_Z \leq \sqrt{|\Omega|} \varepsilon. \quad (5.10)$$

Altogether, (5.8)–(5.10) yield $\|y - y_k\|_Z \leq (1 + 2\sqrt{|\Omega|})\varepsilon$ for all $k \geq \max\{k_1^\varepsilon, k_2^\varepsilon\}$, which along with $y_k \in \mathcal{D}_k$ implies (ii).

To verify (iii), let now $k \in \mathbb{N}$ and $y_k \in \mathcal{D}_k$ be fixed, but arbitrary. Then, by definition of \mathcal{D}_k , there exist $\xi_T^{(k)} \in \mathcal{D}_k^{\text{loc}}$, $T \in \mathcal{T}_{h_k}$, such that $y_k = \sum_{T \in \mathcal{T}_{h_k}} \xi_T^{(k)} \chi_T$ a.e. in Ω . In view of (5.7), for every T , we find $\xi_T \in \mathcal{D}^{\text{loc}}$ such that $\|\xi_T - \xi_T^{(k)}\|_{\mathbb{R}^d \times \mathbb{R}^d} \leq \rho_k$. Therefore, if we define $y \in \mathcal{D}$ by $y := \sum_{T \in \mathcal{T}_{h_k}} \xi_T \chi_T$, then

$$\|y_k - y\|_Z^2 = \sum_{T \in \mathcal{T}_{h_k}} \int_T \|\xi_T - \xi_T^{(k)}\|_{\mathbb{R}^d \times \mathbb{R}^d}^2 dx \leq |\Omega| \rho_k^2,$$

which is (iii) with $t_k = \sqrt{|\Omega|} \rho_k$. \square

Let us denote the number of elements in \mathcal{T}_{h_k} by $N_k := |\mathcal{T}_{h_k}|$. Then \mathcal{D}_k as defined in (5.6) is isomorphic to the finite dimensional set

$$\mathbb{D}_k := \{A \in \mathbb{R}^{N_k \times d \times d} : A_i \in \mathcal{D}_k^{\text{loc}} \forall i = 1, \dots, N_k\},$$

which is clearly compact provided that $\mathcal{D}_k^{\text{loc}}$ is so. This observation immediately implies the following

Proposition 5.4. *Suppose that \mathcal{D} , given as in (5.5), is discretized as in Proposition 5.3 with approximate local material data sets $\mathcal{D}_k^{\text{loc}}$ that are compact for every $k \in \mathbb{N}$. Assume moreover, that the equilibrium constraint set is discretized as in (3.1) with spaces Q_{h_k} and U_{h_k} satisfying Assumption 3.1. Then, for each $k \in \mathbb{N}$, the discretized data-driven problem given by*

$$\left. \begin{array}{l} \min \quad \frac{1}{2} \|y - z\|_Z^2 \\ \text{s.t.} \quad y \in \mathcal{D}_k, z \in \mathcal{E}_k \end{array} \right\} \quad (\text{P}_k)$$

admits a globally optimal solution.

Proof. Throughout the proof, let us suppress the index k in h_k to simplify the notation. First we rewrite (P_k) as

$$(\text{P}_k) \iff \left\{ \begin{array}{l} \min_{(\mathbf{r}, \mathbf{w}) \in \mathcal{D}_k} \min_{(\mathbf{q}, u)} \quad \frac{1}{2} \|(\nabla u_h, \mathbf{q}_h) - (\mathbf{w}, \mathbf{r})\|_{L^2(\Omega; \mathbb{R}^d)^2}^2 \\ \text{s.t.} \quad \mathbf{q}_h \in Q_h, u_h \in U_h, -\text{div}_h \mathbf{q}_h = f. \end{array} \right.$$

By standard arguments, the direct method of calculus of variations yields the existence and uniqueness of a solution to the inner minimization problem. Due to strict convexity, it is uniquely characterized by its necessary and sufficient conditions, which read as follows: Thanks to the surjectivity of $\text{div}_h : Q_h \rightarrow V_h^*$, a tuple (\mathbf{q}_h, u_h) is a solution of the inner minimization problem, iff there exists a Lagrange multiplier $\lambda_h \in V_h$ such that

$$\int_{\Omega} \nabla u_h \cdot \nabla \varphi_h dx = \int_{\Omega} \mathbf{w} \cdot \nabla \varphi_h dx \quad \forall \varphi_h \in U_h \quad (5.11a)$$

$$\int_{\Omega} (\mathbf{q}_h \cdot \mathbf{w}_h + \lambda_h \operatorname{div} \mathbf{w}_h) dx = \int_{\Omega} \mathbf{r} \cdot \mathbf{w}_h dx \quad \forall \mathbf{w}_h \in Q_h \quad (5.11b)$$

$$- \int_{\Omega} v_h \operatorname{div} \mathbf{q}_h dx = \int_{\Omega} f v_h dx \quad \forall v_h \in V_h. \quad (5.11c)$$

By Lemma 3.3(ii), the saddle point system (5.11b)–(5.11c) admits a unique solution $(\mathbf{q}_h, \lambda_h) \in Q_h \times V_h$ for every right hand side $(\mathbf{r}, f) \in L^2(\Omega; \mathbb{R}^d) \times L^2(\Omega)$ and the associated solution operator is linear and continuous by (3.9). Moreover, since U_h is a closed subspace of $H_0^1(\Omega)$, the same holds for the discretized Laplace equation in (5.11a), i.e., for every $\mathbf{w} \in L^2(\Omega; \mathbb{R}^d)$, there is a unique solution $u_h \in U_h$ and the solution mapping is linear and continuous. Thus, there is an affine (due to f) and continuous solution operator of (5.11) denoted by

$$G_h : L^2(\Omega; \mathbb{R}^d)^2 \ni (\mathbf{r}, \mathbf{w}) \mapsto (\mathbf{q}_h, \nabla u_h) \in Q_h \times \nabla U_h. \quad (5.12)$$

With the help of G_h , we can rewrite (P_k) equivalently as

$$(P_k) \quad \Longleftrightarrow \quad \min_{(\mathbf{r}, \mathbf{w}) \in \mathcal{D}_k} \frac{1}{2} \|(G_h - \operatorname{id})(\mathbf{r}, \mathbf{w})\|_{L^2(\Omega; \mathbb{R}^d)^2}^2$$

Therefore, since \mathcal{D}_k is compact as explained above and G_h and thus the whole objective is continuous, the existence of a globally optimal solution follows from the Weierstrass theorem. \square

Remark 5.5. Note that the result of Proposition 5.4 also holds for a problem of the form

$$\left. \begin{array}{l} \min \quad \frac{1}{2} \|y - z\|_Z^2 \\ \text{s.t.} \quad y \in \mathcal{D}_k, z \in \mathcal{E} \end{array} \right\} \quad (\tilde{P}_k)$$

with the continuous equilibrium set \mathcal{E} instead of \mathcal{E}_k . The arguments are completely the same as in the proof of Proposition 5.4, since the solution operators associated with the continuous counterparts to the saddle point problem and the Laplace equation are also linear and continuous. Thus, to ensure the mere existence of optimal solutions, only the discretization of \mathcal{D} is necessary, whereas the numerical computation of optimal solutions of course requires a discretization of \mathcal{E} , too.

