

PROXIMAL GALERKIN: A STRUCTURE-PRESERVING FINITE ELEMENT METHOD FOR POINTWISE BOUND CONSTRAINTS

BRENDAN KEITH* AND THOMAS M. SUROWIEC†

Dedicated with respect and admiration to Leszek Demkowicz on the occasion of his 70th birthday anniversary.

Abstract. The proximal Galerkin finite element method is a high-order, low iteration complexity, nonlinear numerical method that preserves the geometric and algebraic structure of pointwise bound constraints in infinite-dimensional function spaces. This paper introduces the proximal Galerkin method and applies it to solve free boundary problems, enforce discrete maximum principles, and develop a scalable, mesh-independent algorithm for optimal design problems with pointwise bound constraints. This paper also provides a derivation of the latent variable proximal point (LVPP) algorithm, an unconditionally stable alternative to the interior point method. LVPP is an infinite-dimensional optimization algorithm that may be viewed as having an adaptive barrier function that is updated with a new informative prior at each (outer loop) optimization iteration. One of its main benefits is witnessed when analyzing the classical obstacle problem. Therein, we find that the original variational *inequality* can be replaced by a sequence of second-order partial differential *equations* (PDEs) that are readily discretized and solved with, *e.g.*, high-order finite elements. Throughout this work, we arrive at several unexpected contributions that may be of independent interest. These include (1) a semilinear PDE we refer to as the *entropic Poisson equation*; (2) an algebraic/geometric connection between high-order positivity-preserving discretizations and certain infinite-dimensional Lie groups; and (3) a gradient-based, bound-preserving algorithm for two-field, density-based topology optimization. The complete latent variable proximal Galerkin methodology combines ideas from nonlinear programming, functional analysis, tropical algebra, and differential geometry and can potentially lead to new synergies among these areas as well as within variational and numerical analysis. This work is accompanied by open-source implementations of our methods to facilitate reproduction and broader adoption.

1. Introduction. Although the origins of variational analysis can be traced back at least to the seventeenth century [164], its role in the modern study of partial differential equations (PDEs) only truly began to take shape around 1847 once William Thomson introduced what is now known as the Dirichlet principle. In contemporary language, this energy principle states that for all functions $f \in L^2(\Omega)$ and $g \in H^1(\Omega)$, the (weak) solution of Poisson’s equation over a Lipschitz domain $\Omega \subset \mathbb{R}^n$,

$$(1.1) \quad -\Delta u = f \quad \text{in } \Omega, \quad u = g \quad \text{on } \partial\Omega,$$

can be obtained as the $H^1(\Omega)$ -minimizer of the Dirichlet energy,

$$(1.2) \quad E(v) = \frac{1}{2} \int_{\Omega} |\nabla v|^2 dx - \int_{\Omega} v f dx,$$

confined to the constraint set $H_g^1(\Omega) = g + H_0^1(\Omega) = \{v \in H^1(\Omega) \mid v = g \text{ on } \partial\Omega\}$.

Owing to the fact that $H_g^1(\Omega)$ is nonempty, closed, and convex, it is a straightforward consequence of the Lions–Stampacchia theorem [183, 140] that the energy minimizer $u^* \in K = H_g^1(\Omega)$ is the unique solution to the variational inequality (VI)

$$(1.3) \quad \int_{\Omega} \nabla u^* \cdot \nabla v dx \geq \int_{\Omega} f v dx \quad \text{for all } v \in K - u^*.$$

*Division of Applied Mathematics, Brown University, Providence, RI 02912 USA (brendan.keith@brown.edu).

†Simula Research Laboratory, Department of Numerical Analysis and Scientific Computing, Kristian Augusts gate 23, 0164, Oslo, Norway, (thomasms@simula.no)

As it happens, the boundary condition in (1.1) is an equality constraint that induces an *affine* structure on the feasible set. Moreover, it is this particular algebraic structure that can be exploited to show that the minimizer $u^* \in H_g^1(\Omega)$ is also uniquely characterized by a variational *equation*. In the setting above, we have

$$(1.4) \quad \int_{\Omega} \nabla u^* \cdot \nabla w \, dx = \int_{\Omega} f w \, dx \quad \text{for all } w \in H_0^1(\Omega).$$

To arrive at this conclusion from (1.3), the key idea is to notice that $H_g^1(\Omega) + H_0^1(\Omega) = H_g^1(\Omega)$ and, therefore, one may replace v in (1.3) by $u^* \pm w$, for any $w \in H_0^1(\Omega)$. For further details, see, e.g., [38, Proposition 9.22] and [48, Theorem 1.2.2].

When partial differential equations (PDEs) are first written down, their essential boundary conditions are often the only explicit constraints that appear on the space of solutions. As such, all students of the finite element method are taught to derive variational equations like (1.4); cf. [102, Section 1.4], [48, Exercise 1.2.2], and [68, Section 31.2.2]. In turn, viewing the variational equations that originate from essential boundary conditions as structure-exploiting reductions of more general VIs may seem esoteric to some practitioners. Yet, the motive becomes clear when the feasible set has alternative algebraic structures.

To illustrate this point, it is instructive to consider imposing a pointwise non-negativity constraint, $u^* \geq 0$, and setting $g \equiv 0$. Hereafter, let $H_+^1(\Omega) = \{v \in H^1(\Omega) \mid v \geq 0 \text{ a.e.}\}$. Thus, the subset

$$(1.5) \quad K = \{v \in H_0^1(\Omega) \mid v \geq 0 \text{ a.e.}\} = H_0^1(\Omega) \cap H_+^1(\Omega),$$

forms a closed convex cone in $H^1(\Omega)$. It is well-known that the *conic* structure of K allows us to write

$$(1.6) \quad \int_{\Omega} \nabla u^* \cdot \nabla v \, dx \geq \int_{\Omega} f v \, dx \quad \text{for all } v \in K,$$

with equality holding for $v = u^*$; see, e.g., [48, Theorem 1.1.2]. Although we acknowledge this simplified variational formulation, we present it only as a motivating example. Our work pursues a different type of algebraic structure.

Our aim is to provide a *multiplicative* structure-preserving approach to solving bound-constrained optimization problems and variational inequalities in Sobolev spaces. This will lead us to working in *Banach algebras*, which are Banach spaces that are closed under a continuous multiplication operation [51]. Instead of performing Lagrangian relaxation or relying on penalty functions, the key component of our approach is an adaptive form of *entropy regularization*. Through entropy regularization, we will find, e.g., that minimizing the Dirichlet energy over functions in $H_g^1(\Omega) \cap H_+^1(\Omega)$ can be reduced to solving a “Bayesian” sequence of second-order semilinear PDEs where each right-hand side is biased by the prior solution.

Algorithm 1 outlines the meta-algorithm for the pointwise non-negativity constraint $u^* \geq 0$ considered above when $f \in L^\infty(\Omega)$ and $g|_{\partial\Omega} \in C(\partial\Omega)$ satisfies $\text{ess inf}_{\partial\Omega} g > 0$. Note that, unlike, e.g., descent methods [208, Section 3], this algorithm converges for *all* step size values $\alpha > 0$; cf. **Theorem 4.13**. A practical version of the algorithm is readily implementable (see, e.g., [111, 113]) and reduces to the finite element method that gives this paper its name.

Remark 1.1 (Proximal Galerkin). To differentiate between the method proposed in this paper and a significantly different “proximal Galerkin” method proposed in

Algorithm 1: Entropic proximal point algorithm for Dirichlet energy minimization with a non-negativity constraint.

input: Step size parameter $\alpha > 0$ and initial solution guess $w \in H_g^1(\Omega) \cap L^\infty(\Omega)$ satisfying $\text{ess inf } w > 0$.

repeat

Solve the entropic Poisson equation,

$$(1.7) \quad \begin{cases} -\Delta u + \alpha^{-1} \ln u = f + \alpha^{-1} \ln w & \text{in } \Omega, \\ u = g & \text{on } \partial\Omega. \end{cases}$$

Assign $w \leftarrow u$.

until a convergence test is satisfied

[146], we consider the full name of our method to be the “*latent variable proximal Galerkin*” method. Throughout the text, we typically use the abbreviated name “proximal Galerkin” without any cause for confusion.

1.1. Notation. Our notation is rather standard for the finite element literature. Norms are denoted by $\|\cdot\|_X$, inner products by $(\cdot, \cdot)_X$, and duality pairings by $\langle \cdot, \cdot \rangle_{X', X}$ for spaces X and its paired topological dual X' . Whenever it is clear in context, we leave off or abbreviate the subscripts in a natural way. Norm convergence will typically be denoted by \rightarrow or \xrightarrow{X} , for weak convergence we use the standard $\overset{X}{\rightharpoonup}$ or \rightharpoonup . For subsets C of infinite dimensional spaces, we denote the closure by $\text{cl } C$, the boundary by $\text{bd } C$, and the interior by $\text{int } C$. For a mapping F between normed linear spaces X and Y , the Fréchet derivative of F at x is indicated by $F'(x)$.

For an open bounded domain $\Omega \subset \mathbb{R}^n$, $n \in \{1, 2, \dots\}$, $L^p(\Omega)$, $p \in [1, \infty]$, denotes the usual Lebesgue space of (equivalence classes of) p -integrable functions when $p \in [1, \infty)$, and essentially bounded functions when $p = \infty$, respectively. Furthermore, we define

$$L_+^p(\Omega) := \{u \in L^p(\Omega) \mid u \geq 0 \text{ a.e. in } \Omega\}$$

for any $p \in [1, \infty]$.

For domains Ω in \mathbb{R}^d , we denote the boundary by $\partial\Omega$ and the closure by $\overline{\Omega}$. The spaces $L^p(\partial\Omega)$ are defined in the usual way. When needed, we indicate the surface measure for $\partial\Omega$ by $d\mathcal{H}_{n-1}$. The space of continuous functions on $\overline{\Omega}$ is written $C(\overline{\Omega})$. Similarly, $C^m(\overline{\Omega})$, $m \in \mathbb{N} \cup \{\infty\}$, is the space of all m -times continuously differentiable functions. The space of smooth compactly supported test functions on Ω is given by $C_c^\infty(\Omega)$. The Sobolev space of $L^2(\Omega)$ functions with $L^2(\Omega)$ integrable weak derivatives is denoted by $H^1(\Omega)$ and its closed subspace of functions u with trace $\gamma u = 0$ is denoted by $H_0^1(\Omega)$. We use $H^s(\partial\Omega)$, $s \in (0, 1)$, for the usual Sobolev–Slobodeckij space on $\partial\Omega$. We refer the reader to a standard text on function spaces for further details, e.g., [4, 154]. Finally, for simplicity of notation, we adopt the following notational conventions: $0 \ln 0 := 0$ and $\|v\|_{H^{-1}(\Omega)} := \sup_{w \in H_0^1(\Omega)} \frac{(v, w)}{\|\nabla w\|_{L^2(\Omega)}}$.

1.2. Outline. We have attempted to provide a scaffolded presentation of our findings. To this end, Section 2 presents preliminary concepts and provides further motivation for this work. Next, Section 3 reviews the literature and summarizes our main contributions. Sections 4 through 6 present the essential features of proximal

Galerkin methods for the obstacle problem, the advection-diffusion equation, and topology optimization, respectively. Each of these sections contains an algorithm that is designed to be implemented in a production-level finite element code. The reader is encouraged to compare these algorithms with our publicly available implementations [111, 112, 113]. The main paper closes with ??, which contains a small number of concluding remarks, and then proceeds to two technical appendices. [Appendix A](#) contains the continuous-level mathematical analysis and [Appendix B](#) contains the discrete-level, finite element theory. [Appendices A](#) and [B](#) are the most specialized sections of the paper and may be skipped by a casual reader.

2. Preserving multiplicative structure. The proximal Galerkin finite element method is a nonlinear numerical method that preserves the algebraic and geometric structure of bound constraints in infinite-dimensional function spaces. In this section, we study the multiplicative structure of non-negative functions and use the Dirichlet energy (1.2) to illustrate how proximal Galerkin preserves this structure.

2.1. Deconstructing the semiring of non-negative functions. We discuss the natural logarithmic transformation between non-negativity constraints and extended real-valued functions that may take the value $-\infty$. This also first introduces the latent variable ψ . We claim this provides a basis for the use of logarithmic transformations to analyze and solve PDEs, an idea that goes back at least to work by Schrödinger in 1926. Finally, we address the somewhat unnatural conditions this transformation imposes on the solution spaces and variational equations themselves. In turn, we show how a simple regularization of the transformed equations remedies these concerns. We use this discussion to motivate the natural function spaces for pointwise bound constraints in $H^1(\Omega)$ and construct a direct link to entropy regularization.

Let \mathcal{X} be a set equipped with two binary operations: addition $\oplus: \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{X}$ and multiplication $\odot: \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{X}$.

DEFINITION 2.1 (Semiring). *We say that \mathcal{X} is a semiring if the following conditions are satisfied [79, 81]:*

- *Addition \oplus and multiplication \odot are associative;*
- *Addition \oplus is commutative;*
- *Multiplication \odot is distributive with respect to addition \oplus .*

We say that \mathcal{X} is a commutative semiring if the conditions above are satisfied and, moreover, multiplication \odot is commutative.

It is easy to check that the set of non-negative Lebesgue measurable functions,

$$(2.1) \quad M_+(\Omega) = \{v: \Omega \rightarrow \mathbb{R}_+ \mid \{v > c\} \text{ is Lebesgue measurable } \forall c > 0\},$$

forms a commutative semiring under the standard binary operations of pointwise addition and multiplication. In particular, note that for any $u, v \in M_+(\Omega)$, we have

$$(2.2) \quad u + v \in M_+(\Omega), \quad uv \in M_+(\Omega).$$

There is an interesting identification between $M_+(\Omega)$ and the space of extended real-valued measurable functions

$$(2.3) \quad M_{\max}(\Omega) = \{\varphi: \Omega \rightarrow \mathbb{R} \cup \{-\infty\} \mid \{\varphi > c\} \text{ is Lebesgue measurable } \forall c \in \mathbb{R}\},$$

induced by the (pointwise) logarithm and exponential operators. Namely, for all $u \in M_+(\Omega)$, $\psi \in M_{\max}(\Omega)$, and $\alpha > 0$, we have that $\alpha^{-1} \ln u \in M_{\max}(\Omega)$ and

$\exp(\alpha\psi) \in M_+(\Omega)$ under the convention that $\ln 0 = -\infty$ and, likewise, $\exp(-\infty) = 0$. Such logarithmic transformations provide a family of semiring isomorphisms between $M_+(\Omega)$ and $M_{\max}(\Omega)$, where $M_{\max}(\Omega)$ is endowed with the following (generalized) addition and multiplication operations:

$$(2.4) \quad \psi \oplus \varphi = \alpha^{-1} \ln(\exp(\alpha\psi) + \exp(\alpha\varphi)), \quad \psi \odot \varphi = \psi + \varphi,$$

respectively [151, 152]. Moreover, in the limit $\alpha \rightarrow \infty$, the generalized addition operation (2.4) becomes the pointwise maximum operation [142, 147]; namely,

$$(2.5) \quad \psi \oplus \varphi \rightarrow \max\{\psi, \varphi\}.$$

Logarithmic transformations of the above form have been used famously over the last century to analyze differential equations in quantum mechanics [180], fluid flow [98, 50], and electrical engineering [177, 148], and, more recently, to study stochastic PDEs [29, 88]. Given that they appear to capture certain key algebraic properties of the set of non-negative functions, it is tempting to use logarithmic transformations to enforce non-negativity constraints on function spaces. Unfortunately, however, special care is required to apply a logarithmic transformation to a non-negative solution variable in a free-boundary problem.

For illustration, consider minimizing the Dirichlet energy (1.2) over the set of non-negative functions

$$(2.6) \quad K = \{v \in H_g^1(\Omega) \mid v \geq 0 \text{ a.e.}\} = H_g^1(\Omega) \cap H_+^1(\Omega).$$

Assuming that $f \in L^2(\Omega)$ and $u^* \in H^2(\Omega)$, the well-known complementarity conditions for the solution are as follows [116, p. 79]:

$$(2.7) \quad u^* \geq 0, \quad -\Delta u^* - f \geq 0, \quad (\Delta u^* + f) u^* = 0 \text{ a.e. in } \Omega.$$

Another perspective uses a dual variable λ^* , also known as a Lagrange multiplier, to formulate (2.7) as a mixed complementarity problem of the form:

$$(2.8) \quad -\Delta u^* - \lambda^* = f, \quad u^* \geq 0, \quad \lambda^* \geq 0, \quad \langle u^*, \lambda^* \rangle = 0.$$

The Lagrange multiplier exhibits rather low regularity for general domains Ω , so the term “ $\lambda^* \geq 0$ ” is actually understood to mean $\langle \lambda^*, w \rangle \geq 0$ for all $w \in H_0^1(\Omega)$ with $w \geq 0$ a.e. in Ω , i.e., without further regularity assumptions λ^* is merely a nonnegative Radon measure on Ω . See [116, Chap. II, Sec. 6] for details.

If we wish to study this problem under a logarithmic transformation, then a formal computation using the substitution $u^* = \exp \psi^*$ leads to the observation that

$$(2.9) \quad \psi^* = -\infty \quad \text{or} \quad -\operatorname{div}(\exp \psi^* \nabla \psi^*) = f,$$

at almost every point in Ω . Analyzing these equations presents challenges, in part, because it requires moving away from well-studied Sobolev spaces [4] and, instead, working in a space of extended real-valued functions [117] endowed with the metric

$$(2.10) \quad d(\psi, \varphi) = \|\nabla \exp \psi - \nabla \exp \varphi\|_{L^2(\Omega)}.$$

One conclusion of this work is that the above concerns are alleviated by a simple regularization of the degenerate PDE in (2.9). In particular, we show that for all bounded $f \in L^\infty(\Omega)$, the latent solution variable $\psi^* = \ln u^*$ is recovered as the

$\alpha \rightarrow \infty$ limit (with respect to the metric (2.10)) of a family of regularized solutions $\psi \in H^1(\Omega) \cap L^\infty(\Omega)$ satisfying

$$(2.11) \quad -\operatorname{div}(\exp \psi \nabla \psi) + \alpha^{-1} \psi = f.$$

Moreover, the latent variable iteration $\psi^0 \in H^1(\Omega) \cap L^\infty(\Omega)$,

$$(2.12) \quad -\operatorname{div}(\exp \psi^k \nabla \psi^k) + \alpha^{-1} \psi^k = f + \alpha^{-1} \psi^{k-1}, \quad k = 1, 2, \dots,$$

formerly written with primal variables in Algorithm 1, converges to ψ^* for all finite $\alpha > 0$; cf. Theorem 4.13.

The ambient function space for the regularized latent variable ψ is interesting from an algebraic point of view because it is a *Banach algebra*. Indeed, the Sobolev subspace $H^1(\Omega) \cap L^\infty(\Omega)$, whose norm is

$$(2.13) \quad \|v\|_{H^1(\Omega) \cap L^\infty(\Omega)} = \max\{\|v\|_{H^1(\Omega)}, \|v\|_{L^\infty(\Omega)}\},$$

is closed under the standard operations of pointwise addition and multiplication [38, Proposition 9.4]. Maintaining closure under multiplication is desirable, in part, because it often allows one to construct a smooth exponential map [75, 76]. Indeed, of particular interest to this work is the Nemytskii operator

$$(2.14) \quad \exp: H^1(\Omega) \cap L^\infty(\Omega) \rightarrow H^1(\Omega) \cap \operatorname{int} L_+^\infty(\Omega),$$

which is an isomorphism between $H^1(\Omega) \cap L^\infty(\Omega)$ and the *Banach–Lie group*

$$(2.15) \quad H^1(\Omega) \cap \operatorname{int} L_+^\infty(\Omega) = \{w \in H^1(\Omega) \cap L^\infty(\Omega) \mid \operatorname{ess\,inf} w > 0\};$$

cf. Proposition A.9. Since the range of this isomorphism is contained in the $H^1(\Omega) \cap L^\infty(\Omega)$ -interior of the set of essentially bounded non-negative $H^1(\Omega)$ functions, we find that the primal iterates,

$$u^k = \exp \psi^k \in H^1(\Omega) \cap \operatorname{int} L_+^\infty(\Omega) \subset \operatorname{int}(H^1(\Omega) \cap L_+^\infty(\Omega)),$$

will always be *interior points*. In the next motivational subsection, we explain that an identical sequence of interior points $u^k \xrightarrow{H^1(\Omega)} u^*$ can be found by regularizing the Dirichlet energy with an appropriate *entropy* functional.

2.2. Dirichlet free energy. Only special function spaces are endowed with a norm topology that permits a continuous multiplication operator. Indeed, it is well-known that $H^1(\Omega)$ is only closed under multiplication when $n = 1$ [4]. Moreover, it is easy to show that $\operatorname{int} H_+^1(\Omega) = \emptyset$ for all $n \geq 2$, which makes it impossible to define an $H^1(\Omega)$ -interior point in any of its subsets (cf. Remark 4.3). Because $H^1(\Omega) \cap L^\infty(\Omega)$ bypasses both of these topological issues, it is appealing to restrict the feasible set K in (2.6) to essentially bounded functions when minimizing the Dirichlet energy.

Unfortunately, requiring the feasible set to be the intersection of K and $L^\infty(\Omega)$ would cause the direct method of calculus of variations [21] to fail. This is because the Dirichlet energy does not provide control over point-wise values of the solution and $K \cap L^\infty(\Omega)$ is not closed in the $H^1(\Omega)$ norm topology. Therefore, one may conclude that maintaining some important mathematical structures is in conflict with the classical energy principle.

Fortunately, it turns out there is resolution to this conflict that exposes the missing algebraic structure; namely, minimizing the *Dirichlet free energy*,

$$(2.16) \quad A(u) = E(u) + \theta S(u).$$

Here, $\theta = \alpha^{-1} > 0$ is a non-dimensional “temperature” parameter and

$$(2.17) \quad S(u) = \int_{\Omega} u \ln u - u \, dx,$$

is the (negative) entropy functional. As we show in [Theorem 4.7](#), all minimizers of [\(2.16\)](#) lie in $K \cap L^{\infty}(\Omega)$, i.e., for all $\theta > 0$,

$$(2.18) \quad u = \arg \min_{v \in K} A(v) = \arg \min_{v \in K \cap L^{\infty}(\Omega)} A(v),$$

is essentially bounded away from zero in Ω , and $u = u(\theta)$ converge linearly with respect to θ to the unique non-negative minimizer of [\(1.2\)](#), $u^* = \arg \min_{v \in K} E(v)$. More specifically, each u is an interior point and

$$(2.19) \quad \|\nabla u^* - \nabla u\|_{L^2(\Omega)}^2 \leq 2\theta(S(u^*) + |\Omega|),$$

whenever $f \in L^{\infty}(\Omega)$; cf. [Theorem A.15](#).

Finally, and just as importantly, the general VI that characterizes [\(2.18\)](#), i.e.,

$$(2.20) \quad \int_{\Omega} \nabla u \cdot \nabla v \, dx + \theta \int_{\Omega} v \ln u \, dx \geq \int_{\Omega} f v \, dx \quad \text{for all } v \in K - u,$$

can be replaced by a variational equality for the weak form of a semilinear PDE we call the *entropic Poisson equation*, $-\Delta u + \theta \ln u = f$; namely,

$$(2.21) \quad \int_{\Omega} \nabla u \cdot \nabla w \, dx + \theta \int_{\Omega} w \ln u \, dx = \int_{\Omega} f w \, dx \quad \text{for all } w \in H_0^1(\Omega).$$

The entropic Poisson equation is the primal form of [\(2.11\)](#) and has numerous interesting properties that we exploit in this work. We also note that $\theta \ln(1/u)$ approximates the true Lagrange multiplier λ^* introduced in [\(2.8\)](#) above.

The essential idea presented above is expanded on in [Sections 5](#) and [6](#) to accommodate bound constraints for general VIs that do not appear as a result of energy principles, as well as those that appear in topology optimization with a view toward other bound-constrained optimization problems. Crucially, and unlike traditional penalty or barrier methods [[163](#), [30](#), [207](#)], it is *not necessary* to take $\theta \rightarrow 0$ in order to get an arbitrarily accurate approximation of u^* . Indeed, the simple adaptive entropic regularization algorithm given in [Algorithm 1](#) (see also [\(2.12\)](#)), which comes from regularizing the Dirichlet energy with a *relative entropy* functional, is far more appealing and is derived in [Section 4](#). [Figure 2.1](#) provides a diagrammatic reference for the main elements of the continuous-level algorithm in the case $\alpha = 1$.

Remark 2.2 (Dirichlet free energy). We propose the name “Dirichlet free energy” for the functional in [\(2.16\)](#) by analogy with the Helmholtz free energy from statistical mechanics [[168](#)], $A = E - TS$, where E denotes total system energy, T denotes absolute temperature, and S denotes thermodynamic entropy.

2.3. Pointwise-positivity for every polynomial degree. The majority of this paper is based on pursuing the aforementioned observation that the solution of VIs for bound constraints, including [\(2.20\)](#), can be approximated arbitrarily accurately by variational *equations* like [\(2.21\)](#). Leveraging this observation for computational purposes leads to a new class of high-order, nonlinear finite element methods we refer to

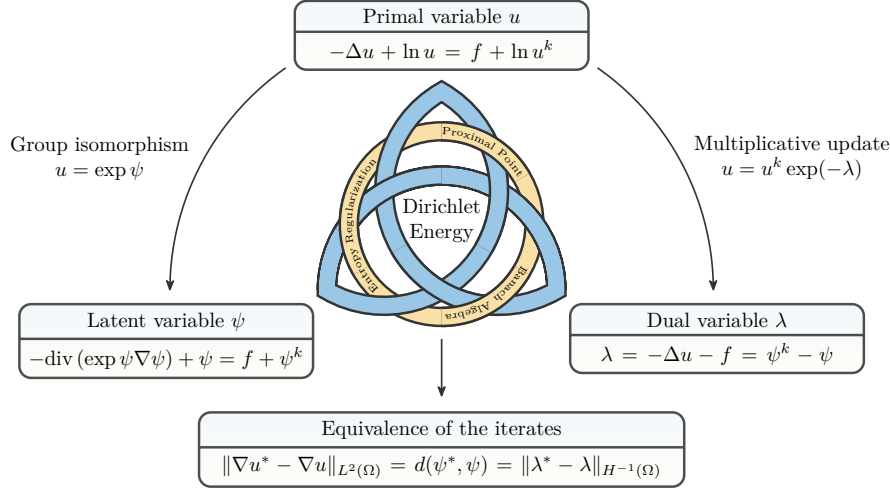


Fig. 2.1: A trinity is formed by the three isomorphic representations of the iterates in the latent variable proximal point method. In this figure, equations for the three representations are given for the problem of minimizing the Dirichlet energy (1.2) over non-negative functions $u \in H_g^1(\Omega) \cap H_+^1(\Omega)$. Note that, for simplicity, the step size here is set to $\alpha = 1$. See Theorems 4.7 and 4.13 for further details and consequences for variable step sizes.

as proximal Galerkin methods. In turn, taking advantage of the multiplicative structure of the solution space leads to non-standard approximation spaces that naturally preserve pointwise positivity at the discrete level.

As we shall also show in Section 4, a very convenient Galerkin discretization of the entropic Poisson equation (1.7) is found by introducing a pair of linear subspaces $V_h \subset H_0^1(\Omega)$ and $W_h \subset L^\infty(\Omega)$ — for instance, spaces of high-degree piecewise polynomials — and simultaneously approximating the solution u in both the primal space and the latent space; namely,

$$(2.22) \quad u \approx u_h \quad \text{and} \quad u \approx \exp \psi_h,$$

where $u_h \in g_h + V_h$, $\psi_h \in W_h$, and $g_h \in H^1(\Omega)$ provides an approximation of the boundary values $g_h|_{\partial\Omega} \approx g|_{\partial\Omega}$. The basic method is outlined in Algorithm 2.

A novelty of the approximate solution $\tilde{u}_h = \exp \psi_h$ is that it is *guaranteed* to deliver *pointwise positivity*. We exploit and extend the above property throughout this work to develop some of the first high-order *bound-preserving* finite element methods for a variety of benchmark problems. Another important property of this exponential discretization is that it preserves the *multiplicative* group structure of the set $\text{int } L_+^\infty(\Omega) = \{w \in L^\infty(\Omega) \mid \text{ess inf } w > 0\}$. More specifically,

$$(2.23) \quad \exp \psi_h \exp \varphi_h = \exp(\psi_h + \varphi_h) \in \exp(W_h) \subset \text{int } L_+^\infty(\Omega),$$

for all ψ_h and $\varphi_h \in W_h$. Before expanding further on this and other topics, we present a comprehensive review of the literature and an itemized list of contributions.

Algorithm 2: Proximal Galerkin method for Dirichlet energy minimization with a pointwise non-negativity constraint.

Input : Step size parameter $\alpha > 0$, linear subspaces $V_h \subset H_0^1(\Omega)$ and $W_h \subset L^\infty(\Omega)$, and initial solution guess $\psi_h \in W_h$.

Output: Two approximate solutions, u_h and $\tilde{u}_h = \exp \psi_h$, and an approximate Lagrange multiplier, $\lambda_h = (\omega_h - \psi_h)/\alpha$.

repeat

 Assign $\omega_h \leftarrow \psi_h$.

 Solve the following (nonlinear) discrete saddle-point problem:

$$\left\{ \begin{array}{l} \text{Find } u_h \in g_h + V_h \text{ and } \psi_h \in W_h \text{ such that} \\ \int_{\Omega} \alpha \nabla u_h \cdot \nabla v \, dx + \int_{\Omega} \psi_h v \, dx = \int_{\Omega} (f + \omega_h) v \, dx \quad \text{for all } v \in V_h, \\ \int_{\Omega} u_h \varphi \, dx - \int_{\Omega} \exp(\psi_h) \varphi \, dx = 0 \quad \text{for all } \varphi \in W_h. \end{array} \right.$$

until a convergence test is satisfied

3. Contributions and related work. The latent variable proximal Galerkin finite element method is as much a finite element method as it is an optimization algorithm. With this in mind, it is important to distinguish proximal Galerkin from the extensive collection of other numerical methods for bound-constrained variational problems. In turn, we choose to survey the optimization literature as well as the numerical PDE literature. The main contributions of this work are highlighted and itemized in [Subsection 3.3](#).

3.1. Optimization methods for pointwise bound constraints. Bound-constrained variational problems arise in many subjects. These include, but are not limited to, contact mechanics [115, 207], financial mathematics [191, Chap. 12], mathematical image processing [13], and the geosciences, such as glaciology [218]. It is here that we are often confronted with the requirement that the solution be pointwise bounded from above or below by some critical threshold over at least a portion of the physical domain or its boundary. In PDE-constrained optimization and optimal control, bounds on the solution of the PDE, i.e., state constraints, naturally arise as a modeling requirement, see the well-known monographs [197, 96] and the references therein, especially [42, 43, 44]. Consequently, a great deal of effort has been spent on treating bound constraints in infinite dimensions.

We mainly restrict our overview to the numerical solution of the obstacle problem (1.3), with $K \subset \{v \in H^1(\Omega) \mid v \geq \phi\}$ for some $\phi \in H^1(\Omega) \cap L^\infty(\Omega)$, since the available solvers capture the main essences of the common techniques for other bound-constrained problems, however, we note that a number of the optimization algorithms listed below are applicable far beyond this setting. Perhaps the most direct approach begins by prescribing a finite-dimensional subspace of $H^1(\Omega)$ for the discrete solution and then solving the associated variational problem by methods of nonlinear programming. In this “first-discretize-then-optimize” class of approaches, the finite-dimensional reformulation typically amounts to a strongly convex quadratic program or a discrete strongly monotone variational inequality. The fact that higher-

order basis functions face numerous challenges when used to enforce pointwise bound constraints limits the benefits of these approaches; for further discussion, see [Subsection 3.2](#). However, a wealth of viable algorithms from nonlinear programming can be applied to lowest-order discretizations; see, e.g., [\[30, 163\]](#). Nevertheless, at least for active set-based approaches, such as in [\[93\]](#), one will almost certainly observe *mesh-dependence*.

Mesh-dependence means that the number of nonlinear solver iterations required to reach a prescribed stopping tolerance (using the appropriate function space norm) will grow without bound on successively finer meshes. Nevertheless, mesh-dependence can be computationally mitigated by appealing to multigrid methods, as was done in the celebrated papers [\[36, 87, 99, 100, 101, 118, 119, 120, 204\]](#); see [\[83\]](#) for a comprehensive review. Despite the favorable behavior of these multigrid methods, there is no proof of mesh-independence in general. In particular, there is no guarantee that a given sequence of meshes will not miss low-dimensional portions (sets of positive capacity [\[116\]](#)) of the active set.

Mesh-dependence in active set methods arises from a lack of generalized differentiability of the (nonsmooth) residual in the function space setting, cf. [\[93, 198\]](#). This has motivated researchers to propose and analyze algorithms for bound-constrained problems in the continuous, i.e., infinite-dimensional setting. If an algorithm can be shown to converge in the continuous setting and the problem of interest exhibits sufficient stability properties around its solution, then this convergence will carry over to perturbed problems. At least for conforming discretizations, the associated finite-dimensional problem can be viewed as such a perturbation provided the discretization is sufficiently fine. For further material on this topic, we refer the reader to the detailed discussions and references to applications in [\[203\]](#) and the pioneering works [\[153, 9, 8\]](#).

Infinite-dimensional algorithms follow their finite-dimensional counterparts and can be roughly split into several categories: penalty methods, barrier methods, augmented Lagrangian methods, and first-order methods of convex optimization. For penalty (approximation) methods, we point the interested reader to the well-known monograph [\[77\]](#), which claims these techniques go back to [\[138, 139\]](#). However, we note that the numerical methods in [\[77\]](#), e.g., coordinate descent, are not seen to be competitive with more recent developments in the subsequent decades after its publication.

Quadratic penalty methods are used widely in PDE-constrained optimization, see, e.g., [\[94, 95, 92\]](#) and readily extendible to numerous applications; see, e.g., [\[122, 114, 3\]](#). These are often referred to as “Moreau–Yosida”-based approaches because the quadratic penalty can be viewed as the Moreau envelope of the indicator function for the bound constraints. The downside of these methods is the requirement to drive the penalty parameter to infinity to restore feasibility. Mirroring their finite-dimensional equivalents, interior point methods have also been investigated in detail for certain classes of PDE-constrained problems, see, e.g., [\[178, 199, 205, 97, 179\]](#). Our method is closer to interior point methods, due to the entropy term [\[187, 169\]](#), and somewhat related to the first-order methods in [\[194, 217\]](#). However, in contrast to traditional interior point methods, the entropy functions employed in the text below do not exclude points from the feasible sets as they are still well-defined for feasible solutions that exhibit contact on sets of positive measure (or capacity). Recently, entropy regularization has become a popular technique to promote exploration in reinforcement learning [\[5, 135, 126\]](#). The same technique is also used in semidefinite programming [\[137\]](#) and optimal transport [\[55\]](#). An early comparison of infinite-

dimensional interior point versus quadratic penalty approaches can be found in [27]. We also point to more recent work [194, 217] on new penalty methods that appear to be mesh-dependent. Finally, though not expressly developed for bound-constrained problems in infinite-dimensional spaces, proximal point methods will play a central role in our method. This is discussed in detail in [Subsection 4.4](#) below.

Augmented Lagrangian approaches have also been developed for variational inequalities and PDE-constrained optimization; see, e.g., early work in [108, 107, 26, 28] along with the monographs [77, 109, 207, 115] and the many references therein. Recent work has extended these methods to more general problems in abstract Banach spaces while simultaneously exploiting advances in matrix-free, inexact subproblem solvers in constrained optimization (such as [91, 121]), see [31, 110, 17]. In finite dimensions, augmented Lagrangian approaches are generally superior to penalty-based methods in the sense that the penalty parameter does not need to be driven to infinity to guarantee feasibility. Moreover, the penalty function in the subproblems is adaptively updated by the dual variables at each iteration. However, as observed in [17, Sec. 5], the situation is more delicate in infinite dimensions, e.g., it may be necessary that some of the penalty parameters need to pass to infinity to guarantee the generation of a sequence of iterates with feasible accumulation points and the dual variables may not be bounded in those function spaces which are more easily treated numerically.

3.2. Numerical methods for pointwise bound constraints. The development of bound-preserving numerical methods for PDEs began in the early days of scientific computing [129, 78] and has remained an important pursuit ever since. Although the present paper focuses on an entirely different category of PDE problems, hyperbolic conservation laws have provided a major source of motivation for research on the topic [53, 90], and have inspired many bound-preserving techniques now applied to other classes of PDEs. In many situations, the challenge lies in the fact that standard high-order numerical methods do not preserve key invariant domain properties of the underlying physics [84], such as pointwise positivity [189], yet, such properties are often required for numerical stability [209].

Some of the earliest attempts to ensure bound constraints involved using artificial viscosity to dampen oscillations that would lead to negativity and other spurious solution features [202, 128]. Later on, more sophisticated “high resolution” flux- and slope-limiting strategies emerged [34, 200, 89, 190]; see also [134] for a classical overview and further references.

One of the most popular approaches to designing high-order bound-preserving methods is flux-corrected transport [34, 213, 123, 124]. The general idea relies on forming a convex combination of a desired high-order solution and a bound-preserving low-order solution. The method then selects the high-order solution wherever the constraint is satisfied and locally transitions to the low-order solution wherever it is necessary to avoid constraint violations. A more recent popular approach [215, 214, 216], which can be traced back to [170], relies on developing high-order schemes with positive cell averages. If such a high-order scheme can be found, the local solution need only be rescaled towards its (positive) mean wherever the constraints are violated.

The majority of high-order bound-preserving numerical methods for PDEs, including the two methods just described for hyperbolic conservation laws, do not constrain the solution to the continuous-level feasible set. This is due, in part, to the fact that checking pointwise bound violations with an arbitrary polynomial is prohibitively expensive [127]. Instead, almost all modern methods involve one of two

common strategies: (1) enlarging the feasible set by only constraining the values of the solution at quadrature or nodal points [215, 189, 136, 64, 20] or (2) diminishing the feasible set by constraining the solution’s basis function coefficients [2, 15, 144, 1]. The former strategy results in a relaxation of the underlying problem that allows for solutions that are not truly positive *pointwise*. The latter strategy typically involves discretizing the solution with a positive basis that guarantees, e.g., that the solution is non-negative *if* its coefficient are non-negative; see [188, 165, 6, 52, 7, 56] for the properties of various choices. If a high-order discretization is used, both strategies lead to basis-dependent solutions, instead of solely approximation space-dependent solutions.

Since limiters tend to have a minimal number of hyperparameters, enforcing bound constraints using many of the techniques above may, at most, reduce to only solving a single-variable optimization problem at each element. Recently, however, optimization-based methods have been explored to enlarge the solution space [210, 32]. In these methods, a nonlinear program is solved at each element. Likewise, global optimization approaches have also been explored, but, possibly owing to the cost, we are only aware of investigations with simple model problems [141, 69].

Finally, logarithm-transformation methods, which date back at least to [104], have been known in the literature for some time [105, 41, 145]. Yet, they have taken on new interest in recent years [155, 143, 73, 201, 63]. Other earlier work of related interest include [60, 39, 156, 131, 70]. The appeal is that discretizing a transformed variable may deliver an approximation that is intrinsically structure-preserving and basis-independent because it encodes geometry of the feasible set. However, as we have already described in detail in [Subsection 2.1](#), naively transforming a PDE variable leads to theoretical concerns when the solution is permitted to reach the boundary of the feasible set. Therefore, implementing these methods in practice can be challenging, and may require ad-hoc assembly rules for the degenerate PDEs that arise, as noted in [201].

3.3. Contributions of the present work. This paper focuses on establishing a mathematical foundation for the proximal Galerkin finite element method and exploring some of its applications. The main technical results are developed specifically for the obstacle problem. Yet, [Sections 5](#) and [6](#) provide further sample applications and suggestions for future work. In order to distinguish our work from previous and parallel efforts described in the literature above, we itemize our primary contributions:

- We introduce a new numerical method to treat infinite-dimensional bound-constrained variational inequality problems. The method hinges on an adaptive entropy regularization technique that was introduced by Nemirovsky and Yudin in [161] for general *reflexive* Banach spaces, but has been primarily explored as an efficient optimization algorithm for finite-dimensional problems [193]. Moreover, the nature of the functionals involved in our approach indicate that we need to work in a non-standard, *non-reflexive* setting that is nevertheless natural for entropy regularization in infinite dimensions.

- The adaptive entropy regularization technique explored in this paper indicates the potential for a broad methodology in which the nonlinearity arising from the variational derivative of the entropy term can be replaced by a slack variable — which we call the “latent” variable — that is isomorphic to the regularized primal variable. This ultimately delivers a greater degree of flexibility in the choice of discretization scheme as the isomorphism naturally facilitates structure-preserving discretizations. We coin this framework the latent variable proximal point (LVPP) methodology.

- We apply the entropy regularization technique to the obstacle problem, and in doing so derive (distributional forms of) the *entropic Poisson equation*,

$$(3.1a) \quad -\Delta u + \ln u = f,$$

and the *binary-entropic Poisson equation*,

$$(3.1b) \quad -\Delta u + \operatorname{arctanh} u = f.$$

When understood as arising from the Euler–Lagrange equations for the regularized energy functionals, these appear to be novel semilinear elliptic PDEs. Though a similar equation to (3.1) has been investigated in [157] and the nonlinearities are, at least when restricted to their *domains*, smooth and monotone, the Nemytskii operators induced by \ln and $\operatorname{arctanh}$ require special care as they have *restricted* domains when defined from the original real-valued functions; cf. [12] and related literature for the analysis of Nemytskii, i.e. nonlinear superposition operators.

- Motivated by the analysis of the entropic Poisson equation, we establish a non-trivial geometric connection between non-negativity-constrained optimization and group theory. Further geometric connections are established via entropy functionals for other bound constraints.

- We present a novel finite element method to solve the entropic Poisson equation and perform preliminary *a priori* error analysis on the resulting nonlinear mixed method. Our numerical experiments indicate that the method is mesh-independent when comparing the number of iterates required to reach a certain solution tolerance; see, e.g., Subsection 4.9.

- We extend the contributions above to arrive at a novel approach to enforce discrete maximum principles on non-symmetric elliptic PDE, e.g., the advection-diffusion equation.

- We introduce *two* different types of stable finite element pairs for proximal Galerkin discretizations of second-order elliptic VIs with pointwise bound constraints. The first type employs a *discontinuous* latent variable ψ_h ; cf. Subsection 4.7. These finite elements lead to a primal solution u_h with a feasible cell average; cf. Remark 4.20. The second type uses a $C^0(\Omega)$ -*continuous* latent variable; cf. Subsection 5.3. In this case, the correct quadrature rule induces a nodally-feasible primal solution; cf. Remark 5.3. Both types of proximal Galerkin discretizations lead to a secondary solution variable \tilde{u}_h that is feasible *pointwise* throughout the domain.

