

**UNIVERSIDAD DE CONCEPCION
ESCUELA DE GRADUADOS
CONCEPCION-CHILE**



**ESTUDIO DE UN PROBLEMA INVERSO PARA UNA ECUACION
PARABOLICA DEGENERADA CON APLICACIONES A LA
TEORIA DE LA SEDIMENTACION**

*Tesis para optar al grado de
Doctor en Ciencias Aplicadas con mención en Ingeniería Matemática*

Aníbal Coronel Pérez

**FACULTAD DE CIENCIAS FISICAS Y MATEMATICAS
DEPARTAMENTO DE INGENIERIA MATEMATICA
2004**

**ESTUDIO DE UN PROBLEMA INVERSO PARA UNA ECUACION
PARABOLICA DEGENERADA CON APLICACIONES A LA
TEORIA DE LA SEDIMENTACION**

Aníbal Coronel Pérez

Director de Tesis: Dr. Mauricio Alejandro Sepúlveda Cortes

Director de Programa: Dr. Mauricio Alejandro Sepúlveda Cortes

COMISION EVALUADORA

Dr. Robert Eymard, Université de Marne-la-Vallée, Francia.

Dr. Nils Henrik Risebro, Universitetet i Oslo, Noruega.

COMISION EXAMINADORA

Firma: _____

Dr. Raimund Bürger

Universität Stuttgart, Alemania

Firma: _____

Dr. Fernando Concha

Universidad de Concepción, Chile

Firma: _____

Dr. François James

Université d'Orléans, Francia

Firma: _____

Dr. Freddy Paiva

Universidad de Concepción, Chile

Firma: _____

Dr. Mauricio Sepúlveda

Universidad de Concepción, Chile

Fecha Examen de Grado: _____

Calificación: _____

Concepción–Septiembre 2004

AGRADECIMIENTOS

Agradezco profundamente a mi entorno familiar, laboral y social, los cuales distintamente estimularon el inicio, continuación y culminación de este trabajo.

Manifiesto mi especial gratitud al Prof. Mauricio Sepúlveda, al Prof. François James y al Prof. Fernando Concha por haberme sugerido el tema de investigación. Además, al Prof. Sepúlveda le debo mi reconocimiento y agradecimiento por haberme apoyado consecuentemente en las distintas etapas de mi formación doctoral, desde mi introducción al estudio de los métodos numéricos para Leyes de Conservación hasta la aceptación y dirección de la tesis; al Prof. James su hospitalidad y conducción durante mi estadía en el MAPMO, la cual marcó exitosamente el inicio de la presente investigación y al Prof. F. Concha su incondicional disponibilidad para compartir su vasta experiencia en el área del modelamiento matemático.

Al Prof. Wolfgang L. Wendland, al Dr. Raimund Bürger y a Stefan Berres, quienes me brindaron una grata estadía en el Instituto de Análisis Aplicado y Simulación Numérica (Institut für Angewandte Analysis und Numerische Simulation, IANS) de la Universidad de Stuttgart que desencadenó en una fructífera producción y especialmente en una latente colaboración.

Al Prof. Georges Guichon y a su grupo de trabajo, por la acogida y apoyo brindado durante mi estadía en la Universidad de Tennessee.

A la Escuela de Graduados de la Universidad de Concepción, al Mecesup UCO9907 y al Fondap en Matemáticas Aplicadas por el financiamiento del desarrollo de la tesis.

A la memoria de Don Rafael Coronel Pérez, mi padre.

A Doña Orfelina Pérez Núñez, mi madre.

A Doña Esperanza del Pilar Lozada Guidichi, mi esposa.

RESUMEN

En esta tesis se estudia un problema inverso para una ecuación parabólica fuertemente degenerada que modela la separación de una mezcla de sólido y fluido por sedimentación. El problema inverso (PI) consiste en la determinación de los coeficientes en la ecuación diferencial que gobierna el proceso a partir de mediciones de la concentración de sólidos que es la variable cuya evolución es descrita por el modelo o problema directo. El PI es formulado como un problema de minimización para una adecuada función de costo que compara la solución del modelo con las observaciones. Se realiza un análisis de PI para dos casos de interés: la sedimentación por efecto de la gravedad (sedimentación) y la sedimentación bajo acción de fuerzas centrípetas (centrifugación). En ambos casos se demuestra un resultado de continuidad de la solución entrópica con respecto a los coeficientes que implica la existencia de soluciones del PI respectivo. La obtención de los puntos estacionarios de la función costo son obtenidos por un método de descenso donde el gradiente es formalmente calculado a través de una formulación Lagrangiana que lleva a la introducción de un estado adjunto dado por un problema retrógrado con valores en la frontera para una ecuación diferencial parabólica lineal fuertemente degenerada y con coeficientes discontinuos. Haciendo un cálculo similar a esta deducción formal del gradiente, se obtiene un método que permite calcular de manera eficiente el gradiente exacto para el modelo discretizado. Este método es utilizado para la identificación de los parámetros que intervienen en el flujo de densidad y el coeficiente de difusión, a través de las relaciones constitutivas propias del modelo.

ABSTRACT

In this thesis we study an inverse problem for a strongly degenerate parabolic equation modelling the separation of a mixture of solid and fluid by sedimentation. The inverse problem (IP) consists in the determination of coefficients of the differential equation governing the processes from concentration measurements, which is the unknown of the model or direct problem. We formulate the IP as a minimization problem for an appropriate cost function which optimize the difference of the solution of the model with the observations. We analyze the IP for two cases of interest: the sedimentation under the acceleration of gravity (sedimentation) and the sedimentation by the influence of centripetal forces (centrifugation). In both cases, we prove a continuous dependence of entropy solution with respect to the coefficients, which implies the existence of solutions for IP. We obtain the stationary points of the cost function by a steepest descent method where the gradient is formally calculated by a Lagrangian formulation which permits the introduction of an adjoint state consisting in a linear backward boundary value problem for a strongly degenerate parabolic equation with discontinuous coefficients. Following a similar computation to the formal deduction of the gradient, we obtain an efficient computation of the exact gradient for the discrete model. This method is used to identify the parameters of the flux density function and diffusion coefficient in the constitutive relations obtained by the modelling assumptions.

Contents

Introducción General	i
1 Numerical identification of parameters for a model of sedimentation processes	1
1.1 Introduction	1
1.2 Statement of the problem	4
1.2.1 The direct problem	4
1.2.2 The inverse problem	5
1.3 Theoretical analysis of the IBVP and the IP	6
1.3.1 The direct problem	6
1.3.2 Existence of solutions to the Inverse Problem	7
1.4 Lagrangian formulation and formal calculus	11
1.5 Numerical schemes and discrete study	12
1.5.1 First and second-order explicit EO scheme	12
1.5.2 Implicit and semi-implicit schemes	16
1.5.3 Numerical tests	18
1.6 Parameters identification results using experimental data as observation.	25
2 Numerical identification of parameters for a strongly degenerate convection-diffusion problem modelling centrifugation of flocculated suspensions	29
2.1 Introduction	30
2.2 Entropy solutions of the direct problem	36
2.2.1 General form of the equations	36
2.2.2 Entropy solutions of the direct problem	37
2.3 Identification as optimization	38
2.3.1 The inverse problem as an optimization problem with PDE constraint	38

2.3.2	Lagrangian formulation	39
2.3.3	Adjoint state	40
2.3.4	Gradient of cost function	40
2.4	Existence of a solution of the inverse problem	41
2.5	Optimization scheme for identification	43
2.5.1	Discrete optimization with PDE as constraint	44
2.5.2	Discrete Lagrangian formulation	44
2.5.3	Discrete adjoint state	45
2.5.4	Discrete gradient of cost function	45
2.6	Derivatives of numerical fluxes	46
2.6.1	Engquist-Osher numerical flux	46
2.6.2	Differentiation with respect to the parameters	46
2.6.3	Differentiation with respect to the unknown	47
2.7	Numerical examples	48
2.7.1	Example 1: Profile of concentration at $t = T$ as observation .	49
2.7.2	Example 2: Profile of concentration at $z = \bar{R} \in [R_0, R]$ as observation	51
2.7.3	Example 3: Profile of concentration with $\sigma = 1$ at $r = \bar{R} \in [R_0, R]$ as observation	53
2.7.4	Example 4: Analytic data as observation	55
2.7.5	Concluding remarks	56
3	Convergence of an upwind scheme for an initial-boundary value problem of a strongly degenerate parabolic equation modelling sedimentation-consolidation processes	59
3.1	Introduction	59
3.2	Preliminaries	62
3.2.1	Definition of entropy solution	62
3.2.2	The difference schemes	65
3.2.3	Compactness criterion	66
3.2.4	Mollifiers and related functions	66
3.3	Estimates on the approximate solutions	67
3.3.1	L^∞ stability	67
3.3.2	BV estimates	69
3.3.3	Global estimates on $A(\phi_\Delta)$	75
3.4	Convergence Analysis	76
3.4.1	The discrete entropy inequality	77
3.4.2	Satisfaction of the entropy inequality	78

3.4.3	Satisfaction of initial and boundary conditions	80
A		85
A.1	Weak and discrete weak formulations	85
	References	89

Introducción General

El compendio de trabajos que se presentan en esta tesis estudian el problema de identificación de las nolinealidades de una ecuación parabólica fuertemente degenerada a partir de un conocimiento de su solución y de sus condiciones iniciales y de frontera. El interés inicial para realizar esta investigación fué el estudio de los métodos numéricos para resolver problemas inversos originados en el modelamiento matemático de los procesos de sedimentación-consolidación a través de la teoría fenomenológica.

La sedimentación es un proceso mecánico para la separación de una mezcla. Su gran utilidad en los procesos industriales lo convierten en un fenómeno relevante para la investigación científica (ver [33, 28]). Los principales aspectos históricos de la evolución de esta teoría aparecen detallados en [25] y [35]. En estos trabajos, se establece como modelo matemático para la descripción del fenómeno de separación sólido-fluido, una ecuación parabólica fuertemente degenerada. Los términos convectivos y difusivos que interviene en este modelo parabólico degenerado, se determinan en la práctica a través de hipótesis constitutivas dadas por expresiones que dependen de un número finito de parámetros.

La teoría fenomenológica de sedimentación está basada en la Teoría de Mezclas de la Mecánica del Medio Continuo y su proceso de modelamiento y puede ser resumido en las siguientes etapas (ver [9, 10, 11, 23, 24, 25, 28]) :

Etapa 1. Se supone que la mezcla, constituida por un sólido finamente dividido inmerso en un fluido, es un sistema particulado que está compuesto de dos medios continuos superpuestos obedeciendo a las siguientes restricciones:

M1 Las partículas sólidas son pequeñas con respecto al recipiente de sedimentación y tienen la misma densidad, tamaño y forma.