As a direct consequence of the previous results, namely Theorem 5.1 and Propositions 5.3 and 5.4, we obtain the following result, which, though it is just a corollary as a consequence of the above findings, can be seen as our main result:

Corollary 5.6. *Let Z , \mathcal{E} , and \mathcal{D} be defined as in (5.1), (5.2), and (5.5). Furthermore, let $\{h_k\}_{k \in \mathbb{N}} \subset \mathbb{R}_{>0}$ be a monotonically decreasing sequence of mesh sizes with $h_k \searrow 0$ as $k \rightarrow \infty$ and suppose the following assumptions:*

- (i) (Transversality) *There are constants $c > 0$ and $b \geq 0$ such that, for all $y \in \mathcal{D}$ and $z \in \mathcal{E}$,*

$$\|y - z\|_Z \geq c(\|y\|_Z + \|z\|_Z) - b.$$

- (ii) (Discrete material data set) *The set \mathcal{D} is approximated as in Proposition 5.3, i.e.,*

$$\mathcal{D}_k := \{y \in Z : y(x) \in \mathcal{D}_k^{\text{loc}} \text{ a.e. in } \Omega, y|_T \in \mathcal{P}_0(T) \forall T \in \mathcal{T}_{h_k}\}$$

and the Hausdorff distance fulfills $d_H(\mathcal{D}^{\text{loc}}, \mathcal{D}_k^{\text{loc}}) \leq \rho_k$ with a sequence $\{\rho_k\}_{k \in \mathbb{N}}$ with $\rho_k \searrow 0$. Moreover, $\mathcal{D}_k^{\text{loc}} \subset \mathbb{R}^d \times \mathbb{R}^d$ is compact for every $k \in \mathbb{N}$.

- (iii) (Conforming discretization) *The finite dimensional spaces Q_k and U_k defining the discrete equilibrium set \mathcal{E}_k satisfy Assumption 3.1.*

Then, for every $k \in \mathbb{N}$, there exists at least one solution of

$$\left. \begin{array}{l} \min \quad \frac{1}{2} \|y - z\|_Z^2 \\ \text{s.t.} \quad y \in \mathcal{D}_k, z \in \mathcal{E}_k \end{array} \right\} \quad (\mathbf{P}_k)$$

and every sequence $\{(y_k, z_k)\} \subset Z \times Z$ of such solutions satisfies the following:

- (a) If $\overline{\mathcal{D} \times \mathcal{E}^\Delta} \neq \emptyset$, then $\|y_k - z_k\|_Z \rightarrow 0$.
 (b) If $\|y_k - z_k\|_Z \rightarrow 0$, then there exists $z \in \overline{\mathcal{D} \times \mathcal{E}^\Delta}$ such that, up to subsequences, $(z, z) = \Delta\text{-}\lim_{k \rightarrow \infty} (y_k, z_k)$.

We again underline that the data closure $\overline{\mathcal{D} \times \mathcal{E}^\Delta}$ is in general not equal to $\mathcal{D} \times \mathcal{E}$, but $\overline{\mathcal{D} \times \mathcal{E}}$ with an enlarged set $\overline{\mathcal{D}} \supset \mathcal{D}$, see Remark 5.2. Corollary 5.6 shows that, under the mentioned assumptions, it is in principle possible to approximate elements of the data closure. The computation of solutions to (\mathbf{P}_k) may however be a very delicate issue, depending on the precise structure of $\mathcal{D}_k^{\text{loc}}$. Before we address this issue in more details in the Section 7 below, let us shortly discuss the question why the use of $H(\text{div})$ -conforming finite elements like the Raviart-Thomas element seems to be indispensable for the discretization of the data-driven problem (\mathbf{DDP}) .

6. WHY $H(\text{div})$ -CONFORMING FINITE ELEMENTS?

As the construction of $H(\text{div})$ -conforming finite elements is rather complicated compared to e.g. classical Lagrangian finite elements, the question arises if their use is really necessary for the discretization of data-driven problems. To answer this question, let us assume that we do not use $H(\text{div})$ -conforming elements, i.e., we discretize the equilibrium set \mathcal{E} by a finite dimensional space $\widehat{Q}_h \subset L^2(\Omega; \mathbb{R}^d)$ with $\widehat{Q}_h \not\subset H(\text{div})$ and define the discrete divergence condition by

$$\mathbf{q}_h \in \widehat{Q}_h, \quad \int_{\Omega} \mathbf{q}_h \cdot \nabla v_h \, dx = \int_{\Omega} f v_h \, dx \quad \forall v_h \in \widehat{V}_h, \quad (6.1)$$

where \widehat{V}_h is a finite dimensional subspace of $H_0^1(\Omega)$. Given a triangulation \mathcal{T}_h of the domain Ω , a classical example for these finite dimensional spaces reads

$$\widehat{Q}_h = \{\mathbf{w} \in L^\infty(\Omega; \mathbb{R}^d) : \mathbf{w}|_T \in \mathcal{P}_0(T) \, \forall T \in \mathcal{T}_h\}, \quad (6.2)$$

$$\widehat{V}_h = \{v \in C(\overline{\Omega}) \cap H_0^1(\Omega) : v|_T \in \mathcal{P}_1(T) \, \forall T \in \mathcal{T}_h\}. \quad (6.3)$$

Revisiting the proof of Theorem 5.1 shows that a central aspect of convergence analysis is that elements from \mathcal{E} can be approximated by elements from \mathcal{E}_h w.r.t. the data topology and vice versa. Let us return to part (a) in the proof of Theorem 5.1. When verifying condition $(\mathbf{DC}2)$ from the definition of data convergence, one considers a sequence $\{(y_h, z_h)\}_{h>0}$, $(y_h, z_h) \in \mathcal{D}_h \times \mathcal{E}_h$, converging in data to (y, z) . To show that (y, z) is an element of the data closure, we need to prove the existence of a sequence $(\hat{y}_h, \hat{z}_h) \in \mathcal{D} \times \mathcal{E}$ with $\hat{y}_h \rightarrow y$, $\hat{z}_h \rightarrow z$, and $\hat{y}_h - \hat{z}_h \rightarrow y - z$, which, in view of the data convergence of $\{(y_h, z_h)\}$ is equivalent to

$$\hat{y}_h - y_h \rightarrow 0, \quad \hat{z}_h - z_h \rightarrow 0, \quad \hat{y}_h - y_h + z_h - \hat{z}_h \rightarrow 0. \quad (6.4)$$

If we assume that the data approximation satisfies the uniform approximation property (iii) from Theorem 5.1, then we already know that a sequence $\{\hat{y}_h\}_{h>0} \subset \mathcal{D}$

exists with $\hat{y}_h - y_h \rightarrow 0$ in $Z = L^2(\Omega; \mathbb{R}^d)^2$. Therefore, if we take this sequence for \hat{y}_h , then the sequence $\{\hat{z}_h\}_{h>0} \subset \mathcal{E}$ must necessarily fulfill

$$\hat{z}_h - z_h \rightarrow 0 \quad \text{in } Z \quad (6.5)$$

in order to guarantee (6.4). This however is not always possible, if nonconforming finite element spaces are used, as we will see in the following when considering the flux component of

$$z_h = (\mathbf{q}_h, \nabla u_h) \in \mathcal{E}_h \subset (\widehat{\mathcal{Q}}_h, \nabla U_h).$$

The best possible choice for the construction of the desired elements from \mathcal{E} is of course to choose the solution $\hat{\mathbf{q}}_h \in H(\text{div})$ of

$$\left. \begin{array}{l} \min_{\mathbf{q} \in L^2(\Omega; \mathbb{R}^d)} \quad \frac{1}{2} \|\mathbf{q} - \mathbf{q}_h\|_{L^2(\Omega; \mathbb{R}^d)}^2 \\ \text{s.t.} \quad -\text{div } \mathbf{q} = f, \end{array} \right\} \quad (6.6)$$

which is uniquely characterized by the existence of $w \in H_0^1(\Omega)$ such that

$$\hat{\mathbf{q}}_h - \mathbf{q}_h + \nabla w = 0, \quad -\text{div } \hat{\mathbf{q}}_h = f.$$

Let us define $\Phi \in H_0^1(\Omega)$ by

$$-\Delta \Phi = f \quad \text{in } H^{-1}(\Omega),$$

as well as the L^2 -projection of \mathbf{q}_h on $\nabla H_0^1(\Omega)$, denoted by $\nabla \Phi^h$ with $\Phi^h \in H_0^1(\Omega)$. Then, we obtain

$$\begin{aligned} \|\hat{\mathbf{q}}_h - \mathbf{q}_h\|_{L^2(\Omega; \mathbb{R}^d)} &= \|\nabla w\|_{L^2(\Omega; \mathbb{R}^d)} \\ &= \sup_{\substack{v \in H_0^1(\Omega) \\ \|\nabla v\|_{L^2(\Omega; \mathbb{R}^d)} \leq 1}} \int_{\Omega} \nabla(\Phi - \Phi^h) \cdot \nabla v \, dx \\ &= \|\nabla \Phi - \nabla \Phi^h\|_{L^2(\Omega; \mathbb{R}^d)}. \end{aligned} \quad (6.7)$$

