

# A numerical approach to study the properties of solutions of the diffusive wave approximation of the shallow water equations

Mauricio Santillana · Clint Dawson

Received: 23 April 2008 / Accepted: 14 January 2009  
© Springer Science + Business Media B.V. 2009

**Abstract** In this paper, we study the properties of approximate solutions to a doubly nonlinear and degenerate diffusion equation, known in the literature as the diffusive wave approximation of the shallow water equations (DSW), using a numerical approach based on the Galerkin finite element method. This equation arises in shallow water flow models when special assumptions are used to simplify the shallow water equations and contains as particular cases the porous medium equation and the p-Laplacian. Diverse numerical schemes have been implemented to approximately solve the DSW equation and have been successfully applied as suitable models to simulate overland flow and water flow in vegetated areas such as wetlands; yet, no formal mathematical analysis has been carried out in order to study the properties of approximate solutions. In this study, we propose a numerical approach as a means to understand some properties of solutions to the DSW equation and, thus, to provide conditions for which the use of the DSW equation may be inappropriate from both the physical and the mathematical points of view, within the context of shallow water modeling. For analysis purposes, we propose a numerical method based on the Galerkin method and we obtain a priori error estimates between the approximate solutions and weak solutions to the DSW equation under physically consistent assumptions. We also present some numerical experiments that provide relevant information about the accuracy of the proposed numerical method

to solve the DSW equation and the applicability of the DSW equation as a model to simulate observed quantities in an experimental setting.

**Keywords** Doubly nonlinear · Degenerate parabolic · Shallow water equations · Doubly degenerate diffusion · Galerkin method nonlinear diffusion · Overland flow · Wetland modeling

## 1 Introduction

In this paper, we study the properties of approximate solutions to a doubly nonlinear and degenerate diffusion equation, known in the literature as the diffusive wave approximation of the shallow water equations (DSW), using a numerical approach based on the Galerkin finite element method. This equation arises in shallow water flow models when special assumptions are used to simplify the shallow water equations, and it gives rise to the following initial/boundary-value problem (IBVP)

$$\left\{ \begin{array}{l} \frac{\partial u}{\partial t} - \nabla \cdot \left( \frac{(u-z)^\alpha}{|\nabla u|^{1-\gamma}} \nabla u \right) = f \quad \text{on } \Omega \times (0, T] \\ u = u_0 \quad \text{on } \Omega \times \{t = 0\} \\ \left( \frac{(u-z)^\alpha}{|\nabla u|^{1-\gamma}} \nabla u \right) \cdot n = B_N \quad \text{on } \partial\Omega \cap \Gamma_N \times (0, T] \\ u = B_D \quad \text{on } \partial\Omega \cap \Gamma_D \times (0, T] \end{array} \right. \quad (1)$$

where  $\Omega$  is an open, bounded subset of  $\mathbb{R}^2$  and  $\Gamma_N$  and  $\Gamma_D$  are subsets of  $\partial\Omega \in C^1$  such that  $\partial\Omega = \Gamma_N + \Gamma_D$ .

M. Santillana (✉) · C. Dawson  
Institute for Computational Engineering and Sciences,  
University of Texas at Austin, Austin, TX, USA  
e-mail: mauricio@ices.utexas.edu

$f : \Omega \times (0, T] \rightarrow \mathbb{R}$ ,  $u_0 : \Omega \rightarrow \mathbb{R}$ ,  $B_N : \partial\Omega \cap \Gamma_N \times (0, T] \rightarrow \mathbb{R}$ , and  $B_D : \partial\Omega \cap \Gamma_D \times (0, T] \rightarrow \mathbb{R}$  are given;  $z : \Omega \rightarrow \mathbb{R}^+$  is a positive time-independent function;  $n$  is the outward normal to  $\Gamma_N$ ;  $0 < \gamma \leq 1$ ;  $1 < \alpha < 2$ ; and  $u : \Omega \times (0, T] \rightarrow \mathbb{R}$  is the unknown. Here,  $\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}$  refers to the Euclidean norm in  $\mathbb{R}^n$  ( $n = 1, 2$  in our work).

Diverse numerical schemes have been implemented to approximately solve the DSW equation appearing in the IBVP (Eq. 1) and have been successfully applied as suitable models to simulate overland flow and water flow in vegetated areas such as wetlands [14, 15, 19, 36, 37]; yet, no formal mathematical analysis has been carried out in order to study the properties of approximate solutions such as well posedness of the discrete problem, error estimates between the true solution (if uniqueness holds) and the numerical approximant, rates of convergence, and so on. This paper is an initial step in that direction.

Shallow water flow over land or in vegetated areas is driven mainly by gravitational forces and dominated by shear stresses, giving rise to uniform and fully developed turbulent flow conditions. These conditions allow the use of scaling arguments to approximate the shallow water system of equations (SWE) by a single doubly nonlinear and degenerate diffusion equation, the DSW equation. Solving the DSW equation computationally requires significantly less work than solving the SWE [19]; thus, it is relevant to explore its applicability and mathematical properties.

Despite the obvious appeal to use the DSW in lieu of the SWE, analyzing the mathematical properties of the DSW equation is not a simple task. Note that the DSW equation contains as particular cases two complicated nonlinear diffusion equations: the porous medium equation (PME) (when  $z = 0$  and  $\gamma = 1$ ) and the  $p$ -Laplacian for  $1 < p < 2$  (when  $\alpha = 0$  and  $p = \gamma + 1$ , this case is not considered in this paper). In fact, to the best of our knowledge, existence, uniqueness, and regularity of solutions of the DSW equation in its general form, as it appears in the IBVP (Eq. 1), have not been studied. However, relevant results and techniques coming from nonlinear diffusion equations and other doubly nonlinear diffusion equations [3, 10, 16, 21, 24, 28] can be applied to the analysis of the DSW equation when topographic effects are ignored (i.e., when  $z = 0$ ). An extensive survey of these results and techniques can be found in [1]. Difficulties arise when *topographic* effects are considered (i.e., when  $z \neq 0$ ) since the mathematical techniques developed to study nonlinear diffusion equations, such as the PME, the  $p$ -Laplacian, and other doubly nonlinear degenerate equations, do not extend directly for this case. Nevertheless, they provide the right setting to start our study.

In this study, we propose a numerical approach as a means to understand some properties of solutions to the DSW equation and, thus, to provide conditions for which the use of the DSW equation may be inappropriate from both the physical and the mathematical points of view, within the context of shallow water modeling. The work presented in this paper is a natural continuation of the study presented in [1]. The outline of this paper is as follows. We begin by presenting a physical motivation and some remarks about the DSW in Sections 1.1 and 1.2. We continue with a review of relevant works existing in the literature, both from the analytical and numerical points of view in Section 1.3. The following sections are devoted to developing our overall strategy, which consists of two steps. First, in Section 2, we will focus our attention on setting up a numerical method based on the Galerkin method to obtain approximate solutions to the IBVP (Eq. 1) and we will obtain a priori error estimates between these approximate solutions and weak solutions to the DSW equation. For this purpose, we will make assumptions about the *true* weak solution to problem 1 based on two criteria, the first one being the physical relevance of solutions to the DSW in the context of shallow water models, and the second one coming from previous analytical results that address the properties of particular cases of the DSW equation in the context of doubly nonlinear and degenerate diffusion equations. Regarding the first criterion, we will restrict our error analysis to approximating the set of nondegenerate solutions to problem 1, whose gradients are bounded. Despite the fact that these solutions may seem very particular, they play an important role in physical applications. The second criterion will ensure uniqueness of solutions from a purely mathematical point of view. Secondly, in Section 3, we will present some numerical experiments that provide relevant information about the accuracy of the proposed numerical method to solve the DSW equation (even for the degenerate case), and the applicability of the DSW equation as a model to simulate observed quantities in an experimental setting. We will also present some numerical experiments aimed at investigating whether some of the qualitative properties, such as the existence of compactly supported solutions or finite speed of propagation of disturbances, found for the solution in [10] for the 1-D case and  $z = 0$ , persist in the more general case for a nonzero and regular topography  $z$ .

### 1.1 Preliminaries

The doubly nonlinear nature of the IBVP (Eq. 1) comes from the fact that the nonlinear behavior appears inside

the divergence term as a product of two nonlinearities involving  $u$  and  $\nabla u$ , namely,  $(u - z)^\alpha$  and  $\nabla u/|\nabla u|^{1-\gamma}$ . For the time being, it will be useful to write the first equation of problem 1 as

$$\frac{\partial u}{\partial t} - \nabla \cdot (a(u, \nabla u) \nabla u) = f, \tag{2}$$

where the diffusion coefficient  $a$  is given by

$$a(u, \nabla u) = \frac{(u - z)^\alpha}{|\nabla u|^{1-\gamma}}, \tag{3}$$

and  $f$  is a given source/sink function independent of  $u$ . Previous work aimed to analyze approximate solutions to nonlinear diffusion equations using the Galerkin finite element method, such as the work of Wheeler [35] and Douglas and Dupont [23], deal with nonlinear diffusion coefficients that only depend on the function  $u$  itself and not on  $\nabla u$ , i.e., diffusion coefficients of the form  $a = a(u)$ . The analysis carried out in such cases requires roughly two assumptions:

$$0 < \mu \leq a(u) \leq M, \quad \text{and} \quad |a'(u)| \leq B \quad \text{for } u \in \mathbb{R} \tag{4}$$

so that  $a$  is uniformly Lipschitz with respect to  $u$  and bounded below by a small constant  $\mu$ . These assumptions ensure, in particular, that one can construct a weak formulation such that, for some Sobolev space  $V$ , one has two fundamental conditions:

$$\begin{aligned} \mu \|u\|_V^2 &\leq (a(u)\nabla u, \nabla u) \quad \text{and} \\ (a(u)\nabla u, \nabla w) &\leq M \|u\|_V \|w\|_V \quad \text{for } u, w \in V, \end{aligned} \tag{5}$$

where  $(\cdot, \cdot)$  represents the appropriate duality pairing. See [31] for a comprehensive study of Galerkin finite element methods for parabolic problems.

The doubly nonlinear nature of the IBVP (Eq. 1) poses new challenges that come from the possible degeneracy of the diffusion coefficient Eq. 3 when  $(u - z) = 0$ , and the nonlinear dependency of Eq. 3 with respect to  $\nabla u$ . In fact, with the condition that  $0 < \gamma \leq 1$ , one can only expect that the diffusion coefficient be uniformly Lipschitz with respect to  $\nabla u$  if  $\gamma = 1$ , that is, when the dependency with respect to  $\nabla u$  disappears and the first equation of problem 1 becomes the PME (for  $z \equiv 0$ ). In general, the diffusion coefficient given by Eq. 3 is, at most, Hölder continuous with respect to  $\nabla u$  and possibly degenerate (i.e.,  $a(u, \nabla u) = 0$ ) in subsets of  $\Omega$ ; thus, one cannot expect that similar expressions such as those shown in Eq. 5 will hold. This fact motivates the need of further assumptions or properties on the type of solutions to be approximated, such as physical consistency, if one is to produce a meaningful numerical method. A natural way to handle the degenerate character of the diffusion coefficient Eq. 3, for

example, is to construct a numerical scheme that approximates nondegenerate problems obtained by substituting the (possibly degenerate) diffusion coefficient Eq. 3,  $a$ , with approximate diffusion coefficients  $a_\epsilon$  in Eq. 1, such that  $0 < \epsilon \leq a_\epsilon(u)$  and with the property that  $a(u) = \lim_{\epsilon \rightarrow 0} a_\epsilon(u)$ , for a small parameter  $\epsilon$ . One then needs to show that solutions of these alternative nondegenerate problems are indeed close, in some sense, to the original (possibly degenerate) solution of problem 3 with  $a$ . This approach has been used previously in the approximation of other degenerate parabolic problems, for example, in [20, 25], and [29]. We will show that, in order to ensure convergence of the proposed method to approximate solutions of problem 1, we will further need to assume that  $\nabla u$  is bounded. This assumption will be shown to be meaningful and physically consistent in the context of shallow water modeling in the next section; see Remark 1.7 in Section 1.3.

### 1.2 Motivation

Even though some of the material of this section has been presented previously in [1], we present it here for completeness. Models for surface water flows are derived from the incompressible, 3-D Navier–Stokes equations, which consist of momentum equations for the three velocity components and a continuity equation. Depending on the physics of the flow, scaling arguments are used in order to obtain effective equations for the problem at hand; see [34]. The IBVP (Eq. 1) is a simplified version of the 2-D shallow water equations called the diffusive wave or zero-inertia approach, which neglects the inertial terms in the horizontal momentum equations.

Recall that, in shallow water theory, the main scaling assumption consists in considering that the vertical scales are small relative to the horizontal ones. This approximation reduces the vertical momentum equation to the hydrostatic pressure relation

$$\frac{\partial p}{\partial z} = \rho g,$$

where  $g$  is the gravitational constant,  $z$  is the vertical coordinate, and  $p$  is the pressure, and leaves us with two effective momentum equations in the horizontal direction. Upon vertical integration, we can obtain the 2-D shallow water momentum equations. In the diffusive wave approximation, the depth-averaged horizontal momentum equations are further approximated using empirical laws, such as Manning’s formula and Chézy’s formula, to find an effective expression for the

horizontal velocity of the fluid in terms of the free water surface slope, given by

$$V = -\frac{(H - z)^{\alpha-1}}{c_f} \frac{\nabla H}{|\nabla H|^{1-\gamma}}, \tag{6}$$

where  $H(t, x)$  is the free water surface elevation or hydraulic head,  $z(x)$  is the bed surface or land elevation,  $0 < \gamma \leq 1$  and  $1 < \alpha < 2$  are non-negative parameters, and  $c_f(x)$  is a friction coefficient.

*Remark 1.1* Manning’s formula in Eq. 6 corresponds to  $\alpha = 5/3$  and  $\gamma = 1/2$ , and  $c_f$  is known as Manning’s coefficient (denoted by  $n$  in the hydraulic literature). Chézy’s formula corresponds to  $\alpha = 3/2$  and  $\gamma = 1/2$ .

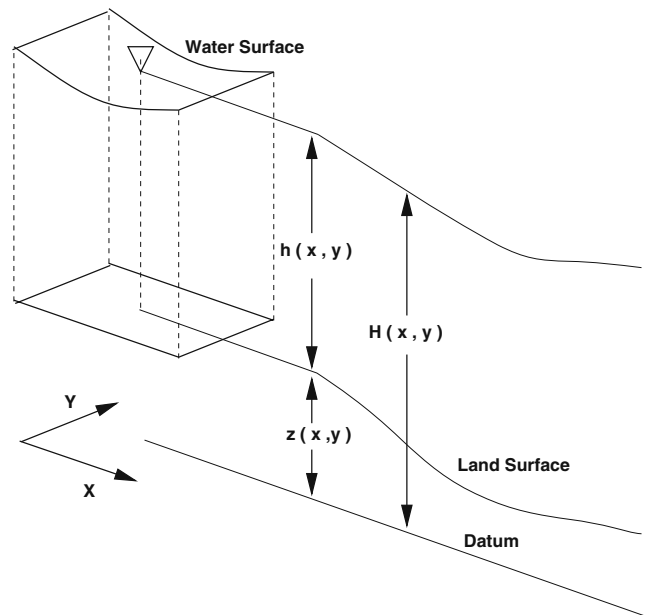
The resulting effective model is given by a doubly nonlinear and degenerate parabolic equation for the water elevation  $H$ , obtained from substituting the particular form of the depth-averaged horizontal velocity, given by Eq. 6, into an equation that arises from combining the depth-averaged continuity equation with the free-surface boundary condition. Equation 1 is a simplification of such model. Furthermore, it is more commonly found in the literature written as

$$\frac{\partial H}{\partial t} - \nabla \cdot \left( \frac{(H - z)^\alpha}{c_f} \frac{\nabla H}{|\nabla H|^{1-\gamma}} \right) = f(t, x), \tag{7}$$

for  $(t, x) \in \mathbb{R}^+ \times \mathbb{R}^2$ ,

where  $f(t, x)$  is a source/sink (such as rain/infiltration). At this point, it is of great importance to mention two of the main requirements for the use of approximation 6, and thus, Eq. 7, to serve as a suitable model to simulate water flow. The first one is that the water depth must be nonnegative,  $(H - z) \geq 0$ , see Fig. 1, and the second one is that the gradient of the water elevation,  $\nabla H$ , needs to be comparable to the gradient of the bathymetry  $\nabla z$ . The latter requirement characterizes water flow regimes not far from uniform flow conditions in open channels, i.e., when the fluid motion is dominated by gravity and balanced by boundary shear stresses, see [8]. In hydrological systems,  $z$  describes the bed surface over which water flows; thus, in physically meaningful situations, one assumes that  $\nabla z$  must be bounded. This in turn implies, in physically meaningful solutions, the boundedness of  $\nabla H$ . This is an extra assumption that will be assumed in the error analysis that aligns well with the physics of the associated problem.