M2 El sólido y el fluido son incompresibles.

M3 No hay transferencia de masa entre el sólido y el fluido.

Etapa 2. Las hipótesis $M1-M3$ permiten aplicar un balance local de masa y momento en el sistema y así obtener

$$\frac{\partial \phi}{\partial t} + \nabla \cdot (\phi \mathbf{v}_s) = 0, \quad (1)$$

$$\frac{\partial(1-\phi)}{\partial t} + \nabla \cdot ((1-\phi)\mathbf{v}_f) = 0, \quad (2)$$

$$\rho_s \phi \left(\frac{\partial \mathbf{v}_s}{\partial t} + (\mathbf{v}_s \cdot \nabla) \mathbf{v}_s \right) = \nabla \cdot \mathbf{T}_s + \rho_s \phi \mathbf{b}_s + \mathbf{m}, \quad (3)$$

$$\rho_f (1-\phi) \left(\frac{\partial \mathbf{v}_f}{\partial t} + (\mathbf{v}_f \cdot \nabla) \mathbf{v}_f \right) = \nabla \cdot \mathbf{T}_f + \rho_f (1-\phi) \mathbf{b}_f - \mathbf{m}, \quad (4)$$

donde ϕ denota la fracción volumétrica local de sólidos (concentración), t el tiempo, \mathbf{m} la fuerza de interacción sólido-fluido por unidad de volúmen y \mathbf{v}_s , ρ_s , \mathbf{T}_s , \mathbf{b}_s , \mathbf{v}_f , ρ_f , \mathbf{T}_f , \mathbf{b}_f denotan la velocidad, la densidad de masa, el tensor de esfuerzo de Cauchy y la fuerza de cuerpo externa para la fase sólida y las cantidades correspondientes para la fase líquida, respectivamente.

Definiendo la velocidad volumétrica promedio como $\mathbf{q} = \phi \mathbf{v}_s + (1-\phi) \mathbf{v}_f$ y al sumar (1) y (2) se deduce la ecuación

$$\nabla \cdot \mathbf{q} = 0. \quad (5)$$

Etapa 3. Se introduce hipótesis constitutivas que permiten simplificar el sistema de ecuaciones (1)-(4).

La hipótesis de separación de fases, un cambio de variables teóricas por variables experimentalmente medibles y un análisis del estado de equilibrio cuando $t \rightarrow \infty$ (ver [26]), implican que \mathbf{T}_s , \mathbf{T}_f y \mathbf{m} se escriban como

$$\mathbf{T}_c = -p_c \mathbf{I}, \quad c \in \{s, f\}, \quad (6)$$

$$p_s = \phi p + \sigma_e, \quad p_f = (1-\phi)p, \quad (7)$$

$$\mathbf{m} = \mathbf{m}_b + \mathbf{m}_d, \quad \mathbf{m}_b = p \nabla \phi, \quad \mathbf{m}_d = \alpha \mathbf{v}_r, \quad (8)$$

donde $\mathbf{v}_r = \mathbf{v}_s - \mathbf{v}_f$ denota la velocidad relativa de la mezcla; p la presión de los poros; α el coeficiente de rozamiento del sedimento; σ_e el esfuerzo efectivo de sólidos y p_c , $c \in \{s, f\}$, las presiones de cada una de las fases. Las relaciones (6)-(8) en conjunto con las ecuaciones (3)-(4) permiten deducir

$$\begin{aligned} \mathbf{v}_r &= \frac{1-\phi}{\alpha} \nabla \sigma_e + \frac{\phi(1-\phi)}{\alpha} \left[\rho_s \left(\frac{\partial \mathbf{v}_s}{\partial t} + (\mathbf{v}_s \cdot \nabla) \mathbf{v}_s \right) - \rho_f \left(\frac{\partial \mathbf{v}_f}{\partial t} + (\mathbf{v}_f \cdot \nabla) \mathbf{v}_f \right) \right] \\ &\quad + \frac{\phi(1-\phi)}{\alpha} (\rho_s \mathbf{b}_s - \rho_f \mathbf{b}_f) \end{aligned} \quad (9)$$

En este punto se hacen dos hipótesis que permiten simplificar el procedimiento general, hecho hasta el momento, a los casos de interés para el presente trabajo: la sedimentación bajo acción de la gravedad (sedimentación) y la sedimentación por efecto de fuerzas centrífugas o centrifugación. Primeramente se consideran que las fuerzas de cuerpo están dadas por $\mathbf{b}_s = \mathbf{b}_f = -g\mathbf{k}$ para la sedimentación y para la centrifugación por

$$\mathbf{b}_c = -g\mathbf{k} - \boldsymbol{\omega} \times \boldsymbol{\omega} \times \mathbf{r} - 2\boldsymbol{\omega} \times \mathbf{v}_c, \quad c \in \{s, f\}, \quad (10)$$

donde g es la constante de aceleración de la gravedad, \mathbf{k} es el vector unitario y $\boldsymbol{\omega} = \omega\mathbf{k}$ es la velocidad angular. En (10), el primer término representa la fuerza gravitacional, el segundo término la fuerza inercial originada por la fuerza centrípeta y el tercer término la aceleración de Coriolis. En segundo lugar, se supone que la mezcla está contenida en un recipiente impermeable con paredes de fricción despreciable. Además, se supone que todas las variables son constantes a lo largo de cualquier sección transversal en el caso de la sedimentación y que existe simetría axial para el caso de la centrifugación. Esta última suposición reduce las ecuaciones multidimensionales consideradas a una sola variable espacial.

La hipótesis constitutiva sobre σ_e y α es que ambas son funciones solamente de la fracción volumétrica de sólidos y satisfacen

$$\sigma_e(\phi) \begin{cases} = 0 & \text{for } \phi \leq \phi_c, \\ > 0 & \text{for } \phi > \phi_c, \end{cases} \quad \sigma'_e(\phi) = \frac{d\sigma_e(\phi)}{d\phi} \begin{cases} = 0 & \text{for } \phi \leq \phi_c, \\ > 0 & \text{for } \phi > \phi_c, \end{cases} \quad (11)$$

$$f_{bk}(\phi) = -\frac{(\rho_s - \rho_f)g\phi^2(1 - \phi)^2}{\alpha(\phi)}, \quad f_{bk}(\phi) < 0, \forall \phi \in]0, \phi_{max}] \subset [0, 1] \quad (12)$$

donde $\phi_c \in]0, \phi_{max}[$ representa la concentración crítica, $f_{bk}(\phi)$ la función de flujo densidad y ϕ_{max} el valor máximo de la concentración.

Etapa 4. De (9) se obtiene la siguiente expresión (ver [28, 26]):

$$v_r = \frac{f_{bk}(\phi)}{\Delta\rho g\phi^2(1 - \phi)} [\sigma'_e(\phi) + \Delta\rho\phi g], \quad \Delta\rho = \rho_s - \rho_f, \quad (13)$$

en el caso de la sedimentación y una expresión similar para el proceso de centrifugación (ver [11]). Sustituyendo (13) en $\phi v_s = \phi q(t) + \phi(1 - \phi)v_r$, donde $q = q(t)$ es consecuencia de (5), se obtiene para la sedimentación el sistema

$$\frac{\partial\phi}{\partial t} + \frac{\partial}{\partial z}(\phi q(t) + f_{bk}(\phi)) = \frac{\partial^2}{\partial z^2}A(\phi), \quad A(\phi) = \int_0^\phi a(s)ds, \quad (14)$$

$$\frac{\partial p_e}{\partial z} = -\frac{\partial\sigma_e(\phi)}{\partial z} - \Delta\rho g\phi, \quad p_e = p - \rho_f g(L - z), \quad (15)$$

y para la centrifugación el sistema

$$\frac{\partial \phi}{\partial t} + \frac{\partial}{\partial r} \left(-f_{\text{bk}}(\phi) \frac{w^2 r}{g} - \sigma \frac{A(\phi)}{r} \right) = \frac{\partial^2}{\partial r^2} A(\phi) + \sigma \left(f_{\text{bk}}(\phi) \frac{\omega^2}{g} + \frac{A(\phi)}{r^2} \right) \quad (16)$$

$$\frac{\partial p_e}{\partial r} = -\frac{\partial \sigma_e(\phi)}{\partial r} + \Delta \rho \phi \omega^2 r, \quad p_e = p - \rho_f \frac{1}{2} \omega^2 r^2, \quad (17)$$

donde z representa la altura, L la altura inicial de la suspensión, r el radio o distancia al eje de rotación, A es la primitiva del coeficiente de difusión

$$a(\phi) = -\frac{f_{\text{bk}}(\phi) \sigma'_e(\phi)}{\Delta \rho g \phi}, \quad (18)$$

y $\sigma \in \{0, 1\}$ es un parámetro introducido de acuerdo con los casos especiales de centrifugación considerados (ver Figura 2.1). De ambos sistemas, (14)-(15) y (16)-(17), se deduce que los procesos de sedimentación y centrifugación quedan modelados completamente por las ecuaciones escalares (14) y (16), respectivamente. Las ecuaciones (15) y (17) son utilizadas para calcular la presión de poros local.

Etapa 5. Para completar el modelo se debe considerar condiciones iniciales y de frontera para (14) y (16).

En el caso de la sedimentación se asume que el flujo de sólidos en $z = 0$ se reduce a su parte convectiva $q(t)\phi(0, t)$ y que la concentración no cambia cuando el sedimento sale del sedimentador, lo que implica la condición de frontera

$$\left(f_{\text{bk}}(\phi) - \frac{\partial A(\phi)}{\partial z} \right) (0, t) = 0. \quad (19)$$

En $z = L$ se considera dos tipos diferentes de condiciones de frontera dadas por

$$\phi(L, t) = \phi_L(t) \quad (20)$$

$$\left(f_{\text{bk}}(\phi) + q(t)\phi - \frac{\partial A(\phi)}{\partial z} \right) (L, t) = f_F(t), \quad (21)$$

donde ϕ_L y f_F son funciones dadas, modelando la concentración de sólidos y el flujo de alimentación de sólidos, respectivamente. La concentración inicial es dada por

$$\phi(z, 0) = \phi_0(z), \quad z \in [0, L]. \quad (22)$$

Los problemas de valores iniciales y de frontera (14), (19), (20), (22) y (14), (19), (21), (22) son llamados “Problema A” y “Problema B”, respectivamente. Note que la diferencia entre ambos problemas es la condición de borde en $x = L$, el Problema

A tiene una condición Dirichlet y el Problema B una condición de flujo. Un caso especial de estos problemas aparece cuando se considera un sedimentador cerrado en $x = 0$, correspondiendo a considerar $q = 0$, llamado proceso de sedimentación batch.