Since \widehat{V}_h is a closed subspace of $H_0^1(\Omega)$, we may decompose $\Phi^h = \Phi_0^h + \Phi_{\perp}^h \in \widehat{V}_h \oplus \widehat{V}_h^{\perp}$, where the orthogonal complement is taken w.r.t. the H_0^1 -scalar product. Due to $-\text{div}_h \mathbf{q}_h = f$, we find for Φ_0^h

$$\int_{\Omega} \nabla \Phi_0^h \cdot \nabla v_h \, dx = \int_{\Omega} \mathbf{q}_h \cdot \nabla v_h \, dx = \langle f, v_h \rangle \quad \forall v_h \in \widehat{V}_h \quad (6.8)$$

and consequently, by the best approximation property of the finite element solution,

$$\|\Phi_0^h - \Phi\|_{H_0^1(\Omega)} \rightarrow 0 \quad \text{as } h \searrow 0, \quad (6.9)$$

follows, provided that $\bigcup_{h>0} \widehat{V}_h$ is dense in $H_0^1(\Omega)$. Therefore, if the sequence $\{\mathbf{q}_h\}_{h>0}$ is such that

$$\liminf_{h \searrow 0} \|\nabla \Phi_{\perp}^h\|_{L^2(\Omega; \mathbb{R}^d)} \geq c > 0, \quad (6.10)$$

then (6.9) implies

$$\begin{aligned} &\liminf_{h \searrow 0} \|\hat{\mathbf{q}}_h - \mathbf{q}_h\|_{L^2(\Omega; \mathbb{R}^d)} \\ &= \liminf_{h \searrow 0} \|\nabla \Phi - \nabla \Phi^h\|_{L^2(\Omega; \mathbb{R}^d)} \\ &\geq \liminf_{h \searrow 0} \left(\|\nabla \Phi_{\perp}^h\|_{L^2(\Omega; \mathbb{R}^d)} - \|\nabla \Phi - \nabla \Phi_0^h\|_{L^2(\Omega; \mathbb{R}^d)} \right) \geq c \end{aligned} \quad (6.11)$$

and hence, (6.5) cannot hold in this case.

Let us give a simple one-dimensional example indicating that (6.10) may well happen for a sequence converging in data. For this purpose, set $\Omega := (0, 2\pi)$ and consider the equidistant triangulation of Ω given by

$$\mathcal{T}_k := \{(0, h_k), (h_k, 2h_k), \dots, ((2k-1)h_k, 2\pi)\}$$

with $h_k := (2k)^{-1} 2\pi$, $k \in \mathbb{N}$. We employ the standard finite element space from (6.3) and denote the nodal basis of $\widehat{V}_k := \widehat{V}_{h_k}$ by φ_i^k , i.e., $\varphi_i^k(j h_k) = \delta_{ij}$, $1 \leq i, j \leq 2k-1$. Note that the meshes and thus the spaces \widehat{V}_k are nested. Furthermore, we abbreviate the solution of (6.8) by $\Phi_0^k := \Phi_0^{h_k} \in \widehat{V}_k$ and set

$$\Psi_k(x) := \frac{1}{\sqrt{2\pi}} \sum_{i=1}^{2k} \varphi_{2i-1}^{2k}(x) \in H_0^1(\Omega), \quad k \in \mathbb{N}. \quad (6.12)$$

Then one easily verifies that

$$\left\| \frac{d}{dx} \Psi_k \right\|_{L^2(0, 2\pi)} = 1, \quad \frac{d}{dx} \Psi_k \rightarrow 0 \quad \text{as } k \rightarrow \infty \quad (6.13)$$

and

$$\int_0^{2\pi} \frac{d}{dx} \Psi_k(x) \frac{d}{dx} \varphi_i^k(x) dx = 0 \quad \forall i = 1, \dots, 2k-1, \quad k \in \mathbb{N}.$$

Since $\widehat{V}_k = \text{span}(\varphi_1^k, \dots, \varphi_{2k-1}^k)$, the latter implies $\Psi_k \in \widehat{V}_k^\perp$. Let us choose the space from (6.2) for \mathbf{q}_h , but with half mesh size, i.e.,

$$\widehat{Q}_k := \text{span}(\chi_{(0, h_{2k})}, \chi_{(h_{2k}, 2h_{2k})}, \dots, \chi_{((4k-1)h_{2k}, 2\pi)})$$

and set

$$\mathbf{q}_k := \frac{d}{dx} \Phi_0^k + \frac{d}{dx} \Psi_k \in \widehat{Q}_k$$

such that $\Phi_\perp^k = \Psi_k$. Then, owing to (6.8) and (6.13), we have

$$-\text{div}_h \mathbf{q}_k = f, \quad \mathbf{q}_k \rightarrow \frac{d}{dx} \Phi \quad \text{in } L^2(0, 2\pi),$$

Together with $\mathbf{q}_k \in \widehat{Q}_k$, this indicates that \mathbf{q}_k can be part of a sequence converging in data, depending on the structure of \mathcal{D}_k . Indeed, if, for instance, there is an $a \in \mathbb{R}$ such that $(a, 1/\sqrt{2\pi}), (-a, -1/\sqrt{2\pi}) \in \mathcal{D}_k^{\text{loc}}$ for all $k \in \mathbb{N}$, then one could choose $U_k := \widehat{V}_{2k}$, $u_k := a \sum_{i=1}^{2k} \varphi_{2i-1}^{2k}$, and set $y_k = (\mathbf{r}_k, \mathbf{w}_k) = (\frac{d}{dx} \Psi_k, \frac{d}{dx} u_k) \in \mathcal{D}_k$. Then $y_k \rightarrow 0$ and

$$z_k - y_k = (\mathbf{q}_k - \mathbf{r}_k, \nabla u_k - \mathbf{w}_k) = (\frac{d}{dx} \Phi_0^k, 0) \rightarrow (\frac{d}{dx} \Phi, 0) \quad \text{in } Z,$$

i.e., $((0, 0), (\frac{d}{dx} \Phi, 0)) = \Delta\text{-}\lim_{k \rightarrow \infty} (y_k, z_k)$. However, due to (6.11) and $\Phi_\perp^k = \Psi_k$, there holds

$$\liminf_{k \rightarrow \infty} \|\hat{\mathbf{q}}_k - \mathbf{q}_k\|_{L^2(0, 2\pi)} \geq \lim_{k \rightarrow \infty} \left\| \frac{d}{dx} \Psi_k \right\|_{L^2(0, 2\pi)} = 1 \quad (6.14)$$

such that (6.5) does not hold in this example.

Altogether we have seen that it is in general not possible to construct a sequence $\{\hat{z}_h\}_{h>0} \subset \mathcal{E}$ such that (6.5) holds, if a finite element space \widehat{Q}_h is used that is not $H(\text{div})$ -conforming. In contrast to this, Lemma 3.5 shows that this is well possible in case of $H(\text{div})$ -conforming finite elements. Nonetheless, this is no proof that no sequences $\{(\hat{y}_h, \hat{z}_h)\}_{h>0} \subset \mathcal{D} \times \mathcal{E}$ exist such that (6.4) holds also in case of nonconforming finite element spaces. Maybe, it is possible to construct a sequence \hat{z}_h by adjusting the sequence \hat{y}_h , but, so far, we have no idea how to do this and this issue gives rise to future research.

7. ALGORITHMS

The last section is devoted to optimization algorithms for the numerical solution of the discretized data-driven problem (\mathbf{P}_k) . In the literature, a projection based fixed-point method is frequently used to solve (\mathbf{P}_k) . We will advance this method by introducing a step size and compare the proximal gradient method arising in this way with two variants of the Douglas-Rachford algorithm applied to (\mathbf{P}_k) . Besides these first-order methods, we also employ the equivalence of (\mathbf{P}_k) to a quadratic semi-assignment problem and develop a heuristic based on this reformulation. The algorithms are tested by means of two examples, a linear and a non-linear material model.