- We present a new algorithm for topology optimization to showcase the breadth of applicability of the geometric optimization techniques developed in this work. The algorithm is efficient and relatively simple to implement. Our results indicate that it is also *mesh-independent*.

- We release our code [111, 112, 113], implemented in part using the finite element software FEniCSx in Python and, otherwise, with the MFEM library in C++ [14], to facilitate broader adoption in the community.

4. The obstacle problem. In Section 2, we surveyed several structural properties that entropy regularization brings to a specific form of the obstacle problem,

$$(4.1) \quad \min_{u \in H_g^1(\Omega)} \frac{1}{2} \int_{\Omega} |\nabla u|^2 dx - \int_{\Omega} f u dx \quad \text{subject to } u \geq \phi \text{ in } \Omega,$$

where $\phi = 0$. In this section, we return to the same motivating example to review these properties in greater detail and extend our conclusions in order to analyze

nonzero obstacle functions $\phi \neq 0$. The main theoretical results in this section are the representation theorem, [Theorem 4.1](#), the characterization theorem, [Theorem 4.7](#), and the convergence theorem, [Theorem 4.13](#). The section closes with a complete proximal Galerkin algorithm to solve the obstacle problem ([Algorithm 4](#)) and a report of our numerical experiments with it ([Subsection 4.9](#)).

4.1. The entropy gradient. Before we can properly investigate entropy regularization and its role in treating the obstacle problem ([4.1](#)), we must closely analyze the regularity of the entropy functional ([2.17](#)) in Lebesgue spaces. Doing so will guide us toward the key geometric structure encoded in this functional. As a pedagogical instrument, we proceed by building an analogy to the finite-dimensional setting.

Let $x \in \mathbb{R}^N$ denote the N -dimensional vector (x_1, \dots, x_N) and denote the non-negative orthant in \mathbb{R}^N by

$$(4.2) \quad \mathbb{R}_+^N = \{(x_1, \dots, x_N) \in \mathbb{R}^N \mid x_i \geq 0 \text{ for all } i = 1, \dots, N\}.$$

Now, consider the corresponding finite-dimensional entropy function $s : \mathbb{R}_+^N \rightarrow \mathbb{R}$ defined by $s(x) = \sum_{i=1}^N x_i \ln x_i - x_i$, wherein we remind the reader of our simplifying convention $0 \ln 0 := 0$. It is easy to see that $s(x)$ is continuous and strictly convex on \mathbb{R}_+^N , but only differentiable on its interior,

$$(4.3) \quad \text{int } \mathbb{R}_+^N = \{(x_1, \dots, x_N) \in \mathbb{R}^N \mid x_i > 0 \text{ for all } i = 1, \dots, N\},$$

due to the logarithmic singularity in the gradient $\nabla s(x) = (\ln x_1, \dots, \ln x_N)$. A careful analysis is required to determine what the effect of the same type of logarithmic singularity will be at the function space level when analyzing the entropy functional S in ([2.17](#)).

As our first key structural result shows, $L_+^\infty(\Omega)$ and $\text{int } L_+^\infty(\Omega)$ reflect the roles played above in finite-dimensions by \mathbb{R}_+^N and $\text{int } \mathbb{R}_+^N$, respectively. The proof is deferred to the outset of [Appendix A.2](#).

THEOREM 4.1 (Gradient representation). *Let $S : L^p(\Omega) \rightarrow \mathbb{R} \cup \{+\infty\}$, $p \in [1, \infty]$, be the negative entropy functional defined by*

$$S(u) = \begin{cases} \int_\Omega u \ln u - u \, dx, & \text{if } u \in L_+^p(\Omega), \\ +\infty, & \text{otherwise.} \end{cases}$$

1. *If $p \in [1, \infty]$, then S is strictly convex and lower semicontinuous.*
2. *If $p \in (1, \infty]$, then S is continuous on $L_+^p(\Omega)$.*
3. *If $p = \infty$, then S is continuously Fréchet differentiable on $\text{int } L_+^\infty(\Omega)$ with respect to the $L^p(\Omega)$ -norm topology. In particular, the $L^\infty(\Omega)$ -Fréchet derivative of S can be uniquely characterized by the variational equation*

$$(4.4) \quad \langle S'(u), v \rangle = \int_\Omega v \ln u \, dx \quad \text{for all } u \in \text{int } L_+^\infty(\Omega) \text{ and } v \in L^\infty(\Omega).$$

Moreover, $\|S'(u)\|_{(L^\infty(\Omega))'} = \|\nabla S(u)\|_{L^1(\Omega)}$, where

$$(4.5) \quad \nabla S(u) = \ln u \in L^\infty(\Omega)$$

is the unique primal representation (i.e., gradient) of $S'(u)$ and is uniquely determined by the variational equation

$$(4.6) \quad (\nabla S(u), v) = \langle S'(u), v \rangle \quad \text{for all } u \in \text{int } L_+^\infty(\Omega) \text{ and } v \in L^1(\Omega).$$

At a first glance, it is tempting to define S from $L^1(\Omega)$ into $\mathbb{R} \cup \{+\infty\}$. This is the perspective taken in much of the literature on infinite-dimensional convex analysis; see, in particular, [35, 22]. In this setting, it is shown that we have strict convexity and lower semicontinuity. However, as noted in [22, Remark 5.7], there are some issues with this viewpoint. For example, S would be nowhere continuous, but it would admit subgradients of the form $\ln u$ whenever $u \in \text{int } L_+^\infty(\Omega)$.

As claimed above, and proven in [Appendix A.2](#), we see that $S : L^p(\Omega) \rightarrow \mathbb{R} \cup \{+\infty\}$ is in fact continuous on $L_+^p(\Omega)$ provided $p > 1$ and even continuously Fréchet differentiable when we take $p = \infty$ and $u \in \text{int } L_+^\infty(\Omega)$. Moreover, the derivative $S'(u)$ admits a “primal” representation of the form $\ln u$, which connects back to the convex analysis literature. Our proof techniques, however, are not based on duality arguments or the properties of subgradients.

Since S will be used to define a Bregman distance below, whose domain needs to fit together with the typical regularity spaces for partial differential operators, we can safely choose any $p \in [1, \infty]$ so that the regularity space is continuously embedded into $L^p(\Omega)$, even if this initially appears to rule out certain functions in the domain of S . For example, if we are working with $u \in H^1(\Omega)$, then we can select $p \in [1, 2]$, regardless of the dimension of Ω or regularity of $\partial\Omega$. On the other hand, if the dimension of Ω is $n = 2$ or higher, then $H^1(\Omega)$ does not continuously embed into $L^\infty(\Omega)$.

Finally, the properties of S given in [Theorem 4.1](#) indicate that $S : L_+^\infty(\Omega) \rightarrow \mathbb{R}$ is part of an important class of *essentially smooth* functions introduced by Rockafellar [172, Section 26] (in finite dimensions) known as *Legendre functions*, which are extended to infinite dimensions in [35, 22]. As discussed in, e.g., [193, Section 2.3], Legendre functions play a crucial role in proximal algorithms for finite-dimensional convex optimization.

To prepare us for non-trivial obstacles $\phi \neq 0$, we have the following corollary to [Theorem 4.1](#) pertaining to the shifted entropy functional $S_\phi(u) = S(u - \phi)$. As with [Theorem 4.1](#), the proof of this result is deferred to [Appendix A.2](#).

COROLLARY 4.2 (Gradient of the shifted entropy functional). *Let $\phi \in L^\infty(\Omega)$. The shifted negative entropy functional $S_\phi(u)$ is strictly convex on*

$$(4.7) \quad L_{\phi,+}^\infty(\Omega) = \{w \in L^\infty(\Omega) \mid w \geq \phi\}.$$

and Fréchet differentiable on

$$(4.8) \quad \text{int } L_{\phi,+}^\infty(\Omega) = \{w \in L_{\phi,+}^\infty(\Omega) \mid \text{ess inf}(w - \phi) > 0\}$$

with respect to the norm topology on $L^\infty(\Omega)$. The Fréchet derivative of S_ϕ can be uniquely characterized by the variational equation

$$(4.9) \quad \langle S'_\phi(u), v \rangle = \int_\Omega v \ln(u - \phi) \, dx \quad \text{for all } u \in \text{int } L_{\phi,+}^\infty(\Omega) \text{ and } v \in L^\infty(\Omega).$$

Moreover, $\|S'_\phi(u)\|_{(L^\infty(\Omega))'} = \|\nabla S_\phi(u)\|_{L^1(\Omega)}$, where

$$(4.10) \quad \nabla S_\phi(u) = \ln(u - \phi) \in L^\infty(\Omega),$$

is the unique primal representation (i.e., gradient) of $S'_\phi : \text{int } L_{\phi,+}^\infty(\Omega) \rightarrow (L^\infty(\Omega))'$ in $L^\infty(\Omega)$, determined by the variational equation

$$(4.11) \quad (\nabla S_\phi(u), v) = \langle S'_\phi(u), v \rangle \quad \text{for all } u \in \text{int } L_{\phi,+}^\infty(\Omega) \text{ and } v \in L^1(\Omega).$$

Remark 4.3 (Empty interior in the $H^1(\Omega)$ topology). We recall that if $K = \{u \in H_g^1(\Omega) \mid u \geq 0\} = H_g^1(\Omega) \cap H_+^1(\Omega)$ and $\Omega \subset \mathbb{R}^n$ with $n > 1$, then $\text{int } K = \emptyset$. This is a simple consequence of the fact that $H^1(\Omega)$ contains unbounded functions, and so we can get arbitrarily close to any $u \in K$ in the H^1 -norm with points outside K .

Remark 4.4 (No Riesz representation theorem). When inspecting [Theorem 4.1](#) and [Corollary 4.2](#), the reader should note that $L^p(\Omega)$ is a Banach algebra only in the case $p = \infty$ and we only prove that $u \mapsto S_\phi(u)$ is Fréchet differentiable with respect to variations in this set; see [\(4.9\)](#). In fact, there is a key step in our proof of [Theorem 4.1](#) that requires all functions u where the functional $S(u)$ is differentiable to have a multiplicative inverse $1/u \in L^\infty(\Omega)$; see [\(A.19\)](#). Based in part on this requirement, we continue to work directly with $L^\infty(\Omega)$, which is a *non-reflexive* Banach space *without* a corresponding Riesz representation theorem [\[4\]](#). It is, therefore, not a trivial consequence of differentiability that the Fréchet derivative $S'_\phi(u) \in (L^\infty(\Omega))'$ has the unique function space representation $\nabla S_\phi(u) \in L^\infty(\Omega)$ given by [\(4.11\)](#). In fact, the derivative of general functionals on $L^\infty(\Omega)$ lie in $(L^\infty(\Omega))'$, which is the space of absolutely continuous, finitely additive set functions of bounded total variation on Ω ; cf. [\[212, p. 118\]](#). Throughout this work, we consciously choose to refer to $\nabla S_\phi: \text{int } L_{\phi,+}^\infty(\Omega) \rightarrow L^\infty(\Omega)$ as the *gradient* of the (shifted) entropy functional, even though we are well aware that the term “gradient” is typically understood as a Hilbert space concept.

4.2. The entropy gradient is an isomorphism. Let us return to the finite-dimensional entropy function $s(x) = \sum_{i=1}^N x_i \ln x_i - x_i$ introduced at the beginning of the previous subsection and focus on its properties in the strictly positive orthant $\text{int } \mathbb{R}_+^N \subset \mathbb{R}^N$. In this case, the reader should note that $x \mapsto \nabla s(x) = (\ln x_1, \dots, \ln x_N)$, is a bijection between the set of component-wise positive vectors $x \in \mathbb{R}_+^N$ and the entire vector space \mathbb{R}^N .

This correspondence has a special algebraic significance if we view \mathbb{R}_+^N as a Lie group under the operation of componentwise multiplication,

$$(4.12) \quad x \otimes y = (x_1 y_1, \dots, x_N y_N),$$

and view \mathbb{R}^N as its associated Lie algebra under addition; cf. [\[132, Example 7.4 \(b\)\]](#). Indeed, the smooth map $\nabla s: \text{int } \mathbb{R}_+^N \rightarrow \mathbb{R}^N$ given above is a Lie group isomorphism because

$$(4.13) \quad \nabla s(x) + \nabla s(y) = (\ln x_1 + \ln y_1, \dots, \ln x_N + \ln y_N) = \nabla s(x \otimes y).$$

It is trivial to see that the same structure is replicated at the infinite-dimensional level between the Banach–Lie algebra $L^\infty(\Omega)$ and its Banach–Lie group $\text{int } L_+^\infty(\Omega)$ since

$$(4.14) \quad \nabla S(u) + \nabla S(v) = \ln u + \ln v = \nabla S(uv).$$

A deeper geometric meaning to this correspondence is revealed if we draw upon the well-known result in differential geometry that all finite-dimensional Lie groups are associated to their Lie algebra by an exponential map [\[132, Proposition 20.8\]](#). In the case of the Lie group $\text{int } \mathbb{R}_+^N$, it may be checked that the inverse of ∇s , defined $(\nabla s)^{-1}(x) = (\exp x_1, \dots, \exp x_N)$, is precisely this map. Conveniently, the finite-dimensional result extends to the Banach–Lie group $\text{int } L_+^\infty(\Omega)$ [\[75, 76\]](#), and we are left with a similar geometric interpretation (cf. [Figure 4.1](#)) of the isomorphism induced by the gradient of the entropy functional $\nabla S: \text{int } L_+^\infty(\Omega) \rightarrow L^\infty(\Omega)$ and its inverse,

$$(4.15) \quad (\nabla S)^{-1}(u) = \exp u.$$

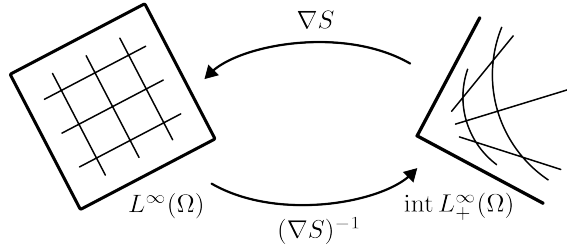


Fig. 4.1: The exponential map $(\nabla S)^{-1}(v) = \exp v$ is an analytic isomorphism between the Banach algebra $L^\infty(\Omega)$ and the Banach-Lie group $\text{int } L_+^\infty(\Omega) = \{v \in L^\infty(\Omega) \mid \text{ess inf } v > 0\}$; see [Proposition A.7](#). Moreover, its restriction to the subalgebra $H^1(\Omega) \cap L^\infty(\Omega)$ forms an isomorphism with the subgroup $H^1(\Omega) \cap \text{int } L_+^\infty(\Omega)$; see [Proposition A.9](#).

Moreover, it can be shown that restricting the exponential map [\(4.15\)](#) to the subalgebra $H^1(\Omega) \cap L^\infty(\Omega)$ induces an isomorphism with the subgroup $H^1(\Omega) \cap \text{int } L_+^\infty(\Omega)$. For further details, see [Propositions A.7](#) and [A.9](#).

Remark 4.5 (Exploiting the geometry of the feasible set). From the optimization point-of-view, there is great value in the isomorphism ∇S residing in the fact that $L^\infty(\Omega)$ is a Banach space and Banach spaces are natural spaces in which to construct additive update formulas (they are complete, normed, and closed under addition). Many competitive algorithms for *unconstrained* optimization problems, such as gradient descent and Newton methods, are additive update formulas that leverage this linear structure in some way [\[163\]](#). Likewise, when dealing with constrained optimization problems, most algorithms appeal to the linear structure of the ambient space containing the feasible set. In [Subsection 4.4](#), we will show how the isomorphism $\nabla S: \text{int } L_+^\infty(\Omega) \rightarrow L^\infty(\Omega)$ allows us to ignore the ambient space the original problem is posed in and work instead with the intrinsic geometry of the constraint set. This, in turn, will allow us to treat constrained optimization problems in Sobolev spaces with methods originally designed only for the unconstrained setting.

4.3. Relative entropy. Entropy not only delivers an isomorphism between the Banach-Lie group $\text{int } L_+^\infty(\Omega)$ and its Banach algebra $L^\infty(\Omega)$. It also induces a valuable distance function called the *relative entropy* or (*extended*) *Kullback-Leibler divergence*.

We assume below that V is a Banach space. For any smooth convex function $G: V \rightarrow \mathbb{R}$, its Bregman divergence is defined by the formula

$$(4.16) \quad D_G(u, v) = G(u) - G(v) - \langle G'(v), u - v \rangle.$$

Encoded in this definition is the important observation that, because G is convex, the graph $\{(u, G(u)) \mid u \in V\}$ will always lie on or above its supporting hyperplanes, $\{(u, G(v) + \langle G'(v), u - v \rangle) \mid u \in V\}$, for every $v \in V$ at which $G'(v)$ exists, see [\[37\]](#) for this and related insights. The Bregman divergence $D_G: \text{dom } G \times \text{dom } G' \rightarrow \mathbb{R}$ measures the vertical distance between these two sets. For nonsmooth convex functionals that are merely subdifferentiable, the definition of subgradients $g' \in \partial G(u)$ implies $G(v) \geq G(u) + \langle g', v - u \rangle$ for all v in $\text{dom } G$. Therefore, Bregman divergences can be defined for nonsmooth functionals using subgradients instead of derivatives, with the caveat that there may be uncountably many g' that describe a supporting hyperplane at points of nonsmoothness; cf. [\[211\]](#).

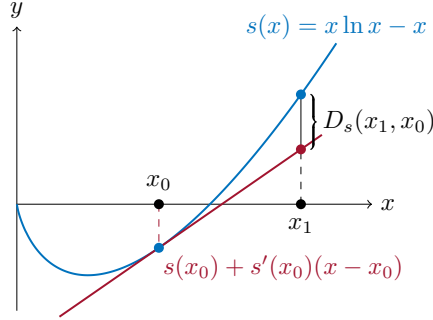


Fig. 4.2: The convex function $s(x) = x \ln x - x$, its supporting hyperplane $\{s'(x_0) + s'(x_0)(x - x_0) \mid x \in \mathbb{R}\}$, and its Bregman divergence $D_s(x_1, x_0) = x_1 \ln(x_1/x_0) - x_1 + x_0$.

Loosely speaking, a Bregman divergence is a generalization of the squared distance between two points in a Hilbert space, and is, therefore, not expected to satisfy a triangle inequality. To see how this interpretation arises, it is a straightforward exercise to check, e.g., that if $G: H_0^1(\Omega) \rightarrow \mathbb{R}$, with $G(u) = \frac{1}{2} \|\nabla u\|_{L^2(\Omega)}^2$, the associated Bregman divergence is

$$(4.17) \quad \begin{aligned} D_G(u, v) &= \frac{1}{2} \|\nabla u\|_{L^2(\Omega)}^2 - \frac{1}{2} \|\nabla v\|_{L^2(\Omega)}^2 - (\nabla v, \nabla u - \nabla v) \\ &= \frac{1}{2} \|\nabla u - \nabla v\|_{L^2(\Omega)}^2. \end{aligned}$$

The relative entropy $D: L_+^p(\Omega) \times \text{int } L_+^\infty(\Omega) \rightarrow \mathbb{R}$, for $p \in [1, \infty]$, is the Bregman divergence induced by the entropy functional S . Given its importance to this work, we neglect to write the subscript- S when working with this measure of distance. In turn, we may select any $u \in L_+^p(\Omega)$ and $v \in \text{int } L_+^\infty(\Omega)$ to explicitly derive the relative entropy as follows,

$$(4.18) \quad D(u, v) = S(u) - S(v) - (\nabla S(v), u - v) = \int_{\Omega} u \ln \frac{u}{v} - u + v \, dx.$$

An illustration of the Bregman divergence of the finite-dimensional entropy function $s(x) = \sum_{i=1}^N x_i \ln x_i - x_i$ is given in Figure 4.2 for the case $N = 1$. We initially use the right-hand side of (4.18) in our study below without requiring its definition as a Bregman divergence. After a careful analysis shows that the relevant solutions are in $L_+^\infty(\Omega)$, we then employ the usual properties of Bregman divergences where required in several convergence proofs. This frees us from the rigid structures of convex analysis, e.g., that often fix the domain V in the beginning and require us to work only in this space and its given topology.

Along with other statistical distances, the relative entropy has a rich history of being used to encode geometric structure in analysis within statistics, probability theory, and information theory [11, 10, 162]. Although a Bregman divergence is not symmetric, i.e., $D_G(u, v) \neq D_G(v, u)$ in general, it will satisfy the following important properties when G is strictly convex [37, 47]:

PROPOSITION 4.6 (Properties of Bregman divergences). *Let $G: V \rightarrow \mathbb{R}$ be smooth and strictly convex. Then the following properties hold:*

Non-negativity. $D_G(u, v) \geq 0$ for all $u \in \text{dom } G$ and $v \in \text{dom } G'$.

Positivity. $D_G(u, v) = 0$ if and only if $u = v$.

Convexity. $D_G(u, v)$ is strictly convex in its first argument. Moreover, if G is strongly convex, then so is $u \mapsto D_G(u, v)$.

Linearity. Let $F: V \rightarrow \mathbb{R}$ be smooth and strictly convex and $\lambda \geq 0$. Then

$$(4.19) \quad D_{G+\lambda F}(u, v) = D_G(u, v) + \lambda D_F(u, v).$$

Three points identity. For all $u \in \text{dom } G$ and $v, w \in \text{dom } G'$, it holds that

$$(4.20) \quad D_G(u, v) - D_G(u, w) + D_G(v, w) = \langle G'(v) - G'(w), v - u \rangle.$$

4.4. Proximal point. Recall that in [Subsection 2.2](#) we proposed the regularized Dirichlet free energy functional

$$(4.21) \quad A(u) = E(u) + \theta S(u),$$

and argued that its minimizer will converge to the solution of the obstacle problem in the limit $\theta \rightarrow 0$; cf. [\(2.19\)](#). Although this approach to solving the non-negative obstacle problem is viable, there is a much more numerically stable alternative. Indeed, it turns out that we can just as readily generate a sequence of positive functions $u^k \rightarrow u^*$ by recursively regularizing the Dirichlet energy $E(u)$ with the Bregman divergence $D(u, u^k)$. The idea is relatively old in finite dimensions [\[45, 192, 47, 193\]](#), and well-explored in reflexive Banach spaces [\[59, 58\]](#). However, given that the algorithm is not well-known in the finite element community, we present a classical description that begins with a Hilbert space framework.

We now introduce the so-called *proximal minimization algorithm* [\[175, 167, 193\]](#), due to Marinet [\[150\]](#). In turn, let H be a Hilbert space and $\alpha > 0$ be a positive step size parameter. The *proximal operator*, introduced in [\[158\]](#) by Moreau, is defined for every proper lower semi-continuous function $F: H \rightarrow \mathbb{R} \cup \{\infty\}$ as follows,

$$(4.22) \quad \text{prox}_{\alpha F}(v) = \arg \min_{u \in H} \left\{ F(u) + \frac{1}{2\alpha} \|u - v\|_H^2 \right\}.$$

The utility of this operator lies largely in the fact that the $\|\cdot\|_H^2$ -regularization term in [\(4.22\)](#) transforms F (which may not be differentiable) into a finite-valued function,

$$(4.23) \quad F_\alpha(v) = \min_{u \in H} \left\{ F(u) + \frac{1}{2\alpha} \|u - v\|_H^2 \right\},$$

with an α^{-1} -Lipschitz continuous gradient [\[175\]](#). Moreover, when F is convex, minimizing either F or F_α is equivalent in the sense that

$$(4.24) \quad \inf_{u \in H} F_\alpha(u) = \inf_{u \in H} F(u).$$

In fact, the set of minimizers, $\arg \min_{u \in H} F(u)$, coincides with the set of fixed points $u^* \in H$ that satisfy $u^* = \text{prox}_{\alpha F}(u^*)$; see, e.g., [\[23, Prop. 12.28\]](#).

Choosing to iterate this fixed point equation with variable step sizes $\alpha_k > 0$ delivers the *proximal minimization algorithm* [150, 174], written explicitly as

$$(4.25) \quad u^0 \in H, \quad u^{k+1} = \text{prox}_{\alpha_{k+1}F}(u^k), \quad k = 0, 1, \dots$$

It is well-known (see, e.g., [85]), that $F(u^k)$ converges to $F(u^*)$ at a rate inversely proportional to the sum of the step sizes. More explicitly, it holds that

$$(4.26) \quad F(u^k) - F(u^*) \leq \frac{1}{2} \frac{\|u^* - u^0\|_H^2}{\sum_{\ell=1}^k \alpha_\ell}.$$

Thus, the function values of proximal iterates (4.25) can converge “arbitrarily” fast (by increasing α_ℓ), and the asymptotic complexity of the iteration (4.25) is determined by the complexity of the method used to solve each subproblem (4.22). Convergence of the function values carries over to convergence of the iterates provided an estimate of the type

$$\sigma(\|u - v\|) \leq F(u) - F(v)$$

holds, where σ is monotone and invertible on \mathbb{R}_+ , e.g., if F is strongly convex.

The potentially arbitrary order of convergence in (4.26) makes the proximal point algorithm an attractive candidate to solve many optimization problems. The drawback, however, is that each iteration of the algorithm requires the solution of a non-smooth optimization problem that may be just as difficult to solve as the original problem; cf. Remark 4.11. At the same time, the proximal operator (4.22) and fixed point iterations (4.25) are fundamental to a broad selection of modern optimization algorithms; see e.g., [23, 24, 193, 125] and the many references therein. They also play a deep role in augmented Lagrangian methods, as recognized at least as early as [173], which have seen a resurgence in interest due, in part, to their applicability for infinite dimensional problems [110, 17].

It turns out many of the most important properties of the proximal minimization algorithm also hold if $\frac{1}{2}\|u - v\|_H^2$ in (4.22) is replaced by a Bregman divergence $D_G(u, v)$ [193]. Indeed, if we assume that $G: V \rightarrow \mathbb{R}$ is a strictly convex functional on a Banach space V , we may define the *Bregman proximal operator*

$$(4.27) \quad \text{prox}_{\alpha F}^G(v) = \arg \min_{u \in \text{dom } F \cap \text{dom } G} \{F(u) + \alpha^{-1}D_G(u, v)\},$$

and the corresponding *Bregman proximal minimization algorithm*

$$(4.28) \quad u^0 \in \text{dom } F \cap \text{dom } G', \quad u^{k+1} = \text{prox}_{\alpha_{k+1}F}^G(u^k), \quad k = 0, 1, \dots$$

Figure 4.2 illustrates the execution of this algorithm for the one-dimensional energy function $e(x) = \frac{1}{2}x^2 + x$ and the relative entropy $D_s(x, y) = x \ln(x/y) - x + y$. Note that under the definitions above, one can show that (4.26) generalizes as follows [47],

$$(4.29) \quad F(u^k) - F(u^*) \leq \frac{D_G(u^*, u^0)}{\sum_{\ell=1}^k \alpha_\ell}.$$

See also Theorem 4.13.

Our contribution is to show that the proximal operator (4.27), with an appropriately defined Bregman divergence, transforms the solution of an infinite-dimensional

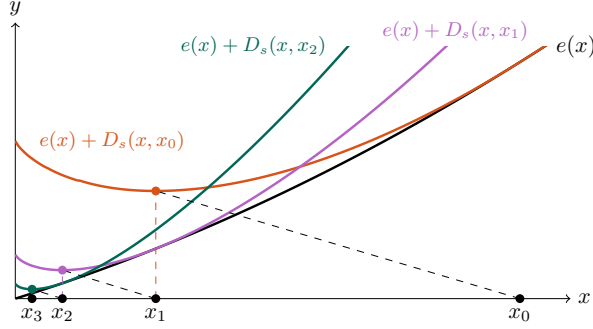


Fig. 4.3: Illustration of convergence to the solution $x^* = 0$ for the constrained minimization problem $\min_{x \in [0, \infty)} e(x)$, where $e(x) = \frac{1}{2}x^2 + x$, by solving the sequence of minimization problems $x_{k+1} = \arg \min_{x \in [0, \infty)} \{e(x) + D_s(x, x_k)\}$ starting at $x_0 = 1$.

constrained optimization problem into a sequence of semi-linear PDEs whose solutions converge to the solution of the underlying VI. In the case of the positive obstacle problem (i.e., $F = E$ and $G = S$), this conclusion hinges on the following result. When combined with (4.28), Theorem 4.7 leads us directly to Algorithm 1, which forms the basis for the proximal Galerkin finite element method. We note that the proof is technical and saved until Appendix A.3.

THEOREM 4.7 (Solution characterization). *Assume $\Omega \subset \mathbb{R}^n$ is an open, bounded Lipschitz domain, $n \geq 1$. Let $K = \{v \in H_g^1(\Omega) \mid v \geq 0\} = H_g^1(\Omega) \cap H_+^1(\Omega)$, where $g \in H^1(\Omega) \cap C(\bar{\Omega})$ satisfies $\min g|_{\partial\Omega} > 0$. Moreover, given $f \in L^\infty(\Omega)$, set*

$$E(v) = \frac{1}{2} \int_{\Omega} |\nabla v|^2 dx - \int_{\Omega} f v dx,$$

and for $w \in \text{int } L_+^\infty(\Omega)$ set $D(v, w) = \int_{\Omega} v \ln(v/w) - v + w dx$. Then, for any step size $\alpha > 0$, the (relative) Dirichlet free energy minimization problem,

$$(4.30) \quad \min_{v \in K} A_\alpha(v) := E(v) + \alpha^{-1} D(v, w),$$

has a unique solution $u \in H_g^1(\Omega) \cap \text{int } L_+^\infty(\Omega)$ that is (uniquely) characterized by the weak form of the entropic Poisson equation; namely,

$$(4.31) \quad (\alpha \nabla u, \nabla v) + (\ln u, v) = (\alpha f, v) + (\ln w, v) \quad \text{for all } v \in H_0^1(\Omega).$$

Remark 4.8 (Adaptive entropy regularization). Similar to the free energy formulation (2.16), where the Lagrange multiplier λ is approximated by $\theta \ln(1/u)$, we see that the subproblems (4.31) give rise to an approximation of the form $\alpha^{-1} \ln(w/u)$. Recalling that $\theta = \alpha^{-1}$ we see that there is fundamental difference in the two approximations given by the inclusion of the function w . Chosen correctly, as with u^k in (4.28), this function can act as an informative prior on the sequence of subproblems. More specifically, $w = u^k$ allows us to view the Bregman divergence $v \mapsto D(v, u^k) = \int_{\Omega} v \ln(v/u^k) - v + u^k dx$ as a Bayesian barrier function that is updated adaptively at each iteration k so that $u^k \rightarrow u^*$ without sending the step size $\alpha \rightarrow \infty$.

From now on, we mainly focus on inhomogeneous obstacle problems; i.e., $\phi \neq 0$. Therefore, we close this subsection with an important corollary for this case. Before we state the result, we note that

$$S_\phi(u) - S_\phi(v) - (\nabla S_\phi(v), u - v) = D(u - \phi, v - \phi)$$

whenever $u \in L_{\phi,+}^p(\Omega)$ and $v \in \text{int } L_{\phi,+}^\infty(\Omega)$. Therefore, $(u, v) \mapsto D(u - \phi, v - \phi)$ is a Bregman divergence on $L_{\phi,+}^p(\Omega) \times \text{int } L_{\phi,+}^\infty(\Omega)$. For technical reasons, we require the obstacles to be in a particular subset of $H^1(\Omega)$ defined by

$$\mathcal{O} := \{ \varphi \in H^1(\Omega) \cap C(\bar{\Omega}) \mid \Delta \varphi \in L^\infty(\Omega) \}.$$

Moreover, like [Theorem 4.7](#), the proof of [Corollary 4.9](#) is delayed until [Appendix A.3](#).

COROLLARY 4.9 (Solution characterization for inhomogeneous obstacles). *In addition to the assumptions of [Theorem 4.7](#), let $\phi \in \mathcal{O}$ such that $\text{ess inf } \gamma(g - \phi) > 0$ on $\partial\Omega$ and define*

$$K_\phi := \{ u \in H_g^1(\Omega) \mid u \geq \phi \text{ a.e. in } \Omega \}.$$

Then for any step size $\alpha > 0$ and $w \in \text{int } L_{\phi,+}^\infty(\Omega)$ the optimization problem

$$(4.32) \quad \min_{v \in K_\phi} E(v) + \alpha^{-1} D(v - \phi, w - \phi),$$

has a unique solution $u \in H_g^1(\Omega) \cap \text{int } L_{\phi,+}^\infty(\Omega)$ that satisfies the weak form of the (generalized) entropic Poisson equation; namely,

$$(4.33) \quad (\alpha \nabla u, \nabla v) + (\ln(u - \phi), v) = (\alpha f, v) + (\ln(w - \phi), v) \quad \text{for all } v \in H_0^1(\Omega).$$

Remark 4.10 (Delicate analysis). Semilinear mixed variational inequalities of obstacle type have been thoroughly studied, as detailed in the famous monograph by J.-F. Rodrigues, [176, Chap. 4.6]. This includes regularity theory and a maximum principle that relates the solution of the VI to the obstacle, forcing term, and boundary values. The techniques go back to the seminal work by Stampacchia [184], Murty and Stampacchia [159] and can also be found in [116]. However, the VI associated with our problem is only valid if we can differentiate the “extra” nonlinearity in the entropy term. This in turn requires the solution u of each subproblem to be essentially bounded and strictly above the obstacle, so we need to resort to a more delicate analysis solely based on the properties of the optimization problem.

Remark 4.11 (Challenges of the Hilbert space setting). Let $\chi_K: H^1(\Omega) \rightarrow \mathbb{R} \cup \{\infty\}$ denote the indicator function $\chi_K(x) = 0$ if $x \in K$ and $\chi_K(x) = \infty$ otherwise. It is interesting to compare the operator

$$(4.34) \quad \text{prox}_{\alpha E + \chi_K}(v) = \arg \min_{u \in K} \left\{ E(u) + \frac{1}{2\alpha} \|u - v\|_{H^1}^2 \right\},$$

to $\text{prox}_{\alpha E}^S(v)$. Indeed, unlike (4.31), the subproblems that (4.34) induces each require the solution of their own VI,

$$(4.35) \quad \int_{\Omega} \nabla((1 + \alpha)u - v) \cdot \nabla w \, dx + \int_{\Omega} (u - v - \alpha f)w \, dx \geq 0 \quad \text{for all } w \in K - u,$$

that is at least as difficult to solve as the original VI defining u^* ; cf. (1.3). Similar issues tend to appear whenever squared norm regularization terms are used to design proximal point algorithms for infinite-dimensional bound constraints.

Alternatively, one can use a penalty method to solve the original problem by considering instead a $C^{1,1}$ -quadratic penalty term of the type

$$\frac{1}{2\alpha} \int_{\Omega} \max\{0, \phi - u\}^2 \, dx.$$

This functional is in fact the prox-operator (in the $L^2(\Omega)$ topology) of the indicator function for the larger feasible set $\{u \in L^2(\Omega) \mid u \geq \phi\}$. See [95, 94] for details including second-order algorithms and an analytical path-following scheme for α . Note that the subproblems using a quadratic penalty would be semismooth semilinear elliptic PDEs. However, since the nonlinearity does not arise from a strictly monotone continuous function, we cannot derive a similar latent variable formulation.

Remark 4.12 (Comparison to the augmented Lagrangian method). It is possible to view classical augmented Lagrangian methods as penalty methods that adaptively change the penalty function and associated penalty parameter via the behavior of the dual variable, i.e., Lagrange multiplier. Aside from identifying an efficient subproblem solver, the challenge is usually to find an appropriate combination of update strategies that allow for inexact subproblem solves and conservative parameter update rules that still exhibit rapid convergence behavior in practice. The method described in this work follows a similar strategy. Indeed, the role of the penalty function is played by a Bregman distance, which is adaptively updated via the primal variable, and the penalty parameter is given by α . Bregman distances allow us to better exploit the geometry of the feasible set and the convergence theory of the proximal point method provides a clear connection to convergence rates that even allows for α to remain constant.

4.5. Latent variable proximal point. An appealing feature of the entropic Poisson equation (4.33) is that its solution permits *two* additional representations; cf. Figure 2.1. In both cases, we take advantage of the entropy gradient $\nabla S(v) = \ln v$ being an isomorphism (cf. Subsection 4.2). First, we may introduce the latent variable representation,

$$(4.36) \quad \psi = \ln(u - \phi) \quad \iff \quad u = \exp \psi + \phi,$$

by simply applying the entropy gradient transformation to the primal solution u . Second, as already noted in Remark 4.8, we may construct a dual variable representation which, for the inhomogeneous obstacle problem, is written as follows:

$$(4.37) \quad \lambda = \alpha^{-1} \ln \frac{w - \phi}{u - \phi} \quad \iff \quad u = (w - \phi) \exp(-\alpha\lambda) + \phi.$$

The utility of these representations is witnessed if we consider how to solve either of the primal subproblems (4.31) or (4.33). Indeed, due to the logarithmic terms, these semi-linear PDEs are only defined if $\text{ess inf}(u - \phi) > 0$, which appears to rule out most efficient root-finding algorithms, such as Newton’s method, and discretization choices, such as the Galerkin method. Fortunately, the alternative solution representations above provide saddle-point relaxations of the entropic Poisson equation that do not suffer from these two drawbacks.

We are now ready to state the final main theoretical result, which also establishes explicit bounds on the optimization error for the latent variable proximal point (LVPP) algorithm, defined via (4.38) below. The proof is given in Appendix A.5.

THEOREM 4.13 (Convergence of LVPP). *Assume $\alpha_{k+1} > 0$, $k = 0, 1, \dots$, is a sequence of positive step size parameters. Furthermore, assume $\Omega \subset \mathbb{R}^n$ is an open, bounded Lipschitz domain, $n \geq 1$, $\phi \in \mathcal{O}$, and let $g \in H^1(\Omega) \cap C(\bar{\Omega})$ such that $\text{ess inf } \gamma(g - \phi) > 0$. Fix $\psi^0 \in H^1(\Omega) \cap L^\infty(\Omega)$ and consider the sequence of functions u^k, ψ^k solving the following coupled system of variational equations:*

$$(4.38) \quad \begin{cases} \text{Find } u^{k+1} \in H_g^1(\Omega) \text{ and } \psi^{k+1} \in L^\infty(\Omega) \text{ such that} \\ (\alpha_{k+1} \nabla u^{k+1}, \nabla v) + (\psi^{k+1}, v) = (\alpha_{k+1} f + \psi^k, v) & \text{for all } v \in H_0^1(\Omega), \\ (u^{k+1}, \varphi) - (\exp \psi^{k+1}, \varphi) = (\phi, \varphi) & \text{for all } \varphi \in L^2(\Omega). \end{cases}$$

Then the Dirichlet energy of the primal iterates is monotonically non-increasing, i.e.,

$$(4.39) \quad E(u^{k+1}) \leq E(u^k).$$

Moreover, if $\sum_{j=1}^k \alpha_j \rightarrow \infty$ as $k \rightarrow \infty$, then the subproblem solutions u^k converge in $H^1(\Omega)$ to

$$(4.40) \quad u^* = \arg \min_{u \in H^1(\Omega)} E(u) \quad \text{subject to } u \geq \phi \text{ in } \Omega \text{ and } u = g \text{ on } \partial\Omega.$$

Furthermore, the functions $\lambda^{k+1} = (\psi^k - \psi^{k+1})/\alpha_{k+1}$ converge strongly in $H^{-1}(\Omega)$ to the Lagrange multiplier $\lambda^* = -\Delta u^* - f$. In fact, the optimization error in both u^k and λ^k are equal and converge at the following arbitrary rate determined by the sequence of step-sizes $\alpha_k > 0$,

$$(4.41) \quad \frac{1}{2} \|\lambda^* - \lambda^k\|_{H^{-1}(\Omega)}^2 = \frac{1}{2} \|\nabla u^* - \nabla u^k\|_{L^2(\Omega)}^2 \leq \frac{D(u^* - \phi, u^0 - \phi)}{\sum_{j=1}^k \alpha_j}.$$

Remark 4.14 (Arbitrary orders of convergence). **Theorem 4.13** shows that the iteration complexity of LVPP depends on the choice of the step sizes α_k . The consequences of different step size choices is summarized in **Corollary A.14**. For example, we find that constant step sizes lead to sublinear convergence and geometrically increasing step sizes lead to first-order convergence. Even faster growing step size sequences will achieve superlinear convergence. See also **Remark 4.18**.

Remark 4.15 (Convergence in the $H^1(\Omega)$ -norm). At first glance, control over the full $H^1(\Omega)$ norm of u^k appears problematic because (4.41) does not include the full norm on $H_g^1(\Omega)$. However, in light of the Poincaré inequality and $u^* - u^k \in H_0^1(\Omega)$, we also obtain

$$\frac{1}{2} \|u^* - u^k\|_{L^2(\Omega)}^2 \leq c \frac{D(u^* - \phi, u^0 - \phi)}{\sum_{j=1}^k \alpha_j},$$

where $c > 0$ is an embedding constant independent of k .

Remark 4.16 (Convergence of the latent variable). If we adopt the conventions $\ln 0 = -\infty$ and $\exp(-\infty) = 0$, we may define $\psi^* = \ln(u^* - \phi)$ as an extended real-valued function on Ω ; cf. **Subsection 2.1**. Likewise, we may understand convergence of the latent variable $\psi \rightarrow \psi^*$ under the metric implied by this transformation. Indeed, consider the metric $d(\psi, \varphi) = \|\nabla \exp \psi - \nabla \exp \varphi\|_{L^2(\Omega)}$, first introduced in (2.10). Clearly,

$$(4.42) \quad d(\psi^*, \psi) = \|\nabla(\exp \psi^* + \phi) - \nabla(\exp \psi + \phi)\|_{L^2(\Omega)} = \|\nabla u^* - \nabla u\|_{L^2(\Omega)},$$

which converges to zero as $k \rightarrow \infty$ by (4.41).

Remark 4.17 (Dual variable mixed formulation). The formulation (4.38) is derived by setting $w = u^k$, $\alpha = \alpha_{k+1}$, and substituting the equation $\psi^{k+1} = \ln(u^k - \phi)$ into (4.33). If, instead, we considered the dual variable substitution $\lambda^{k+1} = \ln((u^k - \phi)/(u^{k+1} - \phi))/\alpha_{k+1}$, we would arrive at the following alternative formulation:

$$(4.43) \quad \begin{cases} \text{Find } u^{k+1} \in H_g^1(\Omega) \text{ and } \lambda^{k+1} \in L^\infty(\Omega) \text{ such that} \\ (\nabla u^{k+1}, \nabla v) - (\lambda^{k+1}, v) = (f, v) \quad \text{for all } v \in H_0^1(\Omega), \\ (u^{k+1}, \varphi) - (u^k \exp(-\alpha_{k+1} \lambda^{k+1}), \varphi) = (\phi, \varphi) \quad \text{for all } \varphi \in L^2(\Omega). \end{cases}$$

Although this is equivalent to (4.38) at the continuous level, it will induce a different Galerkin method; cf. Subsection 4.6. We leave the study of such dual variable proximal Galerkin methods for future research.