En el caso de la centrifugación se considera que la velocidad de la fase sólida se anula en ambos lados $r = R_0$ y $r = R$. Así, las condiciones de frontera para (16) están dadas por

$$\left(f_{\text{bk}}(\phi) \frac{\omega^2 r_b}{g} - \frac{\partial A(\phi)}{\partial r} \right) (r_b, t) = 0, \quad r_b \in \{R_0, R\}. \quad (23)$$

La condición inicial esta dada por $\phi(r, 0) = \hat{\phi}_0(r)$, $r \in [R_0, R]$.

En resumen, de las etapas de modelamiento descritas (Etapa1-Etapa5), se deduce que la forma general de las ecuaciones no lineales involucradas en este estudio está dada por la ecuación parabólica

$$\begin{aligned} \partial_t u + \partial_x F(x, t, u) &= \partial_{xx}^2 A(u) + G(x, t, u), \quad (x, t) \in Q_T =]a, b[\times]0, T[\\ u(x, 0) &= u_0(x) \quad x \in]a, b[\\ \Gamma_\ell(x_\ell, t, u) &= \gamma_\ell(t) \quad t \in]0, T[, \quad x_\ell \in \{a, b\}, \quad \ell = 0, 1, \end{aligned} \quad (24)$$

donde

$$\partial_x A(u) \geq 0, \quad u \in [0, 1]. \quad (25)$$

Los coeficientes toman las formas particulares

$$\begin{aligned} F(x, t, u) &= q(t)u + f_{\text{bk}}(u), \quad G(x, t, u) = 0 \quad \text{y} \\ F(x, t, u) &= -\frac{\omega^2 r}{g} f_{\text{bk}}(u) - \frac{\sigma}{x} A(u), \quad G(x, t, u) = \sigma \left[\frac{\omega^2}{g} f_{\text{bk}}(u) + \frac{A(u)}{x^2} \right], \end{aligned}$$

para la sedimentación gravitatoria y centrifugación, respectivamente. Las condiciones de frontera en esta nueva notación son:

$$\begin{aligned} \Gamma_0(x_0, t, u(x_0, t)) &= F(x_0, t, u(x_0, t)) - \partial_x A(u(x_0, t)), \quad \gamma_0(t) = q(t)u(x_0, t), \\ \Gamma_1(x_1, t, u(x_1, t)) &= u(x_1, t), \quad \gamma_1(t) = \phi_1(t), \end{aligned}$$

para el Problema A,

$$\begin{aligned} \Gamma_\ell(x_\ell, t, u(x_\ell, t)) &= F(x_\ell, t, u(x_\ell, t)) - \partial_x A(u(x_\ell, t)), \quad \ell \in \{0, 1\} \\ \gamma_0(t) &= q(t)u(x_0, t), \quad \gamma_1(t) = f_F(t), \end{aligned}$$

para el Problema B y

$$\begin{aligned}\Gamma_\ell(x_\ell, t, u(x_\ell, t)) &= F(x_\ell, t, u(x_\ell, t)) - \partial_r A(u(x_\ell, t)) + \frac{\sigma}{x_\ell} A(u(x_\ell, t)) \\ \gamma_0(t) &= \gamma_1(t) = 0.\end{aligned}$$

para la centrifugación. La característica fundamental del modelo matemático está dada por el hecho que el término difusivo de la ecuación (24) se *degenera fuertemente* en el sentido dado por la condición (25) y que es consecuencia de (11) y (18).

El análisis matemático de las ecuaciones (24)-(25) ha recibido en las últimas décadas una atención especial ya que su nolinealidad y cambio de comportamiento de parabólico a hiperbólico hacen que la solución se comporte como en las Leyes de Conservación nolineales, es decir se observa la formación de discontinuidades (choques) o una regularización (ondas de rarefacción) en la solución independientemente de la regularidad de las condiciones iniciales y de frontera impuestas (ver [31, 47, 69]). De esta manera las soluciones para (24) deben ser entendidas en el sentido de la entropía de Kružkov.

Un análisis para el problema de Cauchy asociado a la ecuación (24) fue hecho por Carrillo en un extenso y profundo trabajo (ver [31]). Este artículo establece aspectos intrínsecos del comportamiento de la solución bajo hipótesis generales sobre los coeficientes. A pesar de la generalidad en la que se enuncian estos resultados su adaptación al problema con valores en la frontera no es directa, apareciendo complejidades meritorias de un análisis especial. Bürger y coautores en [15], siguiendo el trabajo de Carrillo (ver [31]), realizaron el estudio del modelo de sedimentación (Problema A y Problema B) cuando los coeficientes y datos satisfacen condiciones de regularidad que son detalladas en los capítulos 1 y 3 (ver secciones 1.2 y 3.2). El Problema A necesita una redefinición de la frontera Dirichlet que permita evitar el conflicto con la condición de entropía, esta desigualdad es presentada en [15] y fue obtenida previamente en [23], basándose en los resultados de Volpert y Hudjaev (ver [88, 89]). El estudio de existencia, unicidad y estabilidad con respecto a la condición inicial del problema de centrifugación fue hecho en [20], siguiendo también la técnica de [31, 15].

El método de Volúmenes Finitos es la herramienta natural para la simulación numérica de Leyes de Conservación (ver [48]). Así, en el caso de la ecuación (24) resulta también ser la técnica numérica mas adecuada. La discretización de (24) considerada en esta tesis está dada por tres familias de esquemas: uno formulado de manera implícita, otro de manera semi-implícita y otro de manera explícita. En los tres casos se considera la discretización en el interior del dominio físico y se incorporan adecuadamente las condiciones de frontera. El flujo numérico (para la convección) utilizando en las simulaciones numéricas es el de Engquist-Osher (ver

[42]), cuya evidencia de buen comportamiento y coherencia para reflejar el fenómeno modelado es presentada en los trabajos de R. Bürger y colaboradores, principalmente en [19], donde fue introducido.

En el proceso de simulación es imprescindible determinar valores numéricos para los distintos parámetros físicos que intervienen en los términos convectivo y difusivo de la ecuación parabólica degenerada. A pesar que experimentalmente o mediante tablas se puede obtener aproximaciones, la determinación de estos datos, o valores numéricos de los parámetros, es por lo general un problema difícil. Así, la calibración y validación de estos parámetros se hace a partir de cuán cerca está la solución obtenida mediante simulación numérica del modelo, de una determinada observación experimental. El tratamiento matemático se traduce en un problema de identificación de parámetros que es una situación particular del siguiente problema inverso :

PI. Dadas las funciones $u_0, \gamma_\ell, \ell = 0, 1$ y un conjunto de mediciones experimentales $u^{obs}(x, t)$ encontrar F, G, A de tal manera que la solución entrópica de (24) esté lo “mas cercana posible” de $u^{obs}(x, t)$.

La formulación PI establecida para el problema inverso es equivalente al problema de optimización

$$\min_{F, G, A} \mathcal{J}(u, u^{obs}), \quad \text{sujeto a que } u \text{ sea solución débil de (24)}, \quad (26)$$

donde la función costo \mathcal{J} se escoge para dar precisión matemática (analítica y numérica) del término ambiguo “mas cercana posible”. El funcional natural \mathcal{J} y que será considerado a lo largo de este trabajo es el de mínimos cuadrados, el cual compara las funciones u y u^{obs} en la norma L^2 (ver capítulos 1 y 2). El planteamiento (26) de PI permite estudiar su existencia-unicidad y formular técnicas para su solución.

De una manera general el estudio del problema inverso es posible dividirlo de la siguiente forma

- Análisis de su buen planteamiento.
- Aproximación Numérica.
- Convergencia de los esquemas numéricos.
- Aplicación a datos experimentales.

Todos estos subproblemas, para cada uno de los modelos previamente descritos, son estudiados y reportados en la presente tesis.

La existencia de soluciones para (26) está basado en la dependencia continua de la solución entrópica con respecto a los coeficientes. Sin embargo, la unicidad

de soluciones para (26), es un problema mal-puesto, pudiéndose construir ejemplos explícitos para el caso de las Leyes de Conservación (ver [62]).

La técnica utilizada en esta tesis, para la solución numérica de (26) está basada en el método del gradiente. La nolinealidad de la ecuación (24) implica que su solución no dependa explícitamente de los coeficientes, dificultando de este modo la obtención del gradiente. Sin embargo, haciendo cálculos muy similares a los realizados por James y Sepúlveda para el caso de la Leyes de Conservación (ver [62]) es posible la deducción formal de un gradiente continuo para \mathcal{J} . Este procedimiento se resume en los siguientes pasos :

Paso1 Se introduce una formulación lagrangiana de (26)

$$\mathcal{L}(u, p; \mathbf{e}) = \mathcal{J}(u, u^{obs}) - E(u, p; \mathbf{e}), \quad (27)$$

donde \mathbf{e} denota el vector de los parámetros en las funciones F, G, A , que se van a identificar y E la formulación variacional de (24) dada por

$$\begin{aligned} E(u, p; \mathbf{e}) &= - \int \int_{Q_T} (u \partial_t p + F(x, t, u) \partial_x p + A(u) \partial_x^2 p + G(x, t, u) p) dt dx \\ &\quad + \int_a^b u p \Big|_{t=0}^T dx + \int_0^T ((F(x, t, u) - \partial_x A(u)) p + A(u) \partial_x p) \Big|_{x=a}^b dt, \\ &\quad \forall p \text{ "suficientemente suave"} \end{aligned} \quad (28)$$

y \mathcal{J} la función costo, que se asume de la forma

$$\mathcal{J}(u, u^{obs}) = \frac{1}{2} \int \int_{Q_T} |u(x, t) - u^{obs}(x, t)|^2 dx dt. \quad (29)$$

La función test p puede ser considerado como un multiplicador de Lagrange generalizado (ver [62]). En (28), para cada uno de los casos en estudio se debe sustituir la condición de borde adecuadas.