7.1. Fixed-point methods based on projections. We first consider a class of algorithms that is based on the two L^2 -projections $\pi_{\mathcal{E}_k} : Z \rightarrow \mathcal{E}_k$ and $\pi_{\mathcal{D}_k} : Z \rightarrow \mathcal{D}_k$ with $\mathcal{E}_k := \mathcal{E}_{h_k}$ as in (3.1) and \mathcal{D}_k as in Proposition 5.3, respectively. As the proof of Proposition 5.4 shows, $\pi_{\mathcal{E}_k}$ is simply given by the affine operator G_h defined in (5.12), i.e., given $y = (\mathbf{r}, \mathbf{w}) \in Z$, we have $\pi_{\mathcal{E}_k}(\mathbf{r}, \mathbf{w}) = (\mathbf{q}_{h_k}, \nabla u_{h_k})$, where $(\mathbf{q}_{h_k}, u_{h_k}) \in Q_{h_k} \times U_{h_k}$ is the unique solution to the saddle point system (5.11). To compute the projection $\pi_{\mathcal{D}_k}$, we first notice that we can project an arbitrary function $z \in Z$ onto \mathcal{D}_k by first projecting z onto the space of piecewise constant functions and then projecting the obtained function onto \mathcal{D}_k . To see this, define $\hat{z} := \frac{1}{|T|} \int_T z dx \chi_T$ for a given function $z \in Z$. Then the above assertion follows from

$$\begin{aligned} \arg \min_{y \in \mathcal{D}_k} \|y - z\|_Z^2 &= \arg \min_{y \in \mathcal{D}_k} \|y\|_Z^2 - 2(y, z) + \|z\|_Z^2 \\ &= \arg \min_{y \in \mathcal{D}_k} \|y\|_Z^2 - 2(y, \hat{z}) + \|\hat{z}\|_Z^2 = \arg \min_{y \in \mathcal{D}_k} \|y - \hat{z}\|_Z^2. \end{aligned} \quad (7.1)$$

Moreover, the structure of \mathcal{D}_k yields that the projection can be computed element-wise, for example by using a k -nearest neighbors algorithm. It is to be noted that $\pi_{\mathcal{D}_k}$ is a set-valued map, since the projection on \mathcal{D}_k is clearly non-unique in general due to the non-convexity of \mathcal{D}_k . In the numerical realization of the algorithms introduced below, we pick an arbitrary $y \in \mathcal{D}_k$ attaining the minimum in (7.1).

In [18], a simple projection algorithm was introduced to solve the data driven problem. One iteration step consists of two projections and reads

$$y_{n+1} = \pi_{\mathcal{D}_k}(\pi_{\mathcal{E}_k}(y_n)) \quad (\text{PG})$$

for $n \in \mathbb{N}$ and an arbitrary initial point $y_0 \in Z$. The algorithm terminates if $y_{n+1} = y_n$ for some $n \in \mathbb{N}$. Note that this algorithm is a special case of the proximal gradient method or more specifically of the projected gradient method [3, 8, 20] with step size constant equal to one. To illustrate this, let us shortly recall the concept of the proximal gradient method. Given to mappings $F, G : Z \rightarrow \mathbb{R} \cup \{\infty\}$, where F is smooth, while G may be not, the proximal gradient iteration for the minimization of $F + G$ reads

$$y_{n+1} = \text{prox}_{\gamma G}(y_n - \gamma \nabla F(y_n)), \quad (7.2)$$

where $\text{prox}_{\gamma G}(y) \in \arg \min_{\eta \in Z} \frac{1}{2} \|\eta - y\|_Z^2 + \gamma G(\eta)$ denotes the proximal map and $\gamma > 0$ is a step size. In order to apply the proximal gradient method to the data driven problem, let us rewrite (\mathbf{P}_k) as

$$(\mathbf{P}_k) \quad \iff \quad \min_{y \in Z} \frac{1}{2} \|\pi_{\mathcal{E}_k}(y) - y\|_Z^2 + I_{\mathcal{D}_k}(y), \quad (7.3)$$

where $I_{\mathcal{D}_k}$ denotes the indicator functional of \mathcal{D}_k . If we now set $F(y) := \frac{1}{2}\|\pi_{\mathcal{E}_k}(y) - y\|_Z^2$ and $G := I_{\mathcal{D}_k}$, then (7.2) becomes

$$y_{n+1} = \pi_{\mathcal{D}_k}(y_n - \gamma(y_n - \pi_{\mathcal{E}_k}(y_n))). \quad (\text{PS})$$

If F and G were proper, convex, and lower semi-continuous, then the classical convergence results for proximal gradient methods apply provided that $0 < \gamma_{\min} \leq \gamma < 2L^{-1}$, where L is the Lipschitz constant of ∇F , cf. e.g. [8, Theorems 9.6 and 10.2]. In our case, $\nabla F(y) = y - \pi_{\mathcal{E}_k}(y)$ has the Lipschitz constant $L = 1$ and, accordingly, we have chosen $\gamma \in [1, 2)$ in our computations. We emphasize that $G = I_{\mathcal{D}_k}$ is clearly not convex due to the non-convexity of \mathcal{D}_k and therefore, convergence of the iteration in (PS) to global minimizers cannot be guaranteed. A convergence analysis of the proximal gradient method for non-convex problems in function space similar to (P_k) is presented in [27]. However, due to the lack of convexity, the convergence results are rather limited and one only obtains that weak accumulation points of the sequence of iterates satisfy a rather weak stationarity concept termed L-stationarity, provided that ∇F is completely continuous. One can therefore not expect the iterates produced by the scheme in (PS) to converge to a minimizer (neither global nor local) of (P_k). Nonetheless, the introduction of the step size γ may improve the quality of the numerical results. In our numerical tests, it has turned out that it is favorable to start with $\gamma = 1.4$ and to reduce γ whenever the algorithm circles between two iterates, that is $y_{n+2} = y_n$ for some $n \in \mathbb{N}$, see Table 2 below. The choice $\gamma < 1$ accelerates the convergence of the algorithm to a fixed point with substantially larger objective value and should therefore be avoided.

We compare the performance of the proximal gradient method (with step size $\gamma = 1$ and varying step size) with the Douglas-Rachford algorithm for the computation of elements in the intersection of two sets $A, B \subset Z$, cf., e.g., [1]. This algorithm is also initialized with an arbitrary point $y_0 \in Z$ and the iterates are computed by the formula

$$y_{n+1} = T_{A,B}(y_n) \quad \text{with} \quad T_{A,B} := \frac{\text{id} + R_B \circ R_A}{2},$$

where the reflections R_A and R_B are defined by $R_A := 2\pi_A - \text{id}$ and $R_B := 2\pi_B - \text{id}$, respectively. In the case of the data-driven problem, we have the two alternatives $A := \mathcal{D}_k$ and $B := \mathcal{E}_k$ or $A := \mathcal{E}_k$ and $B := \mathcal{D}_k$, that is

$$y_{n+1} = T_{\mathcal{E}_k, \mathcal{D}_k}(y_n) \quad (\text{DR1})$$

or, respectively,

$$y_{n+1} = T_{\mathcal{D}_k, \mathcal{E}_k}(y_n). \quad (\text{DR2})$$

Since in general $\mathcal{E}_k \cap \mathcal{D}_k = \emptyset$, we can neither expect the existence of a fixed point of $T_{\mathcal{E}_k, \mathcal{D}_k}$ nor $T_{\mathcal{D}_k, \mathcal{E}_k}$, cf. [1, Proposition 9]. Consequently, the iterations in (DR1) and (DR2), respectively, will in general not converge, when applied to the data-driven problem. Though convergence of the iterations in (PG) and (PS) can be expected neither, we observed their convergence to a fixed point in our numerical computations. In contrast to this, the Douglas-Rachford iterations did frequently not converge to a fixed point and therefore, we let the algorithm terminate, if it does not achieve an improvement of the objective value for a fixed number of iterations. Note that the iterates of both variants of the Douglas-Rachford algorithm do not necessarily fulfill $y_{n+1} \in \mathcal{D}_k$ and hence, in order to evaluate the objective of (P_k)

in the form (7.3), we additionally project the current iterate onto \mathcal{D}_k after each iteration to evaluate the objective value.