Remark 4.18 (Strict complementarity). Although Theorem 4.13 allows us to establish arbitrary orders of convergence (see Corollary A.14), it still represents the worst-case iteration complexity. In particular, our numerical experiments in Subsection 4.9.2, suggest that an improved result may be possible if the solution u^* exhibits strict complementarity.

4.6. Proximal Galerkin. Motivated by Theorem 4.13, it is natural to use finite-dimensional subspaces $V_h \subset H_0^1(\Omega)$ and $W_h \subset L^\infty(\Omega)$ in order to form a Galerkin discretization of (4.38). Thus, we arrive at Algorithm 3, which may be seen as a natural extension of Algorithm 2 to the inhomogeneous obstacle problem.

Algorithm 3: Proximal Galerkin method for the obstacle problem.

Input : Linear subspaces $V_h \subset H_0^1(\Omega)$ and $W_h \subset L^\infty(\Omega)$, initial solution guess $\psi_h^0 \in W_h$, unsummable sequence of step sizes $\alpha_k > 0$.
Output: Two approximate solutions, u_h and $\tilde{u}_h = \phi + \exp \psi_h$, and an approximate Lagrange multiplier, $\lambda_h = (\psi_h^{k-1} - \psi_h)/\alpha_k$.

Initialize $k = 0$.

repeat

 Solve the following (nonlinear) discrete saddle-point problem:

$$(4.44) \quad \begin{cases} \text{Find } u_h \in g_h + V_h \text{ and } \tilde{\psi}_h \in W_h \text{ such that} \\ (\alpha_{k+1} \nabla u_h, \nabla v) + (\tilde{\psi}_h, v) = (\alpha_{k+1} f + \psi_h^k, v) \quad \text{for all } v \in V_h, \\ (u_h, \varphi) - (\exp \tilde{\psi}_h, \varphi) = (\phi, \varphi) \quad \text{for all } \varphi \in W_h. \end{cases}$$

 Assign $\psi_h^{k+1} \leftarrow \tilde{\psi}_h$ and $k \leftarrow k + 1$.

until a convergence test is satisfied

Just like Algorithm 2, we find that Algorithm 3 delivers *two* distinct approximations of the exact solution; $u_h \in V_h$ and $\tilde{u}_h \in \phi + \exp(W_h)$. The second of these approximations is unusual because it is guaranteed to satisfy the inequality $\tilde{u}_h > \phi$. Moreover, like the continuous-level algorithm in Theorem 4.13, it also produces an approximate Lagrange multiplier,

$$(4.45) \quad \lambda_h = (\psi_h^{k-1} - \psi_h^k)/\alpha_k,$$

where k denotes the final iterate where the abstract convergence test in Algorithm 3 is satisfied.

Finite element methods typically lead to piece-wise polynomial approximations of the exact solution. Given that $\tilde{u}_h = \phi + \exp \psi_h$ relies on a non-standard type of exponential function approximation, it is natural ask whether \tilde{u}_h can produce an accurate approximation of the continuous-level solution u . The following result provides a partial positive answer to this question. The proof is given in [Appendix B.3](#).

PROPOSITION 4.19 (Approximability). *Let $u \in \text{int } L_+^\infty(\Omega)$ and define $\psi = \ln u$. Moreover, let $\psi_h \in W_h$ and $\tilde{u}_h = \exp \psi_h$. The following identity holds:*

$$(4.46) \quad \|u - \tilde{u}_h\|_{L^\infty(\Omega)} \leq \|u\|_{L^\infty(\Omega)} (\exp \|\psi - \psi_h\|_{L^\infty(\Omega)} - 1).$$

The next ordinary concern would be the stability of the discretization (4.44). In the next subsection, we propose stable pairs of finite elements that can be used to construct V_h and W_h .

4.7. Stable pairs of finite elements I: Discontinuous latent variable.

Subspace pairings determine the stability of finite element methods for saddle-point problems [33]. Thus, it should come as no surprise that the choice of the subspaces V_h and W_h is central to the proximal Galerkin method. For simplicity, we focus on stable pairs of finite elements with discontinuous latent variables ψ_h , as these appear to provide the most efficient conforming approximations per degree of freedom. The elements we propose are defined in (4.48), below. Constructing alternatives using macroelement partitions (e.g., [185]), various non-conforming approximation techniques (e.g., [54, 49]), or even spline-based approximation spaces (cf. [103]) all provide possible alternatives for the design of proximal Galerkin methods. Due to space in this manuscript and the limitations of our present software, we leave these and other possible constructions to future studies.

Here and throughout, \mathcal{T}_h always denotes a shape-regular partition of the domain $\Omega \subset \mathbb{R}^2$ into finitely many open connected triangular or quadrilateral mesh cells T with Lipschitz boundaries ∂T such that Ω is the union of the closure of all mesh cells T in \mathcal{T}_h . Following convention, $h > 0$ denotes the mesh size $h = \max_{T \in \mathcal{T}_h} \text{diam}(T)$. Let $\mathbb{P}_p(T)$ denote the space of polynomials of total order up to and including p on a triangle T . Likewise, let $\mathbb{Q}_p(T)$ denote the space of tensor-product polynomials of order up to and including p on a quadrilateral T [68]. Moreover, for any space $\mathbb{X}(T)$ of polynomials over an element $T \in \mathcal{T}_h$, we abuse notation to denote the corresponding space of “broken” polynomials $\mathbb{X}(\mathcal{T}_h) = \{\varphi \in L^\infty(\Omega) \mid \varphi|_T \in \mathbb{X}(T) \text{ for every } T \in \mathcal{T}_h\}$.

We will require spaces of degree- q polynomials on whose traces on the cell boundary ∂T have lower polynomial degree $p < q$. To this end, define the sets of so-called bubble functions in $\mathbb{P}_q(T)$ and $\mathbb{Q}_q(T)$ to be $\mathring{\mathbb{P}}^q(T) = \{\varphi \in \mathbb{P}_q(T) \mid \varphi|_{\partial T} = 0\}$ and $\mathring{\mathbb{Q}}^q(T) = \{\varphi \in \mathbb{Q}_q(T) \mid \varphi|_{\partial T} = 0\}$, respectively. Accordingly, define $\widehat{\mathbb{P}}_p(T) = \mathbb{P}_p(T) \setminus \mathring{\mathbb{P}}^p(T)$ and $\widehat{\mathbb{Q}}_p(T) = \mathbb{Q}_p(T) \setminus \mathring{\mathbb{Q}}^p(T)$. Finally, let

$$(4.47) \quad \mathbb{P}_p^q(T) = \widehat{\mathbb{P}}_p(T) \oplus \mathring{\mathbb{P}}^q(T) \quad \text{and} \quad \mathbb{Q}_p^q(T) = \widehat{\mathbb{Q}}_p(T) \oplus \mathring{\mathbb{Q}}^q(T).$$

We are now ready to define the finite element spaces, which we chose based on *a priori* analysis of a simple linearization of subproblem (4.44). For further details of the analysis, see [Appendix B](#).

For any integer $p \geq 1$, we define the following two pairs of spaces:

Triangular elements. We refer to the following as the $(\mathbb{P}_p$ -bubble, \mathbb{P}_{p-1} -broken) pairing:

$$(4.48a) \quad V_h = \mathbb{P}_p^{p+2}(\mathcal{T}_h) \cap H_0^1(\Omega), \quad W_h = \mathbb{P}_{p-1}(\mathcal{T}_h).$$

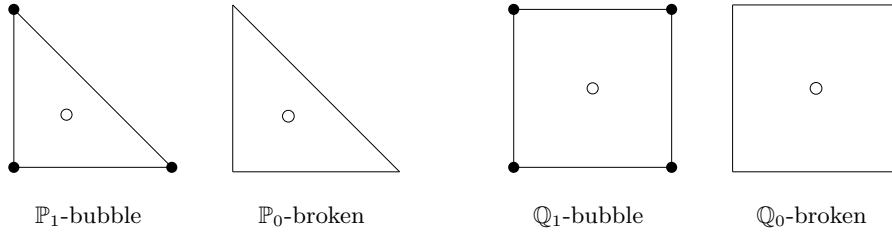


Fig. 4.4: Conventional representation of the $(\mathbb{P}_1\text{-bubble}, \mathbb{P}_0\text{-broken})$ and $(\mathbb{Q}_1\text{-bubble}, \mathbb{Q}_0\text{-broken})$ finite elements in two dimensions. The central degree of freedom in each element can be understood as the average over the individual mesh cell.

Quadrilateral elements. We refer to the following as the $(\mathbb{Q}_p\text{-bubble}, \mathbb{Q}_{p-1}\text{-broken})$ pairing:

$$(4.48b) \quad V_h = \mathbb{Q}_p^{p+1}(\mathcal{T}_h) \cap H_0^1(\Omega), \quad W_h = \mathbb{Q}_{p-1}(\mathcal{T}_h).$$

Figure 4.4 provides a visual representation of the lowest-order versions of these elements.

Remark 4.20 (Positive cell average). Assume $\phi = 0$. Although the piecewise polynomials u_h that arise from solving the subproblems (4.44) can not be guaranteed to preserve pointwise positivity, they are guaranteed to have positive cell averages. Indeed, notice that the subspaces W_h in (4.48) always include piecewise constant functions. Therefore, we may consider the second equation in (4.44) with $\varphi = 1$ in T and $\varphi = 0$ otherwise. Testing with this particular function implies that

$$(4.49) \quad \int_T u_h \, dx = \int_T \exp \psi_h > 0.$$

If $\phi \neq 0$, then a similar argument implies that each cell average of u_h lies above the corresponding cell average of ϕ .

Remark 4.21 (Alternative subspaces). Although variable-order spaces like V_h in (4.48) are supported in some software [61, 74], they may not be available in the preferred software of many users. For this reason, we also recommend the following alternative pairings:

Alternative triangular elements. We refer to the following as the $(\mathbb{P}_{p+2}, \mathbb{P}_{p-1}\text{-broken})$ subspaces:

$$(4.50a) \quad V_h = \mathbb{P}_{p+2}(\mathcal{T}_h) \cap H_0^1(\Omega), \quad W_h = \mathbb{P}_{p-1}(\mathcal{T}_h).$$

Alternative quadrilateral elements. We refer to the following as the $(\mathbb{Q}_{p+1}, \mathbb{Q}_{p-1}\text{-broken})$ subspaces:

$$(4.50b) \quad V_h = \mathbb{Q}_{p+1}(\mathcal{T}_h) \cap H_0^1(\Omega), \quad W_h = \mathbb{Q}_{p-1}(\mathcal{T}_h).$$

Since (4.48) are stable (cf. Lemma B.3), it is a straightforward consequence of the inclusions $\mathbb{P}_p^{p+2}(T) \subset \mathbb{P}_{p+2}(T)$ and $\mathbb{Q}_p^{p+1}(T) \subset \mathbb{Q}_{p+1}(T)$ that (4.50) are also stable. For further details, see Remark 4.21.

4.8. Algorithm. We conducted numerical experiments across two separate codes, FEniCSx [181] and MFEM [14], and have released our implementations to the public [111, 113]. The FEniCSx implementation [113] is Python-based and uses the (\mathbb{P}_p -bubble, \mathbb{P}_{p-1} -broken) triangular elements proposed in Subsection 4.7. The MFEM implementation is written in C++. Because MFEM does not currently support bubble function enrichment, the MFEM implementation [111] uses (\mathbb{Q}_{p+1} , \mathbb{Q}_{p-1} -broken) quadrilateral elements (see Remark 4.21) instead of the (\mathbb{Q}_p -bubble, \mathbb{Q}_{p-1} -broken) elements also proposed in Subsection 4.7.

Algorithm 4 represents the MFEM implementation. Notice that it is a quasi-Newton algorithm that includes a small modification to the local Hessian. We found that for small values of ϵ (we used $\epsilon = 10^{-6}$), this made the built-in GMRES solver more robust and did not strongly affect the convergence rates. Our FEniCSx implementation uses the standard PETSc [19] Newton solver provided by the petsc4py Python package [57]. This implementation does not involve any modifications to the local Hessian. Before stating the algorithm, we introduce notation for the element-wise gradient operator $\nabla_h: \prod_{T \in \mathcal{T}_h} H^1(T) \rightarrow L^2(\Omega)$, which satisfies

$$(4.51) \quad (\nabla_h v, \nabla_h w) = \sum_{T \in \mathcal{T}_h} \int_T \nabla v \cdot \nabla w \, dx.$$

Remark 4.22 (Practical aspects of the implementation: Stopping criteria). We use several metrics to design a meaningful stopping criterion for Algorithm 4. The algorithm uses two loops: an outer loop that updates the parameter α_k and adapts the Bregman term and an inner loop in which the step is calculated using an inexact Newton iteration. In an ideal setting, the outer loop would stop once the residual of the optimality conditions for the original problem has a sufficiently small norm. This is difficult to check in general, e.g., since the H^1 -projection operator onto the feasible set is nontrivial/expensive to evaluate. On the other hand, Theorem 4.13 provides a theoretical convergence rate, which up to a constant could be used to predetermine a maximum number of step sizes to reach a desired tolerance. Furthermore, if we set $\alpha_k = \alpha > 0$ for all k , then we may view the outer loop as a *globally* convergent fixed point iteration. Since these observations are for the ideal setting, we still should check for some indication of convergence. For this reason, we can use the distance between successive iterates as a practical stopping criterion. For the inner iterations, we check either the norm of the residual of the nonlinear term (see [113]) or the length of the proposed step (see [111]). The length of the accepted step informs the next inner iteration tolerance and ensures that more accurate steps are computed as k increases. Since the residual for the first equation is technically exact, we do not include it in the inner stopping criterion.

4.9. Numerical experiments. We performed four sets of numerical experiments in order to validate the proximal Galerkin method. The first experiment involves a smooth biactive manufactured solution that allows us to verify the (mesh-independent) iteration complexity predicted by Theorem 4.13, in addition to high-order convergence rates with respect to the polynomial order of the finite element subspaces. In the second experiment, we check the discrete Karush–Kuhn–Tucker (KKT) conditions on a manufactured solution exhibiting strict complementarity. In this case, we observe *better* iteration complexity than predicted by Theorem 4.13. We conjecture that this improved convergence order holds in general whenever a strict complementarity condition is satisfied; cf. Remark 4.18. The third experiment

Algorithm 4: Quasi-Newton proximal Galerkin method for the obstacle problem.

Input : Piecewise polynomial subspaces $V_h \subset H_0^1(\Omega)$ and $W_h \subset L^\infty(\Omega)$, initial solution guesses $u_h \in V_h$ and $\psi_h \in W_h$, unsummable sequence of step sizes $\alpha_k > 0$, exit tolerance $\text{tol}_{\text{exit}} > 0$, and initial Newton tolerance $\text{tol}_{\text{Newton}} > \text{tol}_{\text{exit}}$.

Output: Two approximate solutions, u_h and $\tilde{u}_h = \phi + \exp \psi_h$, and an approximate Lagrange multiplier, $\lambda_h = (\psi_h^{k-1} - \psi_h)/\alpha_k$.

Initialize $k = 0$.

repeat

Assign $u_h^k \leftarrow u_h$ and $\psi_h^k \leftarrow \psi_h$.

repeat

Assign $w_h \leftarrow u_h$.

Reassign u_h and δ_h by solving the following (linearized) discrete saddle-point problem:

$$\left\{ \begin{array}{l} \text{Find } u_h \in V_h \text{ and } \delta_h \in W_h \text{ such that} \\ (\alpha_{k+1} \nabla u_h, \nabla v) + (\delta_h, v) = (\alpha_{k+1} f + \psi_h^k - \psi_h, v) \quad \text{for all } v \in V_h, \\ (u_h, \varphi) - c_\varepsilon(\psi_h, \delta_h, \varphi) = (\phi + \exp \psi_h, \varphi) \quad \text{for all } \varphi \in W_h, \end{array} \right.$$

where

$$c_\varepsilon(\psi_h, \delta_h, \varphi) = \begin{cases} (\delta_h \exp \psi_h, \varphi) + \varepsilon(\delta_h, \varphi) & \text{if } W_h \in \ker(\nabla_h), \\ (\delta_h \exp \psi_h, \varphi) + \varepsilon(\nabla_h \delta_h, \nabla_h \varphi) & \text{otherwise.} \end{cases}$$

Assign $\psi_h \leftarrow \psi_h + \delta_h$.

until $\|u_h - w_h\|_{L^2(\Omega)}^2 \leq \text{tol}_{\text{Newton}}^2$

Assign $\text{tol}_{\text{Newton}} \leftarrow \|u_h - u_h^k\|_{L^2(\Omega)}$ and $k \leftarrow k + 1$.

until $\text{tol}_{\text{Newton}} < \text{tol}_{\text{exit}}$

involves a non-smooth biactive solution and is included to further stress test the proximal Galerkin method. Finally, in our fourth experiment, we consider a benchmark obstacle problem from the literature and demonstrate our ability solve this problem with the highest order finite elements currently supported in our MFEM code; i.e., we used $p = 12$.

Each of our experiments were conducted on standard sequences of uniformly refined nested meshes $\mathcal{T}_h, \mathcal{T}_{h/2}, \mathcal{T}_{h/4}, \dots$ conforming to unit ball domains in \mathbb{R}^2 . The experiments with the triangular elements (FEniCSx) used an ℓ^∞ -unit ball (i.e., square) domain,

$$(4.52a) \quad \Omega_\infty = \{(x, y) \in \mathbb{R}^2 \mid \max\{|x|, |y|\} < 1\} \subset \mathbb{R}^2,$$

with initial mesh size denoted $h = h_\infty$. Meanwhile, the experiments with the quadrilateral elements (MFEM) used an ℓ^2 -unit ball (i.e., circular) domain,

$$(4.52b) \quad \Omega_2 = \{(x, y) \in \mathbb{R}^2 \mid x^2 + y^2 < 1\} \subset \mathbb{R}^2,$$

with initial mesh size denoted $h = h_2$, which was uniformly refined using a standard

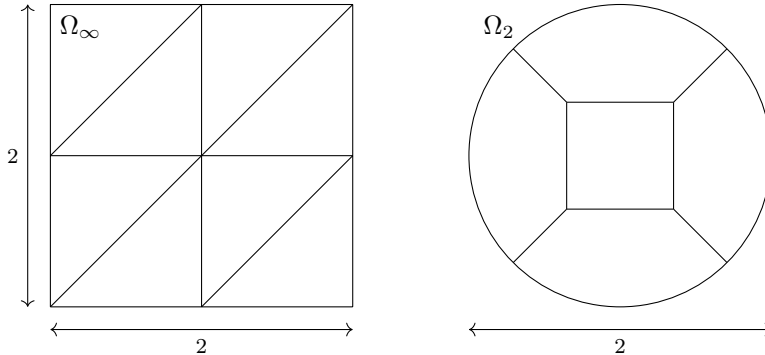


Fig. 4.5: Initial finite element meshes for computational domains Ω_∞ and Ω_2 . We denote their mesh sizes $h = h_\infty$ and $h = h_2$, respectively. Left: Initial triangular mesh for the $(\mathbb{P}_p\text{-bubble}, \mathbb{P}_{p-1}\text{-broken})$ subspace pair on the domain Ω_∞ (FEniCSx experiments). Right: Initial five-element curvilinear quadrilateral mesh for the $(\mathbb{Q}_{p+1}, \mathbb{Q}_{p-1}\text{-broken})$ subspace pair on the domain Ω_2 (MFEM experiments). Across our experiments, we consider various polynomial orders $p \geq 1$ for both of these subspace pairs.

transfinite interpolation rule to handle the curvilinear element mappings [82]. The initial meshes in these sequences are depicted in Figure 4.5.

4.9.1. Experiment 1: Smooth biactive solution. In this experiment, we set $\phi = 0$ and $g = u$, where $u(x, y)$ is the smooth manufactured solution

$$(4.53) \quad u(x, y) = \begin{cases} 0 & \text{if } x < 0, \\ x^4 & \text{otherwise,} \end{cases} \quad \text{implied by } f(x, y) = \begin{cases} 0 & \text{if } x < 0, \\ -12x^2 & \text{otherwise.} \end{cases}$$

See Figures 4.6 and 4.7 for depictions of the exact solution on Ω_∞ and Ω_2 , respectively.

This problem is specifically chosen to exhibit *biactivity*; i.e., both the inequality constraint $u \geq 0$ and the associated Lagrange multiplier are simultaneously equal to zero on a set of positive measure; i.e., on the set $\{(x, y) \mid x < 0\}$. Such problems are notoriously difficult to solve for certain classes of algorithms, such as active set methods. Biactivity, also known as weak activity or lack of strict complementarity, is a notion from nonlinear optimization that indicates a kind of degenerate nonsmoothness of the primal-dual system of equations used to calculate the solution. It is often associated with a lack of stability with respect to perturbations of the data, as well. We refer the interested reader to any standard text of numerical optimization; see, e.g., [163, Definition 12.5].

Our first aim is to use this challenging example to illustrate *mesh-independence* of the proximal Galerkin method. To this end, we use Table 4.1 to record the values of the increments $\|u_h^k - u_h^{k-1}\|_{H^1(\Omega_\infty)}$ taken from a sequence of refined meshes with polynomial orders $p = 1, 2$ from our FEniCSx implementation [113]. The specific step size rule used to generate this data is chosen based on Corollary A.14 to deliver *superlinear* convergence (in iterations), and is given as follows:

$$(4.54a) \quad \alpha_1 = 1, \quad \alpha_k = \min\{\max\{\alpha_1, r^{q^{k-1}} - \alpha_{k-1}\}, 10^{10}\}, \quad k = 2, 3, \dots,$$

where $r = q = 1.5$. Note that, for each iteration k , the increments in Table 4.1

| Progress of the iterates $\ u_h^k - u_h^{k-1}\ _{H^1(\Omega_\infty)}$ for various h and p | | | | | | |
|---|-------------------|--------------------------|-----------------------|-----------------------|--------------------------|-----------------------|
| | | Polynomial order $p = 1$ | | | Polynomial order $p = 2$ | |
| k | α_k | $h_\infty/16$ | $h_\infty/32$ | $h_\infty/64$ | $h_\infty/16$ | $h_\infty/32$ |
| 1 | 1.0 | $2.10 \cdot 10^0$ | $2.10 \cdot 10^0$ | $2.10 \cdot 10^0$ | $2.10 \cdot 10^0$ | $2.10 \cdot 10^0$ |
| 2 | 1.0 | $6.45 \cdot 10^{-1}$ | $6.45 \cdot 10^{-1}$ | $6.45 \cdot 10^{-1}$ | $6.45 \cdot 10^{-1}$ | $6.45 \cdot 10^{-1}$ |
| 3 | 1.49 | $1.73 \cdot 10^{-1}$ | $1.73 \cdot 10^{-1}$ | $1.73 \cdot 10^{-1}$ | $1.73 \cdot 10^{-1}$ | $1.73 \cdot 10^{-1}$ |
| 4 | 2.43 | $1.10 \cdot 10^{-1}$ | $1.10 \cdot 10^{-1}$ | $1.10 \cdot 10^{-1}$ | $1.10 \cdot 10^{-1}$ | $1.10 \cdot 10^{-1}$ |
| 5 | 5.35 | $7.76 \cdot 10^{-2}$ | $7.77 \cdot 10^{-2}$ | $7.77 \cdot 10^{-2}$ | $7.77 \cdot 10^{-2}$ | $7.77 \cdot 10^{-2}$ |
| 6 | $1.64 \cdot 10^1$ | $4.76 \cdot 10^{-2}$ | $4.77 \cdot 10^{-2}$ | $4.77 \cdot 10^{-2}$ | $4.77 \cdot 10^{-2}$ | $4.77 \cdot 10^{-2}$ |
| 7 | $8.50 \cdot 10^1$ | $2.24 \cdot 10^{-2}$ | $2.25 \cdot 10^{-2}$ | $2.25 \cdot 10^{-2}$ | $2.25 \cdot 10^{-2}$ | $2.25 \cdot 10^{-2}$ |
| 8 | $9.35 \cdot 10^2$ | $5.82 \cdot 10^{-3}$ | $5.84 \cdot 10^{-3}$ | $5.85 \cdot 10^{-3}$ | $5.85 \cdot 10^{-3}$ | $5.85 \cdot 10^{-3}$ |
| 9 | $3.17 \cdot 10^4$ | $6.04 \cdot 10^{-4}$ | $6.07 \cdot 10^{-4}$ | $6.07 \cdot 10^{-4}$ | $6.07 \cdot 10^{-4}$ | $6.07 \cdot 10^{-4}$ |
| 10 | $5.85 \cdot 10^6$ | $1.80 \cdot 10^{-5}$ | $1.81 \cdot 10^{-5}$ | $1.81 \cdot 10^{-5}$ | $1.81 \cdot 10^{-5}$ | $1.81 \cdot 10^{-5}$ |
| 11 | $1 \cdot 10^{10}$ | $9.41 \cdot 10^{-8}$ | $9.47 \cdot 10^{-8}$ | $9.49 \cdot 10^{-8}$ | $9.50 \cdot 10^{-8}$ | $9.50 \cdot 10^{-8}$ |
| 12 | $1 \cdot 10^{10}$ | $2.10 \cdot 10^{-12}$ | $2.00 \cdot 10^{-12}$ | $1.96 \cdot 10^{-12}$ | $1.92 \cdot 10^{-12}$ | $1.95 \cdot 10^{-10}$ |
| Tot. linear solves | | 21 | 20 | 19 | 19 | 19 |

Table 4.1: Biactivity. Table of increments $\|u_h^k - u_h^{k-1}\|_{H^1(\Omega_\infty)}$ for various mesh sizes h and polynomial orders p using the triangular element (\mathbb{P}_p -bubble, \mathbb{P}_{p-1} -broken) discretization. The initial degrees of freedom for u_h and ψ_h were set to zero at the beginning of each run. Between eight and ten Newton iterations performed by the PETSc Newton solver used for each initial subproblem solve and then only one Newton solve was used for each of the following subproblems. The convergence of the increments for each fixed k and the boundedness of the number of linear solves indicates *mesh-independence*.

converge to fixed values as the mesh is refined or the polynomial order is raised. Moreover, the total number of linear equation solves remains bounded. Both of these characteristics are emblematic of a *mesh-independent* numerical method.

Our next aim is to verify the convergence orders predicted by [Theorem 4.13](#) and [Corollary A.14](#). In doing so, we consider the double-exponential step size rule [\(4.54a\)](#) alongside the geometric rule

$$(4.54b) \quad \alpha_k = r^{k-1}, \quad k = 1, 2, \dots,$$

with $r = 2$, and the fixed step size rule $\alpha_k = 1$, for all $k = 1, 2, \dots$. The results in [Figure 4.6](#) agree precisely with the predictions made later on in [Corollary A.14](#). In particular, notice that the fixed step rule $\alpha_k = 1$ leads to *sublinear* convergence, the geometric rule [\(4.54b\)](#) induces *linear* convergence, and the double-exponential step size rule [\(4.54a\)](#) delivers *superlinear* convergence.

The final aim of this experiment is to demonstrate that high-order convergence rates (with respect to the mesh size h) can be achieved using polynomial orders $p > 1$. To this end, we use our high-order MFEM implementation to solve for the biactive solution [\(4.53\)](#) on the circular domain $\Omega = \Omega_2$. In [Figure 4.7](#), we plot the approximation errors of the discrete solutions u_h , \tilde{u}_h , and λ_h . From these results, we witness that high-order convergence rates can, indeed, be achieved with the proximal Galerkin method. All results from this experiment can be reproduced by running the FEniCSx code `obstacle.py` or the MFEM code `obstacle.cpp` available at [\[113\]](#).

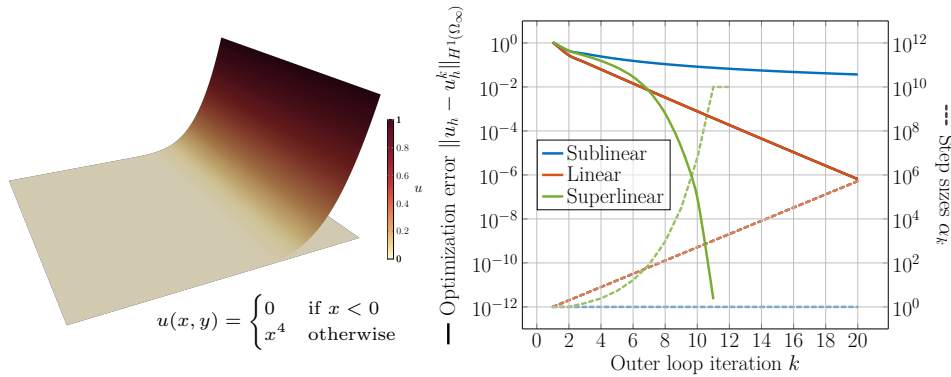


Fig. 4.6: Biactivity. Verifying the convergence orders predicted by [Corollary A.14](#) with the $(\mathbb{P}_1\text{-bubble}, \mathbb{P}_0\text{-broken})$ discretization (FEniCSx). Left: The exact solution. Right: Plots of the optimization error $\|u_h - u_h^k\|_{H^1(\Omega_\infty)}$ and corresponding step size α_k when $h = h_\infty/16$. The blue curve tracks the *sublinear* convergence induced by the fixed step size rule $\alpha_k = 1$. Meanwhile, the step size rules [\(4.54b\)](#) and [\(4.54a\)](#) induce *linear* (red) and *superlinear* (green) convergence, respectively. The results are similar on finer meshes due to mesh-independence; cf. [Table 4.1](#).

4.9.2. Experiment 2: Strict complementarity. In this experiment, we set $\phi = g = 0$ and define

$$(4.55) \quad f(x, y) = 2\pi^2 \sin(\pi x) \sin(\pi y).$$

See [Figure 4.8](#) for a fine mesh ($h = h_\infty/128$) solution u_h as well as the associated Lagrange multiplier λ_h .

When viewed from the perspective of continuum mechanics, the multiplier λ is a resolvent force. It is therefore rare that we would see biactivity of the type in the previous experiments on such large domains as it would correspond to contact without any opposing force resulting from the obstacle.

The first aim of this experiment is to revisit the convergence orders predicted by [Theorem 4.13](#) and [Corollary A.14](#) and demonstrate that they are overly pessimistic for this more typical type of problem. Indeed, as demonstrated in [Figure 4.8](#), we see that linear convergence is achieved using only a fixed step size. In turn, superlinear convergence can be achieved using any unbounded step size rule. For illustration, we have added results using the geometric rule [\(4.54b\)](#) with various growth parameters $r \in \{1.05, 1.1, 2\}$.

The second aim of this experiment is to check convergence of the discrete solution via the KKT conditions. This is useful to assess *a posteriori* the optimality of the discrete solution when the true solution is unknown. To this end, we consider the complementarity condition $|\int_\Omega \lambda u \, dx| = 0$, the primal feasibility condition $\int_\Omega \max\{-u, 0\} \, dx = 0$, and the dual feasibility condition $\int_\Omega \max\{-\lambda, 0\} \, dx = 0$. We note that the discrete solution $\tilde{u}_h = \exp \psi_h$ is always feasible by construction, i.e., $\tilde{u}_h \geq 0$. Therefore, in order to glean more interesting information about the proximal Galerkin solution, we focus on discrete versions of the KKT conditions formulated in terms of the solution variable u_h . The discrete KKT conditions that we

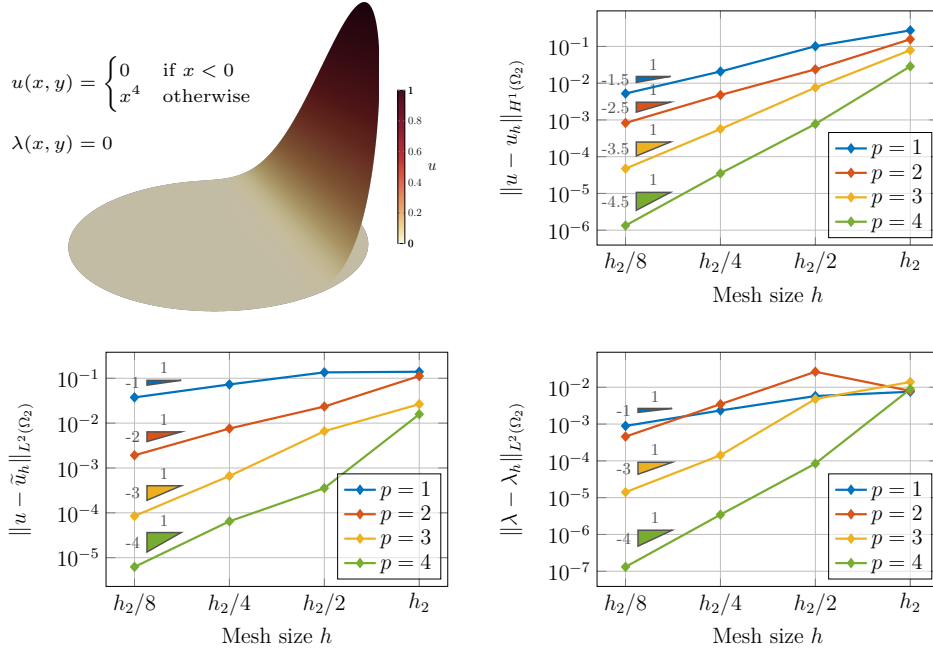


Fig. 4.7: Biactivity. Verifying high-order convergence of the approximation error with various polynomial order ($\mathbb{Q}_{p+1}, \mathbb{Q}_{p-1}$ -broken)-discretizations on $\Omega = \Omega_2$. Top left: The exact solution. Top right: The approximation error of the discrete solution u_h in the H^1 -norm. Bottom left: The approximation error of the discrete solution $\tilde{u}_h = \exp \psi_h$ in the L^2 -norm. Bottom right: The approximation error of the Lagrange multiplier λ_h in the L^2 -norm. Notice that each of the convergence rates for this smooth solution grow with the polynomial order.

checked are recorded in Table 4.2. From the results in this table, we see that discrete primal feasibility, $\int_{\Omega} \max\{-u_h, 0\} dx = 0$, is achieved only in the limit $h \rightarrow 0$. However, discrete complementarity, $|\int_{\Omega} \lambda_h u_h dx| = 0$, and discrete dual feasibility, $\int_{\Omega} \max\{-\lambda_h, 0\} dx = 0$, appear to hold for all mesh sizes. The results of this experiment can be reproduced by running the FEniCSx code `obstacle.py` available at [113].

4.9.3. Experiment 3: Biactive solution, nonsmooth multiplier. In this experiment, we set $\phi = 0$ and $g = u$, where $u(x, y)$ is the smooth manufactured solution

$$(4.56a) \quad u(x, y) = \begin{cases} (1 - 4x^2 - 4y^2)^4 & \text{if } x^2 + y^2 < 1/4, \\ 0 & \text{otherwise,} \end{cases}$$

implied by the forcing function

$$(4.56b) \quad f(x, y) = -\Delta u(x, y) - \begin{cases} 1 & \text{if } x^2 + y^2 > 3/4, \\ 0 & \text{otherwise.} \end{cases}$$

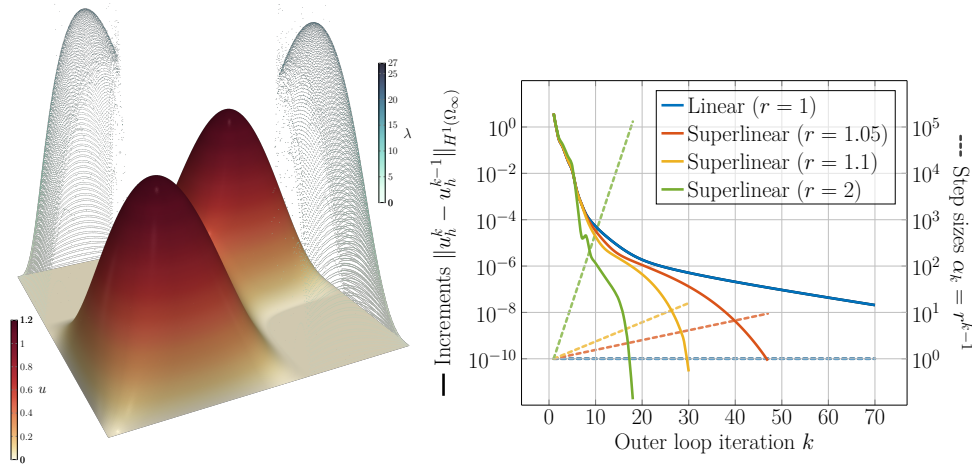


Fig. 4.8: Strict complementarity. Surpassing the convergence orders predicted by [Corollary A.14](#) with the (\mathbb{P}_1 -bubble, \mathbb{P}_0 -broken) discretization (FEniCSx). Left: A high-resolution image of the exact solution u and associated Lagrange multiplier λ . Right: Plots of the primal variable increments $\|u_h^k - u_h^{k-1}\|_{H^1(\Omega_\infty)}$ and corresponding geometric step sizes $\alpha_k = r^{k-1}$ for mesh size $h = h_\infty/16$. The results are similar on all finer meshes due to mesh-independence; cf. [Table 4.1](#). Analogous results can also be obtained for higher-order and quadrilateral element discretizations (not shown).

| h | Complementarity $ \int_{\Omega_\infty} \lambda_h u_h \, dx $ | Primal feasibility $\int_{\Omega_\infty} \max\{-u_h, 0\} \, dx$ | Dual feasibility $\int_{\Omega_\infty} \max\{-\lambda_h, 0\} \, dx$ |
|---------------|---|--|--|
| h_∞ | (all less than 10^{-14}) | $6.97 \cdot 10^{-3}$ | (all less than 10^{-12}) |
| $h_\infty/2$ | | $9.09 \cdot 10^{-3}$ | |
| $h_\infty/4$ | | $1.16 \cdot 10^{-3}$ | |
| $h_\infty/8$ | | $1.69 \cdot 10^{-4}$ | |
| $h_\infty/16$ | | $4.08 \cdot 10^{-5}$ | |
| $h_\infty/32$ | | $4.53 \cdot 10^{-6}$ | |

Table 4.2: Strict complementarity. Checking the discrete KKT conditions for the proximal Galerkin solution owing to [\(4.55\)](#). Here, we see that primal feasibility is achieved in the limit $h \rightarrow 0$. Meanwhile, complementary and dual feasibility holds on all meshes.

Clearly, this is another solution exhibiting biactivity. In this case, however, the multiplier,

$$(4.56c) \quad \lambda(x, y) = \begin{cases} 1 & \text{if } x^2 + y^2 > 3/4, \\ 0 & \text{otherwise,} \end{cases}$$

is discontinuous. See [Figure 4.9](#) for a depiction of the exact solution u as well as the associated Lagrange multiplier λ on the domain $\Omega = \Omega_\infty$.

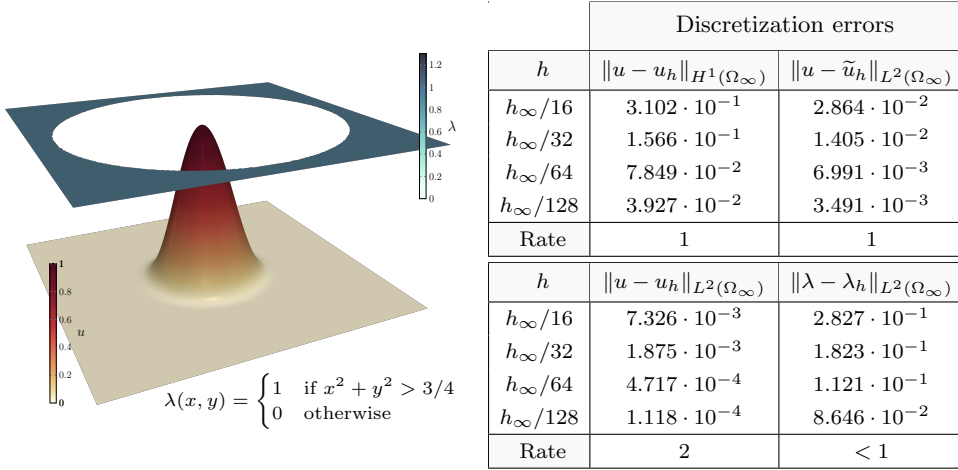


Fig. 4.9: Biactive solution, nonsmooth multiplier. Checking the discretization errors in various standard norms. Notice that the error in the Lagrange multiplier variable does not decay linearly with respect to the L^2 -norm.

We use this experiment to inspect the approximation error of the (\mathbb{P}_1 -bubble, \mathbb{P}_0 -broken) discretization. See Figure 4.9 for our results. As expected, unlike for the biactive solution studied in Subsection 4.9.1, the L^2 -error of the Lagrange multiplier does not decay to zero linearly. We observed no other adverse effects from this nonsmooth manufactured solution. The results of this experiment can be reproduced by running the FEniCSx code `obstacle.py` available at [113].

4.9.4. Experiment 4: Spherical obstacle. Our final experiment is motivated by an exact solution in [86]. Here, we set both $f = 0$ and $g = 0$ and define the obstacle to be the upper surface of a sphere of radius $1/2$, namely

$$(4.57) \quad \phi(x, y) = \sqrt{1/4 - x^2 - y^2},$$

if $\sqrt{x^2 + y^2} \leq 1/2$, and assume that ϕ is sufficiently negative when $\sqrt{x^2 + y^2} > 1/2$ so that no contact happens on that subdomain. Exploiting radial symmetry, the exact solution on the circular domain $\Omega = \Omega_2$ is found to be

$$(4.58) \quad u(x, y) = \begin{cases} A \ln \sqrt{x^2 + y^2} & \text{if } \sqrt{x^2 + y^2} > a, \\ \phi(x, y) & \text{otherwise,} \end{cases}$$

where $a = \exp(W_{-1}(\frac{-1}{2e^2})/2 + 1) \approx 0.34898$, $A = \sqrt{1/4 - a^2}/\ln a \approx -0.34012$, and $W_j(\cdot)$ is the j -th branch of the Lambert W-function.

Figure 4.10 presents the very high order ($p = 12$) proximal Galerkin solutions u_h and \tilde{u}_h to this problem on the coarsest mesh, $h = h_2$, which has only five elements. This is the highest order discretization of the obstacle problem that we have seen in the literature. Table 4.3 compares the subproblem error, $\|u - u_h^k\|_{H^1(\Omega_2)}$, on a sequence of uniformly meshes. From this table, we see that if the number of outer iterations k is held fixed, then the error converges to fixed values as the mesh is refined. This is another hallmark of mesh-independence. This experiment is an official part of MFEM 4.6; in particular, see MFEM Example 36 [111].