Paso2 Se introduce la ecuación adjunta para (26). El gradiente de la función costo en (27) está dado por

$$\nabla_{\mathbf{e}} \mathcal{J}(u(\mathbf{e})) = \langle \partial_u \mathcal{L}(u(\mathbf{e}), p; \mathbf{e}), \nabla_{\mathbf{e}}(u) \rangle + \nabla_{\mathbf{e}} \mathcal{L}(u(\mathbf{e}), p; \mathbf{e}), \quad (30)$$

debido a que $E(u(\mathbf{e}), p; \mathbf{e}) = 0$. En esta derivación formal de la derivada total de la función costo $\nabla_{\mathbf{e}}(u)$ no puede ser calculado, dado que la solución del problema directo u no puede ser calculado explícitamente en términos de \mathbf{e} . Así, se hace

necesaria la introducción de una función test p de modo que $\partial_u \mathcal{L}(u(\mathbf{e}), p; \mathbf{e}) = 0$, es decir que anula el primer término del lado derecho de (30). La derivada de \mathcal{L} en la dirección δu está dada formalmente por

$$\begin{aligned} <\partial_u \mathcal{L}(u(\mathbf{e}), p; \mathbf{e}), \delta u> &= \int \int_{Q_T} (\partial_t p + \partial_u F(x, t, u) \partial_x p + \partial_u A(u) \partial_x^2 p \\ &+ \partial_u G(x, t, u) p + (u - u^{obs})(x, t)) dt dx + \int_a^b u(x, T) p(x, T) dx \\ &+ \int_0^T (\partial_u (F(x, t, u) - \partial_x A(u)) p + \partial_u A(u) \partial_x p) \Big|_{x=a}^b dt. \end{aligned} \quad (31)$$

A fin de anular esta expresión, basta que p satisfaga la denominada ecuación adjunta

$$\begin{aligned} \partial_t p + \partial_u F(x, t, u) \partial_x p &= -a(u) \partial_{xx}^2 p - \partial_u G(x, t, u) \quad (x, t) \in Q_T, \\ p(x, t) &= \partial_u \mathcal{J}, \quad x \in]a, b[, \\ \Gamma'_\ell(x_\ell, t, u, p) &= \gamma'_\ell(t), \quad t \in]0, T[, \end{aligned} \quad (32)$$

donde Γ'_ℓ y γ'_ℓ toman distintas formas de acuerdo a cada uno de los problemas directos en estudio.

Paso3 Se obtiene el gradiente. Siendo p solución del problema adjunto (32) se sigue que el gradiente de la \mathcal{J} viene dado por

$$\begin{aligned} \nabla_{\mathbf{e}} \mathcal{J}(u(\mathbf{e})) &= \nabla_{\mathbf{e}} \mathcal{L}(u(\mathbf{e}), p; \mathbf{e}) = -\nabla_{\mathbf{e}} E(u(\mathbf{e}), p; \mathbf{e}) \\ &= \int \int_{Q_T} (\nabla_{\mathbf{e}} F(x, t, u) \partial_x p + \nabla_{\mathbf{e}} A(u) \partial_x^2 p + \nabla_{\mathbf{e}} G(x, t, u) p) dt dx \\ &- \int_0^T (\nabla_{\mathbf{e}} (F(x, t, u) - \partial_x A(u)) p + \nabla_{\mathbf{e}} A(u) \partial_x p) \Big|_{x=a}^b dt \end{aligned} \quad (33)$$

Los Pasos 1 a 3 describen un cálculo formal que permite resolver el problema de minimización asociado al problema inverso. Sin embargo, aparecen dos dificultades. La primera, se presenta a nivel continuo, y es acerca de la validez de este cálculo. En efecto, debido a las discontinuidades de la solución del problema directo, y de las eventuales discontinuidades del problema adjunto, la diferenciabilidad de la función costo en un sentido clásico, no es un problema fácil de determinar, tal

como se observa por ejemplo para el caso puramente hiperbólico (ver [62]). La segunda dificultad, es acerca de la discretización del problema (33), para el cálculo del gradiente discretizado. El problema (33) corresponde a una ecuación lineal de convección-difusión fuertemente degenerada con coeficientes discontinuos, cuya existencia y unicidad de soluciones es hasta ahora un problema abierto. Si bien se puede hacer una analogía al caso hiperbólico, en que se tiene un conocimiento parcial de existencia y unicidad, a través de las soluciones reversibles, el problema de escoger un método numérico correcto que aproxime a la solución, sigue siendo un problema difícil (ver [27, 55]). En la práctica, subsanamos momentáneamente estos dos problemas, repitiendo el cálculo formal para el caso discreto, es decir calculando un esquema adjunto asociado al esquema del problema directo y luego un gradiente de la función costo discretizada que resulta ser un gradiente exacto. La convergencia del gradiente exacto al gradiente del problema continuo (o quizás a algún subgradiente) sigue siendo hasta ahora un problema abierto.

A nivel discreto, partiendo de una aproximación por Volúmenes Finitos de (24) dada por

$$\begin{aligned} u_j^{n+1} &= u_j^n - \lambda(F_{j+1/2}^n - F_{j-1/2}^n) + \mu(A_{j+1/2}^n - A_{j-1/2}^n) + g_j^n, \quad j = 1, \dots, J-1 \\ \Gamma_\Delta(u_\ell^n, F_{\ell+1/2}^n, A_{\ell+1/2}^n, g_\ell^n) &= \gamma^n, \quad \ell \in \{-1, J\} \end{aligned} \quad (34)$$

donde Γ_ℓ^Δ denota la discretización de la condición de borde que se hace explícitamente para cada una de las condiciones de borde. Se continúa con una formulación Lagrangiana discreta que permite introducir el estado adjunto discreto

$$\begin{aligned} p_j^n &= p_j^{n+1} - \lambda \partial_u \sum_{k=-K}^{K-1} \partial_{u_j^n} F_{j+k+1/2}^n (p_{j+k}^{n+1} - p_{j+k+1}^{n+1}) \\ &\quad + \mu \sum_{k=-\hat{K}}^{\hat{K}-1} \partial_{u_j^n} A_{j+1/2}^n (p_{j+k}^{n+1} - p_{j+k-1}^{n+1}) \\ &\quad + \sum_{k=-\bar{K}}^{\bar{K}-1} g_j^n p_j^{n+1}, \quad j = 1, \dots, J-1, \\ \Gamma'_\Delta(u_\ell^n, F_{\ell+1/2}^n, A_{\ell+1/2}^n, g_\ell^n, p_\ell^{n+1}) &= \gamma'(t^n), \quad \ell \in \{1, J-1\} \end{aligned} \quad (35)$$

se calcula el siguiente gradiente discreto para la función costo

$$\begin{aligned} \nabla_{\mathbf{e}} J &= -\Delta t \Delta x \sum_{(j,n) \in Q_\Delta} \{ \lambda \nabla_{\mathbf{e}} F_{j+1/2}^n (p_j^{n+1} - p_{j+1}^{n+1}) \\ &\quad - \mu \nabla_{\mathbf{e}} A_{j+1/2}^n (p_j^{n+1} - p_{j+1}^{n+1}) - \nabla_{\mathbf{e}} g_j^n p_j^{n+1} \} \end{aligned} \quad (36)$$

Se observa que el gradiente calculado corresponde a la discretización explícita del problema directo, en el caso de los otros esquemas se deduce una expresión equivalente muy similar.

Organización de la Tesis

La tesis es presentada de la siguiente manera:

En el Capítulo 1 se presenta el estudio del problema inverso para el Problema A cuyos resultados fueron publicados en [37]:

A. CORONEL, F. JAMES AND M. SEPULVEDA, Numerical identification of parameters for a model of sedimentation processes. *Inverse Problems*, 19(4):951–972, 2003.

En este Capítulo se presenta un resultado de continuidad de la solución entrópica del problema directo con respecto a los coeficientes, lo que implica la existencia de soluciones del problema inverso. Los resultados numéricos de validación muestran la eficacia del método, a pesar de la carencia de pruebas teóricas. En el artículo [37] no se consideraron datos experimentales reales, ya que no se contaba inicialmente con ellos. Los datos que se muestran al final del Capítulo 1, fueron obtenidos recientemente por Pamela Garrido. Se consideró importante agregar una sección adicional (ver Sección 1.6) respecto del artículo publicado para mostrar la validez del método de identificación de parámetros, al aplicarlo sobre datos observados provenientes de columnas de sedimentación de laboratorio.

En el Capítulo 2 se presenta el estudio del problema inverso para el caso de centrifugación y corresponde a los resultados de los artículos [5] y [6]:

S. BERRES, R. BÜRGER, A. CORONEL, AND M. SEPÚLVEDA. Numerical identification of parameters for a strongly degenerate convection-diffusion problem modelling centrifugation of flocculated suspensions. Technical Report 2003-14, Departamento de Ingeniería Matemática, Universidad de Concepción, Chile, 2003. *To appear in Appl. Numer. Math.*

S. BERRES, R. BÜRGER, A. CORONEL AND M. SEPÚLVEDA, Numerical identification of parameters for a flocculated suspension from concentration measurements during batch centrifugation. Preprint 2004-06, Departamento de Ingeniería Matemática, Universidad de Concepción, Concepción, Chile, *To appear in Chem. Eng. J.*

En [5] se demuestra una condición necesaria para la existencia de soluciones del problema inverso y además se estudia numéricamente la convergencia del método del

lagrangiano. En tanto que en [6] se enfatiza sobre los aspectos físicos del problema, los detalles del algoritmo numérico y su aplicación.

En el Capítulo 3 se presenta el estudio teórico de convergencia del esquema de Volúmenes Finitos (34) y corresponde el trabajo [8]:

R. BÜRGER, A. CORONEL AND M. SEPÚLVEDA. Convergence of upwind schemes for an initial-boundary value problem of a strongly degenerate parabolic equation modelling sedimentation-consolidation processes. Preprint 2004-09, Departamento de Ingeniería Matemática, Universidad de Concepción. *Submitted to Math. Comp.*

Este artículo presenta los resultados de convergencia para un esquema de Volúmenes finitos general, bajo las características de consistencia y conservación. Las pruebas son presentadas para el problema B. Sin embargo, estos resultados son generalizables al problema A y al modelo de centrifugación.

Problemas Abiertos

Dentro de las dificultades no resueltas en la presente tesis se mencionan las siguientes:

- Los cálculos que conducen al gradiente continuo son formales. El gradiente puede no existir, tal como fué observado, para el caso de las Leyes de Conservación, por F. James y M. Sepúlveda en [62]. Esto es consecuencia de que las soluciones del problema directo son generalmente funciones discontinuas. Esta dificultad puede ser solucionada dando al gradiente una interpretación en un sentido débil, partiendo por considerar una diferenciabilidad débil de la solución entrópica con respecto a los coeficientes. Sin embargo, lo único que se conoce es un comportamiento de continuidad de la solución entrópica con relación a los coeficientes (ver teoremas 1.3.2 y 2.4.1) y no algún tipo de diferenciabilidad.

Un avance notable del estudio de la diferenciabilidad es el desarrollado por Stefan Ulbrich para Leyes de Conservación (ver [86, 87]). En [86] se introduce el concepto de “shift-differentiability” con el cual se deduce la diferenciabilidad en el sentido de Fréchet para una clase de funcionales de costo. Esta técnica es basada fuertemente en el trabajo de las características generalizadas de Dafermos (ver [40]) y asumiendo que el flujo es convexo, siendo precisamente estas dos restricciones las que no permiten su directa generalización al problema estudiado en esta tesis.