7.2. Quadratic assignment and local search. In order to design an alternative strategy for the solution of (P_k) , we rewrite the discretized data-driven problem as a *quadratic semi-assignment problem*, where we assign the measured data to the elements of the finite element grid such that the distance to the corresponding projection onto the equilibrium set becomes minimal. For this purpose, let $\mathcal{T}_{h_k} = \{T_1, \dots, T_l\}$, $l = l(h_k) \in \mathbb{N}$, be the considered triangulation of Ω and let \mathcal{D}_k be defined as in (5.6) with $\mathcal{D}_k^{\text{loc}} := \{\hat{y}_1, \dots, \hat{y}_m\}$, $m \in \mathbb{N}$. Then (P_k) can be rewritten as

$$(P_k) \iff \left\{ \begin{array}{l} \min \quad \frac{1}{2} \|\pi_{\mathcal{E}_k}(y) - y\|_Z^2 \\ \text{s.t.} \quad y|_{T_i} = \sum_{j=1}^m x_{i,j} \hat{y}_j \quad \forall i = 1, \dots, n \\ \sum_{j=1}^m x_{i,j} = 1 \quad \forall i = 1, \dots, l \\ x_{i,j} \in \{0, 1\} \quad \forall i = 1, \dots, l, \quad j = 1, \dots, m \end{array} \right\} \quad (\text{QSAP})$$

with $x = (x_{1,1}, \dots, x_{1,m}, x_{2,1}, \dots, x_{2,m}, \dots, x_{l,1}, \dots, x_{l,m})^T \in \mathbb{R}^{lm}$, or equivalently

$$(\text{QSAP}) \iff \left\{ \begin{array}{l} \min \quad x^T A x + b^T x + c \\ \text{s.t.} \quad \sum_{j=1}^m x_{i,j} = 1 \quad \forall i = 1, \dots, l \\ x_{i,j} \in \{0, 1\} \quad \forall i = 1, \dots, l, \quad j = 1, \dots, m \end{array} \right.$$

with a symmetric and positive semidefinite matrix $A \in \mathbb{R}^{lm \times lm}$, $b \in \mathbb{R}^{lm}$, and $c \in \mathbb{R}$ arising from the affine solution operator to the saddle point system (5.11), the mass matrices for the scalar products in Z , and a matrix containing the elements of $\mathcal{D}_k^{\text{loc}}$. Note that it is allowed to assign the same measuring point to more than one element and that the objective function is quadratic.

Remark 7.1. Since quadratic semi-assignment problems are NP-hard [24], the reformulation as (QSAP) shows that it is in general not possible to solve (P_k) efficiently (provided that $NP \neq P$). This makes it practically impossible to solve the discretized data-driven problem exactly, when we are faced with a big amount of measured data and a fine finite element grid.

In view of the above remark, we resort to a heuristic for the solution of (QSAP). Such a heuristic method is the local search, that is we start with an arbitrary initial assignment and change it on only one element of the triangulation while the assignment to the remaining elements is fixed. If the change leads to an improvement, the current solution is updated accordingly and one moves on to the next element in the triangulation, where the same procedure is applied. To test every point in local data set $\mathcal{D}_k^{\text{loc}}$ for the respective element of the triangulation rapidly becomes too costly, even for moderate finite element meshes, since one needs to evaluate the objective in (QSAP) every time, which in turn requires the solution of the saddle point problem (5.11). In order to reduce the effort, one can restrict to the K nearest neighbors of the data point which is currently assigned to the considered element.

But still, for reasonable finite element discretizations, the effort of the method easily becomes too high. A possible remedy is to employ a reduced model-order approach such as proper orthogonal decomposition (POD) for the saddle point problem in (5.11), cf. e.g. [14, 26] and the references therein. Since the right hand side only changes in one single element from one iteration of the local search to the next, only a moderate number of snapshots are needed to achieve a sufficient accuracy of the POD basis. The overall method is sketched in Algorithm 1. Note that we only accept an improvement on the objective if it is computed with the exact finite element model. A crucial issue for the performance of Algorithm 1 is the choice of the initial solution. One possibility is to initialize the local search by a solution obtained with one of the projection algorithms from Section 7.1. This allows to escape from fixed points that are not optimal or just local minimizers. A second option is to solve (QSAP) exactly with a small selection of measured data points on a very coarse finite element mesh. Note that, since (QSAP) is NP-hard, the problem size has to be fairly small to allow for an exact algorithm. We employ the algorithm from [6] on a coarse grid with only a few measured data points and project its solution to the finer mesh to initialize Algorithm 1.

8. NUMERICAL EXPERIMENTS

In our numerical tests, we consider the domain $\Omega = (0, 1)^2$ such that Assumption 2.3(i) is fulfilled. For the finite element discretization, we use an exact triangulation of Friedrich-Keller type and choose the following function spaces

$$Q_h = \mathcal{RT}_0(\mathcal{T}_h), \quad V_h = \mathcal{P}_0(\mathcal{T}_h),$$

and

$$U_h = \{u \in C(\bar{\Omega}) \cap H_0^1(\Omega) : u|_T \in \mathcal{P}_1(T) \forall T \in \mathcal{T}_h\}. \quad (8.1)$$

As seen in Proposition 3.7, these finite element spaces satisfy Assumption 3.1.

We test all previously introduced algorithms with the following parameters and settings: For the projection-based fixed point algorithms (PG), (PS), (DR1), and (DR2), we use $y_0 = 0$ as starting point. If not stated otherwise, we choose the initial step size $\gamma = 1.4$ for (PS) and reduced γ by the factor 0.9 whenever the algorithm circled between two iterates. This choice is a compromise of accuracy and running time motivated by the observations in Table 2 below. As mentioned above, the Douglas-Rachford methods in (DR1) and (DR2) do frequently not converge to a fixed point and hence, we let the algorithms terminate after 50 iterations without an improvement of the objective value.

We also test two variants of Algorithm 1. In the first one, we initialized the algorithm with the best out of ten solutions of the projection algorithm (PS) with random starting points with components in $[-4, 4]^4$ and used the remaining for the computation of the initial POD basis. In the second one, we initialized the algorithm with the exact solution on a coarse grid consisting of eight triangles with mesh size $\frac{1}{\sqrt{2}}$ and 16 selected data points of our measured data set. As indicated above, we computed this solution by means of the exact algorithm from [6]. In both cases, the parameters of Algorithm 1 are set to $\varepsilon_1 = 0.002$, $\varepsilon_2 = 0.001$, $\varepsilon_3 = 0.01$, and $K = 20$. These choices are motivated by numerical experiments based on the linear material law from Section 8.1.