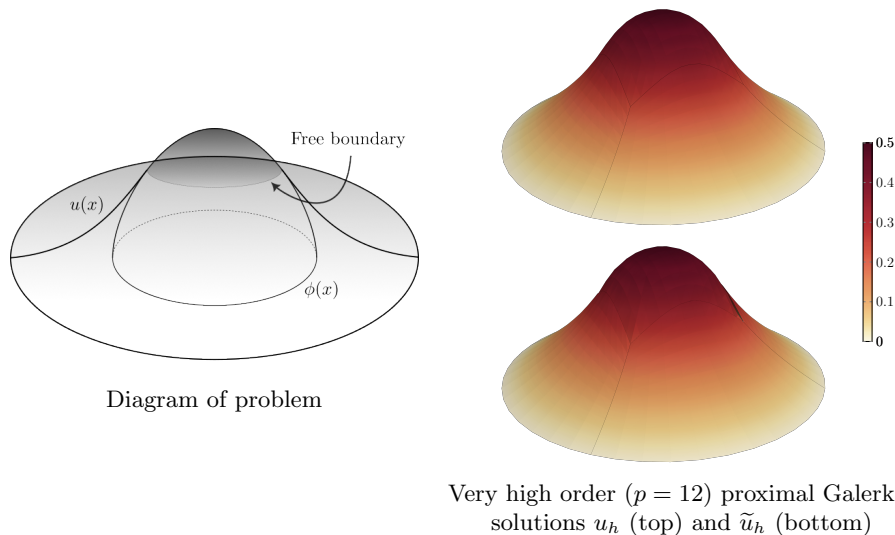


Fig. 4.10: Spherical obstacle. Benchmark obstacle problem from [86]. Left: Diagram of the problem set-up. Right: Five-element proximal Galerkin solutions u_h and \tilde{u}_h .

| k | Linear solves | Primal errors $\ u - u_h^k\ _{H^1(\Omega_2)}$ for $p = 1$ | | | | |
|---------------------|---------------|---|----------------------|----------------------|----------------------|----------------------|
| | | $h_2/8$ | $h_2/16$ | $h_2/32$ | $h_2/64$ | $h_2/128$ |
| 1 | 3 | $2.72 \cdot 10^{-1}$ | $2.70 \cdot 10^{-1}$ | $2.70 \cdot 10^{-1}$ | $2.70 \cdot 10^{-1}$ | $2.70 \cdot 10^{-1}$ |
| 2 | 1 | $1.37 \cdot 10^{-1}$ | $1.38 \cdot 10^{-1}$ | $1.38 \cdot 10^{-1}$ | $1.38 \cdot 10^{-1}$ | $1.38 \cdot 10^{-1}$ |
| 3 | 1 | $3.62 \cdot 10^{-2}$ | $3.33 \cdot 10^{-2}$ | $3.31 \cdot 10^{-2}$ | $3.31 \cdot 10^{-2}$ | $3.31 \cdot 10^{-2}$ |
| \vdots | \vdots | \vdots | \vdots | \vdots | \vdots | \vdots |
| Total iterations | | 11 | 11 | 11 | 11 | 11 |
| Total linear solves | | 13 | 13 | 13 | 13 | 13 |
| Final error | | $1.98 \cdot 10^{-2}$ | $8.73 \cdot 10^{-3}$ | $3.49 \cdot 10^{-3}$ | $1.18 \cdot 10^{-3}$ | $3.85 \cdot 10^{-4}$ |

Table 4.3: Spherical obstacle. Checking the subproblem error, $\|u - u_h^k\|_{H^1(\Omega_2)}$, for various mesh sizes using the $(\mathbb{Q}_2, \mathbb{Q}_0)$ -broken discretization. We used $\alpha_k = 1$ for all $k = 1, \dots$ and stopped the algorithm when $\|u_h^k - u_h^{k-1}\|_{L^2(\Omega_2)} < 10^{-6}$.

5. Extensions I: More general bound constraints and variational inequalities with an application to enforcing discrete maximum principles.

The purpose of this section is to move beyond the proximal framework developed in Section 4 for energy principles with pointwise lower bound constraints. To this end, we aim to answer the following two necessary questions:

1. Can proximal Galerkin be used to *simultaneously* enforce pointwise upper and lower bound constraints?
2. Can proximal Galerkin be applied to variational inequalities that *do not* arise from an energy minimization principle?

The answer to both of these questions is *yes*.

We use our answers to these questions to construct a discrete maximum principle-preserving proximal Galerkin method for the advection-diffusion equation,

$$(5.1) \quad -\epsilon \Delta u + \beta \cdot \nabla u = f \quad \text{in } \Omega, \quad u = g \quad \text{on } \partial\Omega,$$

where $\epsilon > 0$ and $\beta \in \mathbb{R}^d$ are fixed, $f \in L^\infty(\Omega)$, and $g \in H^1(\Omega) \cap C(\bar{\Omega})$. Along the way, we introduce the *binary entropy* (Subsection 5.1), a proximal algorithm for VIs with *non-symmetric* coercive bilinear forms (Subsection 5.2), and an alternative type of proximal Galerkin discretization employing a *continuous* latent variable (Subsection 5.3). The section features an implementable algorithm and closes with a brief survey of numerical experiments.

5.1. Binary entropy. In the previous section, enforcing a pointwise lower bound on the minimizer of the Dirichlet energy led us to consider a sequence of entropy-regularized energy minimization problems. We now consider the situation of enforcing pointwise upper and lower bounds *simultaneously*. For simplicity, we illustrate the approach on the so-called double-obstacle problem written in (5.2) below.

Let $\phi_1, \phi_2 \in H^1(\Omega) \cap L^\infty(\Omega)$ with $\text{ess sup}(\phi_2 - \phi_1) > 0$ and $\text{ess sup} \gamma(\phi_1 - g) < 0 < \text{ess inf} \gamma(\phi_2 - g)$, and consider minimizing the Dirichlet energy under the pointwise bound constraints $\phi_1 \leq v \leq \phi_2$. More specifically,

$$(5.2) \quad u^* = \arg \min_{v \in K} E(v),$$

where $E(v) = \frac{1}{2} \int_{\Omega} |\nabla v|^2 dx - \int_{\Omega} v f dx$ and $K = \{v \in H_g^1(\Omega) \mid \phi_1 \leq v \leq \phi_2\}$. A natural way to apply entropy regularization to (5.2) is revealed if we rewrite the problem with two new variables $v_1 = v - \phi_1$ and $v_2 = \phi_2 - v$. Doing so, we arrive at the equivalent equality-constrained optimization problem

$$(5.3) \quad (u_1^*, u_2^*) = \arg \min_{(v_1, v_2) \in K_1 \times K_2} E(v_1 + \phi_1) \quad \text{subject to } v_1 + v_2 = \phi_2 - \phi_1,$$

where $K_1 = \{v_1 \in H_{g-\phi_1}^1(\Omega) \mid v_1 \geq 0\}$ and $K_2 = \{v_2 \in H_{\phi_2-g}^1(\Omega) \mid v_2 \geq 0\}$. It can be readily verified that $u_1^* = u^* - \phi_1$ and $u_2^* = \phi_2 - u^*$.

Following our treatment of entropy regularization for pointwise non-negativity constraints in Section 4, it stands to consider the sequence $u^k = u_1^k - \phi_1 = \phi_1 - u_2^k \rightarrow u^*$ defined

$$(5.4a) \quad (u_1^k, u_2^k) = \arg \min_{(v_1, v_2) \in K_1 \times K_2} \{E(v_1 + \phi_1) + \alpha_k^{-1} (D(v_1, u_1^{k-1}) + D(v_2, u_2^{k-1}))\}$$

$$(5.4b) \quad \text{subject to } v_1 + v_2 = \phi_2 - \phi_1.$$

We may now resubstitute $v_1 = v - \phi_1$ and $v_2 = \phi_2 - v$ into (5.4a), which leads to

$$(5.5) \quad u^k = \arg \min_{v \in K} \{E(v) + \alpha_k^{-1} D_B(v, u^{k-1})\},$$

where

$$(5.6) \quad D_B(v, w) = \int_{\Omega} (v - \phi_1) \ln \left| \frac{v - \phi_1}{w - \phi_1} \right| + (\phi_2 - v) \ln \left| \frac{\phi_2 - v}{\phi_2 - w} \right| dx,$$

is the Bregman divergence of the (generalized) *binary entropy*

$$(5.7) \quad B(v) = \int_{\Omega} (v - \phi_1) \ln |v - \phi_1| + (\phi_2 - v) \ln |\phi_2 - v| dx.$$

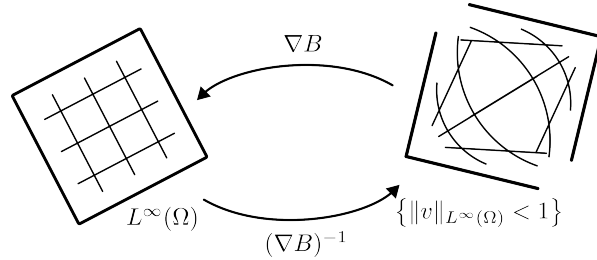


Fig. 5.1: The sigmoid map $(\nabla B)^{-1}(v) = \tanh v$ is a diffeomorphism between the Banach algebra $L^\infty(\Omega)$ and the $L^\infty(\Omega)$ -unit ball.

The cases $(\phi_1, \phi_2) = (0, 1)$ and $(\phi_1, \phi_2) = (-1, 1)$ are somewhat special for the binary entropy functional (5.7). In the first of these, (5.7) is usually referred to as the (negative) Fermi–Dirac or electronic entropy [193]. As these particular upper and lower bounds will appear prominently later on, we choose to adopt special notation for the corresponding entropy gradient and its inverse; namely,

$$(5.8) \quad \nabla B(v) = \text{lnit } v := \ln \frac{v}{1-v} \quad \text{and} \quad (\nabla B)^{-1}(v) = \text{expit } v := \frac{\exp v}{\exp v + 1}.$$

The second case, $(\phi_1, \phi_2) = (-1, 1)$, provides a gradient that is an explicit diffeomorphism between the $L^\infty(\Omega)$ -unit ball, denoted $\mathcal{B}^\infty(\Omega) = \{v \in L^\infty(\Omega) \mid \|v\|_{L^\infty(\Omega)} < 1\}$, and the entire Banach algebra $L^\infty(\Omega)$. More explicitly, we write

$$(5.9) \quad \nabla B(v) = \text{arctanh } v \quad \text{and} \quad (\nabla B)^{-1}(v) = \tanh v,$$

with the diffeomorphism illustrated visually in Figure 5.1. For posterity, we use the latter case of the binary entropy functional to define a canonical *binary-entropic Poisson equation*,

$$(5.10) \quad -\Delta u + \text{arctanh } u = f,$$

which follows from writing the strong form of the first-order optimality condition for (5.5) with $\phi_1 = -1$, $\phi_2 = 1$, $\alpha_k = 1$, and $u^{k-1} = 0$.

5.2. Variational inequalities with non-symmetric bilinear forms. In order to develop a proximal Galerkin method for the advection-diffusion equation (5.1), we first propose a continuous-level proximal algorithm for non-symmetric bilinear forms. The approach is based on the proof technique for the classical theorem of Lions and Stampacchia on the existence of solutions to elliptic variational inequalities with non-symmetric bilinear forms, which can be understood as employing a “linearized” version of algorithm (4.25).

Assume that $H \subset L^2(\Omega)$ is a real separable Hilbert space and $a: H \times H \rightarrow \mathbb{R}$ is bilinear, continuous, and coercive. In particular, assume that there exist constants $C_a, c_a > 0$ such that

$$(5.11) \quad a(w, v) \leq C_a \|w\|_H \|v\|_H \quad \text{and} \quad c_a \|v\|_H^2 \leq a(v, v),$$

for all $w, v \in H$. For any nonempty closed convex set $K \subset H$ and function $f \in L^2(\Omega)$, the Lions–Stampacchia theorem [140] states that the following variational inequality

is well-posed:

$$(5.12) \quad \begin{cases} \text{Find } u^* \in K \text{ such that} \\ a(u^*, v - u^*) \geq (f, v - u^*) \text{ for all } v \in K. \end{cases}$$

The proof works by arguing that for $\rho \in (0, 2c_a/C_a^2)$, the mapping $Q_\rho : H \rightarrow K$, defined as the unique solution of the problem

$$(5.13a) \quad \begin{cases} \text{Given } u \in H, \text{ find } w \in K \text{ such that} \\ (w, v - w)_V \geq (u, v - w)_V - \rho[a(u, v - w) - (f, v - w)] \text{ for all } v \in K, \end{cases}$$

is a contraction. Note that we may equivalently write

$$(5.13b) \quad Q_\rho(u) = \arg \min_{v \in K} \left\{ a(u, v) - (f, v) + \frac{1}{2\rho} \|v - u\|_H^2 \right\},$$

which illustrates the following relationship to the proximal operator introduced in [Section 4.4](#), $Q_\rho(u) = \text{prox}_{\rho[a(u, \cdot) - (f, \cdot)]}(u)$. Clearly, if we find the fixed point $u = Q_\rho(u)$, then [\(5.13a\)](#) reduces to [\(5.12\)](#) and we deduce that $u = u^*$.

This method of successive approximations,

$$(5.13c) \quad u^0 \in H, \quad u^{k+1} = Q_\rho(u^k), \quad k = 0, 1, \dots$$

is well-known and has been analyzed in, e.g., [\[195\]](#). However, it is not exactly amenable to computation because it requires a separate subproblem solver for each of the VIs [\(5.13a\)](#). Given a set K with appropriate structure, we can circumvent this issue using the proximal Galerkin methodology.

We begin by regularizing the continuous-level subproblems [\(5.13a\)](#). A first approach would be to use the Bregman proximal point algorithm [\(4.28\)](#) to solve the subproblems to an iteration-dependent tolerance $\text{tol}_k > 0$. This would result in a sequence of *inexact* successive approximations $\|u^{k+1} - Q_\rho(u^k)\|_V \leq \text{tol}_k$ that could converge to u^* if the sequence of tolerances decays to zero as $k \rightarrow \infty$. The potential drawback of this approach is that it creates an additional nested sequence of iterations. In turn, generating each iterate u^{k+1} may require numerous individual proximal point iterations [\(4.28\)](#) for every inexact fixed point iteration $u^{k+1} \approx Q_\rho(u^k)$.

Instead of using the Bregman proximal point algorithm as a subproblem solver, we propose to modify the original fixed point map [\(5.13b\)](#) by adding an additional Bregman divergence term. More specifically, we propose considering the alternative fixed-point iteration

$$(5.14a) \quad u^0 \in \text{dom } G', \quad u^{k+1} = Q_\rho^{\alpha_{k+1}}(u^k), \quad k = 0, 1, \dots$$

where $G: \text{dom } G \rightarrow \mathbb{R} \cup \{\infty\}$ is a strictly convex entropy functional associated to the feasible set $K \supset \text{int dom } G$ and $Q_\rho^\alpha: \text{int dom } G \rightarrow \text{int dom } G$ is an operator formally defined for all $\rho, \alpha > 0$ as follows:

$$(5.14b) \quad Q_\rho^\alpha(u) = \arg \min_{v \in K} \left\{ a(u, v) - (f, v) + \frac{1}{2\rho} \|v - u\|_H^2 + \frac{1}{\alpha} D_G(v, u) \right\}.$$

Returning to the advection-diffusion problem [\(5.1\)](#), we now assume that $0 \leq f \leq 1$ a.e., $H = H^1(\Omega)$, and $K = \{v \in H_g^1(\Omega) \mid 0 \leq v \leq 1\}$, where g is such that $\text{ess sup } \gamma(0 - g) < 0 < \text{ess inf } \gamma(1 - g)$. Instead of iterating the fixed point

operator $Q_\rho : H \rightarrow K$ for some $\rho > 0$, we propose to generate a sequence of iterates $\{u^k\}$ from (5.14) with $G = B$ set to be the binary entropy considered in (5.8). More explicitly, we choose an appropriate $\rho > 0$ and a sequence $\{\alpha_k\}$ of positive real numbers, and then solve the sequence of resulting subproblems

$$(5.15) \quad \begin{cases} \text{Given } u^{k-1}, \text{ find } u^k \in H_g^1(\Omega) \cap L^\infty(\Omega) \text{ such that} \\ \frac{\alpha_k}{\rho} (\nabla u^k, \nabla v)_{L^2(\Omega)} + (\text{lnit}(u^k), v) = \frac{\alpha_k}{\rho} (\nabla u^{k-1}, \nabla v)_{L^2(\Omega)} + (\text{lnit}(u^{k-1}), v) \\ \quad - \alpha_k [a(u^{k-1}, v) - (f, v)] \quad \text{for all } v \in H_0^1(\Omega). \end{cases}$$

Once discretized with a slack variable ψ_h , we arrive at the algorithm written below.

Algorithm 5: A Proximal Galerkin method for advection-diffusion.

Input : Linear subspaces $V_h \subset H_0^1(\Omega)$ and $W_h \subset L^\infty(\Omega)$, initial solution guesses $u_h^0 \in V_h$, $\psi_h^0 \in W_h$, step sizes $\alpha_k > 0$, $\rho > 0$.

Output: Two approximate solutions, u_h and $\tilde{u}_h = \text{expit } \psi_h$.

Initialize $k = 0$.

repeat

Solve the following (nonlinear) discrete saddle-point problem:

$$(5.16) \quad \begin{cases} \text{Find } u_h \in g_h + V_h \text{ and } \psi_h \in W_h \text{ such that} \\ \frac{\alpha_k}{\rho} (\nabla u_h, \nabla v) + (\psi_h, v) = \alpha_k L(u_h^k, v) + (\psi_h^k, v) \quad \text{for all } v \in V_h, \\ (u_h, \varphi) - (\text{expit } \psi_h, \varphi) = 0 \quad \text{for all } \varphi \in W_h. \end{cases}$$

where

$$L(u, v) = (1/\rho - \epsilon)(\nabla u, \nabla v) - (\beta \cdot \nabla u - f, v).$$

Assign $\psi_h^{k+1} \leftarrow \psi_h$ and $k \leftarrow k + 1$.

until a convergence test is satisfied

5.3. Stable elements II: Continuous latent variable. Constructing a stable finite element discretization for Algorithm 5 has the same challenges we witnessed in solving the obstacle problem with Algorithm 3. Namely, we must construct a *stable* pair of finite element subspaces V_h and W_h . To this end, Subsection 4.7 introduced a class of possible pairings based on the requirement that the latent variable ψ_h be *discontinuous*. We could use the same finite elements here because the saddle-point problem (5.16) has the same structure after linearization as (4.44). Instead, however, we use this subsection to introduce the following alternative class of equal-order finite element pairings where ψ_h is *continuous*.

For any integer $p \geq 1$, we define the following two pairs of spaces:

Triangular elements. We refer to the following as the $(\mathbb{P}_p, \mathbb{P}_p)$ pairing:

$$(5.17a) \quad V_h = \mathbb{P}_p(\mathcal{T}_h) \cap H_0^1(\Omega), \quad W_h = V_h.$$

Quadrilateral elements. We refer to the following as the $(\mathbb{Q}_p, \mathbb{Q}_p)$ pairing:

$$(5.17b) \quad V_h = \mathbb{Q}_p(\mathcal{T}_h) \cap H_0^1(\Omega), \quad W_h = V_h.$$

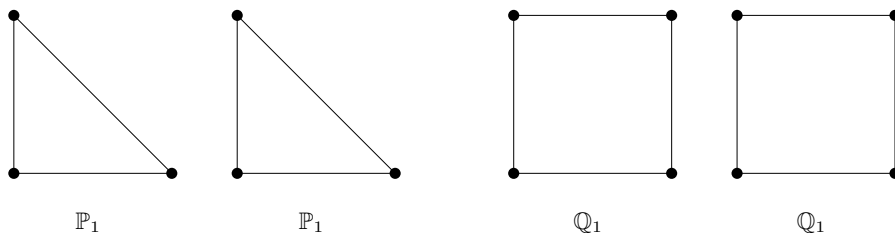


Fig. 5.2: Conventional representation of the $(\mathbb{P}_1, \mathbb{P}_1)$ and $(\mathbb{Q}_1, \mathbb{Q}_1)$ subspace pairs in two dimensions.

Figure 5.2 provides a visual representation of the lowest-order versions of these elements.

The equal-order finite element subspaces in (5.17) are appealing because they can be formed from standard $C^0(\Omega)$ -finite elements that are found in nearly all production codes. However, due to limitations of the stability proof outlined in Remark 5.2, we only advocate for using these elements when the mesh sequence $\{\mathcal{T}_h\}$ is *quasi-uniform* [68, Definition 22.20].

DEFINITION 5.1 (Quasi-uniformity). *A shape-regular sequence of meshes is called quasi-uniform when there exists a mesh-independent constant $c > 0$ such that $h_T \geq ch$ for all $T \in \mathcal{T}_h$ and all $h > 0$ in the index set.*

A careful inspection of Appendix B shows that we do not need to place such a restriction on the mesh when using the non-standard elements proposed in Subsection 4.7.

In addition to the quasi-uniformity assumption above, the reader may notice that the number of degrees of freedom per element pair in (5.17) is larger than that of the analogous order- p pairs proposed in (4.48). We pause to point out that neither of these factors preclude using (5.17) in practical applications. Indeed, many practical applications are solved on quasi-uniform meshes. Moreover, it turns out that the additional computational cost can be mitigated via a mass lumping technique described in Remark 5.3 below.

As also explained in this remark, a particularly interesting consequence of mass lumping is that it induces a *nodally*-feasible primal solution u_h . Moreover, because nodal boundedness extends to pointwise boundedness when $p = 1$, the primal solution u_h will be *pointwise*-feasible for any lowest-order $(\mathbb{P}_1, \mathbb{P}_1)$ or $(\mathbb{Q}_1, \mathbb{Q}_1)$ proximal Galerkin discretization with pointwise box constraints. In contrast, the lowest-order elements in (4.48) only induce a primal discretization with a feasible cell average; cf. Remark 4.20.

Remark 5.2 (Stability). As described in Appendix B, uniform stability of the proximal Galerkin discretization rests on satisfying the Ladyzhenskaya–Babuška–Brezzi (LBB) condition

$$(5.18) \quad \inf_{\varphi \in W_h} \sup_{v \in V_h} \frac{(\varphi, v)}{\|\varphi\|_{H^{-1}(\Omega)} \|\nabla v\|_{L^2(\Omega)}} \geq \beta_0 > 0,$$

with β_0 independent of the mesh size $h > 0$. Verifying this condition is often nontrivial. However, it reduces to a one-line argument given in (5.20) below if the global L^2 -

orthogonal projection $\mathcal{P}_h: H_0^1(\Omega) \rightarrow V_h$, defined

$$(5.19a) \quad (\mathcal{P}_h v, w) = (v, w) \quad \text{for all } w \in V_h,$$

is stable in the (equivalent) $H^1(\Omega)$ -norm; i.e., if there exists a constant $c > 0$, independent of h , such that

$$(5.19b) \quad \|\nabla(\mathcal{P}_h v)\|_{L^2(\Omega)} \leq c \|\nabla v\|_{L^2(\Omega)},$$

for all $v \in H_0^1(\Omega)$.

As shown in, e.g., [68, Proposition 22.21], quasi-uniformity of the mesh sequence implies (5.19) for $V_h = \mathbb{P}_p(\mathcal{T}_h) \cap H_0^1(\Omega)$ and $V_h = \mathbb{Q}_p(\mathcal{T}_h) \cap H_0^1(\Omega)$. Therefore, if we assume $\{\mathcal{T}_h\}$ is quasi-uniform and we are using any of the equal-order pairs in (5.17), then $\mathcal{P}_h(H_0^1(\Omega)) \subset V_h = W_h$ and (5.19) imply that there exists $\beta_0 = 1/c > 0$ such that

$$(5.20) \quad \sup_{v \in V_h} \frac{(\varphi, v)}{\|\nabla v\|_{L^2}} \geq \sup_{v \in H_0^1} \frac{(\varphi, \mathcal{P}_h v)}{\|\nabla(\mathcal{P}_h v)\|_{L^2}} \geq \beta_0 \sup_{v \in H_0^1} \frac{(\varphi, v)}{\|\nabla v\|_{L^2}} = \beta_0 \|\varphi\|_{H^{-1}},$$

for all $\varphi \in V_h$, as necessary.

Remark 5.3 (Nodal feasibility). We choose to focus this remark on the saddle-point problem in Algorithm 5. Yet, similar conclusions could be drawn about Algorithm 3 and, potentially, other future proximal Galerkin algorithms for second-order elliptic VIs.

Consider the equal-order finite element subspaces in (5.17) and let $\{\varphi_i\}_{i=1}^N$ be a basis for $V_h = W_h$. It follows that there exist coefficients \mathbf{c}_j , and \mathbf{d}_j , $j = 1, 2, \dots, N$, such that $u_h(x) = \sum_{j=1}^N \mathbf{c}_j \varphi_j(x)$ and $\psi_h(x) = \sum_{j=1}^N \mathbf{d}_j \varphi_j(x)$. Substituting these expressions into the second variational equation in (5.16), setting $\varphi = \varphi_i$, and replacing the Lebesgue integral $\int_{\Omega} \varphi(x) dx$ with a global quadrature rule $\sum_{l=1}^M w_l \varphi(x_l)$, where $w_l \neq 0$ and $x_l \in \bar{\Omega}$ for $l = 1, 2, \dots, M$, we find that

$$(5.21) \quad \sum_{l=1}^M \sum_{j=1}^N w_l \mathbf{c}_j \varphi_j(x_l) \varphi_i(x_l) = \sum_{l=1}^M w_l \text{expit} \left(\sum_{j=1}^N \mathbf{d}_j \varphi_j(x_l) \right) \varphi_i(x_l)$$

for each index $i = 1, \dots, N$.

We now employ the nodal quadrature technique [72] that is commonly used in, e.g., spectral element methods [62]. Namely, if we assume $M = N$ and that φ_j are formed from a nodal basis with nodes corresponding to the quadrature points x_j (e.g., a Lagrange basis with Gauss–Lobatto nodes [171]), then $\varphi_j(x_l) = \delta_{jl}$ for all $j, l = 1, \dots, N$. In turn, we find that (5.21) reduces to

$$(5.22) \quad \mathbf{c}_i = \text{expit}(\mathbf{d}_i).$$

Finally, we note that $\mathbf{c}_i = u_h(x_i)$ and $\mathbf{d}_i = \psi_h(x_i)$ since we have assumed the basis $\{\varphi_i\}_{i=1}^N$ is nodal. Thus, the primal variable u_h is nodally feasible; i.e., $0 \leq u_h(x_i) \leq 1$ at all nodes $i = 1, \dots, N$.

5.4. Numerical experiments. In this set of experiments, we follow [46, Section 4.1] and consider the exact solution of a model problem attributed to Eriksson and Johnson [65]. In particular, we set $\Omega = (0, 1)^2$, $f = 0$, and $\beta = (1, 0)^\top$ in (5.1) and, therefore, write

$$(5.23) \quad -\epsilon \left(\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right) + \frac{\partial u}{\partial x} = 0 \quad \text{in } \Omega, \quad u = g \quad \text{on } \partial\Omega.$$

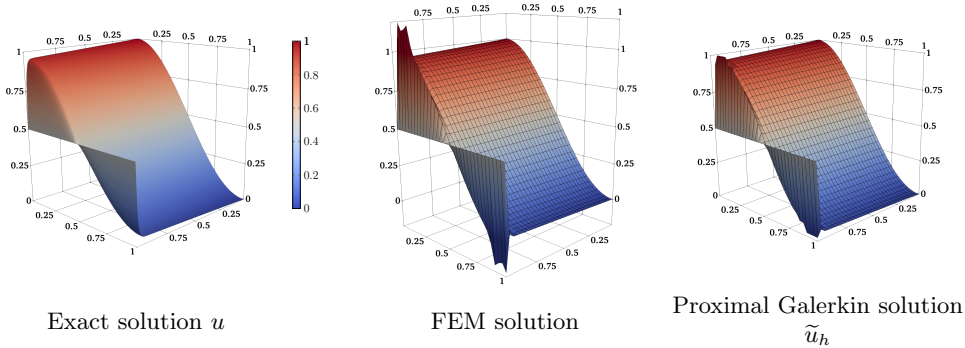


Fig. 5.3: The Eriksson–Johnson problem (5.23) for $\epsilon = 10^{-2}$. Left: The exact solution. Middle: A first-order Bubnov–Galerkin numerical solution that clearly violates the strong maximum principle $0 \leq u(x) \leq 1$. Right: The corresponding $(\mathbb{Q}_1, \mathbb{Q}_1)$ -proximal Galerkin solution $\tilde{u}_h = \text{expit}(\psi_h)$ satisfies the strong maximum principle, by construction. These results can be reproduced by running the MFEM code `advection_diffusion.cpp` available at [113].

Exact solutions of this problem for arbitrary boundary data g can be derived using the separation of variables technique. We choose to isolate the solutions satisfying $u = 0$ where $x = 1$ and $\nabla u \cdot n = 0$ where $y = 0$ and 1 . Doing so generates the following series expansion:

$$(5.24) \quad u = \sum_{n=1}^{\infty} C_n \frac{\exp(r_2(x-1)) - \exp(r_1(x-1))}{r_1 \exp(-r_2) - r_2 \exp(-r_1)} \cos(n\pi y),$$

where $r_{1,2} = \frac{1 \pm \sqrt{1+4\epsilon\lambda_n}}{2\epsilon}$ and $\lambda_n = n^2\pi^2\epsilon$. Note that the constants C_n can be determined from the prescribed values of g on $\{(x, y) \in \bar{\Omega} \mid x = 0\} \subset \partial\Omega$. Since we have not yet prescribed g on this part of the boundary, we define g there to be (5.24) with $C_1 = 1$ and $C_{n \neq 1} = 0$; i.e, we treat the first mode in this series as a manufactured solution for (5.23).

Figure 5.3 places the standard FEM solution of this problem for $\epsilon = 10^{-2}$ alongside the corresponding maximum principle-preserving proximal Galerkin solution \tilde{u}_h . Here, the proximal Galerkin solution can be found by running only two iterations of Algorithm 5 with constant values for ρ and α_k , and using the standard first-order FEM solution to provide an initial guess for u_h^0 and ψ_h^0 . Note that the two discrete solutions are similar, except the proximal Galerkin solution preserves the maximum principle $0 \leq u \leq 1$. We also experimented with using the mass lumping technique described in Subsection 5.3. The results of this experiment are given in Figure 5.4. Here, note that both the the primal solution u_h and the latent variable solution \tilde{u}_h are bound-preserving in this case. For further details, or to reproduce the our experiments, the interested reader is directed to our open-source FEniCSx and MFEM implementations found at [113].

Note that neither of the numerical approximations depicted in Figures 5.3 and 5.4 were obtained with numerical stabilization techniques that are common for this class of singularly-perturbed problems and usually required to consider smaller values of $\epsilon > 0$. Just as conventional stabilized finite element methods do not necessarily

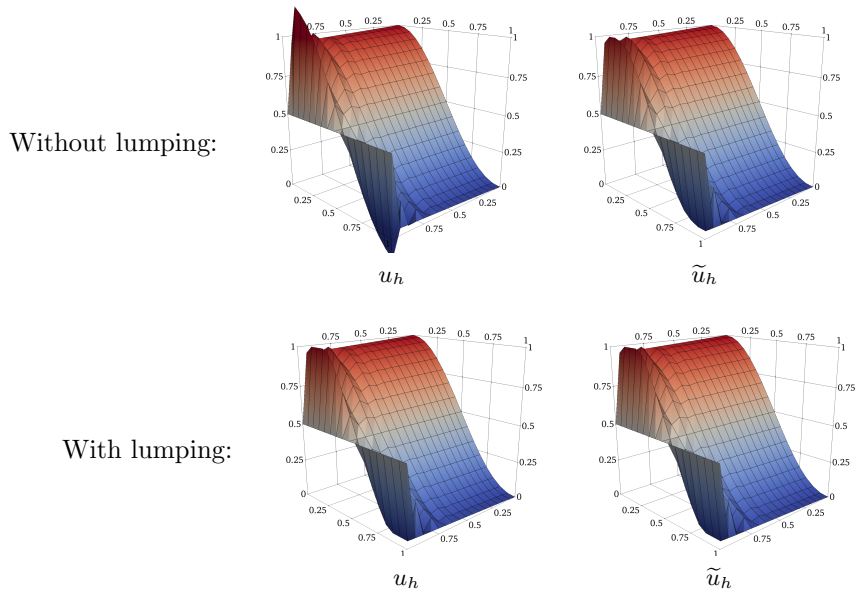


Fig. 5.4: Comparison of two $(\mathbb{Q}_1, \mathbb{Q}_1)$ -proximal Galerkin discretizations of the Eriksson–Johnson problem (5.23) also considered in Figure 5.3. In the first row, we see a pair of solutions corresponding to a proximal Galerkin discretization where standard Gaussian quadrature is used to evaluate every integral. In the second row, we see a similar pair of solutions obtained from a discretization employing the nodal-quadrature mass lumping technique described in Remark 5.3. As argued above, the latter discretization delivers *two* feasible discrete solutions. These results can be reproduced by running the FEniCSx code `advection_diffusion.py` available at [113].

preserve maximum principles [69], we find that entropy regularization does not necessarily induce robustness with respect to the diffusion parameter ϵ . Our conclusion is, therefore, that future work is required to develop robust proximal Galerkin finite element methods for singularly-perturbed PDEs.

6. Extensions II: Non-convex objective functions and a structure-preserving algorithm for topology optimization. The variational problems considered in the sections above share several common features. The most decisive feature is convexity. This raises the question as to whether entropy regularization can be as effective in a non-convex infinite-dimensional setting. We investigate this possibility here by providing a new proximal gradient (entropic mirror descent) framework for possibly non-convex, bounded-constrained optimization in infinite dimensions. Our benchmark problem for this setting is a well-known problem in topology optimization. As before, the section closes with an explicit algorithm and a brief account of numerical experiments. In the interest of completeness, we recall several details from abstract mirror descent methods. Although these methods are widely used in finite-dimensional convex optimization, and much of our treatment is inspired by the more recent works [25, 193], it is important to note that Nemirovski and Yudin did not restrict themselves to finite dimensions in their original works many decades ago

[160, 161].

6.1. Mirror descent. Subsection 4.4 introduced a proximal framework that was applied to solve the obstacle problem. Subsection 5.2 introduced a linearized proximal framework to solve variational inequalities with non-symmetric bilinear forms. The purpose of the present subsection is to combine those two approaches into a general first-order framework for non-convex optimization problems,

$$(6.1) \quad \min_{v \in V} F(v) \quad \text{subject to } v \in K \subset V,$$

where K is a nonempty, closed convex subset of a Banach space V and $F: V \rightarrow \mathbb{R}$ is continuously Fréchet differentiable. We closely follow [25, 193] below to provide intuition for the method. In several places, we are purposely vague. This is particularly the case for the assumption that a Bregman divergence D_G induced by the derivative G' is available or that $\text{int dom } G$ is non-empty with respect to the topology on V .

We begin by introducing the Bregman gradient step operator,

$$(6.2) \quad P_\alpha(w) = \arg \min_{v \in V} \{ \langle F'(w), v \rangle + \alpha^{-1} D_G(v, w) \}, \quad w \in \text{int dom } G,$$

where $G: \text{dom } G \rightarrow \mathbb{R} \cup \{\infty\}$ is strongly convex with derivative $G'(w) \in V'$. When V is a Hilbert space and $G(v) = \frac{1}{2} \|v\|_V^2 = \frac{1}{2} (v, v)_V$, the use of the gradient step operator leads to the standard gradient descent rule. This follows from a straightforward computation of the first-order optimality criteria, which leads to

$$(6.3) \quad P_\alpha(w) = w - \alpha \nabla F(w),$$

where $\nabla F: V \rightarrow V$ is the gradient of F characterized by the variational equation

$$(6.4) \quad (\nabla F(w), v)_V = \langle F'(w), v \rangle \quad \text{for all } v \in V.$$

More generally, assuming the minimizer exists and $G': V \rightarrow V'$ is invertible, (6.2) returns the formula

$$(6.5) \quad P_\alpha(w) = (G')^{-1}(G'(w) - \alpha F'(w)).$$

Recalling the classical steepest descent method, see, e.g., [163], it is not surprisingly that iterating (6.2) can generate a convergent algorithm to solve (6.1) when $K = V$ and an appropriate step size rule for α is available. Indeed, convergence of this algorithm is illustrated in Figure 6.1 for optimizing the scalar objective function $e(x) = \frac{1}{2}x^2 + x$ with the Bregman divergence $D_s(x, x_k)$ from the scalar entropy function $s(x) = x \ln x - x$. This naturally leads to the so-called mirror descent method [161, 25], which, given a sequence of positive step sizes $\{\alpha_k\}$, generates a sequence of iterates $\{u^k\}$ according to the following scheme:

$$u^0 \in \text{int dom } G, \quad u^{k+1} = P_{\alpha_{k+1}}(u^k), \quad k = 0, 1, 2 \dots$$

Nemirovski and Yudin point out that the motion of the iterates $\{u^k\}$, which takes place in the primal space V , is a “shadow” or “image”, of the main motion: $G'(u^k) - \alpha_{k+1} F'(u^k)$, which by definition takes place in the dual space; whence the name “method of mirror descent” [161, p. 88]. This is easily witnessed by introducing a

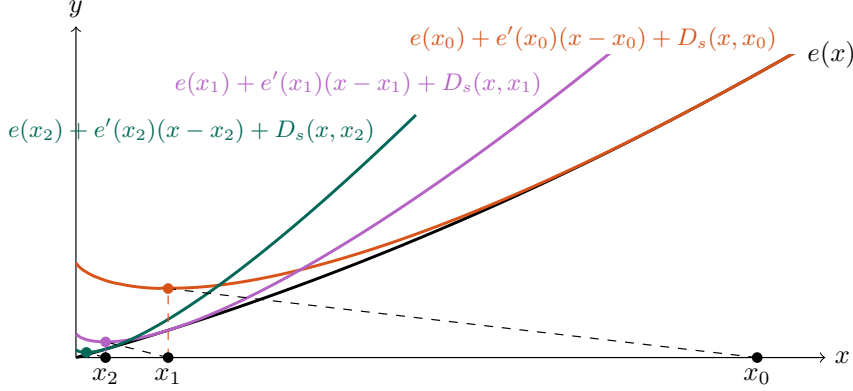


Fig. 6.1: Illustration of convergence to the solution $x^* = 0$ for the constrained minimization problem $\min_{x \in [0, \infty)} e(x)$, where $e(x) = \frac{1}{2}x^2 + x$, by solving the sequence of minimization problems $x_{k+1} = \arg \min_{x \in [0, \infty)} \{e'(x_k)x + D_s(x, x_k)\}$ starting at $x_0 = 1$.

dual variable $\lambda := G'(w)$. Then under the assumptions that G' is invertible, the new step in the dual space takes a somewhat more familiar form:

$$\lambda^{k+1} = \lambda^k - \alpha_{k+1}(F' \circ (G')^{-1})(\lambda^k).$$

This important distinction is often lost in finite dimensions and to some extent in the Hilbert space setting, where G' and F' are often identified with their Riesz representations in V ; i.e., the gradients ∇G and ∇F , respectively.

6.2. Mirror descent with a linear equality constraint. For constrained problems, it is essential that G properly captures the geometry of the feasible set, as was done in the previous sections on the obstacle problem and advection-diffusion equations. Many problems of interest have the following form:

$$(6.6) \quad \min_{v \in K_1 \cap K_2} F(v),$$

where K_1 and K_2 are nonempty, closed convex subsets of V and F is differentiable. For example, suppose that $K = K_1$ is a nonempty, closed convex set and $K_2 := \{v \in V \mid \ell(v) = c\}$ for some linear functional $\ell \in V'$ and constant $c \in \mathbb{R}$; i.e., $K_2 = \ell^{-1}(\{c\})$. Furthermore, suppose that D_G is a Bregman divergence associated with a distance generating function G , which is a Legendre function whose critical domain is linked to the properties of K . In this setting, rather than using (6.2), we fix $\alpha > 0$ and define the operator

$$T_\alpha(w) := \arg \min_{v \in K \cap \ell^{-1}(\{c\})} \{F(w) + \langle F'(w), v - w \rangle + \alpha^{-1} D_G(v, w)\}.$$

We assume here that $D_G(\cdot, w)$ over $K \cap \ell^{-1}(\{c\})$ has all the properties needed to ensure T_α is single-valued. Using standard optimality theory, e.g., [106], we can argue that $u := T_\alpha(w)$ satisfies the inclusion

$$(6.7) \quad 0 \in \alpha F'(w) + G'(u) - G'(w) + \mathcal{N}_{K \cap \ell^{-1}(\{c\})}(u),$$

where $\mathcal{N}_{K \cap \ell^{-1}(\{c\})}(u)$ is the normal cone from convex analysis [106], defined by

$$\mathcal{N}_{K \cap \ell^{-1}(\{c\})}(u) := \{\lambda \in V' \mid \langle \lambda, v - u \rangle \leq 0 \quad \forall v \in K \cap \ell^{-1}(\{c\})\}.$$

Note that if $w = T_\alpha(w)$, then (6.7) reduces to

$$0 \in \alpha F'(w) + \mathcal{N}_{K \cap \ell^{-1}(\{c\})}(w),$$

which indicates that w is a first-order stationary point of (6.6).

If we furthermore assume that K contains a subset \mathcal{B} such that $\ell(\mathcal{B}) \subset (c - \epsilon, c + \epsilon)$, for some $\epsilon > 0$, then $\{c\} - \ell(K)$ contains an open neighborhood of 0. This constraint qualification [106] allows us to rewrite (6.7) as

$$(6.8) \quad 0 \in \alpha F'(w) + G'(u) - G'(w) + \mathcal{N}_K(u) + \mathcal{N}_{\{c\}}(\ell(u)) \circ \ell,$$

where $\mathcal{N}_{\{c\}}(\ell(u)) \circ \ell = \{\mu \ell \in V' \mid \mu \in \mathbb{R}\}$ provided $\ell(u) = c$. Continuing on, we may assume for the sake of argument that the use of D_G forces $\mathcal{N}_K(u) = \{0\}$ and $u \in K$. This happens, for example, if u remains away from the boundary of K . For pointwise bound constraints in L^p -spaces of the type $0 \leq u \leq 1$ considered below, we would also have $\mathcal{N}_K(u) = \{0\}$ when $0 < u < 1$ almost everywhere, even if the set K does not have a non-empty interior. The remaining normal cone is trivial to compute and yields $\mathcal{N}_{\{c\}}(\ell(u)) = \mathbb{R}$.