*Remark 1.2* Note that if one identifies the water elevation  $H$  with the hydrostatic pressure  $p$ , the expression that relates the velocity and the water elevation



**Fig. 1** Water depth diagram,  $h = H - z$

gradient Eq. 6 becomes a modified nonlinear version of the empirical Darcy’s law for gas flow through a porous medium. Indeed, flow in vegetated areas such as wetlands can be understood as a flow through a porous medium.

*Remark 1.3* In this context, Eq. 7 makes sense physically only if  $H - z \geq 0$ . It is with this in mind that we will not pay attention to the approximation of negative solutions of Eq. 1. Note that, in writing Eq. 1, we have assumed that  $c_f(x) \equiv 1$ .

*Remark 1.4* Whenever  $H - z = 0$  (or alternatively,  $u - z = 0$  in Eq. 1), Eq. 7 degenerates, i.e., it is no longer of parabolic type.

*Remark 1.5* Note in particular that, for the case when  $\gamma = 1$ ,  $c_f \equiv 1$ , and  $z \equiv 0$ , Eq. 7 becomes the PME. A comprehensive study of the PME can be found in the book by Vázquez [33], and appropriate error analysis references for this particular case can be found, for example, in [25] and [20].

### 1.3 Literature review

To the best of our knowledge, existence, uniqueness, and regularity of solutions of the DSW equation in its general form, as in the IBVP (Eq. 1), have not been studied. However, when (a) topographic effects are ignored (i.e., when  $z = 0$ ) and (b) zero-Dirichlet

boundary conditions are assumed ( $\partial\Omega = \Gamma_D$ ), the DSW equation can be rewritten in the form:

$$\frac{\partial u}{\partial t} - \nabla \cdot (|\nabla u^m|^{\gamma-1} \nabla u^m) = f, \tag{8}$$

with  $m = 1 + \alpha/\gamma$ . Esteban and Vázquez [10] studied this equation in 1-D for the Cauchy problem ( $\Omega = \mathbb{R}$ ) and established: (1) existence, uniqueness, and regularity of strong solutions of Eq. 8; (2) existence and regularity of free boundaries; and (3) asymptotic behavior of solutions and free boundaries for initial data with compact support. In Section 3, we will investigate whether some of the qualitative properties they found for the solution of Eq. 8 persist in the more general case for a nonzero and regular topography  $z$  in 1-D.

Other relevant works have studied the nonlinear diffusion equation shown in Eq. 9 in higher dimensions. This equation can be immediately related to the DSW, for nonnegative solutions and  $z = 0$ . This is achieved using the change of variables  $u = v^{1/m}$  to transform the DSW into:

$$\frac{\partial v^{1/m}}{\partial t} - m^{-\gamma} \nabla \cdot \left( \frac{\nabla v}{|\nabla v|^{1-\gamma}} \right) = f, \tag{9}$$

where  $0 < 1/m = \gamma/(\alpha + \gamma) < 1$ . This change of variables allows sending the nonlinearity  $u^\alpha$ , inside the divergence term in Eq. 1, to the time derivative term in Eq. 9. This, in turn, moves the difficulty of dealing with the possible degeneracy when  $u = 0$  to dealing with differentiability issues of the function  $v^{1/m}$  at  $v = 0$ . Existence, uniqueness, and some regularity of solutions of Eq. 9 have been studied in higher dimensions considering zero Dirichlet boundary conditions by [3, 16, 21, 28], and [1]. It is worthwhile to mention that uniqueness can only be ensured for particular conditions that will be explained further in this section. We pay particular attention to [1], where we provide a constructive method to prove the existence of weak degenerate solutions of Eq. 9 using the Faedo Galerkin method. This proof offers a natural setting for the current numerical analysis and computational method. In this proof, we propose the use of a sequence of regularized functions  $\{\phi_\epsilon(v_\epsilon)\}$  converging uniformly to  $\phi(v) = v^{1/m}$  with the property that  $\phi'_\epsilon(0) < +\infty$  in order to construct approximate solutions to Eq. 9. If one transforms Eq. 9 using the change of variables  $u^m = v$  into

$$\frac{\partial u}{\partial t} - m^{-\gamma} \nabla \cdot \left( ((\phi^{-1})'(u))^\gamma \frac{\nabla u}{|\nabla u|^{1-\gamma}} \right) = f, \tag{10}$$

or further into the DSW for  $z = 0$ ,

$$\frac{\partial u}{\partial t} - \nabla \cdot \left( |u|^\alpha \frac{\nabla u}{|\nabla u|^{1-\gamma}} \right) = f, \tag{11}$$

by observing that  $(\phi^{-1})'(x) = m x^{m-1}$  and  $m - 1 = \alpha/\gamma$ , then this regularized function strategy naturally suggests the use of a sequence of nondegenerate diffusion coefficients  $a_\epsilon$  such that  $0 < \epsilon \leq a_\epsilon(u, \nabla u)$  and given by

$$a_\epsilon(u, \nabla u) = \frac{(\phi_\epsilon^{-1})'(u)}{|\nabla u|^{1-\gamma}}. \tag{12}$$

as a means to approximate the solutions of problem 1 when  $z = 0$  and for small  $\epsilon$ . This is the case since the property that  $\phi'_\epsilon(0) < +\infty$  implies directly that  $(\phi_\epsilon^{-1})'(0) > 0$ . When proceeding this way, one then needs to show that the solution of these alternative nondegenerate problems are indeed close, in some sense, to the original (possibly degenerate) problem 2. This is done in [1] for the case when  $z = 0$ ; see Lemma 1.1. The overall strategy will be used in the analysis of the current numerical scheme even in the case when  $z \neq 0$ .

Finite difference schemes and finite element techniques have been implemented to approximate the solution of the DSW equation and have been used successfully to simulate water flow in shallow systems in [5, 14, 15, 19, 32, 36, 37]. However, no formal mathematical analysis has been carried out in order to show, for example, that the proposed methods converge in some sense to the *true* solution of the IBVP (Eq. 1). This is not surprising given the complexity of the general formulation of the IBVP (Eq. 1) and the lack of analytical techniques to prove for example uniqueness of solutions in the presence of topographic effects. Some relevant works approximating degenerate parabolic equations include, for example, [12, 13, 17, 20, 25, 29, 30], and [2]. In particular, in [25], the authors present a numerical method to approximate degenerate parabolic problems similar to the one used in this paper. In their study, they analyze equations of the form

$$\frac{\partial u}{\partial t} - \nabla \cdot (\nabla v + b(r(v))) + f(r(v)) = 0, \quad u \in m(v), \tag{13}$$

where  $m(v)$  is a maximal monotone graph in  $\mathbb{R} \times \mathbb{R}$  possibly with a singularity at the origin ( $m'(0) = \infty$ ). Stefan type, nonstationary filtration type, and porous-medium type degenerate parabolic equations can be written in the form Eq. 13. Even though the authors introduce the use of  $r(v)$  to obtain more general results, one may replace  $r(v)$  by  $v$  for clarity. In cases when



singularities in  $m$  appear, the authors use a smoothing procedure similar to the one explained above used in [1]. As a first step of the approximation, they construct a numerical scheme that approximates a regularized problem obtained by replacing  $m$  by a smooth function  $m_\epsilon$  with maximal slope equal to  $1/\epsilon$ , for some regularization parameter  $\epsilon > 0$ . Then, they discretize this regularized problem in space and time to compute the regularized numerical approximation  $U_\epsilon^h$ . Finally, roughly speaking, in order to obtain global error estimates between the solution  $u(t)$  of Eq. 13 and the regularized numerical approximation  $U_\epsilon^h(t)$ , they obtain bounds for the quantities  $\|u(t) - u_\epsilon(t)\|_{L^\infty}$  and  $\|u_\epsilon(t) - U_\epsilon^h(t)\|_{L^\infty}$  for  $0 \leq t \leq T$ , where  $u_\epsilon(t)$  is the true solution of the regularized problem (solving Eq. 13 with  $m_\epsilon$  instead of  $m$ ). From the two  $L^\infty$ -estimates, they can obtain a global estimate of  $\|u(t) - U_\epsilon^h(t)\|_{L^\infty}$  using the triangle inequality. Even though our strategy is similar, in our case, the analysis will not be as complete. We will find bounds similarly, for the quantity  $\|u_\epsilon(t) - U_\epsilon^h(t)\|_{L^2}$  associated to the nondegenerate problem, however, estimating the difference between the solution  $u$  to the IBVP (Eq. 1) and  $u_\epsilon$ , the solutions to the nondegenerate problems (solving Eq. 1 with a nondegenerate diffusion coefficient  $a_\epsilon$  instead of  $a$  as described in the previous paragraphs) is not yet completely understood in the general setting when topographic effects are considered ( $z \neq 0$ ). This is so, in this case, since an appropriate proof of uniqueness of solutions has not been developed yet.

To overcome this difficulty and justify our approach, we rely on the fact that, in [1], we proved that one can approximate a (unique) weak solution  $v$  of Eq. 9 as  $v = \lim_{\epsilon \rightarrow 0} v_\epsilon$  using the Faedo Galerkin method for flat topographies ( $z = 0$ ) and under physically consistent conditions. Recalling that the sequence of  $\{v_\epsilon\}$  is obtained by solving the regularized problem substituting  $v^{1/m}$  by  $\{\phi_\epsilon(v_\epsilon)\}$  as described in the previous paragraphs, one can see that such limit implies that a (unique) solution  $u$  of the IBVP (Eq. 1) can be constructed similarly as  $u = \lim_{\epsilon \rightarrow 0} u_\epsilon$  (where the sequence of  $\{u_\epsilon\}$  is obtained by solving the nondegenerate problems substituting  $a$  in Eq. 3 by  $a_\epsilon$ ) for flat topographies ( $z = 0$ ).

According to Bamberger [3], when  $z = 0$ , a sufficient condition for uniqueness of solutions is that  $u_t \in L^1(0, T; L^1(\Omega))$ ; see Theorem 1.1. This condition implies, since  $u \in L^\infty(0, T; L^\infty(\Omega))$ , that  $u \in C^0(0, T; L^1(\Omega))$ . If one identifies  $u$  (or  $H$  as described in Section 1.2) with the free water surface elevation in a hydrological context, this condition implies that there will be a unique solution if the volume of water in the domain  $\int_\Omega u$  changes continuously in time. This

is a natural and physically consistent condition when modeling hydrologic systems.

*Remark 1.6* Note that the natural norm induced by multiplying the DSW equation by  $u$  and integrating by parts, in the nondegenerate case, is the  $W^{1,1+\gamma}(\Omega)$  norm. Indeed,

$$\int_\Omega \frac{\partial u}{\partial t} u - \int_\Omega \left( \frac{(u-z)^\alpha}{|\nabla u|^{1-\gamma}} \nabla u \right) \cdot \nabla u = \int_\Omega f u$$

implies that, for a sufficiently regular  $u$  and, say, zero Neumann boundary conditions,

$$\frac{1}{2} \frac{\partial}{\partial t} \|u\|_{L^2(\Omega)} + \int_\Omega (u-z)^\alpha |\nabla u|^{1+\gamma} = \int_\Omega f u.$$

After some manipulations on the previous expression along with assumptions on the nondegeneracy  $u - z > \epsilon > 0$ , and  $u_0, f \in L^2(\Omega)$  for all  $t \in [0, T]$ , one obtains the analytic stability result,

$$\|u\|_{L^2(\Omega)} \leq C (\|u_0\|_{L^2(\Omega)}, \|f\|_{L^2(\Omega)})$$

and

$$\|\nabla u\|_{L^{1+\gamma}(\Omega)} \leq C (\|u_0\|_{L^2(\Omega)}, \|f\|_{L^2(\Omega)})$$

In our error analysis, we obtain stability and a priori error estimates for the approximations of  $u$ , and  $\nabla u$ , namely,  $U$  and  $\nabla U$ , in the  $L^2$ -norm, by assuming the appropriate regularity on  $u$ , nondegeneracy and the physical-consistency assumption on the uniform boundedness of  $\nabla u_\epsilon$ . Working with a norm that is not naturally induced by the problem has advantages and disadvantages. The advantages are that the arguments used to prove estimates in our study extend naturally from the classical arguments developed by Wheeler [35], and Douglas and Dupont [23] for nonlinear parabolic problems. The disadvantage is that the error bounds may be too conservative. The previous statement is supported by the numerical findings on the performance of our method presented in Section 3, which show higher convergence rates than those ensured by our analysis.

*Remark 1.7* It is important to mention that the condition that  $\nabla u$  be bounded in the  $L^\infty$  sense does not necessarily hold for all solutions of Eq. 1. In fact, even in the particular case when Eq. 1 becomes the PME ( $z = 0$  and  $\gamma = 1$ ) in two or higher dimensions, there exists a class of solutions called *focusing solutions* that exhibit no local regularity on the gradient in subsets of  $\Omega$ ; see Chapter 19 in [33]. Even though we ignore this class of solutions in our analysis by assuming boundedness of

$\nabla u$ , it is justified to do so in the context of shallow water modeling where, for small beds  $z$ , one expects *small* values of  $\nabla z$  and, thus, *small* values of  $\nabla u$ , as explained in Section 1.2.

*Remark 1.8* In order to prove uniqueness of solutions for problem 1 in its general form, presumably, one has to impose two conditions: one on the regularity of the time-independent function  $z(x)$  describing the topography and another one in the form of an entropy condition as described in [7]. The latter one may provide means to identify unique physically consistent solutions. This in turn may imply regularity properties on the gradient of the solution  $\nabla u$  such as boundedness in some appropriate norm.

### 1.4 Notation

We will use the standard notation introduced in [11]. Let  $X$  be a real Banach space, with norm  $\| \cdot \|$ . The symbol  $L^p(0, T; X)$  will denote the Banach space of all measurable functions  $u : [0, T] \rightarrow X$  such that

- $\|u\|_{L^p(0, T; X)} := \left( \int_0^T \|u(t)\|^p \right)^{1/p} < \infty$ , for  $1 \leq p < \infty$ , and
- $\|u\|_{L^\infty(0, T; X)} := \text{ess sup}_{0 \leq t \leq T} \|u(t)\| < \infty$ .

For  $1 \leq p \leq +\infty$ , we will denote its conjugate as  $p^*$ , i.e.,  $1/p + 1/p^* = 1$ . For any measurable set  $E \subset \Omega$  and real valued vector functions  $u \in L^p(E)$  and  $v \in L^{p^*}(E)$ , we will denote for the duality pairing between  $u$  and  $v$  as

$$(u, v)_E := \int_E u \cdot v.$$

For simplicity, we use  $(u, v) := (u, v)_\Omega$ . Throughout the paper,  $C$  will be a generic positive constant with different values, and the explicit dependence with respect to parameters will be written inside parentheses.