Algorithm 1 Local Search with POD

```

1: Choose triangulation  $\mathcal{T}$  of  $\Omega$ 
2: Choose  $\varepsilon_1, \varepsilon_2, \varepsilon_3 > 0$  and  $K \in \mathbb{N}$ 
3: Compute initial POD-basis
4: Choose initial solution  $y$  with objective value  $v := \frac{1}{2} \|\pi_{\mathcal{E}_k}(y) - y\|_Z^2$ 
5: Set  $\bar{v} := v + 1$ 
6: while  $v \neq \bar{v}$  do
7:   Update:  $\bar{v} \leftarrow v$ 
8:   for all  $\tilde{T} \in \mathcal{T}$  do
9:     Find  $K$  nearest neighbors  $y_1, \dots, y_K \in \mathcal{D}_k^{\text{loc}}$  of  $y|_{\tilde{T}}$ 
10:    for  $j = 1, \dots, K$  do
11:      Define  $\tilde{y}$  on each  $T \in \mathcal{T}$  by

$$\tilde{y}|_T := \begin{cases} y|_T & \text{if } T \neq \tilde{T} \\ y_j, & \text{if } T = \tilde{T} \end{cases}$$

12:      Solve POD-model to approximate  $\pi_{\mathcal{E}_k}(\tilde{y})$  by  $z_a$ 
13:      Set  $v_a := \frac{1}{2} \|\tilde{y} - z_a\|_Z^2$ 
14:      if  $\frac{|v_a - v|}{|v|} > \varepsilon_1$  then
15:        Compute exact solution  $z_e := \pi_{\mathcal{E}_k}(\tilde{y})$  and set

$$v_e := \frac{1}{2} \|\tilde{y} - z_e\|_Z^2$$

16:        if  $\frac{|v_a - v_e|}{|v_e|} > \varepsilon_2$  then
17:          Add  $z_e$  as snapshot to compute new POD-basis
18:        end if
19:        if  $v_e < v$  then
20:          Update:  $y \leftarrow \tilde{y}, v \leftarrow v_e$ 
21:          Add  $z_e$  as snapshot to compute new POD-basis
22:        end if
23:      end if
24:      if  $v_a < (1 + \varepsilon_3)v$  then
25:        Compute exact solution  $z_e := \pi_{\mathcal{E}_k}(\tilde{y})$  and set

$$v_e := \frac{1}{2} \|\tilde{y} - z_e\|_Z^2$$

26:        if  $v_e < v$  then
27:          Update:  $y \leftarrow \tilde{y}, v \leftarrow v_e$ 
28:          Add  $z_e$  as snapshot to compute new POD-basis
29:        end if
30:      end if
31:    end for
32:  end for
33: end while

```

8.1. **Fourier's law.** For the first numerical tests, we considered Fourier's law where the material law is linear. To be more precise, we set $\kappa : \mathbb{R}^2 \rightarrow \mathbb{R}^2$, $\kappa(w) := w$ such that (1.1) simply becomes Poisson's equation, i.e., $-\Delta u = f$ in Ω . Moreover, the right hand side is set to $f(x) = 2\pi^2 \sin(\pi x_1) \sin(\pi x_2)$, which clearly fulfills

the regularity condition in Assumption 2.3. Accordingly, the exact solution reads $u(x) = \sin(\pi x_1) \sin(\pi x_2)$. For the local material data set, we sampled 11025 evenly distributed data points (r, w) fulfilling $r = \kappa(w)$ within $[-4, 4]^4$.

To compare the algorithms, we list the returned objective values and the needed iterations and wall clock times in Table 1. Moreover, we compute the relative distance of the exact solution to the outcome of the respective algorithm measured in the L^2 -norm and in the H_0^1 -seminorm. For this purpose, denote the result of the respective algorithm by \bar{y} and set $(\bar{\mathbf{q}}_{h_k}, \nabla \bar{u}_{h_k}) := \pi_{\varepsilon_k}(\bar{y})$. We then compute

$$\mathbf{err}_{L^2} := \frac{\|u - \bar{u}_{h_k}\|_{L^2(\Omega)}}{\|u\|_{L^2(\Omega)}} \quad \text{and} \quad \mathbf{err}_{H_0^1} := \frac{\|\nabla u - \nabla \bar{u}_{h_k}\|_{L^2(\Omega; \mathbb{R}^d)}}{\|\nabla u\|_{L^2(\Omega; \mathbb{R}^d)}}. \quad (8.2)$$

As one can see, the returned objective values of the algorithms and the distances to

Algorithm	Objective value	\mathbf{err}_{L^2}	$\mathbf{err}_{H_0^1}$	Iterations	Time
Projection (PG)	1.281e-02	1.731e-02	7.973e-02	10	0.477
Projection with stepsize (PS)	1.248e-02	8.869e-03	7.895e-02	17	0.554
Douglas-Rachford (DR1)	1.305e-02	7.878e-03	7.891e-02	99	9.137
Douglas-Rachford (DR2)	1.299e-02	8.292e-03	7.870e-02	52	3.846
Algorithm 1 with initialization by (PS)	1.247e-02	8.800e-03	7.886e-02	4	181.242
Algorithm 1 with exact initialization	1.248e-02	8.800e-03	7.887e-02	21	4351.329

TABLE 1. Fourier's law with $|\mathcal{D}_k^{\text{loc}}| = 11025$ and $h_k := \sqrt{2}/20$.

the exact solution are all similar to each other. It is the computing time that makes the most significant difference. Since the local search in Algorithm 1 considers each element of the triangulation separately, there is much computational effort for only little improvements. However, the algorithm involves several parameters, beside the tolerances ε_1 , ε_2 , and ε_3 for the update of POD-basis and the number K of nearest neighbors, in addition the size of the POD-basis and the choice of the initial point. Moreover, the exact algorithm from [6] can recursively be applied within the local search by fixing the assignment on large parts of the domain, while one applies the exact algorithm on a small amount of elements with only a few selected measured data points. A comprehensive numerical study of Algorithm 1 including the adjustment of all these parameters however goes beyond the focus of this work and gives rise to future research. But, in view of substantial differences w.r.t. the computing time, it is at least doubtful, if the local search could ever be competitive compared to the projection-based methods. For this reason, we left out Algorithm 1 for the following numerical study of a non-linear material law.

8.2. A non-linear material law. In our second numerical test, we considered the non-linear material law $\kappa : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ defined by

$$\kappa(\mathbf{w}) := (2 \tan^{-1}(\|\mathbf{w}\|^2 - 1) + 0.5\pi + 2)\mathbf{w}. \quad (8.3)$$

Note that κ is strongly monotone and globally Lipschitz so that, according to the Browder and Minty theorem, there is a unique solution in $H_0^1(\Omega)$ of (1.1) for every right hand side in $H^{-1}(\Omega)$. This time the right hand side is set to $f(x) = -\text{div}(\kappa(\nabla u(x)))$ with $u(x) = \sin(\pi x_1) \sin(\pi x_2)$, so that the exact solution of (1.1) is given by u .

For the numerical computations, we randomly sampled $|\mathcal{D}_k^{\text{loc}}| = 1000, 10000, 50000, 100000$ uniformly distributed data-points (\mathbf{r}, \mathbf{w}) fulfilling $\mathbf{r} = \kappa(\mathbf{w})$ with $\mathbf{w} \in [-4, 4]^2$. Additionally, a uniformly distributed noise $\mathbf{s}_i \in [-\bar{s}, \bar{s}]^4$, $i = 1, \dots, |\mathcal{D}_k^{\text{loc}}|$,

with $\bar{s} = 0, 0.1, 0.01, 0.001$ is randomly added to each data point. For the finite element discretization, we considered the mesh sizes $h_1 = \frac{\sqrt{2}}{50}$, $h_2 = \frac{\sqrt{2}}{100}$, and $h_3 = \frac{\sqrt{2}}{200}$.

First we investigated the choice of the initial step size in (PS). Recall that the step size is reduced, whenever the algorithm circles between two iterates. The results are shown in Table 2 and illustrate that the initial choice $\gamma = 1.4$ is favorable in terms of the value of the objective, but let the computing time increase in comparison to the original projection algorithm PG, which corresponds to the case $\gamma = 1$. The course of the objective values for $\gamma = 1.4$ produced by (PS) is illustrated in Figure 1.

$ \mathcal{D}_k^{\text{loc}} = 11025, \bar{s} = 0$			$ \mathcal{D}_k^{\text{loc}} = 5000, \bar{s} = 0.1$		
Init. γ	Iterations	Objective value	Init. γ	Iterations	Objective value
0.8	14	3.907e-03	0.8	12	9.839e-02
0.9	12	3.289e-03	0.9	12	8.913e-02
1.0	11	2.899e-03	1.0	11	7.841e-02
1.1	13	2.700e-03	1.1	12	7.420e-02
1.2	13	2.606e-03	1.2	23	6.414e-02
1.3	18	2.573e-03	1.3	47	6.155e-02
1.4	36	2.558e-03	1.4	93	6.034e-02
1.5	40	2.564e-03	1.5	112	6.208e-02
1.6	49	2.566e-03	1.6	137	6.504e-02
1.7	57	2.562e-03	1.7	168	6.816e-02
1.8	58	2.567e-03	1.8	218	7.102e-02
1.9	85	2.564e-03	1.9	317	7.054e-02
2.0	88	2.563e-03	2.0	651	6.944e-02
2.1	560	2.562e-03	2.1	763	8.738e-02
2.2	566	2.567e-03	2.2	1052	6.451e-02

TABLE 2. Testing the projection algorithm with step size (PS) with different initial values for γ with $h = \sqrt{2}/50$.