- El estudio analítico del problema adjunto (32). Esta ecuación, a pesar de ser lineal, posee dos dificultades: la primera, originada debido a que los coeficientes dependen de la solución del problema directo y por ende son generalmente discontinuas, y la segunda, causada por la fuerte degeneración del término difusivo.

El análisis de este tipo de problemas, sin considerar términos difusivos fuertemente degenerados y términos fuente, fue hecho por F. Bouchut y F. James en [27] y recientemente extendido en [86], al caso de términos fuente. Sin embargo, la generalización de la técnica utilizada en [27] al estudio de (32) no es directamente posible debido a que la condición de Lipschitz-continuidad de un solo lado (OSLC) deja de ser válida.

- La consideración anterior acarrea otra dificultad: La falta de resultados de convergencia del estado adjunto discreto. La carencia de una definición precisa de solución para el problema adjunto (32) impide dar sentido a la convergencia del estado adjunto discreto (35) y a cualquier otra aproximación numérica que se considere. En este sentido cabe mencionar que un análisis similar al desarrollado para la convergencia del problema directo (ver Capítulo 3) lleva a establecer las condiciones necesarias utilizadas en [55] para la convergencia de los esquemas de Volúmenes Finitos en las ecuaciones de transporte con coeficientes discontinuos.
- Las hipótesis $M1$ puede ser relajada de modo de obtener un modelo mas realista del fenómeno de sedimentación. Si se supone que las partículas de sólido son de distinto tamaño y densidad lleva al establecimiento de un sistema de ecuaciones diferenciales no lineales de tipo parabólico fuertemente degenerado (ver [4, 7]) y por lo tanto más difícil de resolver.

Chapter 1

Numerical identification of parameters for a model of sedimentation processes

In this chapter we present the identification of parameters in the flux and diffusion functions for a quasilinear strongly degenerate parabolic equation which models the physical phenomenon of flocculated sedimentation. We formulate the identification problem as a minimization of a suitable cost function and we derive its formal gradient by means of adjoint equation which is a backward linear degenerate parabolic equation with discontinuous coefficients. For the numerical approach, we start with the discrete Lagrangian formulation and assuming that the direct problem is discretized by the Engquist-Osher scheme obtain a discrete adjoint state associated to this scheme. The conjugate gradient method permits to find numerically the physical parameters. In particular, it allows to identify as well the critical concentration level at which solid flocs begin to touch each other and determines the change of parabolic to hyperbolic behavior in the model equation.

1.1 Introduction

Batch sedimentation is a classical procedure to separate a flocculated suspension into a concentrated sediment of practical interest, and a clear fluid. It is used, for example in metallurgy and food industry. The experimental setting consists of a vertical column, with a surface feed at the top, and a discharge surface at the bottom. Under the influence of gravity, the suspension separates into a clear fluid and a compressible sediment which is collected at the bottom of the column. See [23, 24, 26, 16, 12, 28] for more complete descriptions and details.

Under several constitutive and simplifying assumptions, it turns out that this mixture of two continuous media (fluid and solid flocs) can be described by a single model equation, namely a partial differential equation of mixed type: hyperbolic and parabolic. The unknown is the volumetric solid concentration ϕ , which is a function of the time t and the height in the column z .

The mixed type nature of the equation results from the different behaviour of the solid flocs on the one hand and the fluid on the other. The former lies in a compression region, where the solid effective stress σ_e is not constant, while in the fluid zone it can be taken as equal to a constant. Therefore, in the solid zone (sediment), the equation will be of parabolic type, and in the fluid zone, it will be hyperbolic.

More precisely, it can be assumed that σ_e is a nonlinear function of ϕ alone, with the following shape:

$$\sigma'_e(\phi) = \begin{cases} = 0 & \text{for } \phi \leq \phi_c \\ > 0 & \text{for } \phi > \phi_c \end{cases}$$

The concentration ϕ_c is the so-called *critical concentration*, above which the flocs get into contact with each other and form a network. The location of the interface $\phi = \phi_c$ is a specific problem in direct simulations.

Another constitutive relation can be obtained through Kynch's kinematical theory of sedimentation. This gives the velocity of the mixture as a function of the concentration ϕ ,

$$f(\phi, t) = q(\phi, t) + f_{bk}(\phi).$$

Here q is the volume average velocity of the mixture, $q = v_s\phi + v_f(1 - \phi)$ (v_s and v_f are the solid and fluid velocities), and f_{bk} the so-called batch flux density function. A constitutive equation should be formulated for this function.

We are interested here in the inverse problem which consists in identifying the constitutive laws f_{bk} and σ_e from experimental data. This kind of problem is in general impossible to solve in its full generality, i.e. considering f_{bk} and σ_e as general functions. Therefore we shall consider more precise constitutive assumptions, which give explicit expressions for both functions, depending on a finite number of parameters. Among these we emphasize the value of the critical concentration ϕ_c , which is of great practical importance, and very difficult to access from experimental data. We mention here that there are several experimental methods to obtain the involved parameters, see [3, 13] for an overview. They employ a set of data of local concentration and initial settling velocities for sedimentation, and concentration gradients for consolidation, to obtain approximate correlations from which the settling and consolidation parameters may be obtained by algebraic manipulations. In contrast, in our approach we need only experimental concentration data and we compute the optimal f_{bk} and σ_e .

There is a large list of authors who have proposed analytical and numerical methods to solve inverse problems for partial differential equations of evolution. For example, parameter identification methods for parabolic PDEs can be found in [30, 29, 57, 70, 71] and references therein. Recently, Yamamoto and Zou [92] reconstructed the radiative coefficient and the initial data for a linear parabolic equation using a piecewise linear finite-element method for the discretization and a nonlinear gradient multigrid method for accelerating the reconstruction process. Several difficulties arise in the case we consider, since we want to reconstruct nonlinear coefficients, but such a strategy is likely to be useful to improve the results. Moreover, the hyperbolic degeneracy gives rise to shocks in the concentration profiles. Both the existence and uniqueness for weak solutions to the direct problem rely on the specific entropy conditions, which has physical relevancy for the model, and must be considered owing to the nonlinearity of the flux and degeneracy of the diffusion, see [23, 24, 19].

The existence of solutions for the inverse problem is a consequence of the continuous dependence of the entropy solutions with respect to the flux and the diffusion (see Theorem 1.3.2 below and [32, 46]). Uniqueness cannot be ensured because of the hyperbolic behaviour, see [62, 61]. Therefore we shall rewrite the identification problem by adapting the technique developed by James and Sepúlveda [62, 61, 63], which is numerically tested as an efficient method to reconstruct the flux of a particular hyperbolic system. The idea is to write the inverse problem as an optimization problem for an appropriate cost function and then to apply the classical conjugate gradient method.

The continuous gradient stems from an adjoint state to the model consisting of a backward linear degenerate convection-diffusion equation with discontinuous coefficients and boundary conditions. As in the purely hyperbolic case, the adjoint equation is ill-posed in uniqueness. We obtain the discrete gradient by computing the exact gradient of the discrete formulation of the optimization problem. This has become a classical technique, instead of computing the discretization of the formal gradient, because identification problems are generally badly conditioned or ill-posed, see [63]. All the computations are based upon the explicit (first and second order), semi-implicit and implicit Engquist-Osher (or generalized upwind) schemes for numerical computation of the solution of the sedimentation model, see Bürger et al. [19]. In each case, the adjoint scheme and the discrete gradient are provided.

The remainder of this paper is organized as follows. In section 1.2 we provide the formulation of the direct and inverse problems. In section 1.3 we analyse the question of the well-posedness. In section 1.4 we present the formal calculus of the gradient. In section 1.5 we introduce the numerical schemes for the identification of

the parameters and we present some numerical results.

1.2 Statement of the problem

1.2.1 The direct problem

Summarizing the results given in [9, 23, 24, 26] for the mathematical model of the sedimentation processes, we have the following IBVP

$$\frac{\partial \phi}{\partial t} + \frac{\partial}{\partial z} (q(t)\phi + f_{bk}(\phi)) = -\frac{\partial}{\partial z} \left(f_{bk}(\phi) \frac{\sigma'_e(\phi)}{\Delta \rho g \phi} \frac{\partial \phi}{\partial z} \right), \quad (z, t) \in Q_T, \quad (1.1)$$

$$\phi(z, 0) = \phi_0(z), \quad z \in [0, L], \quad (1.2)$$

$$\phi(L, t) = \phi_2(t), \quad t \in [0, T], \quad (1.3)$$

$$f_{bk}(\phi) \left(1 + \frac{\sigma'(\phi)}{\Delta \rho g \phi} \frac{\partial \phi}{\partial z} \right) \Big|_{z=0} = 0, \quad t \in [0, T], \quad (1.4)$$

where ϕ is the non-negative unknown function, $Q_T = [0, L] \times [0, T]$, $\Delta \rho$ and g are positive constants, q , f_{bk} , σ_e , ϕ_0 and ϕ_2 are given functions with the following supposed behaviour:

- q is a non-positive Lipschitz function, this is

$$q \in \text{Lip}([0, T]) \text{ and } q(t) \leq 0, \quad \forall t \in [0, T]. \quad (1.5)$$

- $f_{bk}(\phi)$ is a smooth function such that

$$f_{bk}(\phi) = \begin{cases} = 0 & \text{for } \phi \in \mathbb{R} -]0, \phi_{max}[, \\ < 0 & \text{for } \phi \in]0, \phi_{max}[. \end{cases} \quad 0 < \phi_{max} \leq 1 \quad (1.6)$$

- $\sigma_e(\phi)$ is a C^3 function, constant for $\phi < \phi_c$, monotonically increasing for $\phi > \phi_c$ where $\phi_c \in]0, \phi_{max}[$ is a constant, and its first derivative, σ'_e , satisfies

$$\sigma'_e(\phi) = \begin{cases} = 0 & \text{for } \phi \leq \phi_c , \\ > 0 & \text{for } \phi > \phi_c . \end{cases} \quad (1.7)$$

- ϕ_0 and ϕ_2 are piecewise continuous functions such that

$$0 \leq \phi_0(z), \quad \phi_2(t) \leq \phi_{max}, \quad z \in [0, L], \quad t \in [0, T] \quad (1.8)$$

with ϕ_2 changing its monotonicity behaviour a finite number of times.