### 1.5 Auxiliary results

In this section we present the nondegenerate problem we will approximate numerically along with some properties and results that will be used in the analysis carried out in the next sections. We begin by introducing the nondegenerate version of the IBVP (Eq. 1) obtained by replacing the function  $(s - z)^\alpha$  with a sequence of bounded Lipschitz functions  $\{\beta_\epsilon(s)\}$ , with the properties that (1)  $\{\beta_\epsilon(s)\}$  converges uniformly to  $(s - z)^\alpha$  as  $\epsilon \rightarrow 0$ , and (2) for small  $\epsilon > 0$  the following holds  $\beta_\epsilon(s) \geq \epsilon$  for all  $t \in [0, T]$ . To this end, the bathymetry  $z(x)$

will be assumed to be a smooth and bounded time-independent function defined in  $\Omega$ . The nondegenerate IBVP is given by

$$\left\{ \begin{array}{l} \frac{\partial u}{\partial t} - \nabla \cdot \left( \beta_\epsilon(u) \frac{\nabla u}{|\nabla u|^{1-\gamma}} \right) = f \quad \text{on } \Omega \times (0, T] \\ u = u_0 \quad \text{on } \Omega \times \{t = 0\} \\ \left( \frac{\beta_\epsilon(u)}{|\nabla u|^{1-\gamma}} \nabla u \right) \cdot n = B_N \quad \text{on } \partial\Omega \cap \Gamma_N \times (0, T] \\ u = B_D \quad \text{on } \partial\Omega \cap \Gamma_D \times (0, T]. \end{array} \right. \tag{14}$$

In the next section, we develop a numerical scheme to approximate this nondegenerate problem as explained in Section 1.3. The fact that solutions to the nondegenerate problem 14 are close to the original solution to problem 1 as  $\epsilon \rightarrow 0$  will be understood as in [1] for  $z = 0$  and will be assumed for the general case  $z \neq 0$ .

*Remark 1.9* In Eq. 14, for each  $\epsilon$  and, thus, for each  $\beta_\epsilon(u)$ , one has a solution  $u_\epsilon$ .

*Remark 1.10* Frequently, in the actual computational code, one does not need to implement the sequence of  $\{\beta_\epsilon(u)\}$ ; however, it becomes crucial to use this sequence if one wants to find error estimates. For intuition purposes, one could choose, for example, the following sequence:  $\beta_\epsilon(u) = (u - z)^\alpha + \epsilon$ .

We proceed to list some results about the solution of problem 1 under conditions a and b stated in Section 1.3, coming from previous analytical works [3, 10, 16, 21, 28], and [1].

**Lemma 1.1** *We can approximate a nonnegative solution  $v$  of Eq. 9 as:*

$$v = \lim_{\epsilon \rightarrow 0} v_\epsilon,$$

where  $\{v_\epsilon\}$  is a sequence of Faedo–Galerkin solutions of the regularized problems obtained by substituting  $\phi_\epsilon(v_\epsilon)$  for  $\phi(v) = v^{1/m}$  in Eq. 9, with the properties that  $\{\phi_\epsilon(v_\epsilon)\}$  converges uniformly to  $\phi = v^{1/m}$  (or more generally to  $\phi = v|v|^{\frac{1-m}{m}}$ ) as  $\epsilon \rightarrow 0$ , and  $\phi'_\epsilon(0) < +\infty$ . Furthermore,  $u$ , defined as  $u = v^{1/m}$ , is a nonnegative solution of the IBVP (Eq. 1) under conditions a and b, defined in Section 1.3.

*Proof* See step 5 in proof of Theorem 2.1 and Corollary 2.1 in [1]. □

**Lemma 1.2** *If  $u$  is a solution of the IBVP (Eq. 1) under conditions  $a$  and  $b$ , and*

$$u_0 \in L^\infty(\Omega) \quad \text{and} \quad f \in L^\infty(0, T; L^\infty(\Omega)),$$

then

$$\sup_{t \in [0, T]} \|u\|_{L^\infty(\Omega)} \leq C (\|u_0\|_{L^\infty(\Omega)}, \|f\|_{L^\infty(0, T; L^\infty(\Omega))}, T).$$

*Proof* See Corollary 3.2 in [1]. □

**Theorem 1.1** (Bamberger) *Assume  $u$  and  $v$  are weak solutions of the IBVP (Eq. 1) under conditions  $a$  and  $b$  and satisfying*

$$u_t, v_t \in L^1(0, T; L^1(\Omega)), \tag{15}$$

then  $u = v$ .

*Proof* This is a consequence of Theorem 4.1 in [1]. See also [3]. □

**Lemma 1.3** *Let  $u_1$  and  $u_2$  be nonnegative  $L^\infty(\Omega)$  functions; then, for  $\alpha \geq 1$ ,*

$$|u_1^\alpha - u_2^\alpha| \leq \alpha (\max(\|u_1\|_{L^\infty(\Omega)}, \|u_2\|_{L^\infty(\Omega)}))^{\alpha-1} |u_1 - u_2| \tag{16}$$

*Proof* We can express

$$\begin{aligned} |u_1^\alpha - u_2^\alpha| &= \left| \int_0^1 \frac{d}{d\tau} (\tau u_1 + (1-\tau)u_2)^\alpha d\tau \right| \\ &\leq \alpha |u_1 - u_2| \int_0^1 (\tau u_1 + (1-\tau)u_2)^{\alpha-1} d\tau \\ &\leq \alpha |u_1 - u_2| \int_0^1 (\tau \|u_1\|_{L^\infty(\Omega)} \\ &\quad + (1-\tau)\|u_2\|_{L^\infty(\Omega)})^{\alpha-1} d\tau \\ &\leq \alpha |u_1 - u_2| (\max(\|u_1\|_{L^\infty(\Omega)}, \|u_2\|_{L^\infty(\Omega)}))^{\alpha-1} \end{aligned}$$
□

**Lemma 1.4** (Coercivity and continuity) *Let  $u_1$  and  $u_2$  be  $L^\infty(\Omega)$  positive functions with the property that  $\nabla u_1, \nabla u_2 \in L^\infty(\Omega)$ ; then, the following estimates hold true:*

$$\begin{aligned} \gamma \mathcal{A}_0 |\nabla u_1 - \nabla u_2|^2 &\leq \left( \frac{\nabla u_1}{|\nabla u_1|^{1-\gamma}} - \frac{\nabla u_2}{|\nabla u_2|^{1-\gamma}} \right) \\ &\quad \times (\nabla u_1 - \nabla u_2) \end{aligned} \tag{17}$$

and

$$\begin{aligned} \left| \frac{\nabla u_1}{|\nabla u_1|^{1-\gamma}} - \frac{\nabla u_2}{|\nabla u_2|^{1-\gamma}} \right| &\leq \mathcal{A}_0 |\nabla u_1 - \nabla u_2| \\ &\leq \frac{2}{\gamma} |\nabla u_1 - \nabla u_2|^\gamma, \end{aligned} \tag{18}$$

where

$$\mathcal{A}_0 := \int_0^1 |\lambda \nabla u_1 + (1-\lambda) \nabla u_2|^{\gamma-1} d\lambda$$

*Proof* See [9], pp. 348–350. □

See also [4] and the references therein for a more general result.

### 1.5.1 Interpolation theory results

For Lemmas 1.5 and 1.6, we will consider  $\tau$  to be a quasiuniform triangulation of  $\Omega$  into elements  $E_i$ ,  $i = 1, \dots, m$ , with  $\text{diam}(E_i) = h_i$  and  $h = \max_i(h_i)$ .  $\mathcal{M}(= \mathcal{P}^k)$  will denote a finite-dimensional subspace of  $H_0^1(\Omega)$  defined on this triangulation consisting of piecewise polynomials of degree at most  $k$ , and  $K_0$  will denote a constant independent of  $h$  and  $v$ .

**Lemma 1.5** (Interpolation error) *Let  $u \in H^{k+1}(\Omega)$ ; then, there exists  $\hat{u} \in \mathcal{M}$ , interpolant of  $u$ , defined by*

$$\int (\hat{u} - u)v = 0 \quad \forall v \in \mathcal{M}, \tag{19}$$

with the following property:

$$\|\hat{u} - u\|_{H^s(\Omega)} \leq C h^{k+1-s} \|u\|_{H^{k+1}(\Omega)},$$

where  $0 \leq s \leq k$ .

*Proof* See Section 4.4 in [6]. □

**Remark 1.11** Lemma 1.5 implies that, for a subspace  $\mathcal{M} = \mathcal{P}^1$  consisting of piecewise linear polynomials,

$$\begin{aligned} \|\hat{u} - u\|_{L^2(\Omega)} &\leq C h^2 \|u\|_{H^2(\Omega)} \quad \text{and} \\ \|\nabla \hat{u} - \nabla u\|_{L^2(\Omega)} &\leq C \|\hat{u} - u\|_{H^1(\Omega)} \leq C h \|u\|_{H^2(\Omega)} \end{aligned} \tag{20}$$

These inequalities will be useful in Section 3.

**Lemma 1.6** (Inverse inequalities) *Let  $v \in \mathcal{M}$ ; then, there exists a constant  $K_0$  independent of  $h$  and  $v$  such that*

$$\begin{aligned} \|v\|_{L^\infty(\Omega)} &\leq K_0 h^{-1} \|v\|_{L^2(\Omega)} \quad \text{and} \\ \|\nabla v\|_{L^\infty(\Omega)} &\leq K_0 h^{-1} \|\nabla v\|_{L^2(\Omega)} \end{aligned}$$

*Proof* See Section 4.5 in [6]. □



## 2 Discrete approach

In this section, we will use the continuous Galerkin method in order to numerically approximate the solution of the initial/boundary value problem 14. We will provide a priori error estimates between the true solution of Eq. 14,  $u_\epsilon$ , and the Galerkin approximate solutions  $U_\epsilon^h$ , both in the semidiscrete and fully discrete cases for the zero Dirichlet and Newmann boundary conditions. The analysis will be an extension of the techniques presented in [35] and [23] and holds true for any sequence of Lipschitz functions,  $\beta_\epsilon(u)$ , with properties 1 and 2 described in Section 1.5. The (unique) solution  $u_\epsilon$  to the nondegenerate problem 14 and the Galerkin approximate solution  $U_\epsilon^h$  will be denoted with  $u$  and  $U$ , respectively, in the following paragraphs. Based on Lemmas 1.1 and 1.2, the following assumptions will be made in the following paragraphs for the general case when  $z \neq 0$ , as discussed in Section 1.3:

- Solutions to the nondegenerate problem 14 are close to the original solution of problem in some sense as  $\epsilon \rightarrow 0$ . See Lemma 1.1.
- $u \in L^\infty(0, T; L^\infty(\Omega))$ . See Lemma 1.2.
- $\nabla u \in L^\infty(0, T; L^\infty(\Omega))$

The latter assumption restricts our analysis to physically consistent solutions in the context of shallow water modeling. Our numerical analysis is carried out for piecewise polynomial basis functions of order  $k$ . However, the limited regularity of solutions of the DSW calls in general for lower-order approximation spaces.

### 2.1 The semidiscrete case

In the Galerkin method, we seek a differentiable function  $U(\cdot, t) \in \mathcal{M}$ , a finite dimensional subspace of  $H^1(\Omega)$  if the boundary conditions in problem 14 are of Newman-type, or  $H_0^1(\Omega)$  if they are of Dirichlet-type, such that it satisfies the following weak form:

$$\left\{ \begin{aligned} &\left( \frac{\partial U}{\partial t}, v \right) + \left( \beta_\epsilon(U) \frac{\nabla U}{|\nabla U|^{1-\gamma}}, \nabla v \right) = (f, v) \quad t > 0, \forall v \in \mathcal{M}, \\ &\text{and} \quad (U(\cdot, 0), v) = (u_0, v) \quad t = 0, \forall v \in \mathcal{M}, \end{aligned} \right. \tag{21}$$

where  $\mathcal{M}$  denotes the span  $\{v_i\}_{i=1}^M$ , and  $v_1, \dots, v_M$  are linearly independent functions in  $H^1(\Omega)$  and  $\beta_\epsilon(u)$  is a Lipschitz function, with properties 1 and 2 as described in Section 1.5. By construction, we can represent any

function in  $\mathcal{M}$  as a linear combination of the family  $\{v_i\}$ , thus, in particular,

$$U(x, t) = \sum_{i=1}^M \zeta_i(t) v_i(x). \tag{22}$$

Substituting Eq. 22 in Eq. 21, we observe that the semi-discrete problem can be stated: Find coefficients  $\zeta_i(t)$  in Eq. 22 such that

$$\sum_{i=1}^M \zeta_i'(t) (v_i, v_j) + \sum_{i=1}^M \zeta_i(t) (\beta_\epsilon^*(\zeta) \nabla v_i, \nabla v_j) = (f, v_j) \tag{23}$$

for  $j = 1, \dots, M$ ,

with  $\sum_{i=1}^M \zeta_i(0) (v_i, v_j) = (u_0, v_j)$ , and

$$\beta_\epsilon^*(\zeta) := \beta_\epsilon \left( \sum_{i=1}^M \zeta_i(t) v_i(x) \right) \left( \sum_{i=1}^M [\zeta_i(t) \nabla v_i(x)]^2 \right)^{\frac{\gamma-1}{2}}. \tag{24}$$

Equivalently, we can express the previous problem as the initial value problem for the system of nonlinear ordinary differential equations given by

$$\begin{cases} G \zeta'(t) = -B(\zeta) \zeta + F, \\ G \zeta(0) = b, \end{cases} \tag{25}$$

where the entries of the matrix  $G = (G_{ij})$  are given by  $G_{ij} = (v_i, v_j)$ ; the entries of the matrix  $B(\zeta) = (B_{ij}(\zeta))$  are given by

$$B_{ij}(\zeta) = (\beta_\epsilon^*(\zeta) \nabla v_i, \nabla v_j); \tag{26}$$

the components of the vectors  $b = (b_j)$  and  $F = (F_j)$  are given by  $b_j = (u_0, v_j)$  and  $F_j = (f, v_j)$ , respectively; and the vector of unknowns is  $\zeta(t) = (\zeta_j(t))$ .

Whenever  $U(x, t)$  given by Eq. 22 exists, it is called the *continuous in time Galerkin approximation* or *semidiscrete approximation* to the weak solution of problem 14. Though this approximation is never computed in practice, it is easy to understand and gives us insight into our method development. For computational purposes, the variable  $t$  is discretized and a fully algebraic nonlinear system of equations needs to be solved at each time step in order to obtain a fully discrete approximation to the solution. This is studied in Section 2.2.

#### 2.1.1 Existence of the continuous in time Galerkin approximation

**Theorem 2.1** *There exists at least one solution to problem 25.*

*Proof* Since the family  $\{v_i\}$  is linearly independent, the mass matrix  $G$  is a Gram matrix, and thus, in particular, it is a positive definite and invertible matrix. Hence, problem 25 can be written as:

$$\begin{cases} \zeta'(t) = -G^{-1}B(\zeta)\zeta + G^{-1}F, \\ \zeta(0) = G^{-1}b. \end{cases} \tag{27}$$

Note that the mapping  $\Theta(\zeta) : \mathbb{R}^M \rightarrow \mathbb{R}^M$  defined by  $\Theta(\zeta) = -G^{-1}B(\zeta)\zeta + G^{-1}F$  is  $\gamma$ -Hölder continuous, and thus, by *Peano's* existence theorem for ordinary differential equations (see [18] pp. 10), there exists at least one solution to problem 27, which immediately implies the statement of the Theorem.  $\square$

*Remark 2.1* The fact that  $\Theta(\zeta)$  is  $\gamma$ -Hölder continuous is a consequence of  $B(\zeta)\zeta$  being  $\gamma$ -Hölder continuous itself. This can be easily verified recalling definition 26. Indeed, if we let  $Y = \nabla U$ , the diffusion operator  $\beta_\epsilon(U)Y|Y|^{\gamma-1}$  is Lipschitz continuous with respect to  $U$  and  $\gamma$ -Hölder continuous with respect to  $Y$ . Since  $U$  and  $Y$  depend linearly on  $\zeta$ , the diffusion operator will be  $\gamma$ -Hölder continuous with respect to  $\zeta$ , and thus, the discrete diffusion operator defined by Eq. 26 will inherit this property; see Lemma 1 in the Appendix.