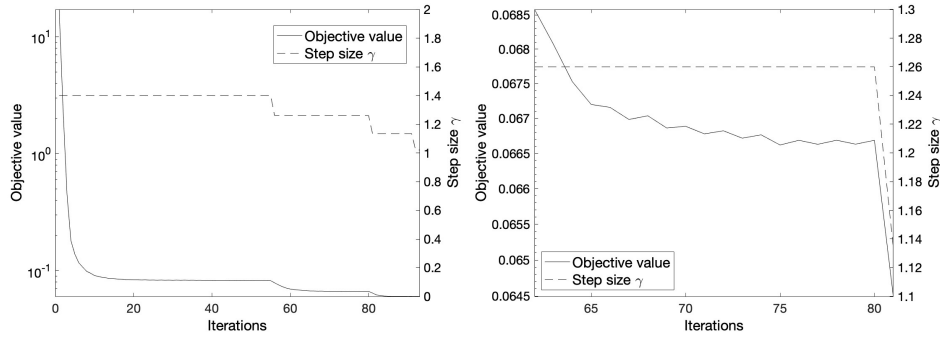


FIGURE 1. Objective values produced by (PS) with initial value $\gamma = 1.4$, $h = \sqrt{2}/50$, $|\mathcal{D}_k^{\text{loc}}| = 5000$, and $\bar{s} = 0.1$.

The errors from (8.2) as well as the objective value and the computing times for this example for the different data sets $\mathcal{D}_k^{\text{loc}}$, various noise levels \bar{s} , and different mesh sizes are shown in Table 3. In addition, the number of iterations and the computing times for the various objectives are listed. The following observations can be made:

With respect to accuracy, all algorithms deliver comparable results with (PS) being slightly superior regarding the objective and $\mathbf{err}_{H_0^1}$ and (DR2) being slightly superior regarding \mathbf{err}_{L^2} . To see this more clearly, we highlight the largest differences (compared to the result of (PS) and (DR2), respectively) concerning objective value and the errors \mathbf{err}_{L^2} and $\mathbf{err}_{H_0^1}$ in boldface.

With regard to the effort, the simple projection method (PG) is nearly always the best. This concerns the number of iterations as well as the consumed computing time. Moreover, the differences are substantial. For instance, in the line highlighted in boldface, (PG) is up to approximately 5, 18, and 8 times faster than (PS), (DR1), and (DR2), respectively.

In accordance with our theoretical findings, the objective value decreases if the noise level \bar{s} is reduced and/or if the data sample set is getting larger. Concerning the noise level, this reduction is however rather moderate. Even the largest reduction indicated with gray background is by less than 10 %, though \bar{s} is reduced from 0.1 to zero. Regarding the errors \mathbf{err}_{L^2} and $\mathbf{err}_{H_0^1}$, the dependency on the noise level is indifferent. The largest reduction, again indicated by gray background, is about 80 % w.r.t. the L^2 -error and 65 % for the H_0^1 -error. On the other hand, there are also instances, two of them marked in light gray, where the errors not only stagnate, but even increase with decreasing noise. In summary, the influence of the noise level appears to be limited and all algorithms behave comparatively robust w.r.t. noisy data.

The situation changes when the sample size is changed. With respect to the sample size, the reduction of the objective is more significant. Here, the largest reduction, when increasing the sample size from 5.000 to 10.000, marked by a box, is about 80 %. In case of the L^2 -error, the largest reduction is about one order of magnitude, while the largest reduction of $\mathbf{err}_{H_0^1}$ is approximately 35 %, both again marked by boxes. In summary, the sample size has a significantly larger impact on the accuracy of the results than the noise level.

The influence of the mesh size to the L^2 - and the H_0^1 -error is similar to that of “classical” finite element simulations. We highlight the smallest L^2 - and H_0^1 -error for zero noise level and maximum $|\mathcal{D}_k^{\text{loc}}|$ in italic type. In case of \mathbf{err}_{L^2} , the best result is obtained by (DR2), in case of $\mathbf{err}_{H_0^1}$ by (PS). In Table 4, these values are compared with a “classical” finite element computation. For the latter, we discretized (1.1) with κ as defined in (8.3) by choosing U_h from (8.1) as trial and test space. The resulting nonlinear system of equation is solved by Newton’s method. The relative errors are denoted by $\mathbf{err}_{L^2}^{\text{FE}}$ and $\mathbf{err}_{H_0^1}^{\text{FE}}$, respectively. We moreover present the experimental order of convergence defined by

$$\text{EOC}_{L^2}^{\text{FE}} := \frac{\log(\mathbf{err}_{L^2}^{\text{FE}}(h_1)) - \log(\mathbf{err}_{L^2}^{\text{FE}}(h_2))}{\log(h_1) - \log(h_2)},$$

where h_1 and h_2 are two mesh sizes and $\mathbf{err}_{L^2}^{\text{FE}}(h_i)$, $i = 1, 2$, the associated relative L^2 -errors. Furthermore, $\text{EOC}_{H_0^1}^{\text{FE}}$ is defined analogously and $\text{EOC}_{L^2}^{\text{DR2}}$ and $\text{EOC}_{H_0^1}^{\text{PS}}$ denote the respective orders of convergence generated by the data driven approach

with (DR2) and (PS), respectively. As the theory predicts, cf. e.g. [5], Table 4

N	$\mathbf{err}_{L^2}^{\text{FE}}$	$\text{EOC}_{L^2}^{\text{FE}}$	$\mathbf{err}_{L^2}^{\text{DR2}}$	$\text{EOC}_{L^2}^{\text{DR2}}$	$\mathbf{err}_{H_0^1}^{\text{FE}}$	$\text{EOC}_{H_0^1}^{\text{FE}}$	$\mathbf{err}_{H_0^1}^{\text{PS}}$	$\text{EOC}_{H_0^1}^{\text{PS}}$
50	1.601e-03	—	1.598e-03	—	3.143e-02	—	3.174e-02	—
100	4.004e-04	1.9995	3.616e-04	2.1438	1.571e-02	1.0005	1.637e-02	0.9552
200	1.001e-04	2.0000	1.343e-04	1.4289	7.854e-03	1.0002	9.078e-03	0.8506

TABLE 4. Relative errors and experimental orders of convergence for the “classical” finite element solution and the best results of the data-driven approach in dependence of the mesh sizes $h = \frac{\sqrt{2}}{N}$.

shows a quadratic and a linear order of convergence for the relative errors $\mathbf{err}_{L^2}^{\text{FE}}$ and $\mathbf{err}_{H_0^1}^{\text{FE}}$, respectively. We moreover observe that the errors produced by the data driven approach are more or less of the same size as the “classical” finite element error except for the smallest mesh with $N = 200$, where the error caused by the sample size becomes predominant compared to the error induced by the mesh.

To summarize, it is to be noted that the accuracy of all algorithms improve, if the noise level is reduced, the data sample set is enlarged, and the mesh is refined, with the noise level having the smallest impact on the accuracy. We moreover observe that all algorithms yield satisfactorily results w.r.t. the accuracy, the best results being even comparable to a “classical” finite element computation. Concerning the performance of the algorithms, the original projection algorithm (PG) turns out to be superior in the sense that it provides an accuracy similar to the other algorithms with the fastest computing time. The proximal gradient method defined by (PS) in average returns the most accurate results (in particular w.r.t. the H_0^1 -seminorm) but with a greater computational effort. Both variants of the Douglas-Rachford algorithm also provide comparable results but one has to be cautious concerning the termination criterion due to the lack of fixed points.