These observations justify the following first-order optimality system that characterizes the map $w \mapsto u := T_\alpha(w)$: Find $(u, \mu) \in K \times \mathbb{R}$ such that

$$(6.9) \quad u = (G')^{-1}(G'(w) - \alpha F'(w) + \mu \ell) \text{ and } \ell(u) = c.$$

In other words, given w and α , compute the increment $\tilde{\lambda} := G'(w) - \alpha F'(w)$ and find $\mu \in \mathbb{R}$ by solving the equation

$$\ell((G')^{-1}(\tilde{\lambda} + \mu \ell)) = c.$$

Repeating the process

$$u^0 \in \text{int dom } G, \quad u^{k+1} = T_{\alpha_{k+1}}(u^k), \quad k = 0, 1, 2 \dots$$

generates a sequence of dual variables. Indeed, given a sequence of positive step sizes $\{\alpha_k\}$, we can generate $\{\lambda^k\}$ according to Algorithm 6.

Algorithm 6: Half-step mirror descent rule in Banach space

Input : Initial dual variable $\lambda^0 \in V'$ and sequence of step sizes $\alpha_k > 0$.

Output: Stationary dual variable $\bar{\lambda}$.

Initialize $k = 0$.

repeat

// Dual space half step (gradient descent)
Assign $\lambda^{k+1/2} \leftarrow \lambda^k - \alpha_{k+1}(F' \circ (G')^{-1})(\lambda^k)$.
// Compute Lagrange multiplier
Solve for $\mu^{k+1} \in \mathbb{R}$ such that $\ell((G')^{-1}(\lambda^{k+1/2} + \mu^{k+1} \ell)) = c$.
// Dual space feasibility correction
Assign $\lambda^{k+1} \leftarrow \lambda^{k+1/2} + \mu^{k+1} \ell$.
Assign $k \leftarrow k + 1$.

until a convergence test is satisfied

The pre-image of G' is tacitly assumed to be contained in K . Therefore, the abstract scheme [Algorithm 6](#) theoretically provides a sequence of feasible primal iterates

$$u^{k+1} := (G')^{-1}(\lambda^{k+1/2} + \mu^{k+1}\ell).$$

Checking for optimality is rather difficult in general, as the evaluation of the residual of first-order optimality conditions may require the computation of a projection operator in non-trivial settings; recall the discussion in [Remark 4.11](#) above. On the other hand, we demonstrated above that fixed points of T_α are stationary for the original problem. This motivates the simple stopping rule: $\|T_{\alpha_{k+1}}(u^k) - u^k\|_V < \text{tol}$, with $\text{tol} > 0$ sufficiently small, or some variant using absolute and relative tolerances. However, in order to remove the influence of α_k , which will typically change with k , we advocate for rescaling the fixed point residual and also consider the relative quantities

$$\eta_k := \|T_{\alpha_{k+1}}(u^k) - u^k\|_V / \alpha_{k+1} = \|u^{k+1} - u^k\|_V / \alpha_{k+1}.$$

In the unconstrained, Hilbert-space setting, we have $\eta_k = \|\nabla F(u^k)\|_V$. Therefore, if $\liminf_k \eta_k \rightarrow 0$, then $\liminf_k \|\nabla F(u_k)\|_V = 0$; i.e., we get a limiting stationarity condition. For example, if $\{u^{k+1} - u^k\}$ is a null sequence in $o(\alpha_k)$ for $\alpha_k \downarrow 0$, then clearly $\eta_k \downarrow 0$. We therefore also use η_k as a heuristic stopping measure in our experiments below.

The abstract derivation above yields an iterative scheme in the dual space V' . Implementing finite-dimensional approximations of negative-order Sobolev spaces can be challenging. However, the bound-constrained variational problem we have considered happens to have a substantial degree of useful structure, and entropy regularization of the associated bound constraints provides us with representations of G' , G'^{-1} and ℓ that lead to a latent space reformulation of [Algorithm 6](#) that is readily treated with finite elements.

6.3. An entropic mirror descent algorithm for topology optimization.

We consider the benchmark topology optimization problem of elastic compliance optimization of a cantilever beam; see, e.g., [\[16\]](#). In particular, we use the two-field filtered density approach to topology optimization [\[182, Section 3.1.2\]](#) to formulate the optimal cantilever beam problem.

The purpose of the problem is to find a material density $0 \leq \rho \leq 1$, where zero indicates no material, and one indicates the complete presence of material, that induces a minimal elastic compliance, $\widehat{F}(\mathbf{u}, \rho) = \int_\Omega \mathbf{u} \cdot \mathbf{f} \, d\mathbf{x}$. In this expression, the displacement $\mathbf{u} = \mathbf{u}(\rho)$ is determined by a variable material density ρ and a fixed body force \mathbf{f} through the classical linear elasticity equation [\[149\]](#), $-\text{Div}(r(\tilde{\rho}) \boldsymbol{\sigma}) = \mathbf{f}$. In this equation, we are meant to understand that

$$(6.10a) \quad \boldsymbol{\sigma} = \lambda \text{div}(\mathbf{u})I + \mu(\nabla \mathbf{u} + (\nabla \mathbf{u})^\top),$$

with Lamé parameters $\lambda, \mu > 0$, is the Cauchy stress of a homogeneous, isotropic material, $\text{Div}(\cdot)$ denotes the row-wise divergence operator, $\tilde{\rho}$ is a regularized (filtered) density function [\[40, 130\]](#), and $r(\tilde{\rho}) > 0$ is a local model for the Young's modulus. For our work, we use the well-known (modified) solid isotropic material penalization (SIMP) model $r(\tilde{\rho}) = \underline{\rho} + \tilde{\rho}^3(1 - \underline{\rho})$, where $0 < \underline{\rho} \ll 1$ is a nominal constant assigned to void regions in order to prevent the stiffness matrix from becoming singular [\[16\]](#).

The full problem formulation is written as follows:

$$(6.10b) \quad \min_{\rho \in L^1(\Omega)} \left\{ \widehat{F}(\mathbf{u}, \rho) = \int_\Omega \mathbf{u} \cdot \mathbf{f} \, d\mathbf{x} \right\},$$

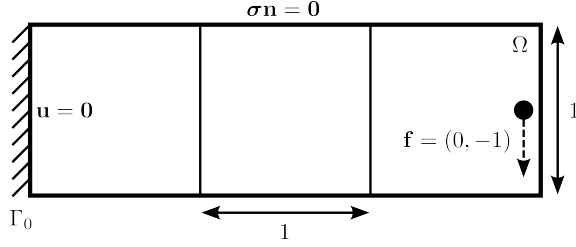


Fig. 6.2: The design domain Ω for the cantilever beam problem (6.10) with corresponding boundary conditions and three-element initial mesh with length $h_0 = 1$.

subject to the constraints

$$(6.10c) \quad \begin{cases} -\text{Div}(r(\tilde{\rho}) \boldsymbol{\sigma}) = \mathbf{f} & \text{in } \Omega \quad \text{with } \mathbf{u} = 0 \text{ on } \Gamma_0, \quad \boldsymbol{\sigma} \mathbf{n} = 0 \text{ on } \partial\Omega \setminus \Gamma_0, \\ -\epsilon^2 \Delta \tilde{\rho} + \tilde{\rho} = \rho & \text{in } \Omega \quad \text{with } \nabla \tilde{\rho} \cdot \mathbf{n} = 0 \text{ on } \partial\Omega, \\ \int_{\Omega} \rho(\mathbf{x}) \, d\mathbf{x} = \theta |\Omega|, \quad 0 \leq \rho \leq 1, & r(\tilde{\rho}) = \underline{\rho} + \tilde{\rho}^3(1 - \underline{\rho}), \end{cases}$$

where $\epsilon > 0$ is a *length scale* and $0 < \theta < 1$ is the desired *volume fraction*, which constrains the amount of the domain Ω occupied by the design. The design domain Ω and associated boundary conditions are depicted in Figure 6.2. We defer a rigorous mathematical discussion to the literature and simply note that it can be shown that \mathbf{u} can be understood, via $\tilde{\rho}$, as a differentiable mapping from ρ into an appropriate regularity space; e.g., a subspace of $[H^1(\Omega)]^2$. Therefore, we replace the objective function in (6.10b) by the reduced functional

$$(6.11a) \quad F(\rho) := \widehat{F}(\mathbf{u}(\rho), \rho)$$

and arrive at the reduced space optimization problem

$$(6.11b) \quad \min_{\rho \in L^1(\Omega)} F(\rho) \quad \text{subject to } 0 \leq \rho \leq 1 \text{ and } \int_{\Omega} \rho \, dx = \theta |\Omega|.$$

We can now solve this problem with a custom version of Algorithm 6 that employs the binary entropy-based Bregman divergence for the pointwise bound constraints found in (6.11b); cf. Subsection 5.1. In particular, the favorable structure of this problem lends itself nicely to a *latent space* representation, given below, that makes use of the transformations

$$\rho^k = \text{expit}(\psi^k) \quad \iff \quad \psi^k = \text{lnit}(\rho^k),$$

as well as the following variational characterization of the gradient $\nabla F(\rho^k)$:

$$(6.12) \quad \begin{cases} \text{Find } \nabla F(\rho^k) := \tilde{w} \in H^1(\Omega) \text{ such that} \\ \epsilon^2 (\nabla \tilde{w}, \nabla v) + (\tilde{w}, v) = -(r'(\tilde{\rho}^k) \boldsymbol{\sigma}(\mathbf{u}^k) : \nabla u^k, v) \quad \text{for all } v \in H^1(\Omega). \end{cases}$$

A visual representation of a single iteration of Algorithm 7 is given in Figure 6.3.

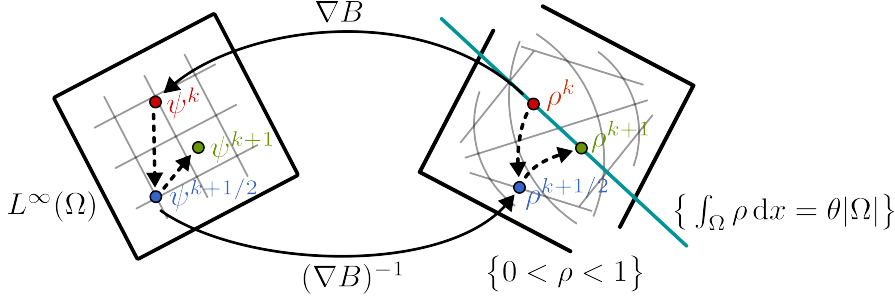


Fig. 6.3: Illustration of motion of the primal and latent iterates in Algorithm 7. When viewed in the primal space, we find that both steps of the progression $\rho^k \mapsto \rho^{k+1/2} \mapsto \rho^{k+1}$ involve *nonlinear* transformations of the primal variables. However, when viewed in the latent space, $L^\infty(\Omega)$, these transformations are simply just *translations* of the latent variables; namely, $\psi^{k+1/2} = \psi^k - \alpha_{k+1} \nabla F(\rho^k)$ and $\psi^{k+1} = \psi^{k+1/2} + c$.

Algorithm 7: Entropic mirror descent for topology optimization.

Input : Initial latent variable $\rho^0 \in L^\infty(\Omega)$, sequence of step sizes $\alpha_k > 0$, increment tolerance `itol.` > 0 , and normalized tolerance `ntol.` > 0 .

Output: Optimized material density $\bar{\rho} = \text{expit}(\psi^k)$.

Initialize $k = 0$.

```

while  $\|\text{expit}(\psi^k) - \text{expit}(\psi^{k-1})\|_{L^1(\Omega)} > \min\{\alpha_k \text{ntol.}, \text{itol.}\}$  do
  // Latent space gradient descent
  Assign  $\psi^{k+1/2} \leftarrow \psi^k - \alpha_{k+1} \nabla F(\text{expit}(\psi^k))$ .
  // Compute Lagrange multiplier
  Solve for  $c \in \mathbb{R}$  such that  $\int_{\Omega} \text{expit}(\psi^{k+1/2} + c) dx = \theta|\Omega|$ .
  // Latent space feasibility correction
  Assign  $\psi^{k+1} \leftarrow \psi^{k+1/2} + c$ .
  Assign  $k \leftarrow k + 1$ .

```

6.4. Numerical experiments. In this set of experiments, we estimate the gradients $\nabla F(\rho^k)$ in Algorithm 7 by discretizing (6.12) with $C^0(\Omega)$ -conforming, quadrilateral finite elements of degree $p \geq 1$. Likewise, the discrete displacements $\mathbf{u}_h^k \approx \mathbf{u}^k$ and filtered densities, $\tilde{\rho}_h^k \approx \tilde{\rho}^k$, are also computed with conforming finite elements of degree p . Finally, unlike the physical variables above, the latent variable ψ^k is approximated by discontinuous piecewise polynomials ψ_h^k of degree $p - 1$. Note that this induces a discontinuous primal variable $\rho_h^k := \text{expit}(\psi_h^k)$ satisfying $0 < \rho_h^k < 1$; see also Remark 6.1. We then apply the resulting discretized version of Algorithm 7 to solve (6.10) with $\rho = 10^{-6}$, $\lambda = \mu = 1$, $\theta = 0.5$, and $\epsilon = 0.02$. This experiment is an official part of MFEM 4.6; in particular, see MFEM Example 37 [112]. For sake of space, we have focused on presenting results with low-order discretizations (i.e., $p = 1, 2$) of the above form and left the exploration of higher-order discretizations to future work.

A sequence of iterates converging to a discrete solution with mesh size $h = h_0/128$ and polynomial degree $p = 1$ are depicted in Figure 6.4. From this figure, we observe typical first-order convergence behavior to a standard truss-like structure. To generate

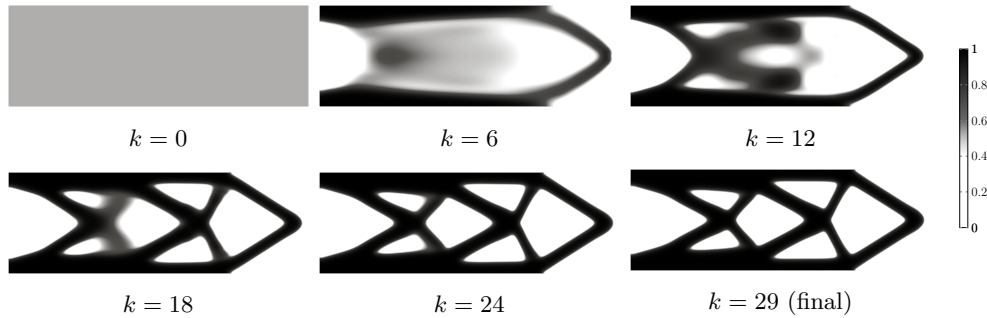


Fig. 6.4: Subsequence of material densities $\tilde{\rho}_h^k$ from [Algorithm 7](#) for selected iterations k . Results obtained with problem parameters $\epsilon = 2 \cdot 10^{-2}$ and $\theta = 0.5$; algorithm parameters $\text{itol.} = 10^{-2}$, $\text{ntol.} = 10^{-5}$, and $\alpha_k = 25k$; and discretization parameters $h = h_0/128$ and $p = 1$.

this figure, we used the heuristic step size sequence $\alpha_k = 25k$ and tolerances $\text{itol.} = 10^{-2}$ and $\text{ntol.} = 10^{-5}$. Although the conventional wisdom from finite-dimensional optimization theory would indicate that α_k should tend to zero, or at the very least be less than the reciprocal of the Lipschitz constant of ∇F , our experiments indicate that we can moderately increase the step sizes and still obtain satisfactory convergence behavior. To be fair, “satisfactory” convergence is based on the heuristic stopping rule given in [Algorithm 7](#).

Future work is needed to develop an adaptive step size selection procedure. In turn, we include [Figure 6.5](#) to show the different effects the step size sequence can have on the final solution. Here, we witness that different sequences — e.g., $\alpha_k = 10k$, $\alpha_k = 25k$, and $\alpha_k = 50k$ — can lead [Algorithm 7](#) to converge to significantly different local optima. This class of non-convex optimization problems is widely known to exhibit multiple local optima, though procedures are available to compute them [166]. In particular, notice from the two top left images that different final designs are possible just by changing the step size rule. The suspicious design on the bottom left is found because the $\alpha_k = 50k$ step size rule is too aggressive in the early iterations. Thereafter, a “design locking” phenomenon that is common in topology optimization problems keeps the design close to its nearly-binary, early state. To generate the results in [Figure 6.5](#), we changed the length scale to $\epsilon = 10^{-2}$ because it invokes a higher parameter sensitivity.

Finally, we return to the case considered in [Figure 6.4](#) (i.e., we again set $\epsilon = 2 \cdot 10^{-2}$ and $\alpha_k = 25k$) to record the sequence of error indicators $\eta_k = \|\rho_h^k - \rho_h^{k-1}\|_{L^1(\Omega)}/\alpha_k$ with different discretization parameters $h \in \{h_0/64, h_0/128, h_0/256\}$ and $p \in \{1, 2\}$. The results are given in [Table 6.1](#). From these results, we see that the number of iterations required to reach the tolerance $\|\rho_h^k - \rho_h^{k-1}\|_{L^1(\Omega)}/\alpha_k < 10^{-5}$ tends to a fixed value as the mesh is refined or the polynomial order is elevated. Moreover, the individual values of η_k appear to stabilize as $h \rightarrow 0$, for both $p = 1, 2$. Both of these properties suggest *mesh-independence* of [Algorithm 7](#).

Remark 6.1 (Preserving pointwise bound constraints at the discrete level). No matter the polynomial degree $p \geq 1$, the discrete primal variable $\rho_h^k = \text{expit}(\psi_h^k)$ is guaranteed to satisfy the pointwise bound constraint $0 \leq \rho_h^k \leq 1$. This is an immediate consequence of the sigmoid map $\text{expit}: \mathbb{R} \rightarrow (0, 1)$ and the decision to discretize the

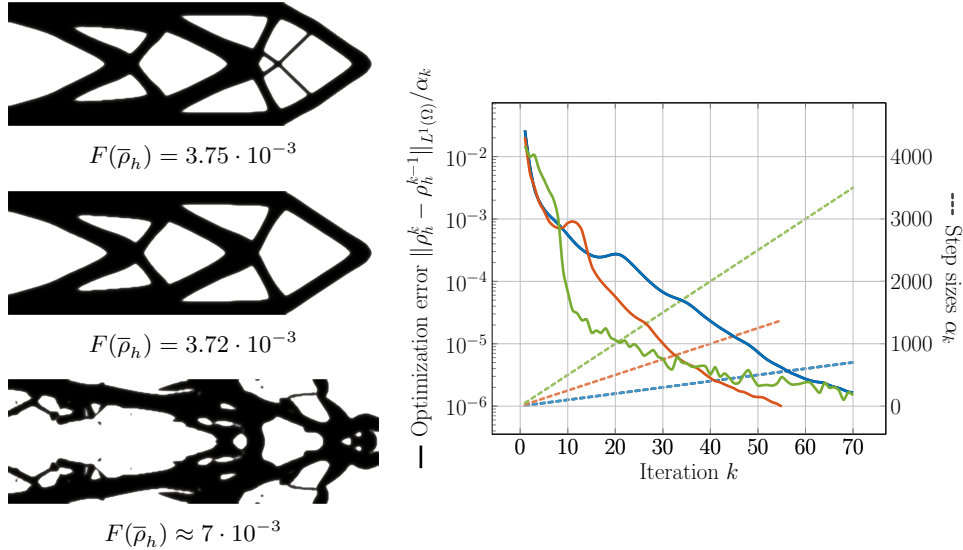


Fig. 6.5: An aggressive step size rule can lead to a better convergence rate. However, if the step size rule is too aggressive, the algorithm may converge to a sub-optimal local minimum or even diverge. Left: The final densities $\tilde{\rho}_h$ and associated compliance values $F(\rho_h)$ for step sizes $\alpha_k = 10k$ (top), $\alpha_k = 25k$ (middle), and $\alpha_k = 50k$ (bottom). Right: The normalized error estimates for the various step size sequences α_k , with $k = 1, 2, \dots$. These results were obtained using the problem parameters $\epsilon = 10^{-2}$ and $\theta = 0.5$ and the discretization parameters $h = h_0/128$ and $p = 1$.

| | | Optimization error $\ \rho_h^k - \rho_h^{k-1}\ _{L^1(\Omega)}/\alpha_k$ for various h and p | | | | |
|--------------------------------------|------------|---|----------------------|----------------------|--------------------------|----------------------|
| | | Polynomial order $p = 1$ | | | Polynomial order $p = 2$ | |
| k | α_k | $h_0/64$ | $h_0/128$ | $h_0/256$ | $h_0/64$ | $h_0/128$ |
| 1 | 25 | $2.00 \cdot 10^{-2}$ | $2.06 \cdot 10^{-2}$ | $2.05 \cdot 10^{-2}$ | $2.07 \cdot 10^{-2}$ | $2.05 \cdot 10^{-2}$ |
| 2 | 50 | $5.42 \cdot 10^{-3}$ | $5.80 \cdot 10^{-3}$ | $5.76 \cdot 10^{-3}$ | $5.88 \cdot 10^{-3}$ | $5.74 \cdot 10^{-3}$ |
| 3 | 75 | $2.97 \cdot 10^{-3}$ | $3.30 \cdot 10^{-3}$ | $3.27 \cdot 10^{-3}$ | $3.38 \cdot 10^{-3}$ | $3.25 \cdot 10^{-3}$ |
| 4 | 100 | $1.61 \cdot 10^{-3}$ | $1.87 \cdot 10^{-3}$ | $1.85 \cdot 10^{-3}$ | $1.94 \cdot 10^{-3}$ | $1.83 \cdot 10^{-3}$ |
| 5 | 125 | $1.12 \cdot 10^{-3}$ | $1.30 \cdot 10^{-3}$ | $1.29 \cdot 10^{-3}$ | $1.36 \cdot 10^{-3}$ | $1.28 \cdot 10^{-3}$ |
| \vdots | \vdots | \vdots | \vdots | \vdots | \vdots | \vdots |
| Total iterations | | 30 | 29 | 29 | 29 | 29 |
| Final compliance $F(\tilde{\rho}_h)$ | | $3.86 \cdot 10^{-3}$ | $4.04 \cdot 10^{-3}$ | $4.02 \cdot 10^{-3}$ | $4.08 \cdot 10^{-3}$ | $4.01 \cdot 10^{-3}$ |

Table 6.1: Table of error estimates $\eta_k = \|\rho_h^k - \rho_h^{k-1}\|_{L^1(\Omega)}/\alpha_k$ for various mesh sizes h and polynomial orders p . The initial density was set to the constant function $\rho_h^0 = \theta$ (i.e., $\psi_h^0 = \text{init } \theta$) at the beginning of each run and each run was stopped once $\eta_k < 10^{-5}$. These results were obtained using the problem parameters $\epsilon = 2 \cdot 10^{-2}$ and $\theta = 0.5$.

latent variable with finite elements. Had we followed the literature and, instead, directly discretized the primal variable ρ^k with finite elements, then the property $0 \leq \rho_h^k \leq 1$ would have to be enforced by introducing discrete-level pointwise bound constraints. This is a common concern in standard topology optimization approaches since the number of discrete-level pointwise bound constraints must grow with the size of the finite element space; cf. [Subsection 3.2](#).

7. Conclusion. We have introduced a new nonlinear finite element method that hinges on a mathematical technique called entropy regularization. We refer to the method as the latent variable proximal Galerkin finite element method or, simply, “proximal Galerkin” for short. The essential feature of the proximal Galerkin method is to provide robust, high-order, and *pointwise* bound-preserving discretizations accompanied by a built-in, low iteration complexity, mesh-independent solution algorithm. We have derived, analyzed, and implemented proximal Galerkin for the obstacle problem and used the advection-diffusion equation and topology optimization to motivate our wider vision for the method. Each of our numerical experiments is accompanied by an open-source implementation to facilitate reproduction of our results and broader adoption of the proximal Galerkin method.

The upshot of this work is that computational techniques for variational inequalities, maximum principles, and bound constraints in optimal design can be unified with a rigorous mathematical framework rooted in nonlinear programming and nonlinear functional analysis. We hope that the proximal Galerkin methods that arise from this framework will lead to new challenges and opportunities in optimization theory, analysis of PDEs, and numerical analysis, as well as provide promising alternatives to the more classical procedures used in industry for industrial-scale problem solving.

Appendix A. Mathematical results I: Isomorphisms, regularity, characterizations, and convergence. This appendix contains proofs and continuous-level structural results supporting the main sections of the paper.

A.1. Structural results on the set $\text{int } L_+^\infty(\Omega)$. The following concepts and results are not commonly used in the finite element literature. Although they can be derived from diverse sources, such as [\[51, 35, 75\]](#), we assemble them here for the reader’s convenience.

DEFINITION A.1 (Group of units). *Let \mathcal{X} be a semiring equipped with two binary operations: addition $\oplus: \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{X}$ and multiplication $\odot: \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{X}$. An element u of \mathcal{X} is called a unit if there exists an inverse element in \mathcal{X} , denoted $\frac{1}{u}$, such that $u \odot \frac{1}{u} = \frac{1}{u} \odot u = 1$. The group of units of \mathcal{X} , denoted \mathcal{X}^\times , is the set of all units in \mathcal{X} .*

This work is largely centered around the group of units $(L_+^\infty(\Omega))^\times$. We prove $(L_+^\infty(\Omega))^\times = \text{int } L_+^\infty(\Omega)$, along with several other algebraic/topological identities, at the end of this subsection; see [Proposition A.7](#).

It is well-known that algebraic and topological structures are often entwined, as the following definition and result shows.

DEFINITION A.2 (Banach algebra). *A Banach algebra is a complete normed vector space that is closed under multiplication.*

PROPOSITION A.3 (Topology of the group of units). *For any Banach algebra \mathcal{X} , its group of units \mathcal{X}^\times is open. Moreover, the inversion map $\mathcal{X}^\times \rightarrow \mathcal{X}^\times: u \mapsto \frac{1}{u}$ is continuous.*

Proof. See [\[51, Theorem 2.2, p. 192\]](#). □

Notably, this result also implies that $\text{int } L_+^\infty(\Omega)$ is a Banach–Lie group.

DEFINITION A.4 (Banach–Lie group). *A Banach manifold is a topological space \mathcal{M} where each point $u \in \mathcal{M}$ has an open neighborhood that is homeomorphic to an open set in a Banach space. A set \mathcal{G} is a Banach–Lie group if it is a Banach manifold that is closed under continuous multiplication and inversion operations.*

An important property of Lie groups is the existence of a smooth exponential map, $\exp: \mathcal{X} \rightarrow \mathcal{G}$, where \mathcal{X} is the associated Lie algebra; cf. [132].

DEFINITION A.5 (Banach–Lie algebra). *A Lie algebra \mathcal{X} is a vector space endowed with an antisymmetric bilinear form called the Lie bracket $[\cdot, \cdot]: \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{X}$ satisfying the Jacobi identity $[\psi, [\varphi, \omega]] + [\varphi, [\omega, \psi]] + [\omega, [\psi, \varphi]] = 0$ for all $\psi, \varphi, \omega \in \mathcal{X}$. A set \mathcal{X} is Banach–Lie algebra if it is both a Lie algebra and a Banach space.*

Using well-known results on Nemytskii operators between Lebesgue spaces, we can argue that $L^\infty(\Omega)$ is the Banach–Lie algebra associated to the Banach–Lie group $\text{int } L_+^\infty(\Omega)$; cf. Proposition A.7. In particular, as a result of Lemma A.6, the Nemytskii operator generated by the standard exponential function on \mathbb{R} provides the exponential map from $L^\infty(\Omega)$ to $\text{int } L_+^\infty(\Omega)$. Moreover, this map is surjective and, thus, the inverse of the entropy gradient $(\nabla S)^{-1} = \exp: L^\infty(\Omega) \rightarrow \text{int } L_+^\infty(\Omega)$ is a group isomorphism. Finally, in our setting, $\text{int } L_+^\infty(\Omega)$ is commutative and so the Lie bracket is trivial; i.e., $[\psi, \varphi] = \psi\varphi - \varphi\psi = 0$ for all $\psi, \varphi \in L^\infty(\Omega)$.

LEMMA A.6. *The Nemytskii operator $\psi \mapsto \exp \psi$ is infinitely continuously Fréchet differentiable on $L^\infty(\Omega)$.*

Proof. We first observe several crucial properties of the exponential function that carry over to the Nemytskii operator. Let $c > 0$ and fix $x \in [-c, c]$. By monotonicity we have $|\exp(x)| \leq \exp(c)$. Hence, for all $c > 0$, there exists $k(c)$ such that $|\exp(x)| \leq k(c)$ for all $x \in [-c, c]$. By [80, Theorem 1 (iv)], the Nemytskii operator maps L^∞ into itself. Similarly, for $c, \varepsilon > 0$ we define $\delta(c, \varepsilon) = \varepsilon / \exp(c)$ and observe that for all $x, y \in \mathbb{R}$ such that $|x| \leq c$, $|y| \leq c$, and $|x - y| < \delta(c, \varepsilon)$ we have $|\exp(x) - \exp(y)| < \varepsilon$. Then by [80, Theorem 5], the Nemytskii operator is continuous from L^∞ into itself. Finally, since $\exp(x)$ is infinitely continuously differentiable with $\exp'(x) = \exp(x)$ for all $x \in \mathbb{R}$ these results carry over to the Nemytskii operators defined by the pointwise derivatives. We may then apply [80, Theorem 7] to argue that the Nemytskii operator is infinitely continuously Fréchet differentiable from L^∞ into itself. \square

Proposition A.7 summarizes various useful interpretations of the set $\text{int } L_+^\infty(\Omega)$. See also Remark A.8.

PROPOSITION A.7. *Nemytskii operator $\psi \mapsto \exp \psi$ is a C^1 -diffeomorphism between $L^\infty(\Omega)$ and $\text{int } L_+^\infty(\Omega)$ for which the following definitions are equivalent:*

(a) *$\text{int } L_+^\infty(\Omega)$ is the set of all positive functions in $L^\infty(\Omega)$ whose reciprocals lie in $L^\infty(\Omega)$,*

$$(A.1a) \quad \text{int } L_+^\infty(\Omega) = \{w \in L^\infty(\Omega) \mid 1/w \in L^\infty(\Omega) \text{ and } w > 0\}.$$

In other words, $\text{int } L_+^\infty(\Omega) = (L_+^\infty(\Omega))^\times$ is the group of units in $L_+^\infty(\Omega)$.

(b) *$\text{int } L_+^\infty(\Omega)$ is the set of all functions in $L^\infty(\Omega)$ whose logarithm is bounded in $L^\infty(\Omega)$,*

$$(A.1b) \quad \text{int } L_+^\infty(\Omega) = \ln^{-1}(L^\infty(\Omega)).$$

(c) *$\text{int } L_+^\infty(\Omega)$ is the image of $L^\infty(\Omega)$ under the exponential map,*

$$(A.1c) \quad \text{int } L_+^\infty(\Omega) = \exp(L^\infty(\Omega)).$$

(d) $\text{int } L_+^\infty(\Omega)$ is the set of all positive functions in $L^\infty(\Omega)$ that are strictly bounded away from zero,

$$(A.1d) \quad \text{int } L_+^\infty(\Omega) = \{w \in L^\infty(\Omega) \mid \text{there exists } \epsilon > 0 \text{ such that } w > \epsilon\}.$$

Proof. We begin by proving the equivalence of (A.1a)–(A.1d). We leave off the dependence on Ω in the following arguments for readability and begin with (A.1d). If $u \in L^\infty$ such that there exists $\epsilon > 0$ with $u > \epsilon$ a.e. on Ω , then for any $w \in L^\infty$ such that $\|u - w\|_{L^\infty} < \epsilon/2$ we have $\epsilon/2 < u - \epsilon/2 < w$. Thus, $w \in L_+^\infty$ and, consequently, $u \in \text{int } L_+^\infty$. Now suppose $u \in L_+^\infty$ and for every $\epsilon > 0$, the set $\mathcal{B}_\epsilon := \{x \in \Omega \mid 0 < u(x) < \epsilon/2\}$ has positive Lebesgue measure. Then for all $\epsilon > 0$, the open ball $\{v \in L^\infty : \|u - v\|_{L^\infty} < \epsilon\}$ contains the function $v = u - \epsilon/2$ on \mathcal{B}_ϵ and $v = u$ on $\Omega \setminus \mathcal{B}_\epsilon$, which is clearly not in L_+^∞ . Hence, every open ball of radius $\epsilon > 0$ around u contains a point outside L_+^∞ , i.e., $u \in \text{bd } L_+^\infty = L_+^\infty \setminus \text{int } L_+^\infty$. This proves (A.1d).

Next, let $u \in \text{int } L_+^\infty$. By (A.1d) there exists $\epsilon > 0$ such that $\infty > \|u\|_{L^\infty} \geq u > \epsilon$. Then, by continuity on \mathbb{R}_{++} , the mapping $u \mapsto 1/u$ preserves measurability and $1/u \in [1/\|u\|_{L^\infty}, 1/\epsilon]$. Hence, $u > 0$ and $1/u \in L^\infty$. Conversely, suppose $u \in L^\infty$ such that $u > 0$ and $1/u \in L^\infty$. Then $0 < 1/u \leq \|1/u\|_{L^\infty}$ and $0 < u \leq \|u\|_{L^\infty}$ imply $u \in [[1/\|u\|_{L^\infty}^{-1}, \|u\|_{L^\infty}]$ a.e. on Ω . It follows from (A.1d) that $u \in \text{int } L_+^\infty$. This proves (A.1a).

Next, suppose that $u \in \text{int } L_+^\infty$. Then by (A.1a), we have $u \in [[1/\|u\|_{L^\infty}^{-1}, \|u\|_{L^\infty}]$ a.e. Consequently, the continuity and monotonicity of the natural logarithm on \mathbb{R}_{++} yields $\ln u \in L^\infty$. In other words, $u \in \ln^{-1}(L^\infty)$. Conversely, suppose we have $u \in \ln^{-1}(L^\infty)$. By definition, $u \in L^\infty$. Thus, we deduce the bounds $\underline{m}, \bar{m} \in \mathbb{R}$ such that $\ln u \in [\underline{m}, \bar{m}]$ a.e. Using the fact that the exponential map is positive and monotone, we infer that $u \in [\exp(\underline{m}), \exp(\bar{m})]$ a.e., $u > 0$, and $1/u \in [\exp(-\bar{m}), \exp(-\underline{m})]$ a.e. Since the reciprocal function is continuous away from zero, $u \in \text{int } L_+^\infty$. This proves (A.1b).

The proof of (A.1c) is similar to that of (A.1b). Let $v \in \exp(L^\infty)$ and $\varphi \in L^\infty$ such that $v = \exp \varphi$. Clearly, we have $v > 0$ a.e. Since φ is essentially bounded, there are independent constants $m, M \in \mathbb{R}$ such that $m \leq \varphi \leq M$ a.e. It follows then that $1/v = 1/\exp \varphi \in [\exp(-M), \exp(-m)]$. Consequently, $1/v$ is bounded. Since v is strictly positive and measurable and $(\cdot)^{-1} : \mathbb{R}_{++} \rightarrow \mathbb{R}$ is continuous, $1/v$ is measurable. Therefore, $v \in \text{int } L_+^\infty$. Conversely, let $v \in \text{int } L_+^\infty$. Then by definition, $v > 0$ a.e. and $1/v \in L^\infty$. This implies $v \in [[\frac{1}{v}\|_{L^\infty}^{-1}, \|v\|_{L^\infty}]$ a.e. It follows that $\varphi := \ln v \in L^\infty$ and $-\ln(\|\frac{1}{v}\|_{L^\infty}) \leq \varphi \leq \ln(\|v\|_{L^\infty})$. As a result, $v = \exp \varphi \in \exp(L^\infty)$, as was to be shown. This completes the proof of (A.1c).

Finally, we prove that $\psi \mapsto \exp \psi$ is a diffeomorphism. This requires us to check that for $\psi, \varphi \in L^\infty$, $\exp(\psi + \varphi) = \exp \psi \exp \varphi$, which holds by well-known properties of the exponential map. Furthermore, for $u, v \in \text{int } L_+^\infty = \exp(L^\infty)$, $\exp^{-1}(uv) = \ln(uv) = \ln u + \ln v = \exp^{-1}(u) + \exp^{-1}(v)$, by well-known properties of logarithms. We know from Lemma A.6 that $\psi \mapsto \exp \psi$ is infinitely differentiable on L^∞ . For the inverse mapping, note that for any $u \in \text{int } L_+^\infty$ and sequence $\{h_k\} \subset L^\infty$ such that $\|h_k\|_{L^\infty} \rightarrow 0$ we have (pointwise a.e.):

$$\begin{aligned} \left| \ln(u + h_k) - \ln(u) - \frac{h_k}{u} \right| &= \left| \ln\left(1 + \frac{h_k}{u}\right) - \ln(1) - \frac{h_k}{u} \right| \\ &\leq \left| \frac{h_k}{u} \right| \left| \left[\int_0^1 (1 + \tau h_k/u)^{-1} d\tau - 1 \right] \right| \end{aligned}$$

$$\begin{aligned}
&\leq \|h_k\|_{L^\infty} \left\| \frac{1}{u} \right\|_{L^\infty} \int_0^1 |(1 + \tau h_k/u)^{-1} - 1| \, d\tau \\
&= \|h_k\|_{L^\infty} \left\| \frac{1}{u} \right\|_{L^\infty} \int_0^1 |(1 + u/(\tau h_k))^{-1}| \, d\tau \\
&\leq \|h_k\|_{L^\infty} \left\| \frac{1}{u} \right\|_{L^\infty} (\|u\|_{L^\infty} \|h_k\|_{L^\infty}^{-1} - 1)^{-1}.
\end{aligned}$$

It follows that

$$\left\| \ln(u + h_k) - \ln(u) - \frac{h_k}{u} \right\| = o(\|h_k\|)$$

and, consequently, that the superposition operator $\ln u$ is Fréchet differentiable on $\text{int } L_+^\infty$ with respect to variations in L^∞ . To see that $1/u$ is continuous on $\text{int } L_+^\infty$, let $u \in \text{int } L_+^\infty$ and $\{u_k\} \subset L^\infty$ such that $u_k \rightarrow u$. Since $u \in \text{int } L_+^\infty$ there exists $\varepsilon > 0$ such that $u > \varepsilon$ a.e. Then, for sufficiently large k , we can argue that $u_k \geq \varepsilon/2$ pointwise a.e. This provides a uniform bound on $\|1/u_k\|_{L^\infty}$. Hence,

$$(A.2) \quad \|1/u_k - 1/u\|_{L^\infty} \leq \|u^{-1}\|_{L^\infty} \|u_k^{-1}\|_{L^\infty} \|u_k - u\|_{L^\infty} \rightarrow 0$$

as $k \rightarrow +\infty$. Thus, $\psi \mapsto \exp \psi$ is a C^1 -diffeomorphism, as necessary. \square

Remark A.8 (Analytic isomorphism). Upon closer inspection, we see that the differentiability of the nonlinearity $\ln: \text{int } L_+^\infty \rightarrow L^\infty$ can be shown to be much higher than C^1 . For example, using the same line of argument as the proof above, we see that

$$\begin{aligned}
\left| (u + h_k)^{-1} - u^{-1} + \frac{h_k}{u^2} \right| &= \left| \frac{-h_k}{u(u + h_k)} + \frac{h_k}{u^2} \right| \\
&\leq \|h_k\|_{L^\infty} \|1/u\|_{L^\infty} \left| \frac{1}{u} - \frac{1}{(u + h_k)} \right|,
\end{aligned}$$

which behaves like $o(\|h_k\|)$, in light of the property shown in (A.2). In fact, if we had defined the original exponential function using its power series, then deeper arguments can be used to illustrate that \exp and \ln are even analytic; cf. [75].

It is well-known that $W^{1,p}(\Omega) \cap L^\infty(\Omega)$ is a Banach algebra for every $1 \leq p \leq \infty$; see, e.g., [38, Proposition 9.4]. The following proposition connects this set to the Banach–Lie group $W^{1,p}(\Omega) \cap \text{int } L_+^\infty(\Omega)$.

PROPOSITION A.9. *Let Ω be an open subset of \mathbb{R}^n and $1 \leq p \leq \infty$. Then*

$$\ln: W^{1,p}(\Omega) \cap \text{int } L_+^\infty(\Omega) \rightarrow W^{1,p}(\Omega) \cap L^\infty(\Omega)$$

and

$$\exp: W^{1,p}(\Omega) \cap L^\infty(\Omega) \rightarrow W^{1,p}(\Omega) \cap \text{int } L_+^\infty(\Omega)$$

are isomorphisms. Moreover,

$$(A.3) \quad \nabla \ln u = \frac{1}{u} \nabla u \quad \text{and} \quad \nabla \exp \psi = \exp \psi \nabla \psi,$$

for all $u \in W^{1,p}(\Omega) \cap \text{int } L_+^\infty(\Omega)$ and $\psi \in W^{1,p}(\Omega) \cap L^\infty(\Omega)$.

Proof. We prove $\ln: W^{1,p}(\Omega) \cap \text{int } L_+^\infty(\Omega) \rightarrow W^{1,p}(\Omega) \cap L^\infty(\Omega)$ and $\nabla \ln u = 1/u \nabla u$ for the case that Ω is bounded. The corresponding statements for the exponential map are treated similarly.

Step 0. Let $u \in W^{1,p}(\Omega) \cap \text{int } L_+^\infty(\Omega)$. By [Proposition A.7](#) we know that $\ln u \in L^\infty(\Omega)$ and, moreover, there exists $\epsilon > 0$ such that $\epsilon \leq u(x) \leq 1/\epsilon$ at a.e. $x \in \Omega$. We now follow the proof technique used for [\[38, Proposition 9.4\]](#) to show that $\ln u \in W^{1,p}(\Omega)$.

Step 1. The first step involves constructing a sequence $u_k \in C_c^\infty(\Omega)$ such that

$$(A.4a) \quad u_k \rightarrow u \quad \text{in } L^p(\Omega) \text{ and pointwise a.e. in } \Omega,$$

$$(A.4b) \quad \nabla u_k \rightarrow \nabla u \quad \text{in } [L^p(\omega)]^n \text{ for all } \omega \subset\subset \Omega.$$

Furthermore,

$$(A.4c) \quad \|u_k\|_{L^\infty(\Omega)} \leq \|u\|_{L^\infty(\Omega)}$$

and, for all $\omega \subset\subset \Omega$, it holds that

$$(A.4d) \quad \|1/u_k\|_{L^\infty(\omega)} \leq \|1/u\|_{L^\infty(\Omega)},$$

once k is sufficiently large. For simplicity, we choose to focus on the case where Ω is bounded. This step may be modified by multiplying u_k with a sequence of smooth cut-off functions to treat the case where Ω is unbounded; cf. [\[38, Proof of Theorem 9.2\]](#).