### 2.1.2 Continuous in time a priori error estimate

In this section, we will study how close (possibly non-unique) solutions to the continuous in time Galerkin approximation problem 25,  $U$ , are to the true weak solution  $u$  of problem 14. We focus our analysis on the case when  $\nabla u \in L^\infty(0, T, L^\infty(\Omega))$ . In the following paragraphs, we will assume that there exists a function  $\hat{u}$  called the *interpolant* of  $u$ , provided that  $u$  belongs to some Banach space with certain regularity. The interpolant could be, for example, the  $L^2$  projection as defined in Lemma 1.5. Further assumptions include that  $\beta_\epsilon(\hat{u}), \beta_\epsilon(U), \nabla \hat{u}, \nabla U \in L^\infty(0, T, L^\infty(\Omega))$ ; see Remark 2.2 below. The error between solutions of problem 25 and the solution of problem 14,  $\|u - U\|_{L^2}$ , will be shown to be bounded by terms that only involve approximation errors between the interpolant and the true solution of problem 14,  $\|u - \hat{u}\|_{L^2}$ . Thus, reducing the global problem to well known results in interpolation theory in Hilbert spaces, such as those presented in Section 1.5.1. To simplify the notation in the analysis, we will denote with  $\beta(x)$  any element of the sequence  $\{\beta_\epsilon(x)\}$  as described in Section 1.5.

*Remark 2.2* The fact that  $\beta_\epsilon(\hat{u}), \nabla \hat{u} \in L^\infty(0, T, L^\infty(\Omega))$ , for particular finite element approximation spaces, is a direct consequence of the assumptions that  $\nabla u \in L^\infty(0, T, L^\infty(\Omega))$ ,  $(u - z) \in L^\infty(0, T, L^\infty(\Omega))$ , and Theorem 4.8.7 in [6]. We will further show that, provided  $\|u - z\|_{L^\infty(0, T, L^\infty(\Omega))} \leq K_1$  and  $\|\nabla u\|_{L^\infty(0, T, L^\infty(\Omega))} \leq K_2$  for some (possibly large) constants  $K_1, K_2 > 0$ , then  $\|U - z\|_{L^\infty(0, T, L^\infty(\Omega))} \leq K_1(1 + \epsilon_{k_1})$  and  $\|\nabla U\|_{L^\infty(0, T, L^\infty(\Omega))} \leq K_2(1 + \epsilon_{k_2})$ , for small parameters  $\epsilon_{k_1}$  and  $\epsilon_{k_2}$ , and for particular finite element approximation spaces; see Lemmas 2.1 and 2.2 below.

**Theorem 2.2** *Let  $u \in W^{1,\infty}(\Omega)$  be the solution of problem 14 and let  $U$  be a solution of problem 25. Let  $\chi = u - \hat{u}$  be the approximation error between the interpolant and the true solution of problem 1. Further, assume that  $\nabla \hat{u}, \nabla U \in L^\infty(0, T, L^\infty(\Omega))$ . Then, for all  $t \in [0, T]$ ,*

$$\begin{aligned} \|u - U\|_{L^2(\Omega)}^2 + \|\nabla u - \nabla U\|_{L^2(0, T, L^2(\Omega))}^2 &\leq \|\chi(t)\|_{H^1(\Omega)}^2 \\ &+ C \left( \|b - u_0\|_{L^2(\Omega)}^2 + \|\chi(0)\|_{L^2(\Omega)}^2 + \int_0^T \|\chi_t\|_{L^2(\Omega)}^2 + \right. \\ &\left. + \int_0^T \|\chi\|_{L^2(\Omega)}^2 + \int_0^T \int_\Omega |\nabla \chi|^2 \gamma \right) \end{aligned} \tag{28}$$

*Proof* Note that a weak solution  $u$  of problem 14 satisfies the weak form given by Eq. 21. Thus, in particular, the following holds:

$$\begin{aligned} \left( \frac{\partial \hat{u}}{\partial t}, v \right) + \left( \beta(U) \frac{\nabla \hat{u}}{|\nabla \hat{u}|^{1-\gamma}}, \nabla v \right) \\ = (f, v) - \left( \frac{\partial(u - \hat{u})}{\partial t}, v \right) \\ + \left( \beta(u) \left( \frac{\nabla \hat{u}}{|\nabla \hat{u}|^{1-\gamma}} - \frac{\nabla u}{|\nabla u|^{1-\gamma}} \right), \nabla v \right) \\ + \left( (\beta(U) - \beta(u)) \frac{\nabla \hat{u}}{|\nabla \hat{u}|^{1-\gamma}}, \nabla v \right) \end{aligned} \tag{29}$$

Subtracting Eq. 29 from Eq. 21, we obtain that

$$\begin{aligned} \left( \frac{\partial(U - \hat{u})}{\partial t}, v \right) + \left( \beta(U) \left( \frac{\nabla U}{|\nabla U|^{1-\gamma}} - \frac{\nabla \hat{u}}{|\nabla \hat{u}|^{1-\gamma}} \right), \nabla v \right) \\ = \left( \frac{\partial(u - \hat{u})}{\partial t}, v \right) + \left( \beta(u) \left( \frac{\nabla u}{|\nabla u|^{1-\gamma}} - \frac{\nabla \hat{u}}{|\nabla \hat{u}|^{1-\gamma}} \right), \nabla v \right) \\ + \left( (\beta(u) - \beta(U)) \frac{\nabla \hat{u}}{|\nabla \hat{u}|^{1-\gamma}}, \nabla v \right) \end{aligned} \tag{30}$$

Let  $\xi = U - \hat{u}$  and  $\chi = u - \hat{u}$ . Set  $v = \xi$  in the previous expression to find

$$\begin{aligned} & \left( \frac{\partial \xi}{\partial t}, \xi \right) + \left( \beta(U) \left( \frac{\nabla U}{|\nabla U|^{1-\gamma}} - \frac{\nabla \hat{u}}{|\nabla \hat{u}|^{1-\gamma}} \right), \nabla \xi \right) \\ &= \left( \frac{\partial \chi}{\partial t}, \xi \right) + \left( \beta(u) \left( \frac{\nabla u}{|\nabla u|^{1-\gamma}} - \frac{\nabla \hat{u}}{|\nabla \hat{u}|^{1-\gamma}} \right), \nabla \xi \right) \\ &+ \left( (\beta(u) - \beta(U)) \frac{\nabla \hat{u}}{|\nabla \hat{u}|^{1-\gamma}}, \nabla \xi \right) \end{aligned} \tag{31}$$

The above expression and estimate 17 lead to the following inequality:

$$\begin{aligned} & \frac{1}{2} \frac{\partial}{\partial t} \|\xi\|_{L^2(\Omega)}^2 + \gamma \epsilon \mathcal{A} \|\nabla \xi\|_{L^2(\Omega)}^2 \\ & \leq \left( \frac{\partial \chi}{\partial t}, \xi \right) + \left( \beta(u) \left( \frac{\nabla u}{|\nabla u|^{1-\gamma}} - \frac{\nabla \hat{u}}{|\nabla \hat{u}|^{1-\gamma}} \right), \nabla \xi \right) \\ & + \left( (\beta(u) - \beta(U)) \frac{\nabla \hat{u}}{|\nabla \hat{u}|^{1-\gamma}}, \nabla \xi \right), \end{aligned} \tag{32}$$

where  $\mathcal{A} := \inf_{(0,T) \times \Omega} (\mathcal{A}_0) = (\sup(\|\nabla U\|_{L^\infty(0,T,L^\infty(\Omega))}, \|\nabla \hat{u}\|_{L^\infty(0,T,L^\infty(\Omega))}))^{\gamma-1}$ .

The terms on the right-hand side can be bounded in the following way: The first one, using Young’s inequality:

$$\left| \int_{\Omega} \frac{\partial \chi}{\partial t} \xi \right| \leq \frac{1}{2} \left( \|\chi_t\|_{L^2(\Omega)}^2 + \|\xi\|_{L^2(\Omega)}^2 \right). \tag{33}$$

The second one, using estimate Eq. 18 and Young’s inequality with  $\epsilon_1$ :

$$\begin{aligned} & \left| \int_{\Omega} \beta(u) \left( \frac{\nabla u}{|\nabla u|^{1-\gamma}} - \frac{\nabla \hat{u}}{|\nabla \hat{u}|^{1-\gamma}} \right) \nabla \xi \right| \\ & \leq \frac{2}{\gamma} M \int_{\Omega} |\nabla \chi|^\gamma |\nabla \xi| \end{aligned} \tag{34}$$

$$\leq \frac{2}{\gamma} M \left( \frac{1}{2\epsilon_1} \int_{\Omega} |\nabla \chi|^{2\gamma} + \frac{\epsilon_1}{2} \|\nabla \xi\|_{L^2(\Omega)}^2 \right), \tag{35}$$

where  $M = \|\beta(u)\|_{L^\infty(\Omega)}$ . The third one, using estimate 16 and Cauchy–Schwarz and Young’s inequalities with  $\epsilon_2$

$$\begin{aligned} & \left| \int_{\Omega} (\beta(u) - \beta(U)) \frac{\nabla \hat{u}}{|\nabla \hat{u}|^{1-\gamma}} \nabla \xi \right| \\ & \leq M^* \int_{\Omega} |u - U| |\nabla \hat{u}|^\gamma |\nabla \xi| \\ & \leq M^* \|\nabla \hat{u}\|_{L^\infty(\Omega)}^\gamma \|u - U\|_{L^2(\Omega)} \|\nabla \xi\|_{L^2(\Omega)} \\ & \leq M^* \|\nabla \hat{u}\|_{L^\infty(\Omega)}^\gamma \\ & \quad \times \left( \frac{1}{2\epsilon_2} \left( \|\chi\|_{L^2(\Omega)}^2 + \|\xi\|_{L^2(\Omega)}^2 \right) + \frac{\epsilon_2}{2} \|\nabla \xi\|_{L^2(\Omega)}^2 \right), \end{aligned} \tag{36}$$

where  $M^* = \alpha \max(\|\beta(u)\|_{L^\infty(\Omega)}, \|\beta(U)\|_{L^\infty(\Omega)})^{\alpha-1}$ .

From estimate Eq. 17, provided  $\nabla \hat{u}, \nabla U \in L^\infty(0, T, L^\infty(\Omega))$ , there exists a constant  $\epsilon_3 > 0$  such that

$$\gamma \epsilon \mathcal{A} \geq \epsilon_3 \quad \text{for all } t \in [0, T] \tag{37}$$

we can combine expressions 32–36 and choose  $\epsilon_1$  and  $\epsilon_2$  small enough to obtain that for some  $\bar{\epsilon} > 0$  and some constants  $C_i > 0$  ( $1 \leq i \leq 4$ )

$$\begin{aligned} & \frac{1}{2} \frac{\partial}{\partial t} \|\xi\|_{L^2(\Omega)}^2 + \bar{\epsilon} \|\nabla \xi\|_{L^2(\Omega)}^2 \\ & \leq C_1 \|\xi\|_{L^2(\Omega)}^2 + C_2 \|\chi_t\|_{L^2(\Omega)}^2 + C_3 \|\chi\|_{L^2(\Omega)}^2 \\ & + C_4 \int_{\Omega} |\nabla \chi|^{2\gamma} \end{aligned} \tag{38}$$

Since the second term on the left-hand side is nonnegative, we can use Gronwall’s Lemma in the previous expression to obtain that, for all  $t \in [0, T]$ ,

$$\begin{aligned} \|\xi(t)\|_{L^2(\Omega)}^2 & \leq C_5(T) \left( \|\xi(0)\|_{L^2(\Omega)}^2 + C_2 \int_0^T \|\chi_t\|_{L^2(\Omega)}^2 \right. \\ & \quad + C_3 \int_0^T \|\chi\|_{L^2(\Omega)}^2 \\ & \quad \left. + C_4 \int_0^T \int_{\Omega} |\nabla \chi|^{2\gamma} \right). \end{aligned} \tag{39}$$

Observe that

$$\begin{aligned} \|u - U\|_{L^2(\Omega)}^2 & \leq \|U - \hat{u}\|_{L^2(\Omega)}^2 + \|u - \hat{u}\|_{L^2(\Omega)}^2 \\ & = \|\xi(t)\|_{L^2(\Omega)}^2 + \|\chi(t)\|_{L^2(\Omega)}^2. \end{aligned} \tag{40}$$

and

$$\begin{aligned} \|\xi(0)\|_{L^2(\Omega)}^2 &= \|b - \hat{u}_0\|_{L^2(\Omega)}^2 \\ &\leq \|b - u_0\|_{L^2(\Omega)}^2 + \|u_0 - \hat{u}_0\|_{L^2(\Omega)}^2 \\ &\leq \|b - u_0\|_{L^2(\Omega)}^2 + \|\chi(0)\|_{L^2(\Omega)}^2 \end{aligned} \tag{41}$$

Thus, combining result 39, the two previous expressions and letting  $C = \max(C_5 C_i)$  for  $i = 2, 3, 4$ , we can establish the first portion of the statement of Theorem 28.

In order to complete the proof, we need to find a bound for the gradients. This is done by integrating expression 38 in time. Note that, on the left-hand side, we have

$$\frac{1}{2} \|\xi(T)\|_{L^2(\Omega)}^2 - \frac{1}{2} \|\xi(0)\|_{L^2(\Omega)}^2 + \epsilon \|\nabla u - \nabla U\|_{L^2(0,T,L^2(\Omega))}^2. \tag{42}$$

Since the first term is nonnegative, the following holds from the time integration of Eq. 38

$$\begin{aligned} \|\nabla \xi\|_{L^2(0,T,L^2(\Omega))}^2 &\leq C \left( \frac{1}{2} \|\xi(0)\|_{L^2(\Omega)}^2 + \int_0^T \|\xi\|_{L^2(\Omega)}^2 + \int_0^T \|\chi_t\|_{L^2(\Omega)}^2 \right. \\ &\quad \left. + \int_0^T \|\chi\|_{L^2(\Omega)}^2 + \int_0^T \int_{\Omega} |\nabla \chi|^{2\gamma} \right). \end{aligned} \tag{43}$$

Estimate 28 follows from the observation that

$$\begin{aligned} \|\nabla u - \nabla U\|_{L^2(0,T,L^2(\Omega))}^2 &\leq \|\nabla U - \nabla \hat{u}\|_{L^2(0,T,L^2(\Omega))}^2 + \|\nabla u - \nabla \hat{u}\|_{L^2(0,T,L^2(\Omega))}^2 \\ &\leq \|\nabla \xi(t)\|_{L^2(0,T,L^2(\Omega))}^2 + \|\nabla \chi(t)\|_{L^2(0,T,L^2(\Omega))}^2, \end{aligned}$$

and the combination of Eqs. 39 and 43. This concludes the proof.  $\square$

*Remark 2.3* Note that the error estimate collapses if condition 37 is not satisfied. This is the reason why we need to use both: the family of nondegenerate  $\{\beta_\epsilon(x)\} > \epsilon$  and the assumptions that  $\nabla u \in L^\infty(0, T, L^\infty(\Omega))$  and  $\nabla U \in L^\infty(0, T, L^\infty(\Omega))$ . The latter assumptions ensure that  $\mathcal{A} > 0$  since  $\nabla u \in L^\infty(0, T, L^\infty(\Omega))$  implies that  $\nabla \hat{u} \in L^\infty(0, T, L^\infty(\Omega))$  for an appropriate finite element space; see Lemma 2.2.