In the setting of sedimentation theory $\phi(z, t)$ denotes the volumetric solid concentration at height z in time t , q the volume-averaged velocity of a suspension, f_{bk} the Kynch batch flux density function, ϕ_{max} the maximum concentration value, σ_e the effective solid stress, L the height of thickener feeding level, T the total time for the process, $\Delta\rho$ the difference of solid and fluid mass densities, g the acceleration of gravity and ϕ_c a critical concentration value or gel point, see [9, 10, 26, 28] for specific details.

The flux density function (briefly, flux) and the diffusion coefficient (briefly, diffusion) associated with equation (1.1) are defined by

$$f(\phi, t) = q(t)\phi + f_{bk}(\phi) \quad \text{and} \quad a(\phi) = -\frac{f_{bk}(\phi)\sigma'_e(\phi)}{\Delta\rho g\phi}, \quad (1.9)$$

respectively. Moreover we define the integrated diffusion coefficient A by

$$A(\phi) = \int_0^\phi a(s)ds. \quad (1.10)$$

Due to (1.6), (1.7) and (1.9), (1.1) is a second order parabolic partial differential equation for $\phi \in [\phi_c, \phi_{max}]$, a nonlinear hyperbolic conservation law for $\phi \in [0, \phi_c]$ and a linear advection equation for $\phi \in \mathbb{R} - [0, \phi_{max}]$. In brief, the IBVP is referenced as a quasilinear strongly degenerate parabolic equation. We remark that it is sufficient to consider the degeneracy on $[0, \phi_c] \cup \{\phi_{max}\}$ because for ϕ solution of the IBVP we have that $\phi \in [0, \phi_{max}]$ almost everywhere, see [24].

1.2.2 The inverse problem

Experimental results obtained in industrial and laboratory processes for flocculated sedimentation suggest we consider in the model (1.1)-(1.9), the dependence of a finite number of parameters for the functions f_{bk} and σ_e , see [9, 10, 26, 28]. The determination of these parameters implies solving an ‘‘Identification Problem’’ (IP). We can formulate this, in general way, as follows:

IP *Given observation data $\phi^{obs}(z)$, at time T , and the functions $q(t)$, $\phi(z, 0)$, $\phi(L, t)$, explicitly, and $\phi(0, t)$, implicitly, satisfying (1.5), (1.2), (1.3) and (1.4) respectively, find the flux f and the diffusion a with f_{bk} and σ_e satisfying (1.6) and (1.7) such that the weak solution $\phi(z, T)$, in time T , of the IBVP (1.1)-(1.7) is as close as possible to $\phi^{obs}(z)$ in some suitable norm.*

We can write the IP, see [62, 61, 63], as the optimization of a cost function J

$$\min_{f,a} J(\phi(\cdot, T)), \quad (1.11)$$

under the constraint that for ϕ satisfy weakly the IBVP (1.1)-(1.4), for some f and a . A natural example of cost function J is

$$J(\phi) = \frac{1}{2} \int_0^L |\phi(z, T) - \phi^{obs}(z)|^2 dz, \quad (1.12)$$

for other examples see [61].

We particularize the general situation by considering the parametric dependent analytic form of the flux and the diffusion. In this paper, we seek f_{bk} by the usual formula of Richardson and Zaki [81], and σ_e by a constitutive law, [19], i.e.

$$f_{bk}(\phi) = u_\infty \phi \left(1 - \frac{\phi}{\phi_{max}}\right)^C, \quad (1.13)$$

$$\sigma_e(\phi) = \begin{cases} \text{Cte.} & \text{for } \phi \leq \phi_c, \\ \alpha ((\phi/\phi_c)^\beta - 1) & \text{for } \phi > \phi_c, \end{cases} \quad (1.14)$$

where u_∞ is the flow velocity of a singular particle in an unbounded medium and $C > 1, \alpha > 0, \beta > 1, \phi_c \in]0, \phi_{max}[$ are parameters (see [13, 28, 50] for other examples). Thus, in this particular case, if we denote by \mathbf{e} the parameters to find, the problem (1.11) can be formulated in an equivalent way by

$$\min_{\mathbf{e}} J(\phi_{\mathbf{e}}(\cdot, T)),$$

where $\phi_{\mathbf{e}}$ is a weak solution of IBVP (1.1)-(1.4) with f_{bk} and σ_e given by (1.13)-(1.14), respectively.

The nonlinearity and degeneracy implies that solutions of the IBVP may become discontinuous in finite time. Thus, we need to interpret the IBVP in a weak way. In addition, it is well known that the IBVP is ill-posed in the classical weak sense because there is no uniqueness. In order to have a well-posed problem we must consider an additional condition or entropy condition, see [23, 24, 15].

1.3 Theoretical analysis of the IBVP and the IP

1.3.1 The direct problem

We adopt in this paper the definition of the weak entropy solution introduced by Bürger et al. [15].

Definition 1.3.1 Let $\phi \in L^\infty(Q_T) \cap BV(Q_T)$. Then ϕ is an entropy solution of the IBVP if the following four conditions are valid:

$$(i) \quad \partial_z A(\phi) \in L^2(Q_T)$$

(ii) For all $k \in \mathbb{R}$ and for all $\varphi \in C^\infty([0, 1] \times [0, T])$ such that $\varphi \geq 0$ and $\text{supp } \varphi \subset]0, 1] \times]0, T[$ the entropy inequality

$$\begin{aligned} & \int \int_{Q_T} \left\{ |\phi - k| \partial_t \varphi + \text{sgn}(\phi - k)(f(\phi, t) - f(k, t) - \partial_z A(\phi)) \partial_z \varphi \right\} dz dt \\ & + \int_0^T \left\{ - \text{sgn}(\phi_2(t) - k)[f(\gamma_1 \phi, t) - f(k, t) - \gamma_1 \partial_z A(\phi)] \varphi(1, t) \right. \\ & \left. + [\text{sgn}(\gamma_1 \phi - k) - \text{sgn}(\phi_2(t) - k)][A(\gamma_1 \phi) - A(k)] \partial_z \varphi(1, t) \right\} dt \geq 0 \end{aligned}$$

is satisfied.

$$(iii) \quad \text{For almost all } t \in]0, T[$$

$$\gamma_0(f_{bk}(\phi) - \partial_z A(\phi)) = 0$$

$$(iv) \quad \text{For almost all } z \in [0, L]$$

$$\lim_{t \rightarrow 0} \phi(z, t) = \phi_0(z)$$

The terms γ_0 and γ_1 denote the traces in $BV(Q_T)$ at $z = 0$ and $z = 1$, respectively. For a precise definition see [24]. Item (i) is a technical regularity condition, (ii) is the entropy condition, while (iii) and (iv) are the weak formulation for the boundary condition (1.4) and the initial condition (1.2), respectively.

Denition 1.3.1 implies that the direct problem is well-posed, this was proved in [15], see also [24]. In those works, following the ideas Kružkov [72] and Carrillo [31], are established several properties for the entropy solution of the direct problem. In particular, on the basis of hypothesis (1.5)-(1.9), was proved the following theorem.

Theorem 1.3.1 *There is almost one entropy solution ϕ for the IBVP (1.1)-(1.4).*

1.3.2 Existence of solutions to the Inverse Problem

In this section we provide a sufficient condition for the existence of at least one solution for our IP. The existence result is a consequence of the continuous dependence of the entropy solution with respect to the flux and diffusion. The non-uniqueness is a consequence of the degeneracy and the hyperbolic behaviour.

The continuous dependence for a Cauchy Problem with spatially dependent flux was studied in [46, 68]. Their ideas, inspired by the works of Carrillo [31] and Cockburn and Gipenberg [32], can be extended to our IBVP with time dependent flux. First, we need the following lemma.

Lemma 1.3.1 Assume that $A' > 0$, then for any $\varphi \geq 0$ in $C_0^\infty(Q_T)$ and $k \in \mathbb{R}$ we have

$$\begin{aligned} \iint_{Q_T} \{|\phi - k| \partial_t \varphi + \operatorname{sgn}(\phi - k)(f(\phi, t) - f(k, t) - \partial_z A(\phi)) \partial_z \varphi\} dt dz \\ = \lim_{\varepsilon \rightarrow 0} \iint_{Q_T} A'(\phi) (\partial_z \phi)^2 \operatorname{sgn}'_\varepsilon(\phi - k) \varphi dt dz, \end{aligned} \quad (1.15)$$

where ϕ is an entropy solution of the IBVP and

$$\operatorname{sgn}_\varepsilon(r) = \begin{cases} -1 & r < -\varepsilon, \\ r/\varepsilon & -\varepsilon \leq r \leq \varepsilon, \\ 1 & r > \varepsilon. \end{cases}$$

Proof. Let $\psi_\varepsilon(z) = -\operatorname{sgn}_\varepsilon(A^{-1}(z) - k)$ and $A_{\psi_\varepsilon}(\phi) = \int_k^\phi \psi_\varepsilon(A(r)) dr$. Then, by a “weak chain rule” (see [31, 15]) we have

$$-\int_0^T \langle \partial_t \phi, -\operatorname{sgn}_\varepsilon(\phi - k) \varphi \rangle dt = \iint_{Q_T} A_{\psi_\varepsilon}(\phi) \partial_t \varphi dt dz, \quad (1.16)$$

where $\langle \cdot, \cdot \rangle$ denotes the usual pairing between $H^{-1}(]0, L[)$ and $H_0^1(]0, L[)$.

On the other hand, by definition 1.3.1, it follows that

$$\begin{aligned} -\int_0^T \langle \partial_t \phi, \operatorname{sgn}_\varepsilon(\phi - k) \varphi \rangle dt \\ + \iint_{Q_T} [f(\phi, t) - f(k, t) - \partial_z A(\phi)] \partial_z (\operatorname{sgn}_\varepsilon(\phi - k) \varphi) dt dz = 0. \end{aligned} \quad (1.17)$$

Combining (1.16) with (1.17) and passing to the limit when $\varepsilon \rightarrow 0$, we obtain (1.15). \square

Theorem 1.3.2 Let u and v be entropy solutions of the following IBVPs:

$$\begin{aligned} \frac{\partial u}{\partial t} + \frac{\partial}{\partial z} (q_1(t)u + f_{bk_1}(u)) &= \frac{\partial^2 A(u)}{\partial z^2} \\ u(z, 0) = u_0(z), \quad u(L, t) = u_L(t), \quad f_{bk_1}(\phi) - A(u) \frac{\partial u}{\partial z} \Big|_{z=0} &= 0, \end{aligned}$$

and

$$\begin{aligned} \frac{\partial v}{\partial t} + \frac{\partial}{\partial z} (q_2(t)v + f_{bk_2}(v)) &= \frac{\partial^2 B(v)}{\partial z^2} \\ v(z, 0) = v_0(z), \quad v(L, t) = v_L(t), \quad f_{bk_2}(v) - B(v) \frac{\partial v}{\partial z} \Big|_{z=0} &= 0, \end{aligned}$$

respectively, where the assumptions (1.5)-(1.10) are fulfilled by each one of its coefficients, initial and boundary conditions. Then for almost all $t \in [0, T]$

$$\begin{aligned} \|u(\cdot, t) - v(\cdot, t)\|_{L^1([0, L])} &\leq \|u_0 - v_0\|_{L^1([0, L])} + tC\|q_1 - q_2\|_{Lip([0, T])} \\ &\quad + tC\|f_{bk_1} - f_{bk_2}\|_{Lip([0, u_{max}])} + \sqrt{t}C\|\sqrt{a} - \sqrt{b}\|_{L^\infty([\phi_c, \phi_{max}])} \end{aligned}$$

where $C > 0$ is a constant and $a(u) = A'(u)$, $b(v) = B'(v)$.