Begin by defining

$$(A.5) \quad \bar{u}(x) = \begin{cases} u(x) & \text{if } x \in \Omega, \\ 0 & \text{if } x \in \Omega \setminus \mathbb{R}^n, \end{cases}$$

and set $u_k = \rho_k * \bar{u}$, where $\rho_k \in C_c^\infty(\mathbb{R}^n)$ is a sequence of mollifier functions satisfying

$$(A.6) \quad \text{supp } \rho_k \subset \overline{B(0, 1/k)}, \quad \int_{\mathbb{R}^n} \rho_k = 1, \quad \rho_k \geq 0 \text{ a.e. in } \mathbb{R}^n.$$

Notice that $u_k(x) \leq 1/\epsilon$ at a.e. $x \in \Omega$ since

$$(A.7) \quad u_k(x) = \int_{\mathbb{R}^n} \bar{u}(x-y)\rho_k(y) \, dy \leq \|u\|_{L^\infty(\Omega)} \int_{\mathbb{R}^n} \rho_k(y) \, dy \leq \frac{1}{\epsilon}.$$

This proves [\(A.4c\)](#).

Now, take $\omega \subset\subset \Omega$ and let $\delta > 0$ be chosen small enough so that the open cover $\bigcup_{x \in \bar{\omega}} B(x, \delta)$ is contained in Ω . Then, for all $k > 1/\delta$ and a.e. $x \in \omega$, we have that

$$(A.8) \quad \epsilon = \int_{B(x,\delta)} \epsilon \rho_k(y) \, dy \leq \int_{B(x,\delta)} \bar{u}(x-y)\rho_k(y) \, dy = u_k(x).$$

We have thus shown [\(A.4d\)](#). Properties [\(A.4a\)](#) and [\(A.4b\)](#) are proven for this sequence in [\[38, Theorem 9.2\]](#).

Step 2. The next step is to consider a test function $\varphi \in C_c^1(\Omega)$. Observe that

$$(A.9) \quad \int_{\Omega} \ln(u_k) \nabla \varphi \, dx = - \int_{\Omega} (1/u_k \nabla u_k) \varphi \, dx.$$

Let $\omega = \text{supp } \varphi \subset\subset \Omega$ denote the support of φ . Clearly, $\ln u_k(x) \rightarrow \ln u(x)$ at a.e. point $x \in \omega$. Moreover, it is a straightforward exercise to show that $|\ln u_k(x)| \leq$

$\max\{\ln \|u\|_{L^\infty(\Omega)}, \ln \|1/u\|_{L^\infty(\Omega)}\}$ at a.e. $x \in \omega$. Therefore, by the dominated convergence theorem, we have that

$$(A.10) \quad \lim_{k \rightarrow \infty} \int_{\Omega} \ln(u_k) \nabla \varphi \, dx = \int_{\Omega} \ln(u) \nabla \varphi \, dx.$$

To treat the right-hand side of (A.9), we apply a converse of the dominated convergence theorem to the sequence $\nabla u_k \rightarrow \nabla u$ given in (A.4b). In particular, by [38, Theorem 4.9], we know that there exists a subsequence $\{\nabla u_{k_l}\}_{l=1}^\infty$ and a function $h \in L^p(\omega)$ such that

$$(A.11) \quad |\nabla u_{k_l}(x)| \leq h(x) \quad \text{for all } l \geq 0 \text{ and a.e. } x \in \omega.$$

Next, we use (A.4d) to conclude that $|1/u_k \nabla u_{k_l}(x)| \leq \|1/u\|_{L^\infty(\Omega)} h(x) \in L^p(\omega)$. In turn, dominated convergence theorem implies that

$$(A.12) \quad \lim_{l \rightarrow \infty} \int_{\Omega} (1/u_{k_l} \nabla u_{k_l}) \varphi \, dx = \int_{\Omega} (1/u \nabla u) \varphi \, dx \quad \square$$

because $1/u_k \nabla u_{k_l}(x) \rightarrow 1/u \nabla u(x)$ as $l \rightarrow \infty$ at a.e. point $x \in \omega$. The identity $\nabla \ln u = 1/u \nabla u$ immediately follows from (A.9), (A.10), and (A.12).

A.2. Regularity of the entropy functional. One of the important facts that arise from Proposition A.7 is that $u \in \text{int } L_+^\infty(\Omega)$ implies $u \geq \|\frac{1}{u}\|_{L^\infty}^{-1}$. Indeed, this property allows us to differentiate the negative entropy function

$$(A.13) \quad S(u) = \begin{cases} \int_{\Omega} u \ln u - u \, dx & u \in L_+^1(\Omega) \\ +\infty & \text{otherwise,} \end{cases}$$

on the open set $\text{int } L_+^\infty(\Omega)$ with variations in $L^\infty(\Omega)$. We proceed now with a proof of Theorem 4.1.

Proof of Theorem 4.1. Case 1: $1 \leq p \leq \infty$. For $p = 1$, the properties of strict convexity and lower semicontinuity can be found in the seminal works [35, 22]. Since $\Omega \subset \mathbb{R}^n$ is bounded, the continuous embedding of $L^p(\Omega)$ into $L^1(\Omega)$ imply these same properties for all $p \in (1, \infty]$.

Case 2: $1 < p \leq \infty$. Our proof continues by considering the Nemytskii operator induced by the real-valued function

$$\widehat{s}(x) := x \ln |x| - x.$$

We will show that it is a continuous map from $L^p(\Omega)$ to $L^1(\Omega)$ when $p > 1$. In doing so, we first note that \widehat{s} is continuous when viewed as a real-valued function $x \mapsto x \ln |x| - x$ with $x \in \mathbb{R}$ and, moreover, for any $p > 1$, there exists a constant $C(p)$ such that

$$(A.14) \quad |\widehat{s}(x)| \leq C(p) + |x|^p.$$

Note that $C(p)$ exists on the one hand since $|\widehat{s}(x)| \leq 1$ for $x \in [-e, e]$. Moreover, for $x \in (e, \infty)$ with $x \rightarrow +\infty$, we have $|\widehat{s}(x)|/x^p \rightarrow 0$ for all $p \in (1, \infty)$. By symmetry, the same argument holds for $x \in (-\infty, e)$ with $x \rightarrow -\infty$. Therefore, by Krasnosel'skii's theorem, see e.g. [12, Theorem 2.2], $\widehat{s}: L^p(\Omega) \rightarrow L^1(\Omega)$ is continuous for $p \in (1, \infty)$.

Clearly, if we restrict \widehat{s} to $L_+^p(\Omega)$, then we have a continuous mapping $\widehat{s}|_{L_+^p(\Omega)}$ on $L_+^p(\Omega)$. This function coincides with

$$(A.15) \quad s(x) = \begin{cases} x \ln x - x, & x > 0, \\ 0, & x = 0, \\ +\infty, & \text{otherwise,} \end{cases}$$

on $L_+^p(\Omega)$. Hence, continuity of S on $L_+^p(\Omega)$ now follows from the continuity of the Lebesgue integral $u \mapsto \int_{\Omega} u \, dx$ and the fact that the composition of two continuous functions is also continuous. Finally, suppose $\{u_k\} \subset L_+^{\infty}(\Omega)$ converges to u in $L^{\infty}(\Omega)$. Then $\{u_k\} \subset L_+^p(\Omega)$ for every $p \in (1, \infty)$ and, moreover, $u_k \rightarrow u$ in $L^p(\Omega)$ because Ω is a bounded domain. Consequently, $S(u_k) \rightarrow S(u)$ as $k \rightarrow +\infty$, as conjectured.

Case 3: $p = \infty$. In order to show that S is Fréchet differentiable on $\text{int } L_+^{\infty}(\Omega)$ with respect to the $L^{\infty}(\Omega)$ topology, we will first prove that S is Gâteaux differentiable on $\text{int } L_+^{\infty}(\Omega)$ and, subsequently, that the Gâteaux derivative S'_G is continuous on $\text{int } L_+^{\infty}(\Omega)$. Fréchet differentiability of S will then follow from [12, Theorem 1.9].

To show that S is Gâteaux differentiable on $\text{int } L_+^{\infty}(\Omega)$, we must show that for any fixed $u \in \text{int } L_+^{\infty}(\Omega)$ and $v \in L^{\infty}(\Omega)$,

$$(A.16) \quad \lim_{\tau \rightarrow 0} \int_{\Omega} \frac{s(u + \tau v) - s(u)}{\tau} \, dx = \int_{\Omega} v \ln(u) \, dx.$$

First observe that for almost every $x \in \Omega$, we have pointwise convergence of the associated integrands, namely,

$$(A.17) \quad \lim_{\tau \rightarrow 0} \frac{s(u(x) + \tau v(x)) - s(u(x))}{\tau} = v(x) \ln(u(x)).$$

Next, we know from the proof of [Proposition A.7](#) that $u \geq \|\frac{1}{u}\|_{L^{\infty}}^{-1}$. This implies that for sufficiently small τ , $u + \tau v > \|\frac{1}{u}\|_{L^{\infty}}^{-1}/2$ holds a.e., and we have

$$(A.18) \quad s(u + \tau v) - s(u) = \int_u^{u+\tau v} \ln \sigma \, d\sigma = \tau v \int_0^1 \ln(u + \sigma \tau v) \, d\sigma.$$

The critical step is to see that for the u and v fixed above, we may find $w \in L^{\infty}(\Omega)$ where

$$(A.19) \quad v = uw.$$

As such, for all sufficiently small τ , we may rewrite

$$(A.20) \quad s(u + \tau v) - s(u) = \tau v \left(\ln u + \int_0^1 \ln(1 + \sigma \tau w) \, d\sigma \right),$$

and, consequently,

$$(A.21) \quad \left| \frac{s(u + \tau v) - s(u)}{\tau} - v \ln u \right| = \left| v \int_0^1 \ln(1 + \sigma \tau w) \, d\sigma \right| \leq |v| |\ln(1 + \tau w)|.$$

To arrive at an upper bound that is independent of τ , we use the following well-known inequality:

$$(A.22) \quad \frac{x}{x+1} \leq \ln(1+x) \leq x.$$

In turn, for all $|\tau| < (2\|w\|_{L^\infty})^{-1}$,

$$(A.23) \quad |\ln(1 + \tau w)| \leq 1.$$

The penultimate argument owes to the function $|v|$ belonging to $L^1(\Omega)$ because Ω is bounded. Indeed, by (A.17), (A.21), and (A.23), the dominated convergence theorem provides us with the following well-defined Gâteaux derivative:

$$(A.24) \quad \langle S'_G(u), v \rangle = \lim_{\tau \rightarrow 0} \int_{\Omega} \frac{s(u + \tau v) - s(u)}{\tau} dx = \int_{\Omega} v \ln u dx.$$

It remains to show that $S'_G: \text{int } L_+^\infty(\Omega) \subset L^\infty(\Omega) \rightarrow [L^\infty(\Omega)]'$ is continuous. To this end, let $u \in \text{int } L_+^\infty(\Omega)$ and consider any sequence $\{u_k\}$ in $\text{int } L_+^\infty(\Omega)$ where $u_k \rightarrow u$ in $L^\infty(\Omega)$. Consequently, we know there exists $C > 0$ such that $\|u_k\|_{L^\infty} \leq C$ and $u_k(x) \rightarrow u(x)$ for almost every $x \in \Omega$. Clearly,

$$(A.25) \quad |\langle S'_G(u) - S'_G(u_k), v \rangle| \leq \int_{\Omega} |v| |\ln |u/u_k|| dx \leq \|v\|_{L^\infty} \int_{\Omega} |\ln |u/u_k|| dx,$$

where $|\ln |u/u_k|| \leq |\ln |u|| + |\ln C| \in L^1(\Omega)$. Therefore,

$$(A.26) \quad \|S'_G(u) - S'_G(u_k)\|_{[L^\infty]'} \leq \int_{\Omega} |\ln |u/u_k|| dx$$

and, by the dominated convergence theorem,

$$(A.27) \quad \lim_{k \rightarrow \infty} \|S'_G(u) - S'_G(u_k)\|_{[L^\infty]'} \leq \int_{\Omega} \lim_{k \rightarrow \infty} |\ln |u/u_k|| dx = 0,$$

as necessary.

Step 3. Let $u \in \text{int } L_+^\infty(\Omega)$. By (A.1b) $\ln u \in L^\infty(\Omega)$. Next, we see that $\|S'(u)\|_{[L^\infty]'} \leq \|\ln u\|_{L^1}$, since

$$(A.28) \quad \langle S'(u), v \rangle = \int_{\Omega} v \ln u dx \leq \|v\|_{L^\infty} \|\ln u\|_{L^1}.$$

Moreover,

$$(A.29) \quad \|S'(u)\|_{[L^\infty]'} \geq \int \mathbb{1}_{\{\ln u > 0\}} \ln u dx - \int \mathbb{1}_{\{\ln u < 0\}} \ln u dx = \|\ln u\|_{L^1},$$

and so $\|S'(u)\|_{[L^\infty]'} = \|\ln u\|_{L^1}$. \square

We complete this subsection with a proof of the gradient representation theorem for the shifted entropy functional, $S_\phi(u) = S(u - \phi)$.

Proof of Corollary 4.2. By Theorem 4.1, S is continuous on $L_+^\infty(\Omega)$. If $\phi \in L^\infty(\Omega)$, then the shift operator $T_\phi u := u - \phi$ is continuous on $L^\infty(\Omega)$ for $u \in L^\infty(\Omega)$ with $u \geq \phi$, as well; continuity of the composition follows.

As argued in Theorem 4.1, S is strictly convex on $L_+^\infty(\Omega)$. Taking $w_i \in L_{\phi,+}^\infty(\Omega)$, with $i = 1, 2$ and $w_1 \neq w_2$, we see that $T_\phi w_i = w_i - \phi \geq 0$ a.e. for $i = 1, 2$. Moreover, $T_\phi w_1 = T_\phi w_2$ iff $w_1 = w_2$ and for $\lambda \in (0, 1)$ we have $T_\phi(\lambda w_1 + (1 - \lambda)w_2) = \lambda(w_1 - \phi) + (1 - \lambda)(w_2 - \phi) = \lambda T_\phi w_1 + (1 - \lambda)T_\phi w_2$. Then since $T_\phi w_i \in L_+^\infty(\Omega)$ for $i = 1, 2$, the strict convexity of the composition follows.

We proceed with the characterization of $\text{int } L_{\phi,+}^{\infty}(\Omega)$. First, we can easily show the elementary properties $L_{\phi,+}^{\infty}(\Omega) = \phi + L_+^{\infty}(\Omega)$ and

$$\text{int } L_{\phi,+}^{\infty}(\Omega) = \text{int}(\phi + L_+^{\infty}(\Omega)) = \phi + \text{int } L_+^{\infty}(\Omega).$$

Therefore, $w \in \text{int } L_{\phi,+}^{\infty}(\Omega)$ implies $w - \phi \in \text{int } L_+^{\infty}(\Omega)$. By [Proposition A.7](#), $w \geq \phi$ and $\text{ess inf}(w - \phi) > 0$. Conversely, if $w \in L^{\infty}(\Omega)$ such that $w \geq \phi$ and $\text{ess inf}(w - \phi) > 0$, then [Proposition A.7](#) implies $w - \phi \in \text{int } L_+^{\infty}(\Omega)$. Hence, $w \in \text{int } L_{\phi,+}^{\infty}(\Omega)$.

Let $w_1 \in \text{int } L_{\phi,+}^{\infty}(\Omega)$. Then $w_1 - \phi \in \text{int } L_+^{\infty}(\Omega)$. It follows from [Theorem 4.1](#) that S_{ϕ} is Fréchet differentiable at w_1 .

The formula for the derivative of S'_{ϕ} can be viewed as an application of the chain rule. Indeed, $S_{\phi} = S \circ T_{\phi}$, S is differentiable with respect to the L^{∞} -norm at $T_{\phi}w$ with $w \in \text{int } L_{\phi,+}^{\infty}(\Omega)$ and T_{ϕ} is differentiable with respect to the L^{∞} -norm at (any) $w \in L^{\infty}$ with derivative $A'_{\phi}(w)$ given by the identity on L^{∞} . Therefore, we have [\(4.9\)](#). Since $u - \phi \in \text{int } L_+^{\infty}(\Omega)$ the rest of the computations for the gradient remain unchanged; in particular, we obtain [\(4.10\)](#) and [\(4.11\)](#). \square

A.3. Deriving the entropic Poisson equation. We devote this subsection to proofs of the characterization theorem and its corollary for non-zero obstacles (i.e., [Corollary 4.9](#)). This section also includes a short remark about a weak maximum principle for the entropic Poisson equation that arises from the first of these proofs.

Proof of [Theorem 4.7](#). The proof proceeds in five steps.

Step 1. Show that there exists a unique solution.

The proof of existence is standard. We sketch the main points here; see, e.g., [\[18, Chap. 3.2\]](#) for details. By [\[68, Lem. 3.30\]](#), we have that

$$(A.30) \quad \|v\|_{L^2} - \int_{\partial\Omega} g \, d\mathcal{H}_{n-1} \leq c \|\nabla v\|_{L^2}, \quad \text{for all } v \in H_g^1(\Omega),$$

for some constant c that depends on Ω . Clearly, A_{α} is finite on K . This yields a minimizing sequence $\{u_k\}$. The form of A_{α} consequently yields the boundedness of $\{\|\nabla u_k\|_{L^2}\}$. Combined with [\(A.30\)](#) we deduce boundedness of $\{u_k\}$ in $H^1(\Omega)$. We can readily show that A_{α} is weakly lower-semicontinuous and K weakly sequentially closed. This yields the existence of a minimizer $u \in K$. The minimizer u is unique because A_{α} is strictly convex on K .

Step 2. Show that $u \leq \max\{\|g\|_{L^{\infty}(\partial\Omega)}, \exp(\|\ln w + \alpha f\|_{L^{\infty}(\Omega)})\}$.

For all $N > 1$, define the set $R_N = \{x \in \Omega \mid u(x) > N\}$. By way of contradiction, we assume that $|R_N| > 0$ for all $N > 1$. Now, consider the following function in L^{∞} :

$$(A.31) \quad u_N(x) = \min\{N, u(x)\} \geq 0.$$

We claim that if $N > \text{ess sup}_{x \in \partial\Omega} g(x)$ then $u_N \in K$. Begin by choosing $\{u_m\} \subset C^1(\bar{\Omega})$ such that $u_m \rightarrow u$ (strongly in $H^1(\Omega)$) and define $u_m^N := \min\{N, u_m\}$. The existence of u_m follows from the assumption that $\partial\Omega$ is Lipschitz; see, e.g., [\[4, 3.22 Theorem\]](#). Next let $v \in L^{\infty}(\partial\Omega)$ and consider that

$$\int_{\partial\Omega} \gamma(u_m^N) v \, d\mathcal{H}_{n-1} = \int_{\{\gamma(u_m) \leq N\}} \gamma(u_m) v \, d\mathcal{H}_{n-1} + N \int_{\{\gamma(u_m) > N\}} v \, d\mathcal{H}_{n-1}.$$

Along a subsequence, denoted still by m , $\gamma(u_m)$ converges pointwise almost everywhere to $\gamma(u) = g$. Then, by hypothesis, the sequence of characteristic functions

$f_m := \chi_{\{\gamma(u_m) > N\}} \rightarrow 0$ pointwise almost everywhere. It follows from Lebesgue's dominated convergence theorem that

$$N \int_{\{\gamma(u_m) > N\}} v \, d\mathcal{H}_{n-1} \rightarrow 0 \quad \text{as } m \rightarrow +\infty.$$

Continuing, we appeal to the proof of [116, Thm A.1] and, e.g., [133, Cor. 18.4], to argue that $\min\{N, u_m\} \rightarrow \min\{N, u\}$ weakly in $H^1(\Omega)$ and $\gamma(\min\{N, u_m\}) \rightarrow \gamma(\min\{N, u\})$ strongly in $L^2(\Omega)$, which in turn yields

$$\lim_{m \rightarrow +\infty} \int_{\{\gamma(u_m) \leq N\}} \gamma(u_m) v \, d\mathcal{H}_{n-1} = \int_{\partial\Omega} \gamma(u_N) v \, d\mathcal{H}_{n-1}.$$

Since v is essentially bounded, we have

$$|\gamma(u_m)v - \gamma(u)v| = |v| |\gamma(u_m) - \gamma(u)| \leq \|v\|_{L^\infty} |\gamma(u_m) - \gamma(u)|.$$

Hence, $\gamma(u_m)v$ converges strongly in $L^2(\partial\Omega)$ to $\gamma(u)v$. Similar to the above, we can argue that $f'_m := \chi_{\{\gamma(u_m) \leq N\}} \rightarrow \chi_{\partial\Omega} = 1$ in $L^2(\partial\Omega)$. It follows that

$$\int_{\partial\Omega} \gamma(u_N)v \, d\mathcal{H}_{n-1} = \int_{\partial\Omega} \gamma(u)v \, d\mathcal{H}_{n-1} = \int_{\partial\Omega} gv \, d\mathcal{H}_{n-1}.$$

By the density of $L^\infty(\partial\Omega)$ in $L^2(\partial\Omega)$ and the fundamental lemma of the calculus of variations [68, Theorem 1.32], we deduce $\gamma(u_N) = g$ a.e. on $\partial\Omega$. Consequently, $u_N \in K$ and for sufficiently large N , it holds that

$$(A.32) \quad D(u, w) + \alpha E(u) < D(u_N, w) + \alpha E(u_N)$$

because u is the unique global minimizer of A_α over K .

Note, however, that

$$\alpha E(u) - \alpha E(u_N) = \alpha \int_{R_N} \frac{1}{2} |\nabla u|^2 - (u - N)f \, dx$$

and

$$\begin{aligned} D(u, w) - D(u_N, w) &= \int_{R_N} u \ln u - N \ln N - (1 + \ln w)(u - N) \, dx \\ &= \int_{R_N} (u - N) \left(\int_0^1 \ln(N + t(u - N)) \, dt - \ln w \right) \, dx. \end{aligned}$$

Combining these observations, we see that

$$\begin{aligned} D(u, w) + \alpha E(u) - D(u_N, w) - \alpha E(u_N) &\geq \\ &\frac{\alpha}{2} \|\nabla u\|_{L^2(R_N)}^2 - (\alpha f + \ln w - \ln N, u - N)_{L^2(R_N)}. \end{aligned}$$

Therefore, for any $N > \max\{\|g\|_{L^\infty(\partial\Omega)}, \exp(\|\alpha f + \ln w\|_{L^\infty(\Omega)})\}$, we have

$$(A.33) \quad D(u, w) + \alpha E(u) > D(u_N, w) + \alpha E(u_N),$$

which contradicts the optimality of u . Hence, there exists some $N_0 > 0$ such that $|R_N| = 0$ for all $N > N_0$, and, in turn, $u \in L^\infty$ with $u \leq \max\{\|g\|_{L^\infty(\partial\Omega)}, \exp(\|\alpha f + \ln w\|_{L^\infty(\Omega)})\}$.

Step 3. Show that $u \geq \min\{\text{ess inf}_{x \in \partial\Omega} g(x), \exp(-\|\ln v + \alpha f\|_{L^\infty(\Omega)})\}$. Thus, $u \in H_g^1(\Omega) \cap \text{int } L_+^\infty(\Omega)$.

For all $\epsilon > 0$, define the set $S_\epsilon = \{x \in \Omega \mid u(x) < \epsilon\}$. By way of contradiction, we assume that $|S_\epsilon| > 0$ for all $\epsilon > 0$. Now, consider the following function in $H^1(\Omega) \cap \text{int } L_+^\infty(\Omega)$:

$$(A.34) \quad u_\epsilon(x) = \max\{\epsilon, u(x)\}.$$

The fact that $u_\epsilon \in \text{int } L_+^\infty$ follows from [Proposition A.7](#).

Continuing, we can emulate the arguments of Step 2. above to show that if $\epsilon < \text{ess inf}_{x \in \partial\Omega} g(x)$, then $u_\epsilon \in H_g^1 \cap \text{int } L_+^\infty$. The steps and justifications are almost identical and are therefore omitted. In turn, for sufficiently small $\epsilon > 0$, it holds that

$$(A.35) \quad D(u, w) + \alpha E(u) \leq D(u_\epsilon, w) + \alpha E(u_\epsilon),$$

As above, we obtain the lower bound

$$\begin{aligned} D(u, w) + \alpha E(u) - D(u_\epsilon, w) - \alpha E(u_\epsilon) \geq \\ \alpha \|\nabla u\|_{L^2(S_\epsilon)}^2 + (\|\alpha f + \ln w\|_{L^\infty} + \ln \epsilon, u - \epsilon)_{L^2(S_\epsilon)} \end{aligned}$$

As a result, once $\epsilon < \exp(-\|\alpha f + \ln w\|_{L^\infty})$, we again contradict the optimality of u . Thus, there exists some $\epsilon_0 > 0$ such that $|S_\epsilon| = 0$ for all $\epsilon < \epsilon_0$ and $u \in H_g^1(\Omega) \cap \text{int } L_+^\infty(\Omega)$ by [Proposition A.7](#).

Step 4. Derive the variational equation.

Let $t > 0$ and $v \in K \cap L^\infty(\Omega)$. Then, by definition,

$$\alpha E(u) + D(u, w) \leq \alpha E(v) + D(v, w).$$

Clearly, $u + t(v - u) \in K$, and consequently,

$$0 \leq \frac{\alpha E(u + t(v - u)) - \alpha E(u)}{t} + \frac{D(u + t(v - u), w) - D(u, w)}{t}$$

Since, $u \in H^1(\Omega) \cap \text{int } L_+^\infty(\Omega)$ and $v - u \in H^1(\Omega) \cap L^\infty(\Omega)$, [Theorem 4.1](#) allows us to expand and pass to the limit as $t \downarrow 0$. This yields the variational inequality:

$$0 \leq \alpha E'(u)(v - u) + S'(u)(v - u) - (\ln w, v - u)_{L^2}$$

for all $v \in H_g^1(\Omega) \cap L_+^\infty(\Omega)$. It is readily observed that $H_g^1(\Omega) \cap \text{int } L_+^\infty(\Omega) \subset \text{int}(H_g^1(\Omega) \cap L_+^\infty(\Omega))$. Therefore, as a result of Steps 2. and 3., u is in the $H^1(\Omega) \cap L^\infty(\Omega)$ interior of the set K . In other words, for sufficiently small $\delta > 0$, $u + v \in H_g^1(\Omega) \cap L^\infty(\Omega)$ and $u + v \geq 0$ for any $v \in H_0^1(\Omega) \cap L^\infty(\Omega)$ with $\|v\|_{H^1 \cap L^\infty} < \delta$. In turn, we obtain the following variational equation with test functions $v \in H_0^1(\Omega) \cap L^\infty(\Omega)$ and $\|v\|_{H^1 \cap L^\infty} < \delta$:

$$0 = \alpha E'(u)v + S'(u)v - (\ln w, v)_{L^2}.$$

The first summand is equivalent to $(\alpha \nabla u, \nabla v)_{L^2} - (\alpha f, v)_{L^2}$ and the second and third summands together have the form $(\ln u - \ln w, v)_{L^2}$. Since $u, w \in \text{int } L_+^\infty(\Omega)$, the map $v \mapsto (\ln u - \ln w, v)_{L^2}$ defines a bounded linear functional on $H_0^1(\Omega)$. Finally, by virtue of the inclusion $C_c^\infty(\Omega) \subset H_0^1(\Omega) \cap L^\infty(\Omega)$, we deduce

$$(A.36) \quad (\alpha \nabla u, \nabla v)_{L^2} + (\ln u, v)_{L^2} = (\alpha f, v)_{L^2} + (\ln w, v)_{L^2} \quad \text{for all } v \in H_0^1(\Omega),$$

as was to be shown.

Step 5. Prove that the entropic Poisson equation has a unique solution in $H_g^1(\Omega) \cap \text{int } L_+^\infty(\Omega)$.

Conversely, suppose $u \in H_g^1(\Omega) \cap \text{int } L_+^\infty(\Omega)$ such that (A.36) holds. Then for any $y \in H_g^1(\Omega) \cap L^\infty(\Omega)$, setting $v = y - u \in H_0^1(\Omega) \cap L^\infty(\Omega) \subset H_0^1(\Omega)$ in (A.36) yields

$$A'_\alpha(u)(y-u) = (\alpha \nabla u, \nabla [y-u])_{L^2} + (\ln u, y-u)_{L^2} - (\alpha f, y-u)_{L^2} - (\ln w, y-u)_{L^2} = 0.$$

Since A_α is differentiable at u with respect to variations in $H^1(\Omega) \cap L^\infty(\Omega)$ and convex on $H_g^1(\Omega)$, we have

$$A_\alpha(y) \geq A_\alpha(u) + A'_\alpha(u)(y-u) = A_\alpha(u) \quad \forall y \in H_g^1(\Omega) \cap L^\infty(\Omega).$$

Taking the closure of $H_g^1(\Omega) \cap L^\infty(\Omega)$ with respect to the H^1 -norm, it follows from the continuity of A_α on $H^1(\Omega)$ that u is a minimizer of A_α over $H_g^1(\Omega)$. Uniqueness follows from strict convexity of A_α . \square

Remark A.10 (Maximum principle). Inspecting the proof above, we see that

$$(A.37a) \quad \min\{g_{\min}, \exp(-\|\ln w + \alpha f\|_{L^\infty})\} \leq u \leq \max\{g_{\max}, \exp(\|\ln w + \alpha f\|_{L^\infty})\},$$

or, equivalently,

$$(A.37b) \quad \min\{\ln g_{\min}, -\|\ln w + \alpha f\|_{L^\infty}\} \leq \ln u \leq \max\{\ln g_{\max}, \|\ln w + \alpha f\|_{L^\infty}\},$$

where $g_{\min} = \text{ess inf}_{x \in \partial\Omega} g(x)$ and $g_{\max} = \text{ess sup}_{x \in \partial\Omega} g(x)$.

Proof of Corollary 4.9. The proof follows that of Theorem 4.7. Here, Corollary 4.2 plays the same role as Theorem 4.1.

Existence and uniqueness follows the homogeneous case in light of the implications of Corollary 4.2. We only need argue that K_ϕ is nonempty. Since $g, \phi \in H^1(\Omega) \cap C(\bar{\Omega})$ the function $v := \max\{g, \phi\} = g + \max\{0, \phi - g\}$ is in $H^1(\Omega) \cap C(\bar{\Omega})$ and satisfies $v \geq \phi$. The trace of w is merely the evaluation on the boundary. Then since $\text{ess inf } \gamma(g - \phi) > 0$ on $\partial\Omega$ by assumption, we have $\gamma(v) = \gamma(g)$ and consequently $v \in K_\phi$.

Setting $\tilde{w} = w - \phi$, we can now readily argue that $u = \tilde{u} + \phi$ where \tilde{u} is the solution of

$$(A.38) \quad \min \frac{1}{2} \|\nabla \tilde{v}\|_{L^2}^2 - (f + \Delta\phi, \tilde{v})_{L^2} + \alpha^{-1} D(\tilde{v}, \tilde{w})$$

over $\tilde{v} \in H_{g-\phi}^1(\Omega)$ subject to $\tilde{v} \geq 0$ in Ω .

Theorem 4.7 then guarantees that \tilde{u} solves

$$(\alpha \nabla \tilde{u}, \nabla v) + (\ln \tilde{u}, v) = (\alpha f + \alpha \Delta\phi, v) + (\ln \tilde{w}, v) \quad \text{for all } v \in H_0^1(\Omega).$$

Substituting $u - \phi = \tilde{u}$ and $w - \phi = \tilde{w}$ yields (4.33). \square

A.4. Towards an entropic Poisson equation with homogeneous boundary conditions. In this subsection, we contemplate the mathematical meaning of an entropic Poisson equation with homogeneous boundary conditions. Theorem 4.7 intentionally avoids this setting due to the strict positivity requirement used to argue that $\max\{\epsilon, u\} \in H_g^1(\Omega) \cap \text{int } L_+^\infty(\Omega)$ in step 3 of the proof. Nevertheless, as the following result demonstrates, we can still provide important insight into the Dirichlet free energy minimization problem (4.30) in the situation where $g \equiv 0$ on $\partial\Omega$ by applying Theorem 4.7 sequentially.

THEOREM A.11. *In addition to the assumptions of [Theorem 4.7](#), suppose that $g \equiv \varepsilon > 0$ and denote the solution of the corresponding entropic Poisson equation by u_ε . Denote the unique solution of [\(4.30\)](#) for $g \equiv 0$ by \bar{u} . For any sequence of scalars $\varepsilon_k \downarrow 0$, the sequence $\{u_k\}$ with $u_k := u_{\varepsilon_k}$ satisfies the following properties:*

1. $u_k \rightarrow \bar{u}$ strongly in $H^1(\Omega)$;
2. $u_k \rightarrow \bar{u}$ weak-* in $L^\infty(\Omega)$;
3. $\gamma(u_k - \bar{u}) \rightarrow 0$ strongly in $L^\infty(\partial\Omega)$.

Proof. According to [Remark A.10](#), we have the bounds

$$\min\{\varepsilon_k, \exp(-\|\ln w + \alpha f\|_{L^\infty})\} \leq u_k \leq \max\{\varepsilon_k, \exp(\|\ln w + \alpha f\|_{L^\infty})\}.$$

Therefore, for sufficiently large k , the bounds reduce to $\varepsilon_k \leq u_k \leq \exp(\|\ln w + \alpha f\|_{L^\infty})$. Since $L^\infty(\Omega)$ is the topological dual of a separable Banach space, there exists a subsequence $\{k_l\}$ and $\tilde{u} \in L^\infty(\Omega)$ such that $u_{k_l} \rightarrow \tilde{u}$ weak-* in $L^\infty(\Omega)$. Clearly, \tilde{u} satisfies $0 \leq \tilde{u} \leq \exp(\|\ln w + \alpha f\|_{L^\infty})$ a.e. in Ω .

Continuing on, we note that $v = u_k - \varepsilon_k \in H^1(\Omega)$ and $\gamma(v) = 0$. Using v as a test function in [\(4.31\)](#), we deduce that

$$\alpha(\nabla u_k, \nabla u_k) \leq |\Omega|/\exp(1) + \|\alpha f + \ln w\|_{L^1(\Omega)} \exp(\|\ln w + \alpha f\|_{L^\infty(\Omega)}).$$

This follows from the fact that $(\nabla u_k, \nabla(u_k - \varepsilon_k)) = (\nabla u_k, \nabla u_k)$,

$$0 \leq u_k - \varepsilon_k \leq \exp(\|\ln w + \alpha f\|_{L^\infty(\Omega)}),$$

and

$$(\ln(u_k), u_k - \varepsilon_k) = (\ln(u_k - \varepsilon_k + \varepsilon_k), u_k - \varepsilon_k) \geq (\ln(u_k - \varepsilon_k), u_k - \varepsilon_k),$$

along with $x \ln x \geq -1/\exp(1)$. Consequently, there exists a further subsequence $\{k_{l_m}\}$ such that $u_{k_{l_m}} \rightarrow \tilde{u}$ weakly in $H^1(\Omega)$ and we deduce that \tilde{u} is feasible to [\(4.30\)](#).

We now demonstrate that $\tilde{u} = \bar{u}$. The fact that \bar{u} is unique follows from the same arguments in step 1 of the proof of [Theorem 4.7](#). Defining $\bar{u}_\varepsilon := \max\{\varepsilon, \bar{u}\}$, we have $\bar{u}_\varepsilon \geq 0$ and $\gamma(\bar{u}_\varepsilon) = \varepsilon$. Consider additionally that

$$\frac{1}{2} \|\nabla \bar{u}_{\varepsilon_k}\|_{L^2(\Omega)}^2 = \frac{1}{2} \int_{\{\bar{u} \geq \varepsilon_k\}} |\nabla \bar{u}|^2 dx \rightarrow \frac{1}{2} \int_{\Omega} |\nabla \bar{u}|^2 dx,$$

by monotone convergence and $\bar{u} \geq 0$. Similarly, we see that $\|\bar{u}_{\varepsilon_k} - \bar{u}\|_{L^2(\Omega)}^2 = O(\varepsilon^2)$. By optimality of u_ε we have

$$E(u_k) + \alpha^{-1}D(u_k, w) \leq E(\bar{u}_{\varepsilon_k}) + \alpha^{-1}D(\bar{u}_{\varepsilon_k}, w).$$

Using the subsequence $\{k_{l_m}\}$, we pass to the limit inferior on both sides. The previous observations along with the continuity/weak lower semicontinuity properties of both the entropy and E imply

$$E(\tilde{u}) + \alpha^{-1}D(\tilde{u}, w) \leq E(\bar{u}) + \alpha^{-1}D(\bar{u}, w),$$

whence we have $\tilde{u} = \bar{u}$. Since \bar{u} is unique, it follows from the Urysohn subsequence property that the entire sequence $\{u_k\}$ converges weakly in $H^1(\Omega)$.

Strong convergence follows by rearranging terms in the optimality statement and considering the limit superior. Indeed, we have

$$(A.39) \quad \frac{1}{2} \|\nabla u_{\varepsilon_k}\|_{L^2}^2 - \frac{1}{2} \|\nabla \bar{u}_{\varepsilon_k}\|_{L^2}^2 \leq (f, u_{\varepsilon_k} - \bar{u}_{\varepsilon_k}) + \alpha^{-1}D(\bar{u}_{\varepsilon_k}, w) - \alpha^{-1}D(u_k, w).$$

This indicates that $\{u_k\}$ converges strongly in $H^1(\Omega)$ by the weak lower semicontinuity of the term $\|\nabla \cdot\|_{L^2}^2$ on $H^1(\Omega)$ and the Kadec–Klee property

Finally, we return to the L^∞ -statements. By the linearity of the trace, we have $\gamma(u_k - \bar{u}) = \varepsilon_k$, which clearly converges to 0 in $L^\infty(\partial\Omega)$. Likewise, if we assume that u_k possess a subsequence that does not weak-* converge to \bar{u} in $L^\infty(\Omega)$, then we can find a further subsequence that does, which leads to a contradiction and completes the proof. \square

Without a deeper analysis of the properties of \bar{u} in [Theorem A.11](#), it is difficult to say whether a type of entropic Poisson equation can be derived for the fully homogeneous case. Nevertheless, we do know that \bar{u} is an optimal solution and is, therefore, feasible. On the active set, $\bar{u} = 0$ is fully determined. Moreover, consider that

$$\begin{aligned} (u_k, \ln u_k) - (\bar{u}, \ln \bar{u}) &= (u_k - \bar{u}, \ln u_k) + (\bar{u}, \ln u_k - \ln \bar{u}) \\ &= \langle u_k - \bar{u}, \alpha f + \ln w + \alpha \Delta u_k \rangle + (\bar{u}, \ln u_k - \ln \bar{u}). \end{aligned}$$

Given $u_k \rightarrow \bar{u}$ in $H^1(\Omega)$ implies $\Delta u_k \rightarrow \Delta \bar{u}$ strongly in $H^1(\Omega)'$ as well as $S(u_k) \rightarrow S(\bar{u})$, it must follow that $(\bar{u}, \ln u_k) \rightarrow (\bar{u}, \ln \bar{u})$. This implies a complementarity relation in the limit; namely,

$$\langle -\alpha \Delta \bar{u} + \ln \bar{u} - \alpha f - \ln w, \bar{u} \rangle = 0.$$

It remains to investigate what happens on the inactive set, i.e., where $\bar{u} > 0$. To this aim, let $I_\delta := \{x \in \Omega \mid \bar{u}(x) \geq \delta > 0\}$. By Egorov's theorem, we can find a subsequence $\{u_{k_l}\}$ with the property that for every $\eta > 0$, there exists a measurable set $\mathcal{B}_{\eta, \delta} \subset I_\delta$ with $|I_\delta \setminus \mathcal{B}_{\eta, \delta}| < \eta$ such that $u_{k_l} \rightarrow \bar{u}$ uniformly on $\mathcal{B}_{\eta, \delta}$. This implies that $\ln u_{k_l} \rightarrow \ln \bar{u}$ uniformly on $\mathcal{B}_{\eta, \delta}$, as well. Consequently, we can localize the entropic Poisson equation with test functions $\varphi \in C_c^\infty(\mathcal{B}_{\eta, \delta})$, leading to

$$\langle -\alpha \Delta u_{k_l} + \ln u_{k_l} - \alpha f - \ln w, \varphi \rangle = 0.$$

Passing to the limit, we see that

$$\langle -\alpha \Delta \bar{u} + \ln \bar{u} - \alpha f - \ln w, \varphi \rangle = 0.$$

This is well-defined due to the restriction to subsets of I_δ . Therefore, up to arbitrarily small subsets of the strict inactive sets I_δ (for all $\delta > 0$) we have recovered the entropic Poisson equation $-\alpha \Delta \bar{u} + \ln \bar{u} = \alpha f + \ln w$ in the sense of distributions.

A.5. Convergence of the latent variable proximal point method. In this subsection, we establish arbitrary convergence rates for the continuous-level proximal point algorithm [\(4.38\)](#) to solve the obstacle problem. We begin by proving the following lemma.

LEMMA A.12. *Under the assumptions of [Theorem 4.7](#), the second-order problem*

$$(A.40a) \quad \text{Find } u \in H_g^1(\Omega) \cap \text{int } L_+^\infty(\Omega) \text{ such that } -\Delta u + \ln u = f \text{ in } H^{-1}(\Omega),$$

is equivalent to the saddle-point problem

$$(A.40b) \quad \text{Find } \tilde{u} \in H_g^1(\Omega) \text{ and } \tilde{\psi} \in L^\infty(\Omega) \text{ such that } \begin{cases} -\Delta \tilde{u} + \tilde{\psi} = f & \text{in } H^{-1}(\Omega), \\ \tilde{u} - \exp \tilde{\psi} = 0 & \text{in } L^2(\Omega). \end{cases}$$

More specifically, both problems admit unique solutions that coincide in the sense that $u = \tilde{u}$ and $\ln u = \tilde{\psi}$ a.e in Ω .

Proof. First of all, we know from [Theorem 4.7](#) that there exists a unique solution to [\(A.40a\)](#). Using [Proposition A.7](#), we know that $\exp: L^\infty(\Omega) \rightarrow \text{int } L_+^\infty(\Omega)$ is an isomorphism. Thus, $\tilde{\psi} = \ln u$ and $\tilde{u} = u$ form a solution to [\(A.40b\)](#). Now, assume that $\tilde{u} \in H_g^1(\Omega)$ and $\tilde{\psi} \in L^\infty(\Omega)$ form an arbitrary solution to [\(A.40b\)](#). By the second equation in [\(A.40b\)](#), we know that

$$(A.41) \quad \tilde{u} = \exp \tilde{\psi} \quad \text{a.e. in } \Omega.$$

Now, by [Proposition A.7](#), we know that $\exp \tilde{\psi} \in \text{int } L_+^\infty(\Omega)$. Thus, $\tilde{u} \in H_g^1(\Omega) \cap \text{int } L_+^\infty(\Omega)$. Moreover, by applying \ln to both sides of [\(A.41\)](#), we find that $\tilde{\psi} = \ln \tilde{u}$. Thus, $\tilde{u} \in H_g^1(\Omega) \cap \text{int } L_+^\infty(\Omega)$ solves [\(A.40a\)](#). Since the solution of [\(A.40a\)](#) is unique, we find that $\tilde{u} = u$ and, in turn, $\tilde{\psi} = \ln u$. \square

We now move on to proving the convergence theorem.