**Corollary 2.1** (Stability) *Under the conditions of Theorem 2.2, the method is stable.*

*Proof* Write  $U(t) = U(t) - u(t) + u(t)$  and use the triangle inequality to find

$$\|U(t)\|_{L^2(\Omega)} \leq \|U(t) - u(t)\|_{L^2(\Omega)} + \|u(t)\|_{L^2(\Omega)}$$

Using the previous theorem and assuming that  $u \in L^2(\Omega)$ , the result is immediate.  $\square$

**Corollary 2.2** *If  $u \in W^{1,\infty}(\Omega) \cap H^{k+1}(\Omega)$  is the solution of problem 14 and  $U$  is a solution of problem 21, constructed with piecewise polynomials of degree at most  $k$ , then, for all  $t \in [0, T]$ ,*

$$\begin{aligned} \|u - U\|_{L^2(\Omega)} + \|\nabla u - \nabla U\|_{L^2(0,T,L^2(\Omega))} &\leq C h^{k\gamma} \left( \int_0^T \|u\|_{H^{k+1}(\Omega)}^{2\gamma} \right)^{\frac{1}{2}}. \end{aligned} \tag{44}$$

*Proof* Assuming such regularity on  $u$ , the estimates given by Lemma 1.5 hold. Thus, by applying the result of Theorem 2.2, and using Hölder’s inequality with  $p = 1/\gamma \geq 1$  and  $p^* = 1/(1 - \gamma) \geq 1$  in the following expression,

$$\int_{\Omega} |\nabla \chi|^{2\gamma} \leq \left( \int_{\Omega} |\nabla \chi|^2 \right)^{\frac{2\gamma}{2}} |\Omega|^{1-\gamma} = \|\nabla \chi\|_{L^2(\Omega)}^{2\gamma} |\Omega|^{1-\gamma} \tag{45}$$

we obtain that, for small  $h$ ,

$$\begin{aligned} \|u - U\|_{L^2(\Omega)}^2 + \|\nabla u - \nabla U\|_{L^2(\Omega)}^2 &\leq h^{2k} \|u\|_{H^{k+1}(\Omega)}^2 \\ &\quad + C \left( h^{2(k+1)} \|u_0\|_{H^{k+1}(\Omega)}^2 + h^{2k\gamma} \int_0^T \|u\|_{H^{k+1}(\Omega)}^{2\gamma} \right) \\ &\leq C h^{2k\gamma} \int_0^T \|u\|_{H^{k+1}(\Omega)}^{2\gamma}, \end{aligned}$$

since  $\|b - u_0\|_{L^2(\Omega)}^2 \leq h^{2(k+1)} \|u_0\|_{H^{k+1}(\Omega)}^2$ . The result of the Corollary comes from the observation that, for  $p, q, s$  positive numbers,  $p^2 + q^2 \leq s^2$  implies that  $p + q \leq \sqrt{2} s$ .  $\square$

*Remark 2.4* Note that the error estimate obtained in Corollary 2.2 is constrained by the value of  $\gamma \in (0, 1)$ . In the hydraulic context,  $\gamma = 1/2$  for both Manning’s and Chézy’s formulas. According to Corollary 2.2, we can ensure  $O(h)$  convergence, for  $\gamma = 1/2$ , by using quadratic basis functions ( $k = 2$ ) to approximate a very regular solution  $u \in W^{1,\infty}(\Omega) \cap H^3(\Omega)$  of problem 14. In Section 3, we present numerical experiments that show that our analysis is conservative. In fact, we show that our method, implemented with piecewise

linear basis functions, can approximate the true solution as  $\mathcal{O}(h^2)$  in nondegenerate scenarios, and as  $\mathcal{O}(h)$  even when degeneracy happens in subregions of the domain  $\Omega$ .

**Lemma 2.1** (Boundedness of approximation) *Under the assumptions of Theorem 2.2, choosing  $\gamma \geq 1/2$ , and provided  $h$  is sufficiently small, if  $(u - z) \in L^\infty(0, T, L^\infty(\Omega)) \cap H^4(\Omega)$  and  $\|u - z\|_{L^\infty(0, T, L^\infty(\Omega))} \leq K_1$ , then  $\|U - z\|_{L^\infty(0, T, L^\infty(\Omega))} \leq K_1(1 + \epsilon_{k_1})$  for a small parameter  $\epsilon_{k_1}$ , and for a finite element approximation space consisting of at least piecewise cubic basis functions.*

*Proof* Write  $(U - z)$  as  $(U - \hat{u} + \hat{u} - z)$  to find

$$\|U - z\|_{L^\infty(0, T, L^\infty(\Omega))} \leq \|U - \hat{u}\|_{L^\infty(0, T, L^\infty(\Omega))} + \|\hat{u} - z\|_{L^\infty(0, T, L^\infty(\Omega))}. \tag{46}$$

Choose  $\hat{u}$  as in the definition Eq. 19, and choose at least  $\mathcal{M} = \mathcal{P}^3$  (piecewise cubic basis functions,  $k = 3$ ). The fact that the interpolant  $(\hat{u} - z) \in L^\infty(0, T, L^\infty(\Omega))$  provided  $(u - z) \in L^\infty(0, T, L^\infty(\Omega)) \cap H^4(\Omega)$ , for a bounded and smooth time-independent function  $z$ , is a consequence of Theorem 4.8.7 in [6]. Thus,  $\|\hat{u} - z\|_{L^\infty(0, T, L^\infty(\Omega))} \leq K_1$  for sufficiently small  $h$ . Now, from Corollary 2.2 and Lemma 1.6, we obtain

$$\begin{aligned} \|U - z\|_{L^\infty(0, T, L^\infty(\Omega))} &\leq \|U - \hat{u}\|_{L^\infty(0, T, L^\infty(\Omega))} + K_1 \\ &\leq K_0 h^{-1} \|U - \hat{u}\|_{L^\infty(0, T, L^2(\Omega))} + K_1 \\ &\leq K_0 h^{3\gamma-1} \left( \int_0^T \|u\|_{H^4(\Omega)}^{2\gamma} \right)^{\frac{1}{2}} + K_1 \end{aligned}$$

Thus, we can choose a sufficiently small  $h$  so that  $K_0 h^{1/2} \left( \int_0^T \|u\|_{H^4(\Omega)}^{2\gamma} \right)^{\frac{1}{2}} \leq \epsilon_{k_1} K_1$ , which implies

$$\|U - z\|_{L^\infty(0, T, L^\infty(\Omega))} \leq K_1(1 + \epsilon_{k_1})$$

This establishes the result of the lemma. □

**Lemma 2.2** (Boundedness of the gradient of the approximation) *Under the assumptions of Theorem 2.2, choosing  $\gamma \geq 1/2$ , and provided  $h$  is sufficiently small, if  $u \in L^\infty(0, T, L^\infty(\Omega)) \cap H^5(\Omega)$ ,  $\nabla u \in L^\infty(0, T, L^\infty(\Omega))$  and  $\|\nabla u\|_{L^\infty(0, T, L^\infty(\Omega))} \leq K_2$ , then  $\|\nabla U\|_{L^\infty(0, T, L^\infty(\Omega))} \leq K_2(1 + \epsilon_{k_2})$  for a finite element approximation space consisting of at least fourth-degree piecewise polynomial basis functions.*

*Proof* Similar to the one in Lemma 2.1. □

## 2.2 Fully discrete approximation

In this section, we will further proceed to consider discretization with respect to time. We will denote with  $dt$  the time step and with  $U^n$  the approximation of  $u(t)$  at time  $t = t_n = n dt$ . In order to write down the method, we will replace the time derivative in Eq. 21 with the quotient

$$\delta U^n = \frac{U^n - U^{n-1}}{dt}, \tag{47}$$

to obtain the following backward Euler scheme:

$$\left\{ \begin{aligned} (\delta U^n, v) + \left( \frac{\beta_\epsilon(U^n)}{|\nabla U^n|^{1-\gamma}} \nabla U^n, \nabla v \right) &= (f^n, v) \\ n dt \in [0, T], \quad \forall v \in \mathcal{M}, \\ \text{and} \quad (U^0, v) &= (u_0, v) \\ \forall v \in \mathcal{M}. \end{aligned} \right. \tag{48}$$

The previous expression defines  $U^n$  implicitly for  $U^{n-1}$  given, and it can be written as follows:

$$\begin{aligned} (U^n, v) + dt \left( \frac{\beta_\epsilon(U^n)}{|\nabla U^n|^{1-\gamma}} \nabla U^n, \nabla v \right) \\ = (U^{n-1} + dt f^n, v) \quad \forall v \in \mathcal{M}, \end{aligned} \tag{49}$$

or in matrix form, using the definitions explained in Section 2.1,

$$(G + dt B(\zeta^n)) \zeta^n = G \zeta^{n-1} + dt F(t_n), \tag{50}$$

where  $f^n$  (alt.  $F(t_n)$ ) is a known function (alt. matrix)

### 2.2.1 Fully discrete a priori error estimate

In this section, we will study how close solutions to the fully discrete Galerkin approximation problem 48 are to the true weak solution of problem 1. The proof follows immediately from the semidiscrete estimates, and it is presented for completeness. Once more, in the analysis, we will denote with  $\beta(x)$  any element of the sequence  $\{\beta_\epsilon(x)\}$  as described in Section 1.5.

**Theorem 2.3** *Let  $u \in W^{1,\infty}(\Omega)$  be the solution of problem 14 and let  $U^n$  be a solution of problem 48 at  $t_n = n dt$ . Let also  $\chi^n = u(t_n) - \hat{u}(t_n)$  be the approximation*



error between the interpolant and the true solution of problem 1. Then, for all  $t_n \in [0, T]$ ,

$$\begin{aligned} & \|u(t_n) - U^n\|_{L^2(\Omega)}^2 + \|\nabla u(t_n) - \nabla U^n\|_{L^2(\Omega)}^2 \\ & \leq C\|b - u_0\|_{L^2(\Omega)}^2 + \|\chi^0\|_{L^2(\Omega)}^2 \\ & \quad + dt C \sum_{j=1}^n \left( \|u_t(t_j) - \delta \hat{u}(t_j)\|_{L^2(\Omega)}^2 \right. \\ & \quad \left. + \|\chi^j\|_{L^2(\Omega)}^2 + \int_{\Omega} |\nabla \chi^j|^{2\gamma} \right) \end{aligned} \tag{51}$$

*Proof* Observe that, identifying  $u^n = u(t_n)$  and  $\hat{u}^n = \hat{u}(t_n)$ , a similar calculation to Eq. 30 yields

$$\begin{aligned} & (\delta(U^n - \hat{u}^n), v) + \left( \beta(U^n) \left( \frac{\nabla U^n}{|\nabla U^n|^{1-\gamma}} - \frac{\nabla \hat{u}^n}{|\nabla \hat{u}^n|^{1-\gamma}} \right), \nabla v \right) \\ & = (u_t^n - \delta \hat{u}^n, v) + \left( \beta(u^n) \left( \frac{\nabla u^n}{|\nabla u^n|^{1-\gamma}} - \frac{\nabla \hat{u}^n}{|\nabla \hat{u}^n|^{1-\gamma}} \right), \nabla v \right) \\ & \quad + \left( (\beta(u^n) - \beta(U^n)) \frac{\nabla \hat{u}^n}{|\nabla \hat{u}^n|^{1-\gamma}}, \nabla v \right) \end{aligned} \tag{52}$$

Now, let  $\xi^n = U^n - \hat{u}(t_n)$  and  $\chi^n = u(t_n) - \hat{u}(t_n)$  and choose  $v = \xi^n$ . Using estimates Eqs. 34 and 36, we obtain

$$\begin{aligned} & (\delta \xi^n, \xi^n) + \gamma(\beta(U^n)) \mathcal{A} \|\nabla \xi^n\|_{L^2(\Omega)}^2 \\ & \leq (u_t^n - \delta \hat{u}^n, \xi^n) + \frac{1}{\epsilon_2} C \|\xi^n\|_{L^2(\Omega)}^2 + C(\epsilon_2 + \epsilon_1) \|\nabla \xi^n\|_{L^2(\Omega)}^2 \\ & \quad + \frac{1}{2\epsilon_2} \|\chi^n\|_{L^2(\Omega)}^2 + \frac{1}{2\epsilon_1} C \int_{\Omega} |\nabla \chi^n|^{2\gamma}. \end{aligned} \tag{53}$$

If condition 37 is satisfied, then for  $\epsilon_1$  and  $\epsilon_2$  small enough, the previous expression implies

$$\begin{aligned} \|\xi^n\|_{L^2(\Omega)}^2 & \leq (\xi^{n-1}, \xi^n) + dt (u_t^n - \delta \hat{u}^n, \xi^n) \\ & \quad + dt C \|\xi^n\|_{L^2(\Omega)}^2 + dt C \|\chi^n\|_{L^2(\Omega)}^2 \\ & \quad + dt C \int_{\Omega} |\nabla \chi^n|^{2\gamma}, \end{aligned} \tag{54}$$

which yields

$$\begin{aligned} & (1 - dt C) \|\xi^n\|_{L^2(\Omega)}^2 \\ & \leq \|\xi^{n-1}\|_{L^2(\Omega)}^2 + dt \|u_t^n - \delta \hat{u}^n\|_{L^2(\Omega)}^2 + dt C \|\chi^n\|_{L^2(\Omega)}^2 \\ & \quad + dt C \int_{\Omega} |\nabla \chi^n|^{2\gamma}. \end{aligned} \tag{55}$$

Using the Taylor expansion for  $(1 - dt C)^{-1}$  around zero, we can rewrite the previous expression for small  $dt$ :

$$\|\xi^n\|_{L^2(\Omega)}^2 \leq (1 + dt C) \|\xi^{n-1}\|_{L^2(\Omega)}^2 + dt C R_n, \tag{56}$$

where

$$R_n = \|u_t^n - \delta \hat{u}^n\|_{L^2(\Omega)}^2 + \|\chi^n\|_{L^2(\Omega)}^2 + \int_{\Omega} |\nabla \chi^n|^{2\gamma}$$

Making use of inequality 56 repeatedly, we find that

$$\begin{aligned} \|\xi^n\|_{L^2(\Omega)}^2 & \leq (1 + dt C)^n \|\xi^0\|_{L^2(\Omega)}^2 \\ & \quad + dt C \sum_{j=1}^n (1 + dt C)^{n-j} R_j \\ & \leq C \|\xi^0\|_{L^2(\Omega)}^2 + dt C \sum_{j=1}^n R_j \\ & \quad \text{for } t_n = dt n \in [0, T], \end{aligned} \tag{57}$$

which, together with

$$\begin{aligned} \|\xi^0\|_{L^2(\Omega)}^2 & = \|b - \hat{u}^0\|_{L^2(\Omega)}^2 \leq \|b - u_0\|_{L^2(\Omega)}^2 + \|u_0 - \hat{u}^0\|_{L^2(\Omega)}^2 \\ & \leq \|b - u_0\|_{L^2(\Omega)}^2 + \|\chi^0\|_{L^2(\Omega)}^2, \end{aligned}$$

establish the first portion of the statement of the Theorem. In order to find the appropriate estimate on the gradients, a similar argument to the one in the proof of the continuous in time a priori error estimate has to be used. Both estimates lead to the result of the theorem.  $\square$

**Corollary 2.3** *If  $u \in W^{1,\infty}(\Omega) \cap H^{k+1}(\Omega)$  is the solution of problem 14,  $U^n$  is a solution of problem 48 at  $t_n = n dt$ , constructed with piecewise polynomials of degree at most  $k$ . Then, for all  $t_n \in [0, T]$ ,*

$$\begin{aligned} & \|u(t_n) - U^n\|_{L^2(\Omega)} + \|\nabla u(t_n) - \nabla U^n\|_{L^2(\Omega)} \\ & \leq C(u, t_n) (dt + h^{k\gamma}) \end{aligned} \tag{58}$$

*Proof* Given the regularity of  $u$ , the quantity  $R_j$  can be bounded using the estimates given by Lemma 1.5,

the fact that  $\|u_t^n - \delta \hat{u}^n\|_{L^2(\Omega)}^2 \leq C^*(u) dt^2$  (see [31]), and applying the results of Theorem 2.3. Thus,

$$\begin{aligned} & \|u(t_n) - U^n\|_{L^2(\Omega)}^2 + \|\nabla u(t_n) - \nabla U^n\|_{L^2(\Omega)}^2 \\ & \leq C \|b - u_0\|_{L^2(\Omega)}^2 + h^{2(k+1)} \|u_0\|_{H^{k+1}(\Omega)}^2 \\ & \quad + dt C \sum_{j=1}^n \left( C^*(u) dt^2 + h^{2k} \|u(t_n)\|_{H^{k+1}(\Omega)}^2 + h^{2k\gamma} \|u(t_n)\|_{H^{k+1}(\Omega)}^{2\gamma} \right) \\ & \leq C(u, t_n) (dt + h^{k\gamma})^2. \end{aligned} \tag{59}$$

Recall that, for  $p, q, s$  positive numbers,  $p^2 + q^2 \leq s^2$  implies that  $p + q \leq \sqrt{2} s$ . Thus, the result of the Corollary follows immediately.  $\square$

### 3 Numerical experiments

A lumped mass continuous Galerkin code with piecewise linear basis functions was implemented in order to perform numerical experiments aimed at investigating: (1) the accuracy and validity of the numerical method previously described to solve the DSW equation for the case when  $z = 0$ , see Section 3.1; (2) the applicability of the DSW equation as a model to simulate observed quantities (such as water discharge and depth profile) in an experimental setting for a prescribed bathymetry  $z = z(x) \neq 0$ , see Section 3.2; and (3) the qualitative properties of solutions to the DSW in its general form (Eq. 1).