9. CONCLUSION AND OUTLOOK

In this paper, we studied a finite element discretization of the data driven approach as introduced [18], where we focused on a stationary scalar diffusion problem. We showed that the finite element error analysis can be incorporated into the data convergence analysis of [10] as long as finite elements are used where a vanishing discretized divergence implies that the continuous divergence vanishes, too, cf. Lemma 3.3(i). In the conductivity example, i.e., our stationary scalar diffusion problem, the construction of such elements is comparatively simple and, for instance, Raviart-Thomas elements do the job, see Proposition 3.7. The situation changes, if one turns to problems in elasticity, where the symmetry of the stress tensor significantly complicates the construction of such elements, see for instance [4, Section 4] and [2]. The considerations in Section 6 however indicate that the use of “classical” piecewise linear and continuous finite elements might not be feasible in context of the data driven approach, an aspect that gives rise to future research. Nevertheless, if $H(\text{div})$ -conforming finite elements fulfilling Assumptions 3.1 are used, then, under suitable assumptions on the approximation of the local data set, see (5.7), we obtain the same approximation results as in [10] so that (subsequences of) minimizers of the finite dimensional problems (\mathbf{P}_k) converge in data to elements of the intersection $\mathcal{D} \cap \mathcal{E}$.

Another issue concerns the computation of minimizers of (P_k) . As seen in Section 7.2, the discretized data driven problem (P_k) is equivalent to a quadratic semi-assignment problem and, as such, NP-hard, see Remark 7.1. There is thus no hope to find an efficient algorithm for the computation of minimizers. We therefore presented two heuristic approaches, one based on projections on \mathcal{E}_h and \mathcal{D}_h , the other one based on a local search in combination with model order reduction. It turns out that, the local search is not competitive, at least for the example of Fourier’s law. There are however plenty of parameters to adjust, such as the size of the POD basis and the neighborhood within the local search, and it requires further investigations to analyze if a smart choice of the parameters could result in a competitive algorithm. With regard to the projection-type methods, the most simple alternating projection algorithm according to [18] turns out to be advantageous with respect to the ratio of accuracy and computational time. With regard to accuracy only, the modified projection method (PS) and a variant of the Douglas-Rachford algorithm delivered the best results. If the noise level is zero and the sample size is large enough, then these results are comparable to a “classical” finite element simulation with known material law. There is however no theoretical evidence for a convergence of the projection-based methods, even not to any kind of fixed-point, not to mention a minimizer of (P_k) , and the examples in [17] show that this is indeed an issue. The robustification of projection-based methods for a reliable solution of (P_k) is probably one of the most important open problems in the context of data driven numerics.

REFERENCES

- [1] F. J. Aragón Artacho, R. Campoy, and M. K. Tam. The douglas–rachford algorithm for convex and nonconvex feasibility problems. *Mathematical Methods of Operations Research*, 91(2), 2020.
- [2] D.N. Arnold, F. Brezzi, and J. Douglas. PEERS: A new mixed finite element for plane elasticity. *Japan J. Appl. Math.*, 1:347–367, 1984.
- [3] A. Beck. *First-Order Methods in Optimization*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2017.
- [4] D. Braess. *Finite Elements: Theory, Fast Solvers, and Applications in Solid Mechanics*. Cambridge University, 3. edition, 2007.
- [5] C.S. Brenner and R.L. Scott. *The Mathematical Theory of Finite Element Methods*. Springer-Verlag, New York, 1994.
- [6] C. Buchheim and E. Traversi. Quadratic combinatorial optimization using separable underestimators. *INFROMS Journal on Computing*, 30(3):424–437, 2018.
- [7] P. Carrara, L. De Lorenzis, L. Stainier, and M. Ortiz. Data-driven fracture mechanics. *Computer Methods in Applied Mechanics and Engineering*, 372:113390, 2020.
- [8] C. Clason and T. Valkonen. Introduction to nonsmooth analysis and optimization, 2020.
- [9] S. Conti, S. Müller, and M. Ortiz. Data-driven problems in elasticity. *Archive for Rational Mechanics and Analysis*, 229(1):79–123, Jan 2018.
- [10] S. Conti, S. Müller, and M. Ortiz. Data-driven finite elasticity. *Archive for Rational Mechanics and Analysis*, 237(1):1–33, 2020.
- [11] R. G. Durán. Mixed finite element methods. In *Mixed Finite Elements, Compatibility Conditions, and Applications: Lectures given at the C.I.M.E. Summer School held in Cetraro, Italy June 26–July 1, 2006*, pages 1–44. Springer Berlin Heidelberg, 2008.
- [12] R. Eggersmann, T. Kirchdoerfer, S. Reese, L. Stainier, and M. Ortiz. Model-free data-driven inelasticity. *Computer Methods in Applied Mechanics and Engineering*, 350:81–99, 2019.
- [13] R. Eggersmann, L. Stainier, M. Ortiz, and S. Reese. Model-free data-driven computational mechanics enhanced by tensor voting. *Computer Methods in Applied Mechanics and Engineering*, 373:113499, 2021.

- [14] Martin Kahlbacher and Stefan Volkwein. Galerkin proper orthogonal decomposition methods for parameter dependent elliptic systems. *Discussiones Mathematicae, Differential Inclusions, Control and Optimization*, 27(1):95–117, 2007.
- [15] Y. Kanno. Data-driven computing in elasticity via kernel regression. *Theoretical and Applied Mechanics Letters*, 8:361–365, Dec 2018.
- [16] Y. Kanno. Simple heuristic for data-driven computational elasticity with material data involving noise and outliers: a local robust regression approach. *Japan Journal of Industrial and Applied Mathematics*, 3(1085–1101), 2018.
- [17] Y. Kanno. Mixed-integer programming formulation of a data-driven solver in computational elasticity. *Optimization Letters*, 13(7):1505–1514, 2019.
- [18] T. Kirchdoerfer and M. Ortiz. Data-driven computational mechanics. *Computer Methods in Applied Mechanics and Engineering*, 304:81–101, 2016.
- [19] T. Kirchdoerfer and M. Ortiz. Data-driven computing in dynamics. *International Journal for Numerical Methods in Engineering*, 113, Jun 2017.
- [20] C. Natemeyer and D. Wachsmuth. A proximal gradient method for control problems with non-smooth and non-convex control cost. *Computational Optimization and Applications*, 80(2):639–677, 2021.
- [21] L.T.K. Nguyen and M.-A. Kneip. A data-driven approach to nonlinear elasticity. *Computers & Structures*, 194:97–115, 2018.
- [22] K. Poelstra, T. Bartel, and B. Schweizer. A data driven framework for evolutionary problems in solid mechanics. No. 2021-02, Preprints der Fakultät für Mathematik, TU Dortmund, 2021.
- [23] M. Röger and B. Schweizer. Relaxation analysis in a data driven problem with a single outlier. *Calc. Var.*, 59(119), 2020.
- [24] S. Sahni and T. Gonzalez. P-complete approximation problems. *Journal of the ACM (JACM)*, 23(3):555–565, 1976.
- [25] R. Temam. *Mathematical Problems in Plasticity*. Gauthier-Villars, Paris, 1985.
- [26] Sebastian Ullmann, Marko Rotkvic, and Jens Lang. Pod-galerkin reduced-order modeling with adaptive finite element snapshots. *Journal of Computational Physics*, 325:244–258, 2016.
- [27] D. Wachsmuth. Iterative hard-thresholding applied to optimal control problems with $L^0(\Omega)$ control cost. *SIAM J. Control Optim.*, 57(2):854–879, 2019.

TECHNISCHE UNIVERSITÄT DORTMUND, FAKULTÄT FÜR MATHEMATIK, LEHRSTUHL LSX, VOGELPOTHSWEG 87, 44227 DORTMUND, GERMANY
Email address: annika.mueller@tu-dortmund.de

TECHNISCHE UNIVERSITÄT DORTMUND, FAKULTÄT FÜR MATHEMATIK, LEHRSTUHL LSX, VOGELPOTHSWEG 87, 44227 DORTMUND, GERMANY
Email address: christian2.meyer@tu-dortmund.de