Proof of Theorem 4.13. The proof has three main steps, the first two of which build off of [Lemma A.12](#). Without loss of generality, we focus on the case where $\phi = 0$. The statement for general obstacles $\phi \neq 0$ can be recovered by making the change of variables $u - \phi = \tilde{u}$ and $w - \phi = \tilde{w}$ used in the proof of [Corollary 4.9](#).

Step 0. By [Lemma A.12](#), the sequence of iterates u^k coming from [\(4.38\)](#) and

$$(A.42) \quad (\alpha_{k+1} \nabla u^{k+1}, \nabla v) + (\ln u^{k+1}, v) = (\alpha_{k+1} f + \ln u^k, v) \quad \text{for all } v \in H_0^1(\Omega)$$

are equal a.e. in Ω . We take advantage of this fact throughout the proof below. In particular, given $u^0 = \exp \psi^0$ as in the hypotheses, we suppose that sequence $\{u^k\}$ is generated by the proximal point method, where each u^k solves [\(A.42\)](#), $k = 1, 2, \dots$. By [Theorem 4.7](#), $u^k \in H_g^1(\Omega) \cap \text{int } L_+^\infty(\Omega)$ for all $k = 0, 1, 2, \dots$.

Step 1. Inequality [\(4.39\)](#) is proved by exploiting the fact that $D(u, w) \geq 0$ with equality if and only if $u = w$. In particular,

$$(A.43) \quad \begin{aligned} E(u^{k+1}) &\leq E(u^{k+1}) + D(u^{k+1}, u^k) / \alpha_{k+1} \\ &\leq E(u^k) + D(u^k, u^k) / \alpha_{k+1} = E(u^k). \end{aligned}$$

Step 3. Since $u^j \in H_g^1(\Omega) \cap \text{int } L_+^\infty(\Omega)$ for all $j = 1, 2, \dots, k$, the definition of D as a true Bregman distance according to [\(4.16\)](#) is justified and consequently, the three-points identity [\(4.20\)](#) is as well. This leads to

$$(A.44) \quad D(w, u^j) - D(w, u^{j-1}) + D(u^j, u^{j-1}) = \langle S'(u^j) - S'(u^{j-1}), u^j - w \rangle,$$

where $w \in H_g^1(\Omega)$ such that $w \geq 0$ a.e. Next, notice that [\(A.42\)](#) is equivalent to

$$(A.45) \quad \langle S'(u^j), v \rangle - \langle S'(u^{j-1}), v \rangle = -\alpha_j E'(u^j) v \quad \text{for all } v \in H_0^1(\Omega).$$

Clearly, we have $u^j - w \in H_0^1(\Omega)$. Therefore, [\(A.45\)](#) and the subgradient inequality for E at u^j imply

$$(A.46) \quad \langle S'(u^j) - S'(u^{j-1}), u^j - w \rangle = \langle \alpha_j E'(u^j), w - u^j \rangle \leq \alpha_j E(w) - \alpha_j E(u^j).$$

Together, [\(A.44\)](#) and [\(A.46\)](#) imply that

$$(A.47) \quad D(w, u^k) + \sum_{j=1}^k D(u^j, u^{j-1}) + \sum_{j=1}^k \alpha_j [E(u^j) - E(w)] \leq D(w, u^0).$$

Now, given $D(w, u^k) \geq 0$, $D(u^{k+1}, u^k) \geq 0$ for all k and $E(u^k) \leq E(u^j)$ for all $j \leq k$, by (4.39), along with (A.47), we deduce the bound

$$(A.48) \quad E(u^k) \leq E(w) + \frac{D(w, u^0)}{\sum_{j=1}^k \alpha_j},$$

for all $w \in H_g^1(\Omega)$ satisfying $w \geq 0$. This step is completed by setting $w = u^*$ in the inequality above and using strong convexity of $E: H_0^1(\Omega) \rightarrow \mathbb{R}$. In particular, observe that

$$(A.49) \quad \begin{aligned} \frac{D(u^*, u^0)}{\sum_{j=1}^k \alpha_j} &\geq E(u^k) - E(u^*) \geq \langle E'(u^*), u^k - u^* \rangle + \frac{1}{2} \|\nabla u^* - \nabla u^k\|_{L^2}^2 \\ &\geq \frac{1}{2} \|\nabla u^* - \nabla u^k\|_{L^2}^2, \end{aligned}$$

where, we have used the first-order optimality condition $\langle E'(u^*), v - u^* \rangle \geq 0$ for all $v \in K$ in the final inequality.

Step 4. Finally, we prove the first equality in (4.41). To this end, consider the two equations

$$(A.50a) \quad (\nabla u^k, \nabla v) - (f, v) = (\lambda^k, v) \text{ for all } v \in H_0^1(\Omega),$$

and

$$(A.50b) \quad (\nabla u^*, \nabla v) - (f, v) = (\lambda^*, v) \text{ for all } v \in H_0^1(\Omega).$$

Combining these two equations, we find that

$$(A.51) \quad \|\lambda^* - \lambda^k\|_{H^{-1}(\Omega)} = \sup_{v \in H_0^1(\Omega)} \frac{(\lambda^* - \lambda^k, v)}{\|\nabla v\|_{L^2(\Omega)}} = \sup_{v \in H_0^1(\Omega)} \frac{(\nabla u^* - \nabla u^k, \nabla v)}{\|\nabla v\|_{L^2(\Omega)}}.$$

We now find $\|\lambda^* - \lambda^k\|_{H^{-1}(\Omega)} \leq \|\nabla u^* - \nabla u^k\|_{L^2(\Omega)}$ by applying the triangle inequality to the numerator of the third expression above. Likewise, we find $\|\nabla u^* - \nabla u^k\|_{L^2(\Omega)} \leq \|\lambda^* - \lambda^k\|_{H^{-1}(\Omega)}$ by considering the candidate function $v = u^* - u^k$. \square

We now turn to studying the iteration complexity of LVPP for various step size sequences. To this end, we first recall the standard definitions of Q- and R-convergence.

DEFINITION A.13 (Convergence orders and rates). *Let \mathcal{X} be a Banach space with norm $\|\cdot\|_{\mathcal{X}}$. We say that a sequence $\{x_k\}_{k=0}^{\infty} \subset \mathcal{X}$ converges to $x^* \in \mathcal{X}$ with order $q \geq 1$ and rate $r \geq 0$ if*

$$(A.52) \quad \lim_{k \rightarrow \infty} \frac{\|x_{k+1} - x^*\|_{\mathcal{X}}}{\|x_k - x^*\|_{\mathcal{X}}^q} = r.$$

If $q = 1$ and $r = 1$, then we say x_k converges Q-sublinearly to x^ . If $q = 1$ and $r \in (0, 1)$, then we say x_k converges Q-linearly to x^* . If $q > 0$ or $q = 1$ and $r = 0$, then we say x_k converges Q-superlinearly to x^* .*

If $\|x_k - x^\|_{\mathcal{X}} \leq \epsilon_k$ for all k , where ϵ_k converges Q-sublinearly (linearly, superlinearly) to zero, then we say that x_k converges R-sublinearly (linearly, superlinearly) to x^* .*

The following corollary establishes convergence orders associated to various step size sequences.

COROLLARY A.14 (Prescribed convergence orders). *Fix $C > 0$. Under the assumptions of [Theorem 4.13](#), consider the following candidate sequences of step sizes:*

Case 1: Fix $m \in \mathbb{N}$ and set

$$(A.53a) \quad \alpha_k = Ck(k+1) \cdots (k+m) \quad \text{for all } k = 1, 2, \dots$$

Case 2: Fix $\mu > 1$ and set

$$(A.53b) \quad \alpha_k = C\mu^{k-1} \quad \text{for all } k = 1, 2, \dots$$

Case 3: Set $\alpha_1 = C$ and

$$(A.53c) \quad \alpha_{k+1} = Ckk! \quad \text{for all } k = 1, 2, \dots$$

Case 4: Fix $\mu, q, r > 1$ and set $\alpha_1 = r^{1/(q-1)}\mu$ and

$$(A.53d) \quad \alpha_{k+1} = r^{1/(q-1)}\mu^{q^k} - \alpha_k \quad \text{for all } k = 2, 3, \dots$$

Then sequence [\(A.53a\)](#) delivers a sublinear R-convergence; sequence [\(A.53b\)](#) delivers R-linear convergence with rate $1/\mu$; sequence [\(A.53c\)](#) delivers R-superlinear convergence with order 1 and rate 0; and sequence [\(A.53d\)](#) delivers R-superlinear convergence with order q and rate r .

Proof. Throughout the proof, we use the definition $\epsilon_k = D(u^*, u^0) / \sum_{j=1}^k \alpha_j$.

Case 1: For this case, we can use the hockey-stick identity to show that

$$(A.54) \quad (m+2) \sum_{j=1}^k j(j+1) \cdots (j+m) = k(k+1) \cdots (k+m+1).$$

Thus, if $\alpha_k = Ck(k+1) \cdots (k+m)$ for all k , then

$$(A.55) \quad \frac{\epsilon_{k+1}}{\epsilon_k} = \frac{k(k+1) \cdots (k+m)}{(k+1)(k+2) \cdots (k+m+1)} = \frac{k}{k+m+1} \rightarrow 1 \quad \text{as } k \rightarrow \infty.$$

Case 2: For this case, we use the identity

$$(A.56) \quad \sum_{j=1}^k \mu^{j-1} = \frac{\mu^k - 1}{\mu - 1}.$$

Thus, if $\alpha_k = C\mu^{k-1}$ for all k , then

$$(A.57) \quad \frac{\epsilon_{k+1}}{\epsilon_k} = \frac{\mu^k - 1}{\mu^{k+1} - 1} \rightarrow \frac{1}{\mu} \quad \text{as } k \rightarrow \infty.$$

Case 3: For this case, we use the identity

$$(A.58) \quad \sum_{j=1}^k jj! = (k+1)! - 1.$$

Thus, if $\alpha_k = C(k-1)(k-1)!$ for all $k \geq 2$ and $\alpha_1 = C$, then

$$(A.59) \quad \frac{\epsilon_{k+1}}{\epsilon_k} = \frac{k!}{(k+1)!} = \frac{1}{k+1} \rightarrow 0 \text{ as } k \rightarrow \infty.$$

Case 4: In this case, we use the fact that $\sum_{j=1}^k \alpha_{j-1} = \alpha_{k-1}$ is a telescoping sum by design. Thus, if $\alpha_{k+1} = r^{1/(q-1)}\mu^{q^k} - \alpha_k$ for all $k \geq 1$ and $\alpha_1 = r^{1/(q-1)}\mu$, then

$$(A.60) \quad \frac{\epsilon_{k+1}}{\epsilon_k^q} = \frac{r^{q/(q-1)}(\mu^{q^{k-1}})^q}{r^{1/(q-1)}\mu^{q^k}} = r \text{ for all } k \geq 1. \quad \square$$

A.6. The entropic Poisson equation in the zero-temperature limit. We close this section by showing that the one-parameter family of solutions to the entropic Poisson equation with “temperature” $\theta = \alpha^{-1}$ converge strongly (in $H^1(\Omega)$) to the solution of the obstacle problem as $\theta \rightarrow 0$.

THEOREM A.15. *Assume $\Omega \subset \mathbb{R}^n$ is an open, bounded Lipschitz domain, $n \geq 1$, and let $g \in H^1(\Omega) \cap C(\bar{\Omega})$ such that $\min g|_{\partial\Omega} > 0$. Let $u_\theta \in H_g^1(\Omega) \cap \text{int } L_+^\infty(\Omega)$ denote the solution of the entropic Poisson equation,*

$$(A.61) \quad (\nabla u_\theta, \nabla w) + \theta(\ln u_\theta, w) = (f, v) \text{ for all } w \in H_0^1(\Omega),$$

and let

$$(A.62) \quad u^* = \arg \min_{u \in H^1(\Omega)} E(u) \text{ subject to } u \geq 0 \text{ in } \Omega \text{ and } u = g \text{ on } \partial\Omega.$$

Then $u_\theta \rightarrow u^*$ in $H^1(\Omega)$ linearly with respect to θ . In particular,

$$(A.63) \quad \frac{1}{2} \|\nabla u^* - \nabla u_\theta\|_{L^2(\Omega)}^2 \leq \theta(S(u^*) + |\Omega|).$$

Proof. For all functions $u \in L_+^\infty(\Omega)$, $v, w \in \text{int } L_+^\infty(\Omega)$, the representation theorem, [Theorem 4.1](#), and the three-points identity [\(4.20\)](#), together, give us

$$(A.64) \quad D(u, v) - D(u, w) + D(v, w) = (\nabla S(v) - \nabla S(w), v - u).$$

Moreover, the characterization theorem, [Theorem 4.7](#), tell us that [\(A.61\)](#) is equivalent to

$$(A.65) \quad S'(u_\theta) = -\frac{1}{\theta} E'(u_\theta).$$

Next, notice that $\nabla S(1) = \ln 1 = 0$ and, provided $u \in H^1(\Omega)$,

$$(A.66) \quad \langle E'(u_\theta), u - u_\theta \rangle \leq E(u) - E(u_\theta),$$

by convexity. Thus, taking $u \in H_g^1(\Omega) \cap L_+^\infty(\Omega)$ and setting $v = u_\theta$ and $w = 1$ in [\(A.64\)](#) leads to

$$\begin{aligned} D(u, u_\theta) - D(u, 1) + D(u_\theta, 1) &= (\nabla S(u_\theta) - \nabla S(1), u_\theta - u) \\ &= \langle S'(u_\theta), u_\theta - u \rangle \\ &= \frac{1}{\theta} \langle E'(u_\theta), u - u_\theta \rangle \end{aligned}$$

$$\leq \frac{1}{\theta}(E(u) - E(u_\theta)).$$

Rerranging the inequality above and invoking [Proposition 4.6](#), we find that

$$(A.67) \quad E(u_\theta) - E(u) \leq \theta(D(u, 1) - D(u, u_\theta) - D(u_\theta, 1)) \leq \theta D(u, 1),$$

where the second inequality arises because both $D(u, u_\theta)$ and $D(u_\theta, 1)$ are non-negative. The proof is completed by setting $u = u^*$ and exploiting the strong convexity of E , as done in [\(A.49\)](#), and noting that $D(u, 1) = S(u) + |\Omega|$. \square

Appendix B. Mathematical results II: Elementary finite element error analysis. The purpose of this appendix is to establish certain minor *a priori* error analysis results related to the stability of the linearized subproblems encountered in [Algorithm 4](#). These results are necessary to motivate the finite elements proposed in [\(4.48a\)](#) and [\(4.48b\)](#). In short, the main outcome of this appendix is that the finite elements are stable and, therefore, we expect optimal high-order convergence rates for the solutions of certain *linearized* subproblems. We intentionally stop short of providing a full *a priori* error analysis of the *nonlinear* subproblems or the complete proximal Galerkin method. Such analysis is planned for a forthcoming paper in which we additionally aim to prove that the Proximal Galerkin method is *mesh-independent*.

B.1. Stability of the linearized subproblems. The first result of this section is that the linearized subproblems in [Algorithm 4](#) are stable at the continuous level. In particular, they are uniformly stable with respect to $H^1(\Omega)$ and $H^{-1}(\Omega)$ when we replace V_h and W_h by $H_0^1(\Omega)$ and $L^2(\Omega)$, respectively. The proof uses standard Hilbert space arguments for saddle-point problems; cf. [\[33, Section 4.3\]](#).

THEOREM B.1. *Let $\psi \in L^\infty(\Omega)$. Then, for every $f \in L^2(\Omega)$ and $\phi \in H^1(\Omega)$, the saddle-point problem*

$$(B.1) \quad \begin{cases} (\nabla u, \nabla v) + (\delta, v) = (f, v) & \text{for all } v \in H_0^1(\Omega), \\ (u, \varphi) - (\delta \exp \psi, \varphi) = (\phi, \varphi) & \text{for all } \varphi \in L^2(\Omega), \end{cases}$$

has a unique solution $u \in H_0^1(\Omega)$, $\delta \in L^2(\Omega)$ that satisfies

$$(B.2a) \quad \|\nabla u\|_{L^2(\Omega)} \leq \|f\|_{H^{-1}(\Omega)} + \|\nabla \phi\|_{L^2(\Omega)},$$

and

$$(B.2b) \quad \|\delta\|_{H^{-1}(\Omega)} \leq 2\|f\|_{H^{-1}(\Omega)} + \|\nabla \phi\|_{L^2(\Omega)}.$$

Remark B.2 (Choice of norms). Notice that [\(B.2b\)](#) is independent of ψ . The corresponding bound of the $L^2(\Omega)$ -norm of ψ degenerates as $\text{ess inf } \psi \rightarrow -\infty$; cf. [\(B.4\)](#). Thus, we choose to interpret [\(B.1\)](#) as a singularly-perturbed saddle-point problem and focus mainly on convergence of the incremental latent variable ψ in the $H^{-1}(\Omega)$ norm.

Proof of Theorem B.1. Existence and uniqueness of $u \in H_0^1(\Omega)$, $\delta \in L^2(\Omega)$ follows readily from the Lax–Milgram theorem. Indeed, notice that [\(B.1\)](#) may be rewritten as

$$(B.3) \quad B((u, \delta), (v, \varphi)) = (f, v) - (\phi, \varphi) \quad \text{for all } v \in H_0^1(\Omega) \text{ and } \varphi \in L^2(\Omega),$$

where $B((u, \delta), (v, \varphi)) = (\nabla u, \nabla v) + (\delta, v) - (u, \varphi) + (\delta \exp \psi, \varphi)$. Moreover, it is a straightforward exercise to check that

$$(B.4) \quad \|\nabla u\|_{L^2(\Omega)}^2 + \|\exp(-\psi)\|_{L^\infty(\Omega)}^{-1} \|\delta\|_{L^2(\Omega)}^2 \leq B((u, \delta), (u, \delta)),$$

for all $(u, \delta) \in H_0^1(\Omega) \times L^2(\Omega)$, which establishes the coercivity condition necessary to apply the theorem.

We now turn to proving (B.2). From the first equation in (B.1), notice that

$$\begin{aligned} \|\delta\|_{H^{-1}(\Omega)} &= \sup_{v \in H_0^1(\Omega)} \frac{(\delta, v)}{\|\nabla v\|_{L^2(\Omega)}} \leq \sup_{v \in H_0^1(\Omega)} \frac{(\nabla u, \nabla v)}{\|\nabla v\|_{L^2(\Omega)}} + \sup_{v \in H_0^1(\Omega)} \frac{(f, v)}{\|\nabla v\|_{L^2(\Omega)}} \\ &= \|\nabla u\|_{L^2(\Omega)} + \|f\|_{H^{-1}(\Omega)}. \end{aligned}$$

Thus, we may focus the remainder of the proof on controlling $\|\nabla u\|_{L^2(\Omega)}$. The remaining arguments center on the equation

$$(B.5) \quad \|\nabla u\|_{L^2(\Omega)}^2 + (\delta \exp \psi, \delta) = (f, u) - (\phi, \delta),$$

which follows from setting $v = u$ and $\varphi = \delta$ in (B.3). We now consider the cases $\phi = 0$ and $f = 0$ separately.

Case 1: $\phi = 0$. It is straightforward to see that

$$(B.6) \quad \|\nabla u\|_{L^2(\Omega)}^2 \leq (\nabla u, \nabla u) + (\delta \exp \psi, \delta) = (f, u) \leq \|f\|_{H^{-1}(\Omega)} \|\nabla u\|_{L^2(\Omega)}.$$

Thus, $\|\nabla u\|_{L^2(\Omega)} \leq \|f\|_{H^{-1}(\Omega)}$.

Case 2: $f = 0$. Notice that

$$(B.7) \quad \|\nabla u\|_{L^2(\Omega)}^2 \leq (\nabla u, \nabla u) + (\delta \exp \psi, \delta) \leq \|\nabla \phi\|_{L^2(\Omega)} \|\delta\|_{H^{-1}(\Omega)}.$$

Moreover, observe that

$$(B.8) \quad \|\nabla u\|_{L^2(\Omega)} = \sup_{v \in H_0^1(\Omega)} \frac{(\nabla u, \nabla v)}{\|\nabla v\|_{L^2(\Omega)}} = \sup_{v \in H_0^1(\Omega)} \frac{(\delta, v)}{\|\nabla v\|_{L^2(\Omega)}} = \|\delta\|_{H^{-1}(\Omega)},$$

where the second equality follows from the first equation in (B.1). Therefore,

$$(B.9) \quad \|\delta\|_{H^{-1}(\Omega)} = \|\nabla u\|_{L^2(\Omega)} \leq \|\nabla \phi\|_{L^2(\Omega)}.$$

Collecting the inequalities above leads to (B.2). \square

Our next result is that the finite elements (4.48) are uniformly stable in $H^1(\Omega) \times H^{-1}(\Omega)$; i.e., they satisfy the Ladyzhenskaya–Babuška–Brezzi (LBB) stability condition

$$(B.10) \quad \beta_h := \inf_{\varphi \in W_h} \sup_{v \in V_h} \frac{(\varphi, v)}{\|\varphi\|_{H^{-1}(\Omega)} \|\nabla v\|_{L^2(\Omega)}} > 0,$$

and, furthermore, β_h is strictly bounded away from zero for all mesh sizes $h > 0$ and (clearly) independent of ψ . Note that here and throughout, we typically treat the symbol $C > 0$ as a generic mesh-independent constant.

LEMMA B.3. *Assume that \mathcal{T}_h is a shape-regular sequence of affine meshes covering $\bar{\Omega} = \bigcup_{T \in \mathcal{T}_h} T$. Let V_h and W_h be the finite element spaces defined in (4.48). Then there is a constant β_0 such that for all $h > 0$,*

$$(B.11) \quad \inf_{\varphi \in W_h} \sup_{v \in V_h} \frac{(\varphi, v)}{\|\varphi\|_{H^{-1}(\Omega)} \|\nabla v\|_{L^2(\Omega)}} \geq \beta_0 > 0.$$

Remark B.4 (Idea of the proof). The proof proceeds independently for each finite element pairing by constructing a so-called Fortin operator $\Pi_h: H_0^1(\Omega) \rightarrow V_h$ satisfying

$$(B.12a) \quad \|\Pi_h v\|_{H^1(\Omega)} \leq C \|v\|_{H^1(\Omega)},$$

for some h -independent constant $C > 0$, and

$$(B.12b) \quad (\Pi_h v, \varphi) = (v, \varphi) \quad \text{for all } v \in H_0^1(\Omega) \text{ and } \varphi \in W_h.$$

It is well-known that the existence of such an operator on a fixed mesh \mathcal{T}_h implies the LBB stability condition (B.10); see, e.g., [71] and [33, Proposition 5.4.3]. See also [66, Theorem 1] for the converse. Likewise, h -independence of the constant C in (B.12a) implies the existence of the uniform discrete stability constant β_0 in (B.11).

We employ a standard technique to construct our Fortin operators that involves splitting the operator into two terms; cf. [33, Section 5.4.4]. In particular, for each pair of subspaces (V_h, W_h) , we define $\Pi_h = \tilde{\mathcal{I}}_h + \tilde{\Pi}_h(I - \tilde{\mathcal{I}}_h)$, where $\tilde{\mathcal{I}}_h: H_0^1(\Omega) \rightarrow V_h$ is a quasi-interpolation operator (see, e.g., [68, Section 22.4]) and $\tilde{\Pi}_h: L^2(\Omega) \rightarrow V_h$ is a linear operator satisfying

$$(B.12c) \quad (\tilde{\Pi}_h v, \varphi) = (v, \varphi) \quad \text{for all } v \in L^2(\Omega) \text{ and } \varphi \in W_h,$$

and $\|\tilde{\Pi}_h(I - \tilde{\mathcal{I}}_h)v\|_{H^1(\Omega)} \leq C \|v\|_{H^1(\Omega)}$ for all $v \in H_0^1(\Omega)$. It is easy to check that such an operator satisfies (B.12a) and (B.12b). Indeed,

$$(B.13) \quad \|\Pi_h v\|_{H^1(\Omega)} \leq \|\tilde{\mathcal{I}}_h v\|_{H^1(\Omega)} + \|\tilde{\Pi}_h(I - \tilde{\mathcal{I}}_h)v\|_{H^1(\Omega)} \leq C \|v\|_{H^1(\Omega)},$$

and, moreover,

$$(B.14) \quad (\Pi_h v, \varphi) = (\tilde{\mathcal{I}}_h v, \varphi) + (\tilde{\Pi}_h(I - \tilde{\mathcal{I}}_h)v, \varphi) = (\tilde{\mathcal{I}}_h v, \varphi) + ((I - \tilde{\mathcal{I}}_h)v, \varphi) = (v, \varphi),$$

where the second equality follows from (B.12c).

Proof of Lemma B.3. Case 1. We first consider the $(\mathbb{P}_p$ -bubble, \mathbb{P}_{p-1} -broken) finite elements defined in (4.48a). In this setting, we define $\tilde{\Pi}_h$ satisfying (B.12c) element-wise by solving the following local variational problem at each element $T \in \mathcal{T}_h$:

$$(B.15) \quad \begin{cases} \text{Find } (\tilde{\Pi}_h v)|_T := v_T \in \mathring{\mathbb{P}}_{p+2}(T) \text{ such that} \\ (v_T, \varphi)_T = (v, \varphi)_T \quad \text{for all } \varphi \in \mathbb{P}_{p-1}(T). \end{cases}$$

Notice that $|\mathring{\mathbb{P}}_{p+2}(T)| = p(p+1)/2 = |\mathbb{P}_{p-1}(T)|$ and, in turn, each function $v_T \in \mathring{\mathbb{P}}_{p+2}(T)$ is well-defined. It is straightforward to see that the zero extension of v_T — i.e., the function $\bar{v}_T = v_T$ on T and 0 otherwise — is a member of

$$(B.16) \quad V_h = \{v \in H_0^1(\Omega) \mid v \in \widehat{\mathbb{P}}_p(T) \oplus \mathring{\mathbb{P}}_{p+2}(T) \text{ for all } T \in \mathcal{T}_h\}.$$

Thus, we conclude that the operator $\tilde{\Pi}_h: H_0^1(\Omega) \rightarrow V_h$, given by $v \mapsto \sum_{T \in \mathcal{T}_h} \bar{v}_T$, is well-defined.

We now analyze $\tilde{\Pi}_h$ more closely element-wise. To this end, let $\{b_i\}$ be a basis for $\mathring{\mathbb{P}}_{p+2}(T)$ and $\{\varphi_i\}$ be a basis for $\mathbb{P}_{p-1}(T)$. Upon writing $v_T = \sum_{j=1}^{p(p+1)/2} c_j b_j$, we see that the variational problem (B.15) is equivalent to the invertible linear system

$$(B.17) \quad M_{ij} c_j = d_i, \quad i = 1, 2, \dots, p(p+1)/2,$$

where $\mathbf{M}_{ij} = (b_j, \varphi_i)_T$, and $\mathbf{d}_i = (v, \varphi_i)_T$. Standard scaling arguments (see, e.g., [68, Proposition 28.5]) can be used to show that $\|\mathbf{M}^{-1}\|_{\ell^2} \leq C|T|^{-1}$. Meanwhile, Hölder's inequality can be used to show that

$$(B.18) \quad \mathbf{d}_i = \int_T v \varphi_i \, dx \leq \max_i \|\varphi_i\|_{L^\infty(T)} |T|^{1/2} \|v\|_{L^2(T)},$$

for each $i = 1, 2, \dots, p(p+1)/2$. Thus, we conclude that

$$(B.19) \quad \|\mathbf{c}\|_{\ell^2} \leq \|\mathbf{M}^{-1}\|_{\ell^2} \|\mathbf{d}\|_{\ell^2} \leq C|T|^{-1/2} \|v\|_{L^2(T)}.$$

A similar scaling argument (see, e.g., [68, Lemma 11.7]) implies that

$$(B.20) \quad |b_i|_{H^1(T)} \leq Ch_T^{-1} |T|^{1/2}.$$

Combining (B.19) and (B.20), we find that

$$(B.21) \quad |\tilde{\Pi}_h v|_{H^1(T)} = |v_T|_{H^1(T)} \leq Ch_T^{-1} \|v\|_{L^2(T)}.$$

The next step is to specify the quasi-interpolation operator $\tilde{\mathcal{I}}_h: H_0^1(\Omega) \rightarrow V_h$. We choose to use the operator defined in [67, Equation 6.10], which, by [67, Theorem 6.3], has the property that

$$(B.22) \quad \|(I - \tilde{\mathcal{I}}_h)v\|_{L^2(T)} \leq Ch_T \|\nabla v\|_{L^2(\Omega_T)} \quad \text{for all } v \in H_0^1(\Omega),$$

where $\Omega_T \subset \Omega$ is the union of mesh cells neighboring T . We now find that

$$(B.23) \quad |\tilde{\Pi}_h(I - \tilde{\mathcal{I}}_h)v|_{H^1(T)} \leq Ch_T^{-1} \|(I - \tilde{\mathcal{I}}_h)v\|_{L^2(T)} \leq C \|\nabla v\|_{L^2(\Omega_T)}.$$

Note that the maximum number of elements in Ω_T is bounded uniformly in h owing to the regularity of the mesh sequence. Likewise, we find that

$$(B.24) \quad |\tilde{\Pi}_h(I - \tilde{\mathcal{I}}_h)v|_{H^1(\Omega)}^2 \leq C \sum_{T \in \mathcal{T}} \|\nabla v\|_{L^2(\Omega_T)}^2 \leq C \|\nabla v\|_{L^2(\Omega)}^2.$$

We have succeeded in checking the conditions outlined in Remark B.4 and the proof is complete.

Case 2. We now consider the $(\mathbb{Q}_p$ -bubble, \mathbb{Q}_{p-1} -broken) finite elements defined in (4.48b). In this case, we define $\tilde{\Pi}_h v := v_h \in V_h$ element-wise by solving the local variational problems

$$(B.25) \quad \begin{cases} \text{Find } v_h|_T \in \mathring{\mathbb{Q}}_{p+1}(T) \text{ such that} \\ (v_h, \varphi)_T = (v, \varphi)_T \quad \text{for all } \varphi \in \mathbb{Q}_{p-1}(T). \end{cases}$$

Here, we notice that $|\mathring{\mathbb{Q}}_{p+1}(T)| = p^2 = |\mathbb{Q}_{p-1}(T)|$. In particular, for every element $T \in \mathcal{T}_h$, the variational problem (B.25) is equivalent to an invertible $p^2 \times p^2$ linear system, and so $\tilde{\Pi}_h v := v_h \in V_h$ is well-posed. The remainder of the proof proceeds as done in *Case 1*. \square

Remark B.5 (Alternative subspaces). Notice that Lemma B.3 implies that the pairing (\tilde{V}_h, W_h) is uniformly stable for any subspace $\tilde{V}_h \subset H_0^1(\Omega)$ containing V_h . Indeed, observe that

$$(B.26) \quad \sup_{v \in \tilde{V}_h} \frac{(\varphi, v)}{\|\nabla v\|_{L^2(\Omega)}} \geq \sup_{v \in V_h} \frac{(\varphi, v)}{\|\nabla v\|_{L^2(\Omega)}} \geq \beta_0 \|\varphi\|_{H^{-1}(\Omega)},$$

for all $\varphi \in W_h$. Thus, other elements such as the $(\mathbb{Q}_{p+1}, \mathbb{Q}_{p-1}$ -broken) pair proposed in Remark 4.21, are also stable owing to the embedding $\mathbb{Q}_p^{p+1}(T) = \hat{\mathbb{Q}}_p(T) \oplus \mathring{\mathbb{Q}}_{p+1}(T) \subset \mathbb{Q}_{p+1}(T)$.

B.2. Convergence of the linearized subproblems. This subsection is devoted to a proof that the $(\mathbb{P}_p\text{-bubble}, \mathbb{P}_{p-1}\text{-broken})$ and $(\mathbb{Q}_p\text{-bubble}, \mathbb{Q}_{p-1}\text{-broken})$ elements defined in (4.48) converge optimally toward the solutions of the linearized subproblems (B.1).

THEOREM B.6. *Let u_h and δ_h to be the discrete solutions of the saddle-point problem*

$$(B.27) \quad \begin{cases} \text{Find } u_h \in V_h \text{ and } \delta_h \in W_h \text{ such that} \\ (\nabla u_h, \nabla v) + (\delta_h, v) = (f, v) \quad \text{for all } v \in V_h, \\ (u_h, w) - (\delta_h \exp \psi, w) = (\phi, w) \quad \text{for all } w \in W_h, \end{cases}$$

where V_h and W_h are the $(\mathbb{P}_p\text{-bubble}, \mathbb{P}_{p-1}\text{-broken})$ and $(\mathbb{Q}_p\text{-bubble}, \mathbb{Q}_{p-1}\text{-broken})$ finite element spaces defined in (4.48). Likewise, let $r \geq 1$ be an integer and assume that the unique solutions of the continuous-level variational problem

$$(B.28) \quad \begin{cases} \text{Find } u \in H_0^1(\Omega) \text{ and } \delta \in L^2(\Omega) \text{ such that} \\ (\nabla u, \nabla v) + (\delta, v) = (f, v) \quad \text{for all } v \in H_0^1(\Omega), \\ (u, w) - (\delta \exp \psi, w) = (\phi, w) \quad \text{for all } w \in L^2(\Omega), \end{cases}$$

are sufficiently regular that $u \in H^{r+1}(\Omega)$ and $\delta \in H^{r-1}(\Omega)$. Then, if \mathcal{T}_h is a shape-regular sequence of affine meshes, it holds that

$$(B.29) \quad \|u - u_h\|_{H^1(\Omega)} + \|\delta - \delta_h\|_{H^{-1}(\Omega)} \leq C_1 h^s (|u|_{H^{s+1}(\Omega)} + |\delta|_{H^{s-1}(\Omega)}),$$

for all $1 \leq s \leq \min\{r, p\}$, where C_1 is a mesh-independent constant that remains bounded as $\text{ess inf } \psi \rightarrow -\infty$. Moreover, if in addition $\delta \in H^r(\Omega)$, then there exists a mesh-independent constant C_2 such that

$$(B.30) \quad \|\delta - \delta_h\|_{L^2(\Omega)} \leq C_2 h^s (|u|_{H^{s+1}(\Omega)} + |\delta|_{H^s(\Omega)}),$$

for each $1 \leq s \leq \min\{r, p\}$ above. However, $C_2 \rightarrow \infty$ as $\text{ess inf } \psi \rightarrow -\infty$.

Proof. Most elements of the proof are standard, so we only give a sketch in the case of the $(\mathbb{P}_p\text{-bubble}, \mathbb{P}_{p-1}\text{-broken})$ finite elements defined in (4.48a). We refer the interested reader to [33, Section 5.5] for further details. Let $a: H_0^1(\Omega) \times H_0^1(\Omega) \rightarrow \mathbb{R}$ denote the bilinear form $a(u, v) = (\nabla u, \nabla v)$. Likewise, define denote $b(u, \varphi) = (u, \varphi)$ and $c(\delta, \varphi) = (\delta \exp \psi, \varphi)$ for all $u \in H_0^1(\Omega)$ and $\delta, \varphi \in L^2(\Omega)$. Due to the stability of the linearized discrete subproblem (B.27) given to us by Lemma B.3 and coercivity of a , we appeal to the standard theory of mixed methods to arrive at the following *a priori* estimate:

$$(B.31) \quad \|u - u_h\|_{H^1(\Omega)} + \|\delta - \delta_h\|_{H^{-1}(\Omega)} \leq C \left(\inf_{v \in V_h} \|u - v\|_{H^1(\Omega)} + \inf_{\varphi \in W_h} \|\delta - \varphi\|_{H^{-1}(\Omega)} \right).$$

For any Banach spaces V and W , let

$$(B.32) \quad \|d\| = \sup_{v \in V} \sup_{w \in W} \frac{|d(v, w)|}{\|v\|_V \|w\|_W}$$

denote the norm of a continuous bilinear form $d: V \times W \rightarrow \mathbb{R}$. A classical analysis given in, e.g., [33, Section 4.3], shows that C in (B.31) depends at most on $\|a\|$, $\|b\|$,

$\|c\|$, the coercivity constant for a , and β_0 , and that C remains bounded as $\|c\| \rightarrow 0$. It is straightforward to check that $\|a\|, \|b\| \leq 1$ and $\|c\| = \|\exp \psi\|_{L^\infty(\Omega)}$. Finally, we note that the coercivity constant for a depends only on the Poincaré constant of the domain Ω and β_0 is naturally independent of ψ .

We now bound the two terms on the right-hand side of (B.31). Given that $\mathbb{P}_p(\mathcal{T}_h) \cap H_0^1(\Omega) \subset V_h$, we may consider the standard order- p nodal interpolation operator $\mathcal{I}_h: H_0^{r+1}(\Omega) \rightarrow \mathbb{P}_p(\mathcal{T}_h) \cap H_0^1(\Omega)$; see, e.g., [68, Section 19.3]. By shape-regularity of the mesh sequence and [68, Corollary 19.8], we have that

$$(B.33) \quad \|v - \mathcal{I}_h v\|_{H^1(\Omega)} \leq Ch^s |v|_{H^{s+1}(\Omega)}, \quad 1 \leq s \leq \min\{r, p\}.$$

Thus, we find that

$$(B.34) \quad \inf_{v \in V_h} \|u - v\|_{H^1(\Omega)} \leq \|u - \mathcal{I}_h u\|_{H^1(\Omega)} \leq ch^s |u|_{H^{s+1}}.$$

The second term on the right-hand side of (B.31) is treated with the order- $(p-1)$ $L^2(\Omega)$ -orthogonal projection operator $\mathcal{P}_h: L^2(\Omega) \rightarrow W_h = \mathbb{P}_{p-1}(\mathcal{T}_h)$, which has the well-known property

$$(B.35) \quad \|\varphi - \mathcal{P}_h \varphi\|_{L^2(\Omega)} \leq Ch^t |\varphi|_{H^t(\Omega)}, \quad 0 \leq t \leq \min\{r-1, p\}.$$

For further details, see, e.g., [68, Section 18.4]. We now note that $\delta - \mathcal{P}_h \delta \in L^2(\Omega)$ and

$$(B.36) \quad (\delta - \mathcal{P}_h \delta, \varphi_h) = 0 \quad \text{for all } \varphi_h \in W_h.$$

Therefore, for any $\varphi_h \in W_h$, we may write

$$\|\delta - \mathcal{P}_h \delta\|_{H^{-1}(\Omega)} = \sup_{\varphi \in H_0^1(\Omega)} \frac{(\delta - \mathcal{P}_h \delta, \varphi)}{\|\nabla \varphi\|_{L^2(\Omega)}} = \sup_{\varphi \in H_0^1(\Omega)} \frac{(\delta - \mathcal{P}_h \delta, \varphi - \varphi_h)}{\|\nabla \varphi\|_{L^2(\Omega)}}.$$

Taking $\varphi_h = \mathcal{P}_h \varphi$ and invoking (B.35), we deduce that

$$\begin{aligned} \|\delta - \mathcal{P}_h \delta\|_{H^{-1}(\Omega)} &\leq \|\delta - \mathcal{P}_h \delta\|_{L^2(\Omega)} \sup_{\varphi \in H_0^1(\Omega)} \frac{\|\varphi - \mathcal{P}_h \varphi\|_{L^2(\Omega)}}{\|\nabla \varphi\|_{L^2(\Omega)}} \\ &\leq Ch^s |\delta|_{H^{s-1}(\Omega)}, \end{aligned}$$

since $\|\delta - \mathcal{P}_h \delta\|_{L^2(\Omega)} \leq Ch^{s-1} |\delta|_{H^{s-1}(\Omega)}$ and $\|\varphi - \mathcal{P}_h \varphi\|_{L^2(\Omega)} \leq Ch \|\nabla \varphi\|_{L^2(\Omega)}$. Inequality (B.29) now follows by collecting the above bounds.

To prove (B.30), we appeal to Cea's lemma applied to the bilinear form B in (B.3). In particular, notice that the coercivity constant in (B.4) vanishes as $\text{ess inf } \psi \rightarrow -\infty$. Therefore,

$$(B.37) \quad \|\delta - \delta_h\|_{L^2(\Omega)} \leq C \left(\inf_{v \in V_h} \|u - v\|_{H^1(\Omega)} + \inf_{\varphi \in W_h} \|\delta - \varphi\|_{L^2(\Omega)} \right),$$

where the stability constant blows up, i.e., $C \rightarrow \infty$, as $\text{ess inf } \psi \rightarrow -\infty$. The remainder of the inequality (B.30) is determined by invoking (B.34) and noting that $\|\delta - \mathcal{P}_h \delta\|_{L^2(\Omega)} \leq Ch^s |\delta|_{H^s(\Omega)}$ by allowing $0 \leq t \leq \min\{r, p\}$ in (B.35) because of the additional regularity assumed on δ . \square

B.3. Nonlinear approximability. We finish this appendix with a proof of [Proposition 4.19](#).

Proof of Proposition 4.19. Recall that $u = \exp \psi$ and, therefore,

$$(B.38) \quad \tilde{u}_h - u = \exp \psi_h - \exp \psi = \int_{\psi}^{\psi_h} \exp s \, ds$$

$$(B.39) \quad = (\psi_h - \psi) \int_0^1 \exp(\psi + s(\psi_h - \psi)) \, ds$$

$$(B.40) \quad = (\psi_h - \psi) \exp \psi \int_0^1 (\exp(\psi_h - \psi))^s \, ds.$$

As such, it holds that

$$(B.41) \quad \|u - \tilde{u}_h\|_{L^\infty} \leq \|\psi - \psi_h\|_{L^\infty} \|\exp \psi\|_{L^\infty} \int_0^1 (\exp \|\psi - \psi_h\|_{L^\infty})^s \, ds$$

$$(B.42) \quad = \|\exp \psi\|_{L^\infty} (\exp \|\psi - \psi_h\|_{L^\infty} - 1),$$

where the last line follows from the identity $\int_0^1 a^s \, ds = (a - 1) / \ln a$. \square

Extended dedication from B. Keith. *Feynman once said that calculus is “the language God talks” [206]. Expanding on this mystical statement, Strogatz has suggested that all physical laws are “sentences” in this “language of the universe” [186]. From my perspective, if the above is true, it must be the case that God has written his sentences in variational form.*

The present work deals centrally with variational methods and a seemingly divine entropy functional that has never failed to surprise me since the day I began this project with Thomas. In turn, I have frequently been reminded of von Neumann’s famous quote to Shannon: “No one knows what entropy really is” [196]. Reflecting back on Feynman and Strogatz’s perspectives, it is helpful to think that at least God knows, even if we do not. As such, and on the occasion of his 70th birthday, it feels only fitting that I dedicate this work to Leszek Demkowicz; the kind and deeply religious man who not so long ago taught me a variational perspective on the universe.