Even though Lemmas 2.1 and 2.2 suggest that one should use at least fourth-order polynomial basis functions in order to ensure the boundedness of both  $U - z$  and  $\nabla U$  and, thus, convergence of the numerical scheme, in practice, we found that the use of piecewise linear basis functions was appropriate. Furthermore, our numerical experiments showed that the implementation of the regularizing Lipschitz functional  $\beta_\epsilon(u) > \epsilon$  (instead of  $u - z$ ), as described in Section 1.5, was not necessary. In fact, we found that the stability of the code behaved similarly with or without the implementation of the  $\beta_\epsilon(u) > 0$ , and generally,  $\|U_\epsilon - U\|_{L^2(\Omega)} \sim \mathcal{O}(\epsilon)$ , with  $U_\epsilon$ , the numerical solution of the nondegenerate problem 14, and  $U$  the numerical solution of the possibly degenerate problem 1.

A lumped mass approach was chosen since, for well behaved nonlinear parabolic equations, such schemes satisfy a maximum principle and, thus, provide a monotone and physically consistent way to approach the

solution, see [31]. For computational purposes, problem 50 was approximated by the semi-implicit scheme

$$(G + dt B(\zeta_{(l)})) \zeta^n = G \zeta^{n-1} + dt F(t_n). \tag{60}$$

A Picard iteration approach

$$\zeta_{(l)} = (G + dt B(\zeta_{(l-1)}))^{-1} (G \zeta^{n-1} + dt F(t_n)),$$

with an initial guess  $\zeta_{(0)} = \zeta^{n-1}$  was used in order to resolve the nonlinearity, with the assumption that

$$\lim_{l \rightarrow \infty} (G + dt B(\zeta_{(l)})) = (G + dt B(\zeta^n)). \tag{61}$$

In practice, the iteration process was stopped when  $\|\zeta_{(l)} - \zeta_{(l-1)}\|_{L^2(\Omega)}^2 \leq \tau$  for a prescribed tolerance  $\tau$ , and  $\zeta^n$  was set equal to the value of  $\zeta_{(l)}$  in the last iteration. In all our experiments, we chose  $\alpha = 5/3$  and  $\gamma = 1/2$  (as in [37] and [14]). These values correspond to Manning’s formula. Numerical studies addressing the case when  $z = 0$ ,  $\alpha \geq 1$  and  $\gamma = 1$ , which corresponds to the PME, include, for example, [25] and [20].

#### 3.1 Convergence analysis

As suggested in the previous sections, convergence of the numerical method proposed to approximate the DSW equation may fail if the depth  $u - z$  is zero or if its gradient  $\nabla u$  is unbounded. In order to investigate the accuracy and validity of the numerical method in different circumstances, we chose two approaches: the first one consisted of reproducing a Lipschitz continuous compactly supported solution of Eq. 1 presented in [10] for the 1-D case, for  $z = 0$ ,  $f = 0$ , and  $\Omega = \mathbb{R}$ . The second one consisted of choosing a simple function  $u(x, t) \geq 0$  with unbounded gradient at  $x = 0$  that we used to create a right-hand-side  $f(x, t)$ . This was done by applying the differential operator defined by the left-hand side of Eq. 1. We then used our method to approximate this  $u$ . For both cases, we obtained convergence rates for a variety of scenarios and present them in the following paragraphs.

*Remark 3.1* Despite the fact that Corollary 2.3 would only ensure convergence results of the type,  $\|u(t_n) - U^n\|_{L^2(\Omega)} \leq C(u, t_n) (dt + h^{1/2})$  for  $\gamma = 1/2$ , and non-degenerate solutions  $u \in H^2(\Omega)$  using piecewise linear basis functions ( $k = 1$ ), we chose the time step comparable to or smaller than the square of the grid diameter,  $dt \lesssim h^2 = dx^2$ , in our convergence experiments. This was done in order to investigate if optimality in the convergence rates (i.e.,  $\|u(t_n) - U^n\|_{L^2(\Omega)} \leq C(u, t_n) (dt + h^2)$ ) could be achieved for, say, nondegenerate solutions with bounded gradients. Intuitively, under these

conditions, Eq. 1 should resemble a well behaved non-linear parabolic problem.

### 3.1.1 Compactly supported solution

The explicit expression of a Barenblatt solution for Eq. 1 is presented in [10] and given by

$$u(x, t) = t^{-\frac{1}{\gamma(m+1)}} \left[ C - k(m, \gamma) |\Phi|^{\frac{\gamma+1}{\gamma}} \right]_+^{\frac{\gamma}{m\gamma-1}}, \tag{62}$$

where  $[s(x)]_+$  denotes the positive part of  $s(x)$ ,  $m = 1 + \alpha/\gamma$ ,  $C$  is a positive function related to the initial mass  $M$ , given by

$$M = \int_{-\infty}^{\infty} u(x, t) dx,$$

and

$$k(m, \gamma) = \frac{m\gamma - 1}{m(\gamma + 1)} \left( \frac{1}{\gamma(m+1)} \right)^{\frac{1}{\gamma}}, \quad \text{and} \quad \Phi = x t^{-\frac{1}{\gamma(m+1)}}$$

The function given by Eq. 62 is Lipschitz continuous and compactly supported; thus, it is almost everywhere differentiable and its gradient is bounded wherever it exists. By changing the numerical domain  $\Omega$ , we managed to verify how the numerical method approximates the solution both when it is globally not degenerate,  $u > 0$ , and when it degenerates in some subsets of  $\Omega$ . For our calculations in the degenerate case, we chose a numerical interval  $[-L, L]$  big enough so that the free boundary would always be inside the domain  $\Omega$  for  $t \in [t_0, t_f]$ . The results are shown in Table 1, panel a, for

**Table 1** Convergence rates to approximate Barenblatt solutions for  $\alpha = 5/3$  and  $\gamma = 1/2$  using  $t_0 = 2$  and  $t_f = 2.1$

$dt$	$dx$	$\ U - u\ _{L^2(\Omega)}$	Conv. rate
(a) Nondegenerate case			
1/50	1	$6.34 \times 10^{-3}$	.
1/50	1/2	$1.79 \times 10^{-3}$	1.82
1/100	1/4	$4.87 \times 10^{-4}$	1.88
1/200	1/8	$1.21 \times 10^{-4}$	2.00
1/400	1/16	$2.88 \times 10^{-5}$	2.08
1/1000	1/32	$7.37 \times 10^{-6}$	1.97
1/4000	1/64	$1.87 \times 10^{-6}$	1.98
(b) Degenerate case			
1/50	1	$1.32 \times 10^{-1}$	.
1/50	1/2	$8.39 \times 10^{-2}$	0.65
1/100	1/4	$3.97 \times 10^{-2}$	1.08
1/200	1/8	$2.47 \times 10^{-2}$	0.69
1/400	1/16	$1.36 \times 10^{-2}$	0.85
1/1000	1/32	$7.83 \times 10^{-3}$	0.80
1/4000	1/64	$4.74 \times 10^{-3}$	0.72

(a)  $\Omega = [-5, 5]$  degenerate case and (b)  $\Omega = [-2, 2]$  nondegenerate case

the nondegenerate case, and in Table 1, panel b, for the degenerate case. In all cases,  $t_0 = 2$  and  $t_f = 2.1$ , and the Picard iteration scheme was run until the tolerance value met the condition  $\tau \leq 10^{-10}$ . We computed the numerical solutions and compared them to the true solution using  $dt = 1/10$  and  $dx = 1/20$  to produce the plots in Fig. 2. The computed mass  $M$  of the numerical solution was calculated as a function of time  $t$ , and it was observed to be close to the constant  $M = 5.8465$ , which corresponds to the value of the mass of the true solution Eq. 62 in the time interval  $t \in [1, 2.5]$ .

### 3.1.2 Artificial right-hand side

The function  $u(x, t) \geq 0$  with unbounded gradient at  $x = 0$  that we used to create a right-hand-side  $f(x, t)$  by applying the differential operator defined by the left-hand side of Eq. 1 was:

$$u = \begin{cases} (100 - t^2)\sqrt{x} & \text{if } x \geq 0, \\ 0 & \text{if } x < 0. \end{cases} \tag{63}$$

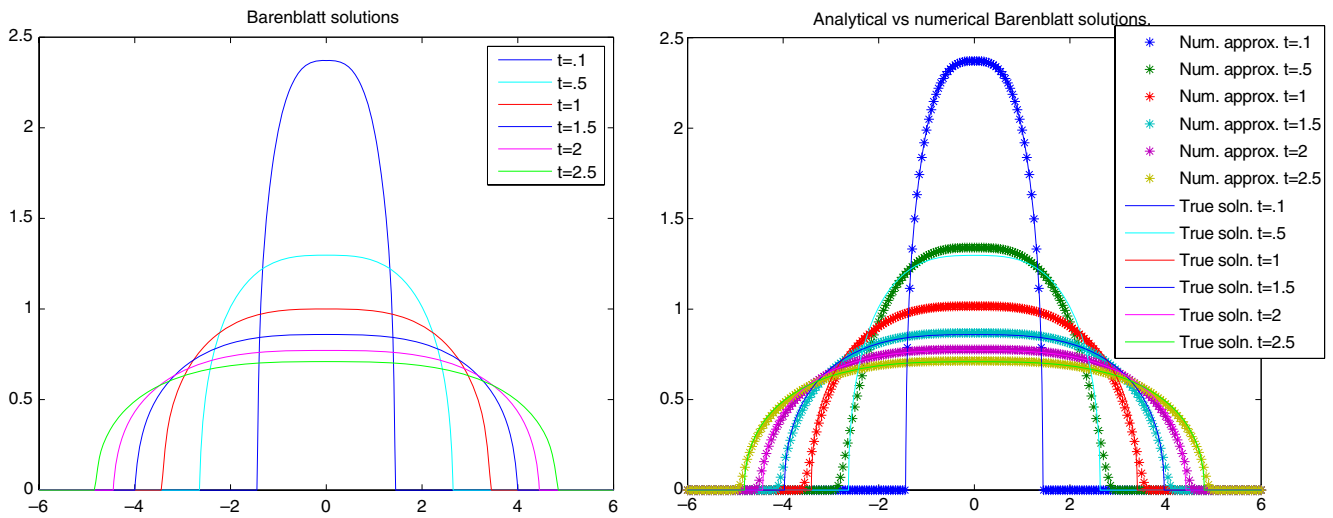
with  $\alpha = \frac{5}{3}$  and  $\gamma = \frac{1}{2}$ , the right-hand-side  $f$  is given by

$$f = \frac{\partial u}{\partial t} - \frac{\partial}{\partial x} \left( \frac{u^{\frac{5}{3}}}{|\frac{\partial u}{\partial x}|^{\frac{1}{2}}} \frac{\partial u}{\partial x} \right) = \begin{cases} -2t\sqrt{x} - \frac{7}{12\sqrt{2}} (100 - t^2)^{\frac{13}{6}} x^{-\frac{5}{12}} & \text{if } x \geq 0, \\ 0 & \text{if } x < 0, \end{cases}$$

Examples of convergence rates between the numerical solution and the true solution Eq. 63 in three different domains  $\Omega$  can be found in Table 2, panels a, b, and c. All results were calculated with  $t_0 = 9$  and  $t_f = 9.1$ , and the Picard iteration scheme was run until the tolerance value met the condition  $\tau \leq 10^{-10}$ . Table 2, panel a, shows results for  $\Omega = [5, 10]$ , where the solution  $u$  is not degenerate and its gradient is bounded. Table 2, panel b, shows results corresponding to  $\Omega = [0, 5]$ . In this case, the gradient of the solution  $u$  is unbounded and  $u = 0$  at  $x = 0$ . Finally, for  $\Omega = [-0.5, 4.5]$ , the results are shown in Table 2, panel c. Here, the solution  $u$  degenerates in the interval  $[-0.5, 0]$  and the gradient of  $u$  is unbounded at  $x = 0$ .

### 3.1.3 Convergence results discussion

Corollary 2.3 establishes that the numerical scheme will approximate the true solution as  $\mathcal{O}(h^{k\gamma})$ , whenever piecewise polynomials of degree at most  $k$  are used and the true solution  $u \in H^{k+1}(\Omega)$  is nondegenerate and its gradient is bounded. In Tables 1, panel a, and 2, panel a, we show the convergence rates of the numerical



**Fig. 2** Comparison of numerically simulated compactly supported solutions for  $\alpha = 5/3$  and  $\gamma = 1/2$ . (Left) Analytical Barenblatt solutions ( $z = 0$ ), (right) numerical vs analytical solutions ( $z = 0$ )

method in cases when the true solution is nondegenerate and its gradient is bounded. Note that the numerical solution converges quadratically to the true solution ( $\mathcal{O}(h^2)$ ) for piecewise linear basis functions ( $k = 1$ ), suggesting that, under these conditions, Eq. 1 becomes a *well behaved* nonlinear parabolic problem and the expected convergence rates from Corollary 2.3 appear to be too conservative. On the other hand, when the

solution degenerates at one point and the gradient is unbounded (thus, the conditions to ensure convergence according to Corollary 2.3 are not satisfied), optimality is lost but the numerical scheme still converges linearly to the true solution ( $\mathcal{O}(h)$ ), see Table 2, panel b. Close to linear convergence is observed when degeneracy happens and the gradient of the solution is discontinuous, see Table 1, panel b. The worst-case scenario takes place when degeneracy happens in an interval and the gradient of the solution is unbounded, see Table 2, panel c. In this case the order of convergence seems to behave close to  $\mathcal{O}(h^{2/5})$ . It is important to mention at this point that the regions of the domain where the numerical approximation differs mostly from the true solution correspond to regions where the gradient is discontinuous, which generally takes place near the free boundary. We show this fact in Fig. 3.