Proof. It is sufficient to prove this theorem when the equations are uniformly parabolic since the degenerate parabolic case is a consequence of the convergence of the vanishing viscosity method, see [24, 15]. The proof will be done by the generalized doubling of variables technique, see [31].

Let us introduce the test function $\varphi \in C_0^\infty(Q_T \times Q_T)$

$$\varphi(z, t, y, s) = \left\{ \int_{-\infty}^{t-\nu} \rho_\theta(s) ds - \int_{-\infty}^{t-\tau} \rho_\theta(s) ds \right\} \rho_\eta(z-y) \rho_\delta(t-s),$$

where $\nu, \tau \in]0, T[$ are Lebesgue points of $\|u(\cdot, t) - v(\cdot, t)\|_{L^1([0, L])}$, $\theta \in]0, \min\{\nu, T - \tau\}[$, $\eta, \delta > 0$ and $\rho_\varepsilon(x) = (1/\varepsilon)\rho(x/\varepsilon)$ for $\varepsilon > 0$ with $\rho \in C_0^\infty(\mathbb{R})$ such that ρ is a even function, $\rho(r) = 0$ for $|r| > 1$ and $\int \rho(r) dr = 1$.

We apply lemma 1.3.1 twice; Firstly with φ as a test function in (z, t) , $\phi = u(z, t)$, $k = v(y, s)$ and integrate with respect to $(y, s) \in Q_T$. Then we apply the lemma with φ as a test function in (y, s) , $\phi = v(y, s)$ and $k = u(z, t)$ and integrate with respect to $(z, t) \in Q_T$. By summing up the results we obtain

$$\begin{aligned} &- \iiint_{Q_T \times Q_T} \left\{ |u - v|(\partial_t \varphi + \partial_s \varphi) \right. \\ &\quad \left. + \text{sgn}(u - v)[(f(u, t) - f(v, t))\partial_z \varphi - (g(v, s) - g(u, s))\partial_y \varphi] \right. \\ &\quad \left. - [\text{sgn}(u - v)\partial_z A(u)\partial_z \varphi + \text{sgn}(v - u)\partial_y B(v)\partial_y \varphi] \right\} dz dt dy ds \\ &= - \lim_{\varepsilon \rightarrow 0} \iiint_{Q_T \times Q_T} \left\{ A'(u)(\partial_z u)^2 + B'(v)(\partial_y v)^2 \right\} \text{sgn}'_\varepsilon(v - u) \varphi dz dt dy ds \\ &\leq - \lim_{\varepsilon \rightarrow 0} \iiint_{Q_T \times Q_T} 2\sqrt{A'(u)}\sqrt{B'(v)}\partial_z u \partial_y v \text{sgn}'_\varepsilon(v - u) \varphi dz dt dy ds = S_1. \end{aligned}$$

Now, by triangle inequality we get

$$I_1 + I_2 + I_3 \leq - \iiint_{Q_T \times Q_T} |u - v|(\partial_t \varphi + \partial_s \varphi) dz dt dy ds \leq S_1 + S_2 + S_3$$

where

$$I_1 = - \iiint_{Q_T \times Q_T} |u(y, t) - v(y, t)|(\partial_t \varphi + \partial_s \varphi) dz dt dy ds$$

$$\begin{aligned}
I_2 &= - \iiint_{Q_T \times Q_T} |v(y, t) - v(y, s)| (\partial_t \varphi + \partial_s \varphi) dz dt dy ds \\
I_3 &= - \iiint_{Q_T \times Q_T} |u(z, t) - u(y, t)| (\partial_t \varphi + \partial_s \varphi) dz dt dy ds \\
S_2 &= \iiint_{Q_T \times Q_T} \operatorname{sgn}(u - v) [(f(u, t) - f(v, t)) \partial_z \varphi - (g(v, s) - g(u, s)) \partial_y \varphi] dz dt dy ds \\
S_3 &= - \iiint_{Q_T \times Q_T} [\operatorname{sgn}(u - v) \partial_z A(u) \partial_z \varphi + \operatorname{sgn}(v - u) \partial_y B(v) \partial_y \varphi] dz dt dy ds.
\end{aligned}$$

Taking into account that $\partial_t \varphi + \partial_s \varphi = [\rho_\theta(t - \nu) - \rho_\theta(t - \tau)] \rho_\eta(z - y) \rho_\delta(t - s)$ and $\partial_z \varphi + \partial_y \varphi = 0$, we obtain that

$$\begin{aligned}
\lim_{\theta \rightarrow 0} I_1 &= \|u(\cdot, \tau) - v(\cdot, \tau)\|_{L^1([0, L])} - \|u(\cdot, \nu) - v(\cdot, \nu)\|_{L^1([0, L])}, \\
\lim_{\theta \rightarrow 0} I_2 &= 0, \\
\lim_{\theta \rightarrow 0} I_3 &\geq -2\eta \sup_{t \in (\nu, \tau)} |u(\cdot, t)|_{BV([0, L])}, \\
\lim_{\theta, \eta, \delta \rightarrow 0} S_2 &\leq (\|q_1 - q_2\|_{Lip} + \|f_{bk1} - f_{bk2}\|_{Lip}) \sup_{t \in (\nu, \tau)} |u(\cdot, t)|_{BV([0, L])}
\end{aligned}$$

and

$$\lim_{\theta \rightarrow 0} S_1 + S_3 \leq (\tau - \nu) \sup_{t \in (\nu, \tau)} |u(\cdot, t)|_{BV([0, L])} \frac{2}{\eta} \|\sqrt{B'} - \sqrt{A'}\|_{L^\infty([\phi_c, \phi_{max}])}^2.$$

This estimates yields the result. \square

Corollary 1.3.1 *Let $\mathcal{M} = Lip([0, T]) \times Lip([0, u_{max}]) \times L^\infty([\phi_c, \phi_{max}])$. The mapping $\tilde{J} : (q, f_{bk}, a) \mapsto J(\phi_{q, f_{bk}, a})$ defined from \mathcal{M} to \mathbb{R}_+ is continuous. Then, if $(q, f_{bk}, a) \in \mathcal{F}$, where \mathcal{F} is a compact subset of \mathcal{M} , there exist at least one solution for the IP.*

Because the square root is not a Lipschitz function it is interesting to note that if we consider the mapping $\hat{J} : (q, f_{bk}, \sqrt{a}) \mapsto J(\phi_{q, f_{bk}, a})$ from \mathcal{M} to \mathbb{R}_+ is Lipschitz continuous, obtaining a strong result of the continuity of the cost function with respect to “ \sqrt{a} ” rather than “ a ”.

It is known that IP is ill-posed in uniqueness if we consider for instance the identification of the parameters of the flux only a shock observation, when $A = 0$ (see [39] for more details). We hope to solve this problem of uniqueness following the idea of [39]: considering several experimental observations with rarefaction waves and a limited number of real parameters to identify.

1.4 Lagrangian formulation and formal calculus

We define a Lagrangian for the problem (1.11) by setting

$$\mathcal{L}(\phi, \psi; f, a) = J(\phi) - E(\phi, \psi; f, a),$$

where ψ is a smooth test function, ϕ is the state variable and

$$\begin{aligned} E(\phi, \psi; f, a) &= - \int_{Q_T} \left(\phi \frac{\partial \psi}{\partial t} + f(\phi) \frac{\partial \psi}{\partial z} + A(\phi) \frac{\partial^2 \psi}{\partial z^2} \right) \\ &\quad + \int_{z=L} \left(\psi f(\phi_2) - \psi \frac{\partial A}{\partial z}(\phi_2) + A(\phi_2) \frac{\partial \psi}{\partial z} \right) \\ &\quad - \int_{z=0} \left(\psi q(t) \phi + A(\phi) \frac{\partial \psi}{\partial z} \right) + \int_{t=T} \psi \phi - \int_{t=0} \psi \phi_0. \end{aligned}$$

When we formally take the derivative of \mathcal{L} in the direction $\delta\phi$ we obtain

$$\begin{aligned} \left\langle \frac{\partial \mathcal{L}}{\partial \phi}, \delta\phi \right\rangle &= \int_{Q_T} \delta\phi \left(\frac{\partial \psi}{\partial t} + f'(\phi) \frac{\partial \psi}{\partial z} + a(\phi) \frac{\partial^2 \psi}{\partial z^2} \right) \\ &\quad + \int_{z=0} \delta\phi \left(\psi q(t) + a(\phi) \frac{\partial \psi}{\partial z} \right) + \int_{t=T} \delta\phi (\phi - \phi^{obs} - \psi), \end{aligned}$$

where we used the fact that ϕ_0, ϕ_2 and ϕ^{obs} are fixed. Now, we are interested in cancelling $\partial\mathcal{L}/\partial\phi$, then ψ should be a solution of

$$\frac{\partial \psi}{\partial t} + (q(t) + f'_{bk}(\phi)) \frac{\partial \psi}{\partial z} = -a(\phi) \frac{\partial^2 \psi}{\partial z^2}, \quad (z, t) \in Q_T, \quad (1.18)$$

$$\psi(z, T) = \phi(z, T) - \phi^{obs}(z), \quad z \in [0, L], \quad (1.19)$$

$$\psi q(t) - a(\phi) \frac{\partial \psi}{\partial z} \Big|_{z=0} = 0, \quad t \in [0, T]. \quad (1.20)$$

This problem is a backward boundary value problem for a linear parabolic degenerate equation with discontinuous coefficients. The end condition (1.19) depends on the cost function, in this case it corresponds to J . In more general case we have the relation $\int_{t=T} \delta\phi \psi = \langle \partial J / \partial \phi, \delta\phi \rangle$, for all $\delta\phi$ smooth. No boundary condition is needed at $z = L$, for the case of batch sedimentation $q = \phi_2 = 0$, since the characteristics are in the domain for the direct problem.