Acknowledgements. The authors gratefully thank Jesse Chan, Jérôme Darbon, Alexandre Ern, Patrick Farrell, Caroline Geiersbach, Dohyun Kim, Tarik Dzanic, Boyan Lazarov, Michael Hintermüller, Thomas J.R. Hughes, Karl-Hermann Neeb, Chi-Wang Shu, and N. Sukumar for their time and helpful discussions during the development of this work. We also thank Dohyun Kim, Tzanio Kolev, Boyan Lazarov, Socratis Petrides, and Jingyi Wang for peer-reviewing our MFEM implementations in order to merge them into the source code as official MFEM examples. Finally, we thank Jørgen Dokken for generously responding to our many inquiries about FEniCSx, as well as his help refactoring and preparing our FEniCSx implementations for public release.

REFERENCES

- [1] R. ABGRALL, P. ÖFFNER, AND H. RANOCHA, *Reinterpretation and extension of entropy correction terms for residual distribution and discontinuous Galerkin schemes: Application to structure preserving discretization*, *Journal of Computational Physics*, 453 (2022), p. 110955.

- [2] R. ABGRALL AND J. TREFILIK, *An example of high order residual distribution scheme using non-Lagrange elements*, Journal of Scientific Computing, 45 (2010), pp. 3–25.
- [3] L. ADAM, M. HINTERMÜLLER, AND T. M. SUROWIEC, *A semismooth Newton method with analytical path-following for the H^1 -projection onto the Gibbs simplex*, IMA Journal of Numerical Analysis, 39 (2018), pp. 1276–1295.
- [4] R. A. ADAMS AND J. J. FOURNIER, *Sobolev spaces*, Elsevier, 2003.
- [5] Z. AHMED, N. LE ROUX, M. NOROUZI, AND D. SCHUURMANS, *Understanding the impact of entropy on policy optimization*, in International conference on machine learning, PMLR, 2019, pp. 151–160.
- [6] M. AINSWORTH, G. ANDRIAMARO, AND O. DAVYDOV, *Bernstein–Bézier finite elements of arbitrary order and optimal assembly procedures*, SIAM Journal on Scientific Computing, 33 (2011), pp. 3087–3109.
- [7] L. ALLEN AND R. C. KIRBY, *Bounds-constrained polynomial approximation using the Bernstein basis*, Numerische Mathematik, 152 (2022), pp. 101–126.
- [8] E. L. ALLGOWER AND K. BÖHMER, *Application of the mesh independence principle to mesh refinement strategies*, SIAM Journal on Numerical Analysis, 24 (1987), pp. 1335–1351.
- [9] E. L. ALLGOWER, K. BOHMER, F. A. POTRA, AND W. C. RHEINBOLDT, *A mesh-independence principle for operator equations and their discretizations*, SIAM Journal on Numerical Analysis, 23 (1986), pp. 160–169.
- [10] S.-I. AMARI, *Information geometry and its applications*, vol. 194, Springer, 2016.
- [11] S.-I. AMARI AND H. NAGAOKA, *Methods of information geometry*, vol. 191, American Mathematical Soc., 2000.
- [12] A. AMBROSETTI AND G. PRODI, *A primer of nonlinear analysis*, 34, Cambridge University Press, 1995.
- [13] L. AMBROSIO AND V. M. TORTORELLI, *Approximation of functionals depending on jumps by elliptical functionals via Γ -convergence*, Commun. Pure Appl. Math., 43 (1990), pp. 999–1036.
- [14] R. ANDERSON, J. ANDREJ, A. BARKER, J. BRAMWELL, J.-S. CAMIER, J. CERVENY, V. DOBREV, Y. DUDOUIT, A. FISHER, T. KOLEV, ET AL., *MFEM: A modular finite element methods library*, Computers & Mathematics with Applications, 81 (2021), pp. 42–74.
- [15] R. ANDERSON, V. DOBREV, T. KOLEV, D. KUZMIN, M. Q. DE LUNA, R. RIEBEN, AND V. TOMOV, *High-order local maximum principle preserving (MPP) discontinuous Galerkin finite element method for the transport equation*, Journal of Computational Physics, 334 (2017), pp. 102–124.
- [16] E. ANDREASSEN, A. CLAUSEN, M. SCHEVENELS, B. S. LAZAROV, AND O. SIGMUND, *Efficient topology optimization in MATLAB using 88 lines of code*, Structural and Multidisciplinary Optimization, 43 (2011), pp. 1–16.
- [17] H. ANTIL, D. P. KOURI, AND D. RIDZAL, *ALESQP: An augmented Lagrangian equality-constrained SQP method for optimization with general constraints*, SIAM Journal on Optimization, 33 (2023), pp. 237–266.
- [18] H. ATTOUCH, G. BUTTAZZO, AND G. MICHAILLE, *Variational Analysis in Sobolev and BV Spaces: Applications to PDEs and Optimization*, MOS-SIAM Series on Optimization, Society for Industrial and Applied Mathematics, second ed., 2014.
- [19] S. BALAY, S. ABHYANKAR, M. F. ADAMS, S. BENSON, J. BROWN, P. BRUNE, K. BUSCHELMAN, E. CONSTANTINESCU, L. DALCIN, A. DENER, V. EIJKHOUT, J. FAIBUSSOWITSCH, W. D. GROPP, V. HAPLA, T. ISAAC, P. JOLIVET, D. KARPEEV, D. KAUSHIK, M. G. KNEPLEY, F. KONG, S. KRUGER, D. A. MAY, L. C. MCINNES, R. T. MILLS, L. MITCHELL, T. MUNSON, J. E. ROMAN, K. RUPP, P. SANAN, J. SARICH, B. F. SMITH, S. ZAMPINI, H. ZHANG, H. ZHANG, AND J. ZHANG, *PETSc/TAO users manual*, Tech. Report ANL-21/39 - Revision 3.19, Argonne National Laboratory, 2023, <https://doi.org/10.2172/1968587>.
- [20] G. BARRENECHEA, E. GEORGIOULIS, T. PRYER, AND A. VEESER, *A nodally bound-preserving finite element method*, arXiv preprint arXiv:2304.01067, (2023).
- [21] S. BARTELS, *Numerical methods for nonlinear partial differential equations*, vol. 47, Springer, 2015.
- [22] H. H. BAUSCHKE, J. M. BORWEIN, AND P. L. COMBETTES, *Essential smoothness, essential strict convexity, and Legendre functions in Banach spaces*, Communications in Contemporary Mathematics, 3 (2001), pp. 615–647.
- [23] H. H. BAUSCHKE AND P. L. COMBETTES, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, Springer New York, NY, 2011.
- [24] A. BECK, *First-Order Methods in Optimization*, Society for Industrial and Applied Mathematics, Philadelphia, PA, 2017.
- [25] A. BECK AND M. TEBoulLE, *Mirror descent and nonlinear projected subgradient methods for*

- convex optimization*, Operations Research Letters, 31 (2003), pp. 167–175.
- [26] M. BERGOUNIOUX, *Augmented Lagrangian method for distributed optimal control problems with state constraints*, Journal of Optimization Theory and Applications, 78 (1993), pp. 493–521.
- [27] M. BERGOUNIOUX, M. HADDOU, M. HINTERMÜLLER, AND K. KUNISCH, *A comparison of a Moreau–Yosida-based active set strategy and interior point methods for constrained optimal control problems*, SIAM Journal on Optimization, 11 (2000), pp. 495–521.
- [28] M. BERGOUNIOUX AND K. KUNISCH, *Augmented Lagrangian techniques for elliptic state constrained optimal control problems*, SIAM Journal on Control and Optimization, 35 (1997), pp. 1524–1543.
- [29] L. BERTINI AND G. GIACOMIN, *Stochastic Burgers and KPZ equations from particle systems*, Communications in mathematical physics, 183 (1997), pp. 571–607.
- [30] D. P. BERTSEKAS, *Nonlinear Programming*, Athena Scientific, 1999.
- [31] E. G. BIRGIN AND J. M. MARTÍNEZ, *Practical augmented Lagrangian methods for constrained optimization*, SIAM, 2014.
- [32] P. BOCHEV, D. RIDZAL, M. D’ELIA, M. PEREGO, AND K. PETERSON, *Optimization-based, property-preserving finite element methods for scalar advection equations and their connection to algebraic flux correction*, Computer Methods in Applied Mechanics and Engineering, 367 (2020), p. 112982.
- [33] D. BOFFI, F. BREZZI, AND M. FORTIN, *Mixed finite element methods and applications*, vol. 44, Springer, 2013.
- [34] J. P. BORIS AND D. L. BOOK, *Flux-corrected transport. I. SHASTA, a fluid transport algorithm that works*, Journal of computational physics, 11 (1973), pp. 38–69.
- [35] J. M. BORWEIN AND A. S. LEWIS, *Strong rotundity and optimization*, SIAM Journal on Optimization, 4 (1994), pp. 146–158.
- [36] A. BRANDT AND C. W. CRYER, *Multigrid algorithms for the solution of linear complementarity problems arising from free boundary problems*, SIAM Journal on Scientific and Statistical Computing, 4 (1983), pp. 655–684.
- [37] L. M. BREGMAN, *The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming*, USSR Computational Mathematics and Mathematical Physics, 7 (1967), pp. 200–217.
- [38] H. BRÉZIS, *Functional analysis, Sobolev spaces and partial differential equations*, vol. 2, Springer, 2011.
- [39] F. BREZZI, L. D. MARINI, AND P. PIETRA, *Two-dimensional exponential fitting and applications to drift-diffusion models*, SIAM Journal on Numerical Analysis, 26 (1989), pp. 1342–1355.
- [40] T. E. BRUNS AND D. A. TORTORELLI, *Topology optimization of non-linear elastic structures and compliant mechanisms*, Computer methods in applied mechanics and engineering, 190 (2001), pp. 3443–3459.
- [41] J. A. CARRILLO, A. JÜNGEL, AND S. TANG, *Positive entropic schemes for a nonlinear fourth-order parabolic equation*, Discrete and Continuous Dynamical Systems - B, 3 (2003), pp. 1–20.
- [42] E. CASAS, *Control of an elliptic problem with pointwise state constraints*, SIAM Journal on Control and Optimization, 24 (1986), pp. 1309–1318.
- [43] E. CASAS, *Boundary control of semilinear elliptic equations with pointwise state constraints*, SIAM Journal on Control and Optimization, 31 (1993), pp. 993–1006.
- [44] E. CASAS, *Pontryagin’s principle for state-constrained boundary control problems of semilinear parabolic equations*, SIAM Journal on Control and Optimization, 35 (1997), pp. 1297–1327.
- [45] Y. CENSOR AND S. A. ZENIOS, *Proximal minimization algorithm with d -functions*, Journal of Optimization Theory and Applications, 73 (1992), pp. 451–464.
- [46] J. CHAN, N. HEUER, T. BUI-THANH, AND L. DEMKOWICZ, *A robust DPG method for convection-dominated diffusion problems II: Adjoint boundary conditions and mesh-dependent test norms*, Computers & Mathematics with Applications, 67 (2014), pp. 771–795.
- [47] G. CHEN AND M. TEBoulLE, *Convergence analysis of a proximal-like minimization algorithm using Bregman functions*, SIAM Journal on Optimization, 3 (1993), pp. 538–543.
- [48] P. G. CIARLET, *The finite element method for elliptic problems*, SIAM, 2002.
- [49] B. COCKBURN, J. GOPALAKRISHNAN, AND R. LAZAROV, *Unified hybridization of discontinuous Galerkin, mixed, and continuous Galerkin methods for second order elliptic problems*, SIAM Journal on Numerical Analysis, 47 (2009), pp. 1319–1365.
- [50] J. D. COLE, *On a quasi-linear parabolic equation occurring in aerodynamics*, Quarterly of

- applied mathematics, 9 (1951), pp. 225–236.
- [51] J. B. CONWAY, *A course in functional analysis*, vol. 96, Springer, 2007.
- [52] J. A. COTTRELL, T. J. HUGHES, AND Y. BAZILEVS, *Isogeometric analysis: Toward integration of CAD and FEA*, John Wiley & Sons, 2009.
- [53] M. G. CRANDALL AND A. MAJDA, *Monotone difference approximations for scalar conservation laws*, Mathematics of Computation, 34 (1980), pp. 1–21.
- [54] M. CROUZEIX AND P.-A. RAVIART, *Conforming and nonconforming finite element methods for solving the stationary Stokes equations I*, Revue française d’automatique informatique recherche opérationnelle. Mathématique, 7 (1973), pp. 33–75.
- [55] M. CUTURI, *Sinkhorn distances: Lightspeed computation of optimal transport*, Advances in neural information processing systems, 26 (2013).
- [56] S. DAHLKE AND T. M. SUROWIEC, *Wavelet-based approximations of pointwise bound constraints in Lebesgue and Sobolev spaces*, IMA Journal of Numerical Analysis, 42 (2022), pp. 417–439.
- [57] L. D. DALCIN, R. R. PAZ, P. A. KLER, AND A. COSIMO, *Parallel distributed computing using python*, Advances in Water Resources, 34 (2011), pp. 1124 – 1139, <https://doi.org/10.1016/j.advwatres.2011.04.013>. New Computational Methods and Software Tools.
- [58] J. DARBON AND G. P. LANGLOIS, *Accelerated nonlinear primal-dual hybrid gradient methods with applications to supervised machine learning*, arXiv preprint arXiv:2109.12222, (2021).
- [59] J. DARBON AND G. P. LANGLOIS, *Efficient and robust high-dimensional sparse logistic regression via nonlinear primal-dual hybrid gradient algorithms*, arXiv preprint arXiv:2111.15426, (2021).
- [60] P. DEGROEN AND P. W. HEMKER, *Error bounds for exponentially fitted Galerkin methods applied to stiff two-point boundary value problems*, Numerical Analysis of singular perturbation problems (PW Hemker & JJH Miller eds), (1979), pp. 217–249.
- [61] L. DEMKOWICZ, *Computing with hp-adaptive finite elements: Volume 1. One and two dimensional elliptic and Maxwell problems*, CRC press, 2006.
- [62] S. DUCZEK AND H. GRAVENKAMP, *Mass lumping techniques in the spectral element method: On the equivalence of the row-sum, nodal quadrature, and diagonal scaling methods*, Computer Methods in Applied Mechanics and Engineering, 353 (2019), pp. 516–569.
- [63] T. DZANIC, W. TROJAK, AND F. D. WITHERDEN, *Bounds preserving temporal integration methods for hyperbolic conservation laws*, Computers & Mathematics with Applications, 135 (2023), pp. 6–18.
- [64] T. DZANIC AND F. D. WITHERDEN, *Positivity-preserving entropy-based adaptive filtering for discontinuous spectral element methods*, Journal of Computational Physics, 468 (2022), p. 111501.
- [65] K. ERIKSSON AND C. JOHNSON, *Adaptive streamline diffusion finite element methods for stationary convection-diffusion problems*, mathematics of computation, 60 (1993), pp. 167–188.
- [66] A. ERN AND J.-L. GUERMOND, *A converse to Fortin’s lemma in Banach spaces*, Comptes Rendus Mathématique, 354 (2016), pp. 1092–1095.
- [67] A. ERN AND J.-L. GUERMOND, *Finite element quasi-interpolation and best approximation*, ESAIM: Mathematical Modelling and Numerical Analysis, 51 (2017), pp. 1367–1385.
- [68] A. ERN AND J.-L. GUERMOND, *Finite Elements I*, Texts in Applied Mathematics, Springer, 2021.
- [69] J. A. EVANS, T. J. HUGHES, AND G. SANGALLI, *Enforcement of constraints and maximum principles in the variational multiscale method*, Computer Methods in Applied Mechanics and Engineering, 199 (2009), pp. 61–76.
- [70] R. FATTAL AND R. KUPFERMAN, *Constitutive laws for the matrix-logarithm of the conformation tensor*, Journal of Non-Newtonian Fluid Mechanics, 123 (2004), pp. 281–285.
- [71] M. FORTIN, *An analysis of the convergence of mixed finite element methods*, RAIRO. Analyse numérique, 11 (1977), pp. 341–354.
- [72] I. FRIED AND D. S. MALKUS, *Finite element mass matrix lumping by numerical integration with no convergence rate loss*, International Journal of Solids and Structures, 11 (1975), pp. 461–466.
- [73] G. FU AND Z. XU, *High-order space-time finite element methods for the Poisson–Nernst–Planck equations: Positivity and unconditional energy stability*, Computer Methods in Applied Mechanics and Engineering, 395 (2022), p. 115031.
- [74] F. FUENTES, B. KEITH, L. DEMKOWICZ, AND S. NAGARAJ, *Orientation embedded high order shape functions for the exact sequence elements of all shapes*, Computers & Mathematics with applications, 70 (2015), pp. 353–458.

- [75] H. GLÖCKNER, *Algebras whose groups of units are Lie groups*, *Studia Mathematica*, 153 (2002), pp. 147–177.
- [76] H. GLÖCKNER, *Lie groups of measurable mappings*, *Canadian Journal of Mathematics*, 55 (2003), pp. 969–999.
- [77] R. GLOWINSKI, *Numerical Methods for Nonlinear Variational Problems*, Springer Berlin Heidelberg, 1984.
- [78] S. K. GODUNOV, *Finite difference method for numerical computation of discontinuous solutions of the equations of fluid dynamics*, *Matematičeskij sbornik*, 47 (1959), pp. 271–306.
- [79] J. S. GOLAN, *Semirings and their Applications*, Springer Science & Business Media, 2013.
- [80] H. GOLDBERG, W. KAMPOWSKY, AND F. TRÖLTZSCH, *On Nemytskij operators in L_p -spaces of abstract functions*, *Mathematische Nachrichten*, 155 (1992), pp. 127–140.
- [81] M. GONDRAN AND M. MINOUX, *Graphs, dioids and semirings: New models and algorithms*, vol. 41, Springer Science & Business Media, 2008.
- [82] W. J. GORDON AND C. A. HALL, *Construction of curvilinear co-ordinate systems and applications to mesh generation*, *International Journal for Numerical Methods in Engineering*, 7 (1973), pp. 461–477.
- [83] C. GRÄSER AND R. KORNHUBER, *Multigrid methods for obstacle problems*, *Journal of Computational Mathematics*, 27 (2009), pp. 1–44.
- [84] J.-L. GUERMOND AND B. POPOV, *Invariant domains and first-order continuous finite element approximation for hyperbolic systems*, *SIAM Journal on Numerical Analysis*, 54 (2016), pp. 2466–2489.
- [85] O. GÜLER, *On the convergence of the proximal point algorithm for convex minimization*, *SIAM journal on control and optimization*, 29 (1991), pp. 403–419.
- [86] T. GUSTAFSSON, R. STENBERG, AND J. VIDEMAN, *On finite element formulations for the obstacle problem—mixed and stabilised methods*, *Computational Methods in Applied Mathematics*, 17 (2017), pp. 413–429.
- [87] W. HACKBUSCH AND H. MITTELMANN, *On multi-grid methods for variational inequalities*, *Numerische Mathematik*, 42 (1983), pp. 65–76.
- [88] M. HAIRER, *Solving the KPZ equation*, *Annals of mathematics*, (2013), pp. 559–664.
- [89] A. HARTEN, *High resolution schemes for hyperbolic conservation laws*, *Journal of Computational Physics*, 49 (1983), pp. 357–393.
- [90] A. HARTEN, P. D. LAX, AND B. V. LEER, *On upstream differencing and Godunov-type schemes for hyperbolic conservation laws*, *SIAM review*, 25 (1983), pp. 35–61.
- [91] M. HEINKENSCHLOSS AND D. RIDZAL, *A matrix-free trust-region SQP method for equality constrained optimization*, *SIAM Journal on Optimization*, 24 (2014), pp. 1507–1541.
- [92] M. HINTERMÜLLER AND M. HINZE, *Moreau–Yosida regularization in state constrained elliptic control problems: Error estimates and parameter adjustment*, *SIAM Journal on Numerical Analysis*, 47 (2009), pp. 1666–1683.
- [93] M. HINTERMÜLLER, K. ITO, AND K. KUNISCH, *The primal-dual active set strategy as a semismooth Newton method*, *SIAM Journal on Optimization*, 13 (2002), pp. 865–888.
- [94] M. HINTERMÜLLER AND K. KUNISCH, *Feasible and noninterior path-following in constrained minimization with low multiplier regularity*, *SIAM Journal on Control and Optimization*, 45 (2006), pp. 1198–1221.
- [95] M. HINTERMÜLLER AND K. KUNISCH, *Path-following methods for a class of constrained minimization problems in function space*, *SIAM Journal on Optimization*, 17 (2006), pp. 159–187.
- [96] M. HINZE, R. PINNAU, M. ULBRICH, AND S. ULBRICH, *Optimization with PDE constraints*, vol. 23 of *Mathematical Modelling: Theory and Applications*, Springer, New York, 2009.
- [97] M. HINZE AND A. SCHIELA, *Discretization of interior point methods for state constrained elliptic optimal control problems: Optimal error estimates and parameter adjustment*, *Computational Optimization and Applications*, 48 (2009), pp. 581–600.
- [98] E. HOPF, *The partial differential equation $u_t + uu_x = u_{xx}$* , *Communications on Pure and Applied Mathematics*, 3 (1950), pp. 201–230.
- [99] R. H. W. HOPPE, *Multigrid algorithms for variational inequalities*, *SIAM J. Numer. Anal.*, 24 (1987), pp. 1046–1065.
- [100] R. H. W. HOPPE, *Une méthode multigrille pour la solution des problèmes d’obstacle*, *ESAIM: Mathematical Modelling and Numerical Analysis - Modélisation Mathématique et Analyse Numérique*, 24 (1990), pp. 711–735.
- [101] R. H. W. HOPPE AND R. KORNHUBER, *Adaptive multilevel methods for obstacle problems*, *SIAM Journal on Numerical Analysis*, 31 (1994), pp. 301–323.
- [102] T. J. HUGHES, *The finite element method: Linear static and dynamic finite element analysis*, Courier Corporation, 2012.

- [103] T. J. HUGHES, J. A. COTTRELL, AND Y. BAZILEVS, *Isogeometric analysis: CAD, finite elements, NURBS, exact geometry and mesh refinement*, Computer methods in applied mechanics and engineering, 194 (2005), pp. 4135–4195.
- [104] F. ILINCA, *Méthodes d'éléments finis adaptatives pour les écoulements turbulents*, PhD thesis, École polytechnique de Montréal, 1996.
- [105] F. ILINCA, D. PELLETIER, ET AL., *A unified finite element algorithm for two-equation models of turbulence*, Computers & fluids, 27 (1998), pp. 291–310.
- [106] A. IOFFE AND V. TIHOMIROV, *Theory of Extremal Problems: Theory of Extremal Problems*, ISSN, Elsevier Science, 2009.
- [107] K. ITO AND K. KUNISCH, *The augmented Lagrangian method for equality and inequality constraints in Hilbert spaces*, Mathematical Programming, 46 (1990), pp. 341–360.
- [108] K. ITO AND K. KUNISCH, *An augmented Lagrangian technique for variational inequalities*, Applied Mathematics and Optimization, 21 (1990), pp. 223–241.
- [109] K. ITO AND K. KUNISCH, *Lagrange Multiplier Approach to Variational Problems and Applications*, Society for Industrial and Applied Mathematics, 2008.
- [110] C. KANZOW, D. STECK, AND D. WACHSMUTH, *An augmented Lagrangian method for optimization problems in Banach spaces*, SIAM Journal on Control and Optimization, 56 (2018), pp. 272–291.
- [111] B. KEITH AND T. M. SUROWIEC, *MFEM Example 36: A peer-reviewed MFEM implementation of proximal Galerkin for the obstacle problem*. <https://github.com/mfem/mfem/pull/3398>, 2023.
- [112] B. KEITH AND T. M. SUROWIEC, *MFEM Example 37: A peer-reviewed MFEM implementation of proximal Galerkin for topology optimization*. <https://github.com/mfem/mfem/pull/3400>, 2023.
- [113] B. KEITH, T. M. SUROWIEC, AND J. S. DOKKEN, *Examples for the Proximal Galerkin Method*, July 2023, <https://github.com/thomas-surowiec/proximal-galerkin-examples>.
- [114] M. KEUTHEN AND M. ULBRICH, *Moreau–Yosida regularization in shape optimization with geometric constraints*, Computational Optimization and Applications, 62 (2014), pp. 181–216.
- [115] N. KIKUCHI AND J. T. ODEN, *Contact problems in elasticity: A study of variational inequalities and finite element methods*, SIAM, 1988.
- [116] D. KINDERLEHRER AND G. STAMPACCHIA, *An introduction to variational inequalities and their applications*, SIAM, 2000.
- [117] V. KOLOKOLTSOV AND V. P. MASLOV, *Idempotent analysis and its applications*, vol. 401, Springer Science & Business Media, 1997.
- [118] R. KORNUBER, *Monotone multigrid methods for elliptic variational inequalities I*, Numerische Mathematik, 69 (1994), pp. 167–184.
- [119] R. KORNUBER, *Monotone multigrid methods for elliptic variational inequalities II*, Numerische Mathematik, 72 (1996), pp. 481–499.
- [120] R. KORNUBER AND R. KRAUSE, *Adaptive multigrid methods for Signorini’s problem in linear elasticity*, Computing and Visualization in Science, 4 (2001), pp. 9–20.
- [121] D. P. KOURI AND D. RIDZAL, *Inexact trust-region methods for PDE-constrained optimization*, Frontiers in PDE-Constrained Optimization, (2018), pp. 83–121.
- [122] K. KUNISCH AND X. LU, *Optimal control for an elliptic system with convex polygonal control constraints*, IMA Journal of Numerical Analysis, 33 (2012), pp. 875–897.
- [123] D. KUZMIN, *Positive finite element schemes based on the flux-corrected transport procedure*, Computational Fluid and Solid Mechanics, Elsevier, (2001), pp. 887–888.
- [124] D. KUZMIN, R. LÖHNER, AND S. TUREK, *Flux-corrected transport: Principles, algorithms, and applications*, Springer, 2012.
- [125] G. LAN, *First-order and Stochastic Optimization Methods for Machine Learning*, Springer Cham, 2020.
- [126] M. LANDAJUELA, B. K. PETERSEN, S. K. KIM, C. P. SANTIAGO, R. GLATT, T. N. MUNDHENK, J. F. PETTIT, AND D. M. FAISSOL, *Improving exploration in policy gradient search: Application to symbolic optimization*, arXiv preprint arXiv:2107.09158, (2021).
- [127] J. B. LASSERRE, *A sum of squares approximation of nonnegative polynomials*, SIAM review, 49 (2007), pp. 651–669.
- [128] P. LAX AND B. WENDROFF, *Systems of conservation laws*, Communications on Pure and Applied Mathematics, 13 (1960), pp. 217–237.
- [129] P. D. LAX, *Weak solutions of nonlinear hyperbolic equations and their numerical computation*, Communications on pure and applied mathematics, 7 (1954), pp. 159–193.
- [130] B. S. LAZAROV AND O. SIGMUND, *Filters in topology optimization based on Helmholtz-type differential equations*, International Journal for Numerical Methods in Engineering, 86

- (2011), pp. 765–781.
- [131] S. LAZGHAB, T. AUKRUST, AND K. HOLTHE, *Adaptive exponential finite elements for the shear boundary layer in the bearing channel during extrusion*, Computer Methods in Applied Mechanics and Engineering, 191 (2002), pp. 1113–1128.
 - [132] J. LEE, *Introduction to Smooth Manifolds*, vol. 218, Springer Science & Business Media, second ed., 2013.
 - [133] G. LEONI, *A First Course in Sobolev Spaces*, Graduate studies in mathematics, American Mathematical Soc., 2009.
 - [134] R. J. LEVEQUE, *Numerical methods for conservation laws*, vol. 214, Springer, 1992.
 - [135] H. LI, S. GUPTA, H. YU, L. YING, AND I. DHILLON, *Quasi-Newton policy gradient algorithms*, arXiv preprint arXiv:2110.02398, (2021).
 - [136] Y. LIN, J. CHAN, AND I. TOMAS, *A positivity preserving strategy for entropy stable discontinuous Galerkin discretizations of the compressible Euler and Navier–Stokes equations*, Journal of Computational Physics, 475 (2023), p. 111850.
 - [137] M. LINDSEY, *Fast randomized entropically regularized semidefinite programming*, arXiv preprint arXiv:2303.12133, (2023).
 - [138] J.-L. LIONS, *Contrôle optimal des systèmes gouvernés par des équations aux dérivées partielles*, Dunod-Gauthier-Villars, Paris, 1968.
 - [139] J.-L. LIONS, *Quelques Méthodes De Résolution Des Problèmes Aux Limites Non Linéaires*, Dunod Paris, 1969.
 - [140] J.-L. LIONS AND G. STAMPACCHIA, *Variational inequalities*, Communications on pure and applied mathematics, 20 (1967), pp. 493–519.
 - [141] R. LISKA AND M. SHASHKOV, *Enforcing the discrete maximum principle for linear finite element solutions of second-order elliptic problems*, Commun. Comput. Phys., 3 (2008), pp. 852–877.
 - [142] G. L. LITVINOV, *Maslov dequantization, idempotent and tropical mathematics: A brief introduction*, Journal of Mathematical Sciences, 140 (2007), pp. 426–444.
 - [143] Z. LIU AND X. LI, *The exponential scalar auxiliary variable (E-SAV) approach for phase field models and its explicit computing*, SIAM Journal on Scientific Computing, 42 (2020), pp. B630–B655.
 - [144] C. LOHMANN, D. KUZMIN, J. N. SHADID, AND S. MABUZA, *Flux-corrected transport algorithms for continuous Galerkin methods based on high order Bernstein finite elements*, Journal of Computational Physics, 344 (2017), pp. 151–186.
 - [145] H. LUO, J. D. BAUM, AND R. LÖHNER, *Computation of compressible flows using a two-equation turbulence model on unstructured grids*, International Journal of Computational Fluid Dynamics, 17 (2003), pp. 87–93.
 - [146] A. MACKEY, H. SCHAEFFER, AND S. OSHER, *On the compressive spectral method*, Multiscale Modeling & Simulation, 12 (2014), pp. 1800–1827.
 - [147] D. MACLAGAN AND B. STURMFELS, *Introduction to tropical geometry*, vol. 161, American Mathematical Society, 2021.
 - [148] P. A. MARKOWICH, *The stationary semiconductor device equations*, Springer Science & Business Media, 1985.
 - [149] J. E. MARSDEN AND T. J. HUGHES, *Mathematical foundations of elasticity*, Courier Corporation, 1994.
 - [150] B. MARTINET, *Regularisation, d’inéquations variationnelles par approximations succesives*, Revue Française d’informatique et de Recherche opérationnelle, (1970).
 - [151] V. P. MASLOV, *A new approach to generalized solutions of nonlinear systems*, in Doklady Akademii Nauk, vol. 292, Russian Academy of Sciences, 1987, pp. 37–41.
 - [152] V. P. MASLOV, *On a new principle of superposition for optimization problems*, Uspekhi Matematicheskikh Nauk, 42 (1987), pp. 39–48.
 - [153] S. F. MCCORMICK, *A revised mesh refinement strategy for Newton’s method applied to nonlinear two-point boundary value problems*, in Lecture Notes in Mathematics, Springer Berlin Heidelberg, 1978, pp. 15–23.
 - [154] W. C. H. MCLEAN, *Strongly elliptic systems and boundary integral equations*, Cambridge university press, 2000.
 - [155] M. S. METTI, J. XU, AND C. LIU, *Energetically stable discretizations for charge transport and electrokinetic models*, Journal of Computational Physics, 306 (2016), pp. 1–18.
 - [156] B. MINOR, *Exponential characteristic spatial quadrature for discrete ordinates neutral particle transport in two-dimensional Cartesian coordinates*, PhD thesis, 1993.
 - [157] M. MONTENEGRO AND O. SANTANA DE QUEIROZ, *Existence and regularity to an elliptic equation with logarithmic nonlinearity*, Journal of Differential Equations, 256 (2009), pp. 482–511.

- [158] J.-J. MOREAU, *Proximité et dualité dans un espace hilbertien*, Bulletin de la Société mathématique de France, 93 (1965), pp. 273–299.
- [159] M. K. V. MURTHY AND G. STAMPACCHIA, *A variational inequality with mixed boundary conditions*, Israel Journal of Mathematics, 13 (1972), pp. 188–224.
- [160] A. S. NEMIROVSKI AND D. B. YUDIN, *Effective methods for the solution of convex programming problems of large dimensions*, Ekonom. i Mat. Metody, 15 (1979), pp. 135–152.
- [161] A. S. NEMIROVSKIJ AND D. B. YUDIN, *Problem complexity and method efficiency in optimization*, Wiley-Interscience, 1983.
- [162] F. NIELSEN, *An elementary introduction to information geometry*, Entropy, 22 (2020), p. 1100.
- [163] J. NOCEDAL AND S. J. WRIGHT, *Numerical optimization*, Springer, 1999.
- [164] J. T. ODEN AND J. N. REDDY, *Variational methods in theoretical mechanics*, Springer Science & Business Media, second ed., 1982.
- [165] A. ORTIZ, M. PUSO, AND N. SUKUMAR, *Maximum-entropy meshfree method for compressible and near-incompressible elasticity*, Computer Methods in Applied Mechanics and Engineering, 199 (2010), pp. 1859–1871.
- [166] I. P. PAPAPOULOS, P. E. FARRELL, AND T. M. SUROWIEC, *Computing multiple solutions of topology optimization problems*, SIAM Journal on Scientific Computing, 43 (2021), pp. A1555–A1582.
- [167] N. PARIKH AND S. BOYD, *Proximal algorithms*, Foundations and trends[®] in Optimization, 1 (2014), pp. 127–239.
- [168] R. K. PATHRIA, *Statistical mechanics*, Elsevier, 2016.
- [169] D. PAVLOV, B. STURMFELS, AND S. TELEN, *Gibbs manifolds*, arXiv preprint arXiv:2211.15490, (2022).
- [170] B. PERTHAME AND C.-W. SHU, *On positivity preserving finite volume schemes for Euler equations*, Numerische Mathematik, 73 (1996), pp. 119–130.
- [171] C. POZRIKIDIS, *Introduction to finite and spectral element methods using MATLAB*, CRC press, 2005.
- [172] R. T. ROCKAFELLAR, *Convex analysis*, vol. 18, Princeton university press, 1970.
- [173] R. T. ROCKAFELLAR, *Augmented Lagrangians and applications of the proximal point algorithm in convex programming*, Mathematics of Operations Research, 1 (1976), pp. 97–116.
- [174] R. T. ROCKAFELLAR, *Monotone operators and the proximal point algorithm*, SIAM journal on control and optimization, 14 (1976), pp. 877–898.
- [175] R. T. ROCKAFELLAR AND R. J.-B. WETS, *Variational analysis*, vol. 317, Springer Science & Business Media, 2009.
- [176] J. RODRIGUES, *Obstacle Problems in Mathematical Physics*, Elsevier Science, 1987.
- [177] D. L. SCHARFETTER AND H. K. GUMMEL, *Large-signal analysis of a silicon read diode oscillator*, IEEE Transactions on electron devices, 16 (1969), pp. 64–77.
- [178] A. SCHIELA, *The Control Reduced Interior Point Method. A Function Space Oriented Algorithmic Approach*, PhD thesis, Freie Universität Berlin, 2006.
- [179] A. SCHIELA AND A. GÜNTHER, *An interior point algorithm with inexact step computation in function space for state constrained optimal control*, Numerische Mathematik, 119 (2011), pp. 373–407.
- [180] E. SCHRÖDINGER, *Quantisierung als eigenwertproblem*, Annalen der physik, 385 (1926), pp. 437–490.
- [181] M. W. SCROGGS, J. S. DOKKEN, C. N. RICHARDSON, AND G. N. WELLS, *Construction of arbitrary order finite element degree-of-freedom maps on polygonal and polyhedral cell meshes*, ACM Transactions on Mathematical Software (TOMS), 48 (2022), pp. 1–23.
- [182] O. SIGMUND AND K. MAUTE, *Topology optimization approaches: A comparative review*, Structural and multidisciplinary optimization, 48 (2013), pp. 1031–1055.
- [183] G. STAMPACCHIA, *Equations elliptiques du second ordre à coefficients discontinus*, Séminaire Jean Leray, (1963), pp. 1–77.
- [184] G. STAMPACCHIA, *Le problème de Dirichlet pour les équations elliptiques du second ordre à coefficients discontinus*, in Annales de l’institut Fourier, vol. 15, 1965, pp. 189–257.
- [185] R. STENBERG, *Analysis of mixed finite elements methods for the Stokes problem: A unified approach*, Mathematics of computation, 42 (1984), pp. 9–23.
- [186] S. STROGATZ, *Infinite powers: How calculus reveals the secrets of the universe*, Eamon Dolan Books, 2019.
- [187] B. STURMFELS, S. TELEN, F.-X. VIALARD, AND M. VON RENESSE, *Toric geometry of entropic regularization*, arXiv preprint arXiv:2202.01571, (2022).
- [188] N. SUKUMAR, *Construction of polygonal interpolants: A maximum entropy approach*, International journal for numerical methods in engineering, 61 (2004), pp. 2159–2181.
- [189] Z. SUN, J. A. CARRILLO, AND C.-W. SHU, *A discontinuous Galerkin method for nonlinear*

- parabolic equations and gradient flow problems with interaction potentials*, Journal of Computational Physics, 352 (2018), pp. 76–104.
- [190] P. K. SWEBY, *High resolution schemes using flux limiters for hyperbolic conservation laws*, SIAM journal on numerical analysis, 21 (1984), pp. 995–1011.
- [191] P. TANKOV AND R. CONT, *Financial Modelling with Jump Processes, Second Edition*, Chapman and Hall/CRC Financial Mathematics Series, Taylor & Francis, 2015.
- [192] M. TEBoulLE, *Entropic proximal mappings with applications to nonlinear programming*, Mathematics of Operations Research, 17 (1992), pp. 670–690.
- [193] M. TEBoulLE, *A simplified view of first order methods for optimization*, Mathematical Programming, 170 (2018), pp. 67–96.
- [194] G. TRAN, H. SCHAEFFER, W. M. FELDMAN, AND S. J. OSHER, *An l^1 penalty method for general obstacle problems*, SIAM Journal on Applied Mathematics, 75 (2015), pp. 1424–1444.
- [195] R. TRÉMOLIÈRES, J. LIONS, AND R. GLOWINSKI, *Numerical Analysis of Variational Inequalities*, North-Holland, Amsterdam, 1981.
- [196] M. TRIBUS AND E. C. MCLRVINE, *Energy and information*, Scientific American, 225 (1971), pp. 179–190.
- [197] F. TRÖLTZSCH, *Optimal Control of Partial Differential Equations*, Graduate Studies in Mathematics, American Mathematical Society, April 2010.
- [198] M. ULBRICH, *Semismooth Newton methods for operator equations in function spaces*, SIAM Journal on Optimization, 13 (2002), pp. 805–841.
- [199] M. ULBRICH AND S. ULBRICH, *Primal-dual interior-point methods for PDE-constrained optimization*, Mathematical Programming, 117 (2007), pp. 435–485.
- [200] B. VAN LEER, *Towards the ultimate conservative difference scheme. V. A second-order sequel to Godunov’s method*, Journal of computational Physics, 32 (1979), pp. 101–136.
- [201] A. VIJAYWARGIYA AND G. FU, *Two finite element approaches for the porous medium equation that are positivity preserving and energy stable*, arXiv preprint arXiv:2303.14216, (2023).
- [202] J. VON NEUMANN AND R. D. RICHTMYER, *A method for the numerical calculation of hydrodynamic shocks*, Journal of applied physics, 21 (1950), pp. 232–237.
- [203] M. WEISER, A. SCHIELA, AND P. DEUFLHARD, *Asymptotic mesh independence of Newton’s method revisited*, SIAM journal on numerical analysis, 42 (2005), pp. 1830–1845.
- [204] B. I. WOHLMUTH AND R. H. KRAUSE, *Monotone multigrid methods on nonmatching grids for nonlinear multibody contact problems*, SIAM Journal on Scientific Computing, 25 (2003), pp. 324–347.
- [205] W. WOLLNER, *A posteriori error estimates for a finite element discretization of interior point methods for an elliptic optimization problem with state constraints*, Computational Optimization and Applications, 47 (2008), pp. 133–159.
- [206] H. WOUK, *The language God talks: On science and religion*, Hachette UK, 2010.
- [207] P. WRIGGERS, *Computational Contact Mechanics*, vol. 2, Springer Berlin Heidelberg, 2006.
- [208] S. J. WRIGHT AND B. RECHT, *Optimization for data analysis*, Cambridge University Press, 2022.
- [209] K. WU AND C.-W. SHU, *Geometric quasilinearization framework for analysis and design of bound-preserving schemes*, arXiv preprint arXiv:2111.04722, (2021).
- [210] B. C. YEE, S. S. OLIVIER, T. S. HAUT, M. HOLEC, V. Z. TOMOV, AND P. G. MAGINOT, *A quadratic programming flux correction method for high-order DG discretizations of SN transport*, Journal of Computational Physics, 419 (2020), p. 109696.
- [211] W. YIN, S. OSHER, D. GOLDFARB, AND J. DARBON, *Bregman iterative algorithms for ℓ_1 -minimization with applications to compressed sensing*, SIAM Journal on Imaging sciences, 1 (2008), pp. 143–168.
- [212] K. YOSIDA, *Functional analysis*, Springer Science & Business Media, 2012.
- [213] S. T. ZALESAK, *Fully multidimensional flux-corrected transport algorithms for fluids*, Journal of computational physics, 31 (1979), pp. 335–362.
- [214] X. ZHANG AND C.-W. SHU, *On maximum-principle-satisfying high order schemes for scalar conservation laws*, Journal of Computational Physics, 229 (2010), pp. 3091–3120.
- [215] X. ZHANG AND C.-W. SHU, *On positivity-preserving high order discontinuous Galerkin schemes for compressible Euler equations on rectangular meshes*, Journal of Computational Physics, 229 (2010), pp. 8918–8934.
- [216] X. ZHANG AND C.-W. SHU, *Maximum-principle-satisfying and positivity-preserving high-order schemes for conservation laws: Survey and new developments*, Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences, 467 (2011), pp. 2752–2776.
- [217] D. ZOZZO, B. OSTING, M. XIA, AND S. J. OSHER, *An efficient primal-dual method for the obstacle problem*, Journal of Scientific Computing, 73 (2017), pp. 416–437.
- [218] T. ZWINGER AND J. C. MOORE, *Diagnostic and prognostic simulations with a full Stokes*

model accounting for superimposed ice of Midtre Lovénbreen, Svalbard, *The Cryosphere*,
3 (2009), pp. 217–229.