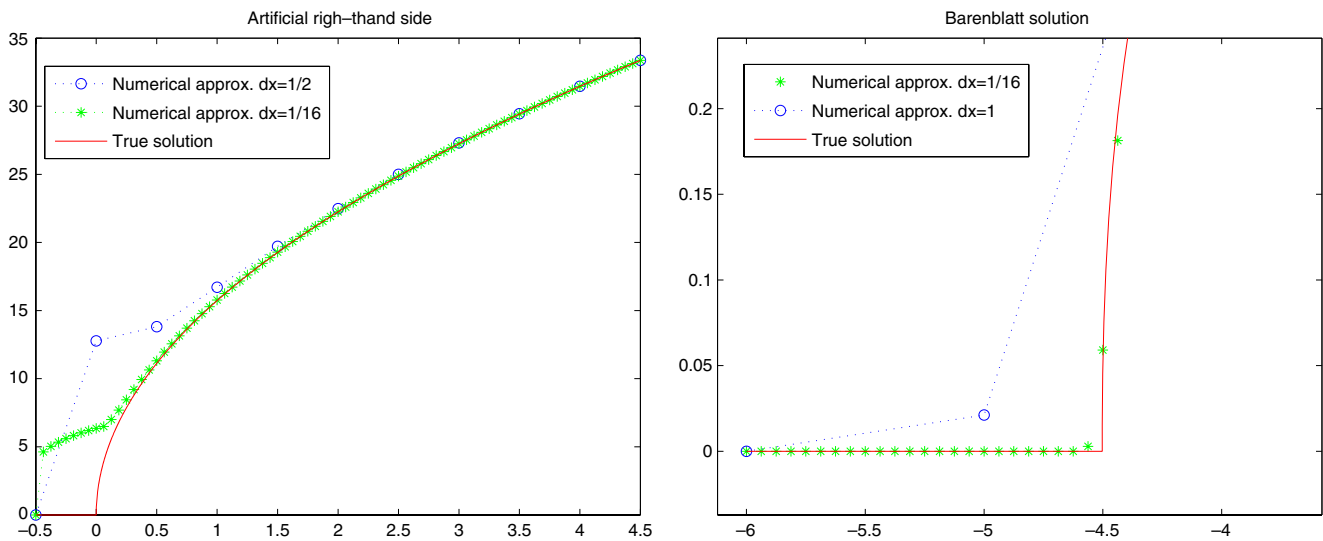
**Table 2** Convergence rates, with  $t_0 = 9$  and  $t_f = 9.1$

$dt$	$dx$	$\ U - u\ _{L^2(\Omega)}$	Conv. rate
(a) Nondeg. & bdd gradient, $\Omega=(5,10)$			
1/50	1	$1.28 \times 10^{-1}$	.
1/50	1/2	$3.52 \times 10^{-2}$	1.86
1/100	1/4	$9.57 \times 10^{-3}$	1.88
1/200	1/8	$2.62 \times 10^{-3}$	1.87
1/400	1/16	$7.53 \times 10^{-4}$	1.80
1/1000	1/32	$2.16 \times 10^{-4}$	1.80
1/4000	1/64	$5.44 \times 10^{-5}$	1.99
(b) Unbdd gradient, $\Omega=(0,5)$			
1/50	1	2.83	.
1/50	1/2	1.48	0.93
1/100	1/4	$7.85 \times 10^{-1}$	0.92
1/200	1/8	$4.05 \times 10^{-1}$	0.95
1/400	1/16	$2.04 \times 10^{-1}$	0.98
1/1000	1/32	$1.03 \times 10^{-1}$	1.00
1/4000	1/64	$5.14 \times 10^{-2}$	1.00
(c) Deg. & unbdd gradient, $\Omega = (-0.5, 4.5)$			
1/50	1/2	4.99	.
1/100	1/4	3.61	0.47
1/200	1/8	2.64	0.45
1/400	1/16	1.96	0.43
1/1000	1/32	1.48	0.41
1/4000	1/64	1.11	0.41

### 3.2 Validation: Iwagaki's experiment

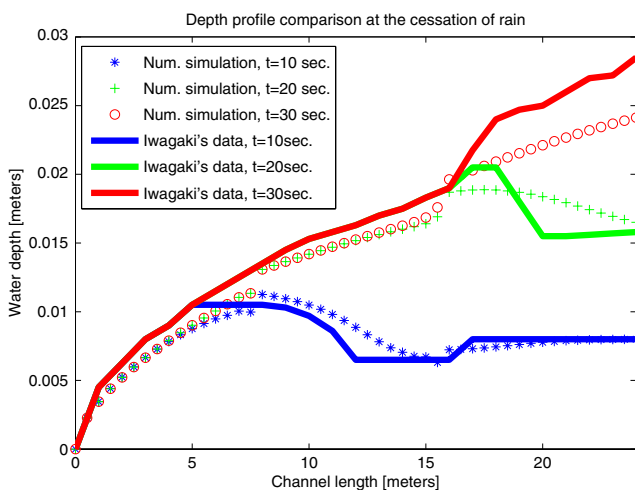
In order to validate our code in the case when  $z \neq 0$ , we chose to numerically reproduce a set of laboratory experiments conducted by Iwagaki [22]. This approach was followed by Zhang and Cundy [37] and by Feng and Molz [14]. We used the friction parameters reported by Iwagaki and no calibration was pursued. The numerical simulations were obtained using a regular mesh with element diameter  $\Delta x = 1/2$  m and a time step  $\Delta t = 1/2$  s.

Iwagaki's experiments were designed to produce unsteady flows in a channel 24 m long with a cross section of 19.6 cm. The channel was divided into three sections of equal length (8 m) and different slopes



**Fig. 3** Comparison of numerically computed solutions (*left*) (with initial condition  $u_0 = 19\sqrt{x}$  at  $t = 9$ ) vs true solutions  $u = (100 - t^2)\sqrt{x}$  at time  $t = 9.2$  using an artificial right-hand-side  $\Omega = [-0.5, 4.5]$  (*right*) Numerical approx. of Barenblatt  $\Omega = [-6, 6]$

( $\theta = 0.02, 0.015, 0.01\%$ ) each. During experiments, three different rainfall intensities ( $f = 0.108, 0.064,$  and  $0.80$  cm/s) were simultaneously applied to each section for different time periods ( $t = 10, 20$  and  $30$  s). Figure 4 shows snapshots of the depth profiles along the domain, both measured and numerically simulated, at the cessation of three different rain events lasting  $t = 10, 20,$  and  $30$  s, respectively. We can see that, overall, the relevant qualitative nature of the phenomena is captured in the numerical simulations. In Fig. 5,

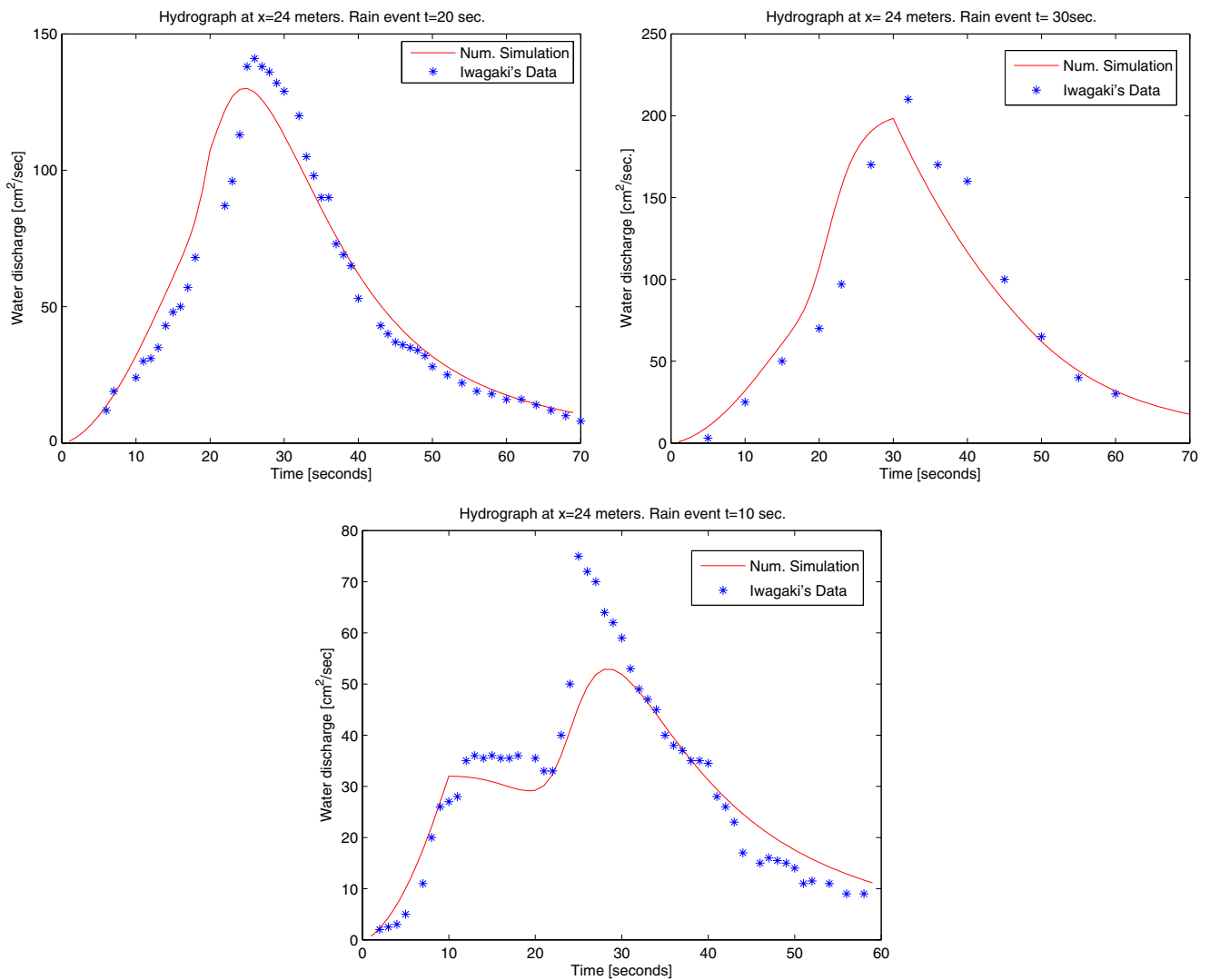


**Fig. 4** Comparison of the numerically calculated depth profiles at the cessation of rainfall with Iwagaki's experimental results on a three plane cascade. *Solid lines* are experimental results and *dashed points* are numerical results

the simulated water discharge  $q = Vh$  as a function of time  $t$ , at the lowest end of the domain,  $x = 24$  m, was plotted and compared to the experimental data for three different rain events.<sup>1</sup> In Fig. 5, top left, the agreement of the hydrograph with the experimental data is very good. The inability to get the full maximum of the curve and the extra spread in the curve by the numerical results is clearly one of the limitations of the diffusive wave approximation. Nevertheless, the time when a maximum discharge is achieved is matched well by the simulation. In Fig. 5, top right, the agreement of the hydrograph with the experimental data is good. In fact, we can observe that the breakthrough time in the simulation is smaller; however, the area under the curve corresponding to the water discharge volume seems to be in good agreement with the experimental data. Figure 5, bottom, shows the hydrograph for a 10-s rain event. Breakthrough times, as well as the main qualitative behavior, are captured by the simulation even though the maximum values are not accurately reproduced due to the diffusive nature of the approximation. It must be emphasized that, in a flooding event, the most relevant pieces of information obtained from hydrographs are the breakthrough time and the overall discharged water volume, which are nicely modeled by the numerical simulation. Moreover, it is interesting to

<sup>1</sup>It is important to note that the discharge  $q$  is calculated using the expression inside the divergence term in the first equation of the IBVP (Eq. 1).





**Fig. 5** Comparison of the numerically calculated hydrographs at  $x = 24$  m with Iwagaki’s experimental results on a three-plane cascade, during a rain event of (top left) 20 s, (top right) 30 s, and

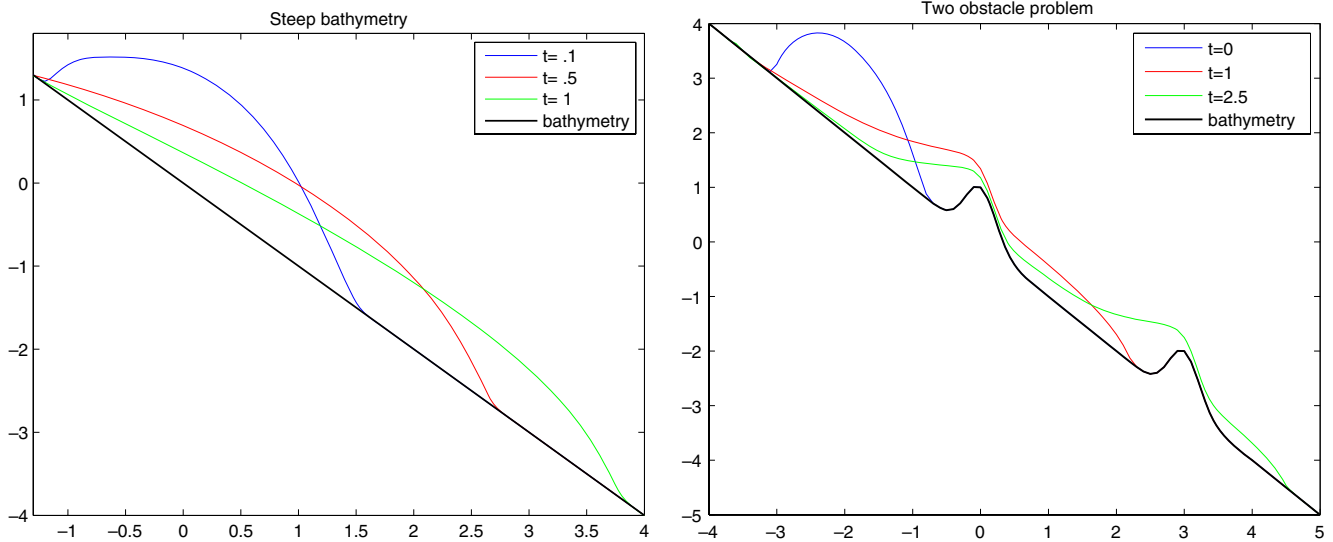
(bottom) 10 s. Solid lines are numerical results and dashed points are experimental results

observe that, despite the fact that, in the derivation of the DSW equation, one assumes uniform flow conditions, the numerical results match reasonably well with the unsteady flow experimental measurements.

### 3.3 Qualitative properties of solutions

In this section, we present qualitative properties of solutions to the DSW when the topographic effects are not neglected ( $z \neq 0$ ). Our findings are based on numerically simulated solutions obtained with our code. Our aim was to investigate if the properties of solutions of the DSW equation found in the 1-D case when  $z = 0$  in [3, 10], and [1] persist in a more general setting

when the bathymetry  $z$  is a smooth and bounded time-independent function. Properties such as boundedness and existence of compactly supported solutions, finite speed of propagation of disturbances, and extinction in finite time were found to persist for regular  $z \neq 0$ . We did not pursue any convergence analysis for this case due to the lack of an analytic expression for true solutions of the IBVP (Eq. 1). We present snapshots of solutions at different times for different bathymetries  $z(x) \neq 0$  in Figs. 6 and 7. In all our numerical experiments, we used  $\Delta x = 1/20$  and  $\Delta t = 1/20$ . Observe in particular that, in Fig. 7 (right), the water depth reaches equilibrium without the appearance of sloshing. This is a consequence of the diffusive wave approximation



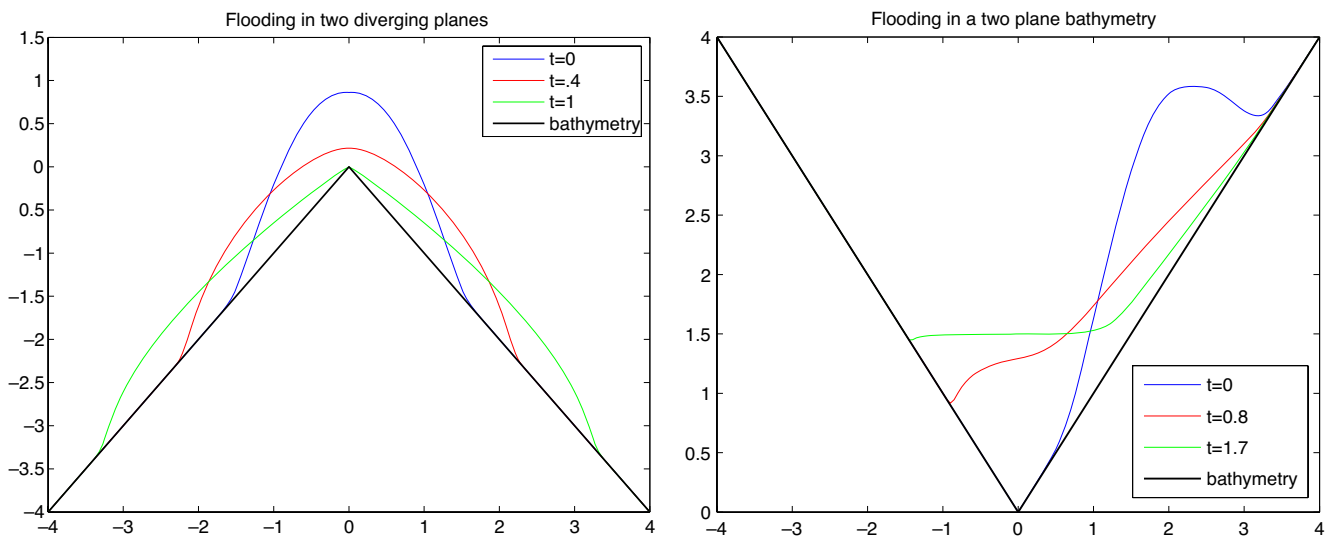
**Fig. 6** Evolution of a compactly supported initial condition simulating a well localized rain event. (Left) Initial condition  $u_0 = z + [-2(x - 1)(x + 1)]_+$  on a steep bathymetry, (right) initial con-

dition  $u_0 = z + [-2(x + 3)(x + 1)]_+$  on an inclined bathymetry with two obstacles

and would model the flow of a very viscous fluid flow or a flow dominated by bottom friction. In Figs. 6 and 7, we observe that the DSW equation has some regularization effect on the solution in regions where the problem is not degenerate. However, the regularity of the solution  $u(x)$  corresponds to that of  $z(x)$  whenever the problem becomes degenerate (i.e., when  $u - z = 0$ ).

### 4 Conclusions

In this study, we have shown the results of a numerical approach to study the properties of solutions of the IBVP (Eq. 1). Our emphasis was placed in analyzing the mathematical properties of the partial differential equation appearing in the IBVP (Eq. 1), the DSW equation, in order to find estimates for the



**Fig. 7** Evolution of a compactly supported initial condition simulating a well localized rain event. (Left) Initial condition  $u_0 = z + 2 \exp(-2x^2)$  on two diverging planes, (right) initial condition  $u_0 = z + 2 \exp(-2(x - 2)^2)$  on a two-plane bathymetry

error between numerical solutions constructed using the Galerkin method and true solutions of this equation. We proved that the numerical solutions converge to the true solution of the DSW equation under certain physically consistent conditions. These consist of requiring that the solution  $u \in L^\infty(0, T; L^\infty(\Omega))$  and  $\nabla u \in L^\infty(0, T; L^\infty(\Omega))$ . Furthermore, we showed in Lemmas 2.1 and 2.2 that these conditions ensure the boundedness of the discrete solution  $U$  and its gradient  $\nabla U$  for particular finite element approximation spaces.

Our analytical a priori error estimates are not optimal. In particular, the absence of appropriate conditions leading to a proof to ensure uniqueness of solutions of the DSW equation, in its general form (Eq. 1), imposes restrictions in our analysis. We presented numerical evidence that shows that the proposed numerical method converges to the true weak solution even when the conditions for Theorem 2.2 to hold are not met, for example, when the true solution degenerates and its gradient is unbounded. Furthermore, we found that, in regions where the solution does not degenerate ( $u - z \neq 0$ ), the method reaches optimal convergence rates. We showed that, despite the fact that the IBVP (Eq. 1) has not been fully studied analytically when  $z \neq 0$ , properties such as boundedness and existence of compactly supported solutions, finite speed of propagation of disturbances, and extinction in finite time found in the 1-D case when  $z = 0$  in [3, 10], and [1], persist for a bounded and smooth bathymetry  $z \neq 0$ , based on our numerical solutions.

For studies addressing the applicability of the DSW equation as a model to simulate shallow water flow, instead of the full Saint Venant (or shallow water) equations in experimental and real life settings, we refer the reader, for example, to the works of Ponce et al. [27] and [26] in the 1-D case and to the references mentioned in Section 1.3 in 2-D cases. In our study, we show solutions of the DSW equation that locally violate some of the essential assumptions in the flow regimes used to derive the DSW from the Navier–Stokes equations. Such is the case for the family of compactly supported Barenblatt solutions exhibited when  $z = 0$  in Section 3.1. The gradient of these solutions (water surface slope) is not comparable to the gradient of the bathymetry  $\nabla z$  close to the free boundary. A more extreme case of solutions that violate the uniform flow conditions happen even when the DSW equation becomes the PME ( $z = 0$  and  $\gamma = 1$ ) in two or higher dimensions. In Chapter 19 in [33], the author shows that there exists a class of solutions called *focusing solutions* that exhibit no local regularity on the gradient

in subsets of  $\Omega$ . The existence of these kinds of solutions serves as a reminder of the limitations of using the DSW equation as a hydrodynamical model.

One interesting fact about solutions of the DSW is the following: When one sees the DSW equation as a conservation law with respect to the depth  $u^* = u - z$ , it becomes

$$\frac{\partial u^*}{\partial t} - \nabla \cdot (u^* V) = f,$$

where the horizontal velocity  $V$  is given by Eq. 6, and its magnitude is

$$|V| = \frac{|u^*|^{\alpha-1}}{c_f} |\nabla u|^\gamma,$$

which indicates that, at the free boundary (interface between regions where  $u^* > 0$  and  $u^* = 0$ ) or any place in the domain where the depth of the water  $u^*$  is zero, the magnitude of the velocity is zero since  $\alpha > 1$  (and provided  $\nabla u$  does not tend to infinity faster than  $u^*$  tends to zero).

In Section 3.2, we show that, despite the limitations of the diffusive wave approximation nature of the DSW equation, the main qualitative behavior of water flow in an experimental setting designed to produce unsteady flows, such as breakthrough time and discharge, were captured by the simulation using the DSW as a model. However, the maximum values of the water depth are not accurately reproduced due to the diffusive nature of the approximation.

#### 4.1 Future work

From a purely mathematical point of view, the inadequacy of the continuous Galerkin method as an approximation technique is seen in locations close to the free boundary; see Fig. 3. Future studies should include the investigation of alternative techniques that are capable of capturing sharper fronts such as the discontinuous Galerkin approach or other stabilized methods.

**Acknowledgements** This work was supported in part by the National Science Foundation, Project No. DMS-0411413, and DMS-0620697, and *Centro de Investigación en Geografía y Geomática, “Ing. Jorge L. Tamayo”, A.C.*

Appendix

**Lemma 1** *The operator  $\mathcal{A}(x) : \mathbb{R}^n \rightarrow \mathbb{R}^n$  defined by*

$$\mathcal{A}(x) = \frac{x}{|x|^{1-\gamma}} \tag{64}$$

*is monotone, i.e., for any  $x, y \in \mathbb{R}^n$*

$$(\mathcal{A}(x) - \mathcal{A}(y)) \cdot (x - y) \geq 0.$$

*Proof* Define the function  $\mathcal{B}(x) : \mathbb{R}^n \rightarrow \mathbb{R}$  as

$$\mathcal{B}(x) = |x|^{\gamma+1} \quad \text{where} \quad |x| = \left( \sum_{j=1}^n x_j^2 \right)^{\frac{1}{2}}$$

and note that

$$\frac{\partial}{\partial x_i} |x|^{\gamma+1} = (\gamma+1)|x|^{\gamma-1}x_i \implies \frac{1}{\gamma+1} \nabla \mathcal{B}(x) = \mathcal{A}(x).$$

Since  $\gamma + 1 > 1$ , the function  $\mathcal{B}(x)$  is strictly convex. The gradient of a convex function is strictly increasing in each and all of its components; thus, the result of the lemma holds true. □

**Lemma 2** *Let  $x \in \mathbb{R}^M$  and  $f(x), g(x)$  be  $L^\infty$  functions. If  $f(x)$  is Lipschitz continuous and  $g(x)$  is  $\gamma$ -Hölder continuous, with  $0 < \gamma < 1$ , then the product  $f(x)g(x)$  is  $\gamma$ -Hölder continuous for  $x_1$  and  $x_2$  in a bounded domain  $\Omega$ .*

*Proof* Observe that

$$\begin{aligned} &|f(x_1)g(x_1) - f(x_2)g(x_2)| \\ &\leq |f(x_1)(g(x_1) - g(x_2))| + |g(x_2)(f(x_1) - f(x_2))| \\ &\leq \|f(x)\|_{L^\infty} |x_1 - x_2|^\gamma + \|g(x)\|_{L^\infty} |x_1 - x_2| \\ &\leq (\|f(x)\|_{L^\infty} + C\|g(x)\|_{L^\infty}) |x_1 - x_2|^\gamma \end{aligned}$$

□

References

1. Alonso, R., Santillana, M., Dawson, C.: On the diffusive wave approximation of the shallow water equations. *Eur. J. Appl. Math.* **19**(5), 575–606 (2008)

2. Arbogast, T., Wheeler, M.F., Zhang, N.Y.: A nonlinear mixed finite element method for a degenerate parabolic equation arising in flow in porous media. *SIAM J. Numer. Anal.* **33**(4), 1669–1687 (1996)

3. Bamberger, A.: Étude d’une équation doublement non linéaire. *J. Funct. Anal.* **24**, 148–155 (1977)

4. Barrett, J.W., Liu, W.B.: Finite element approximation of the parabolic p-laplacian. *SIAM J. Numer. Anal.* **31**(2), 413–428 (1994)

5. Bauer, P., Gumbrecht, T., Kinzelbach, W.: Aregional coupled surface water/groundwater model of the Okavango Delta, Botswana. *Water Resour. Res.* **42** (2006)

6. Brenner, S.C., Scott, R.: *The Mathematical Theory of Finite Element Methods*. Springer, New York (1994)

7. Carrillo, J.: Entropy solutions for nonlinear degenerate problems. *Arch. Ration. Mech. Anal.* **147**(4), 269–361 (1999)

8. Daugherty, R., Franzini, J., Finnemore, J.: *Fluid Mechanics with Engineering Applications*. McGraw-Hill, New York (1985)

9. DiBenedetto, E.: *Degenerate Parabolic Equations*. Springer, New York (1993)

10. Esteban, J.R., Vázquez, J.L.: Homogeneous diffusion in  $\mathbb{R}$  with power-like nonlinear diffusivity. *Arch. Rat. Mech. Anal.* **103**, 39–80 (1988)

11. Evans, L.C.: *Partial Differential Equations*. American Mathematical Society, Providence (2002)

12. Evje, S., Karlsen, K.: Discrete approximations of bv solutions to doubly nonlinear degenerate parabolic equations. *Numer. Math.* **86**(3), 377–417 (2000)

13. Fadimba, K.B., Sharpley, R.C.: Galerkin finite element method for a class of porous medium equations. *Nonlinear Anal. Real World Appl.* **5**, 355–387 (2004)

14. Feng, K., Molz, F.J.: A 2-d diffusion based, wetland flow model. *J. Hydrol.* **196**, 230–250 (1997)

15. Di Giammarco, P., Todini, E., Lamberti, P.: A conservative finite elements approach to overland flow: the control volume finite element formulation. *J. Hydrol.* **175**, 267–291 (1996)

16. Grange, O., Mignot, F.: Sur la résolution d’une équation et d’une inéquation paraboliques non linéaires. *J. Funct. Anal.* **11**, 77–92 (1972)

17. Hansen, E., Ostermann, A.: Finite element runge–kutta discretizations of porous medium-type equations. *SIAM J. Numer. Anal.* **46**(4), 1769–1779 (2008)

18. Hartman, P.: *Ordinary Differential Equations*. Society for Industrial and Applied Mathematics, Philadelphia (2002)

19. Hromadka, T.V., Berenbrock, C.E., Freckleton, J.R., Guymon, G.L.: A two-dimensional dam-break flood plain model. *Adv. Water Resour.* **8** (1985)

20. Yong, W., Pop, I.S.: A numerical approach to degenerate parabolic equations. *Numer. Math.* **92**(2), 357–381 (2002)

21. Ishige, K.: On the existence of solutions of the Cauchy problem for a doubly nonlinear parabolic equation. *SIAM J. Math. Anal.* **27**(5), 1235–1260 (1996)

22. Iwagaki, Y.: Fundamental studies of runoff analysis by characteristics. In: *Bulletin*, vol. 10, p. 25. Disaster Prevention Res. Inst., Kyoto University, Kyoto, Japan (1955)

23. Douglas, J. Jr., Dupont, T.F.: Galerkin methods for parabolic equations. *SIAM J. Numer. Anal.* **7**, 575–626 (1970)

24. Lions, J.L.: *Quelques Méthodes de Résolution des Problèmes aux Limites Non Linéaires*. Dunod Gauthier-Villars, Paris (1969)

25. Nochetto, R.H., Verdi, C.: Approximation of degenerate parabolic problems using numerical integration. *SIAM J. Numer. Anal.* **25**, 784–814 (1988)

26. Ponce, V.M., Li, R.M., Simons, D.B.: Applicability of kinematic and diffusion models. *J. Hydraul. Div.* **104**, 353–360 (1978)
27. Ponce, V.M., Simons, D.B.: Shallow wave propagation in open channel flow. *J. Hydraul. Div.* **103**, 1461–1476 (1977)
28. Raviart, P.A.: Sur la résolution de certaines équations paraboliques non linéaires. *J. Funct. Anal.* **5**, 299–328 (1970)
29. Rose, M.E.: Numerical methods for flows through porous media. *I. Math. Comput.* **40**(162), 435–467 (1983)
30. Rulla, J., Walkington, N.J.: Optimal rates of convergence for degenerate parabolic problems in two dimensions. *SIAM J. Numer. Anal.* **33**(1), 56–67 (1996)
31. Thomée, V.: *Galerkin Finite Element Methods for Parabolic Problems*. Springer Series in Computational Mathematics, no. 25. Springer, New York (1997)
32. Turner, A.K., Chanmeesri, N.: Shallow flow of water through non-submerged vegetation. *Agric. Water Manag.* **8**, 375–385 (1984)
33. Vázquez, J.L.: *The Porous Medium Equation. Mathematical Theory*. Oxford University Press, Oxford (2006)
34. Vreugdenhil, C.B.: *Numerical Methods for Shallow-Water Flow*. Kluwer Academic, Dordrecht (1998)
35. Wheeler, M.F.: A priori  $L^2$  error estimates for galerkin approximations to parabolic partial differential equations. *SIAM J. Numer. Anal.* **10**(4), 723–759 (1973)
36. Xanthopoulos, Th., Koutitas, Ch.: Numerical simulation of a two dimensional flood wave propagation due to dam failure. *J. Hydraul. Res.* **14**(4), 321–331 (1976)
37. Zhang, W., Cundy, T.W.: Modeling of two-dimensional overland flow. *Water Resour. Res.* **25**, 2019–2035 (1989)