Let $\phi_{f,a}$ be solution of the IBVP. If $\phi_{f,a}$ is smooth, then we have $E(\phi_{f,a}, \psi; f, a) = 0$ for each ψ smooth. In this way, $\mathcal{L}(\phi_{f,a}, \psi; f, a) = J(\phi_{f,a}) = \tilde{J}(f, a)$. Thus

$$\begin{aligned} \nabla \tilde{J}(f, a) &= \left(\left\langle \frac{\partial \mathcal{L}}{\partial \phi_{f,a}}, \phi_{f,a} \right\rangle + \frac{\partial \mathcal{L}}{\partial f}, \left\langle \frac{\partial \mathcal{L}}{\partial \phi_{f,a}}, \phi_{f,a} \right\rangle + \frac{\partial \mathcal{L}}{\partial a} \right) \\ &= \left\langle \frac{\partial \mathcal{L}}{\partial \phi_{f,a}}, \phi_{f,a} \right\rangle (1, 1) + \left(\frac{\partial \mathcal{L}}{\partial f}, \frac{\partial \mathcal{L}}{\partial a} \right). \end{aligned}$$

In this equality, for the ψ solution of (1.18)-(1.20) the first term in the right-hand side is zero, so that we obtain

$$\begin{aligned} \langle \nabla \tilde{J}, (\delta f, \delta a) \rangle &= \langle \nabla \mathcal{L}, (\delta f, \delta a) \rangle = -\langle \nabla E, (\delta f, \delta a) \rangle \\ &= \int_{Q_T} \left(\delta f(\phi_{f,a}) \frac{\partial \psi}{\partial z}, a(\phi_{f,a}) \frac{\partial^2 \psi}{\partial z^2} \right) \\ &\quad + \int_{z=L} \left(\delta f(\phi_2) \psi, a(\phi_2) \frac{\partial \psi}{\partial z} \right) - \int_{z=0} \left(0, a(\phi_1) \frac{\partial \psi}{\partial z} \right). \end{aligned} \quad (1.21)$$

1.5 Numerical schemes and discrete study

We divide the interval $(0, L)$ into M subintervals of length $\Delta z = L/M$ and the interval $(0, T)$ into N subintervals of length $\Delta t = T/N$. For $n = 0, \dots, N$ and $j = 0, \dots, M$ we will denote by ϕ_j^n the value of the numerical solution at $(j\Delta z, n\Delta t)$ and by $\phi_j^0, \phi_j^n, \phi_j^{obs}$ the corresponding approximation of $\phi_0, \phi_2, \phi^{obs}$, respectively.

At the discrete level the minimization corresponds to the following problem

$$\min_{\mathbf{e}} J_\Delta(\phi_j^n(\mathbf{e}))$$

where J_Δ is the discrete form of the cost function. In the case of (1.12) this is given by

$$J_\Delta(\phi_j^n(\mathbf{e})) = \frac{1}{2} \sum_{j=0}^M \Delta z |\phi_j^N - \phi_j^{obs}|^2.$$

In this way we define the discrete associated Lagrangian by

$$\mathcal{L}_\cdot(\phi_j^n, \psi_j^n; \mathbf{e}) = \frac{1}{\Delta z} J_\Delta(\phi_j^n) - E_\Delta(\phi_j^n, \psi_j^n; \mathbf{e}),$$

where $E_\Delta(\phi_j^n, \psi_j^n; \mathbf{e})$ denotes a discrete weak formulation, this is obtained by summation by parts of the numerical scheme by the same means as we calculated the numerical solution of the IBVP and ψ_j^n will be chosen as the solution of a discrete adjoint problem.

1.5.1 First and second-order explicit EO scheme

In [19] the IBVP was well discretized by means of the first- and second-order Engquist-Osher scheme. The first-order EO scheme is

$$\begin{aligned} \frac{\phi_j^{n+1} - \phi_j^n}{\Delta t} + q(n\Delta t) \frac{\phi_{j+1}^n - \phi_j^n}{\Delta z} + \frac{f_{bk}^{EO}(\phi_j^n, \phi_{j+1}^n) - f_{bk}^{EO}(\phi_{j-1}^n, \phi_j^n)}{\Delta z} \\ = \frac{A(\phi_{j+1}^n) - 2A(\phi_j^n) + A(\phi_{j-1}^n)}{(\Delta z)^2}, \quad j = 1, \dots, M-1, \end{aligned} \quad (1.22)$$

where

$$\begin{aligned} f_{bk}^{EO}(\phi_j^n, \phi_{j+1}^n) &= f_{bk}^+(\phi_j^n) + f_{bk}^-(\phi_{j+1}^n) \\ &= \left[f_{bk}(0) + \int_0^{\phi_j^n} \max(f'_{bk}(s), 0) ds \right] + \left[\int_0^{\phi_{j+1}^n} \min(f'_{bk}(s), 0) ds \right]. \end{aligned}$$

The second order EO scheme is defined by

$$\begin{aligned} \frac{\phi_j^{n+1} - \phi_j^n}{\Delta t} &+ q(n\Delta t) \frac{\phi_{j+1}^L - \phi_j^R}{\Delta z} + \frac{f_{bk}^{EO}(\phi_j^R, \phi_{j+1}^L) - f_{bk}^{EO}(\phi_{j-1}^R, \phi_j^L)}{\Delta z} \\ &= \frac{A(\phi_{j+1}^n) - 2A(\phi_j^n) + A(\phi_{j-1}^n)}{(\Delta z)^2}, \quad j = 1, \dots, M-1, \quad (1.23) \end{aligned}$$

where

$$\phi_j^L = \phi_j^n - \frac{\Delta z}{2} s_j^n \quad \text{and} \quad \phi_j^R = \phi_j^n + \frac{\Delta z}{2} s_j^n.$$

Here $s_1^n = s_{M-1}^n = 0$ and for $j = 2, \dots, M-2$ we have

$$s_j^n = \text{mm}\left(\theta \frac{\phi_j^n - \phi_{j-1}^n}{\Delta z}, \frac{\phi_{j+1}^n - \phi_{j-1}^n}{2\Delta z}, \theta \frac{\phi_{j+1}^n - \phi_j^n}{\Delta z}\right), \quad \theta \in [0, 2],$$

with mm the minmod function

$$\text{mm}(a, b, c) = \begin{cases} \min(a, b, c), & \text{if } a, b, c > 0, \\ \max(a, b, c), & \text{if } a, b, c < 0, \\ 0, & \text{otherwise.} \end{cases}$$

In both cases, first and second order, ϕ_M^n is obtained by $\phi_M^n = \phi_2(n\Delta t)$ and ϕ_0^n by the following formula:

$$\frac{\phi_0^{n+1} - \phi_0^n}{\Delta t} + q(n\Delta t) \frac{\phi_1^n - \phi_0^n}{\Delta z} + \frac{f_{bk}^{EO}(\phi_0^n, \phi_1^n)}{\Delta z} = \frac{A(\phi_1^n) - A(\phi_0^n)}{(\Delta z)^2}.$$

Denoting by $\lambda = \Delta t/\Delta z$ and $\nu = \Delta t/(\Delta z)^2$ the discrete weak formulation $E_\Delta = E_\Delta(\phi_j^n, \psi_j^n, \mathbf{e})$ for (1.22) is given by

$$\begin{aligned} E_\Delta &= \sum_{n,j} \left\{ \phi_j^{n+1} - \phi_j^n + \lambda q(n\Delta t)(\phi_{j+1}^n - \phi_j^n) + \lambda \left(f_{bk}^{EO}(\phi_j^n, \phi_{j+1}^n, \mathbf{e}) \right. \right. \\ &\quad \left. \left. - f_{bk}^{EO}(\phi_{j-1}^n, \phi_j^n, \mathbf{e}) \right) - \nu(A(\phi_{j+1}^n, \mathbf{e}) - 2A(\phi_j^n, \mathbf{e}) + A(\phi_{j-1}^n, \mathbf{e})) \right\} \psi_j^{n+1} \\ &= \sum_{n,j} \left\{ \phi_j^n \left[\psi_j^n - \psi_j^{n+1} + \lambda q(n\Delta t)(\psi_{j-1}^{n+1} - \psi_j^{n+1}) \right] \right. \\ &\quad \left. + \lambda f_{bk}^{EO}(\phi_j^n, \phi_{j+1}^n, \mathbf{e})(\psi_j^{n+1} - \psi_{j+1}^{n+1}) - \nu A(\phi_j^n, \mathbf{e})(\psi_{j-1}^{n+1} - 2\psi_j^{n+1} + \psi_{j+1}^{n+1}) \right\} \end{aligned}$$

$$\begin{aligned}
& + \sum_{j=0}^{M-1} \left\{ \phi_j^N \psi_j^N - \phi_j^0 \psi_j^0 \right\} \\
& + \sum_{n=0}^{N-1} \left\{ \lambda \left[q(n\Delta t) \phi_M^n \psi_{M-1}^{n+1} - q(n\Delta t) \phi_0^n \psi_{-1}^{n+1} + f_{bk}^{EO}(\phi_{M-1}^n, \phi_M^n, \mathbf{e}) \psi_M^{n+1} \right] \right. \\
& \quad \left. - \nu \left[A(\phi_M^n, \mathbf{e}) \psi_{M-1}^{n+1} - A(\phi_{M-1}^n, \mathbf{e}) \psi_M^{n+1} - A(\phi_0^n, \mathbf{e}) (\psi_{-1}^{n+1} - \psi_0^{n+1}) \right] \right\},
\end{aligned}$$

where we have approximated the boundary condition at $z = 0$ by

$$f_{bk}(\phi) \left(1 + \frac{\sigma'(\phi)}{\Delta \rho g \phi} \frac{\partial \phi}{\partial z} \right) \Big|_{z=0} \approx f_{bk}^{EO}(\phi_{-1}^n, \phi_0^n) - \frac{A(\phi_0^n) - A(\phi_{-1}^n)}{\Delta z} = 0.$$

Taking the derivative of E_Δ with respect to ϕ_j^n we obtain

$$\begin{aligned}
\frac{\partial E_\Delta}{\partial \phi_j^n} &= \psi_j^n - \psi_j^{n+1} + \lambda q(n\Delta t) (\psi_{j-1}^{n+1} - \psi_j^{n+1}) \\
&\quad + \lambda \left\{ \min(f'_{bk}(\phi_j^n), 0) \psi_{j-1}^{n+1} + |f'_{bk}(\phi_j^n)| \psi_j^{n+1} - \max(f'_{bk}(\phi_j^n), 0) \psi_{j+1}^{n+1} \right\} \\
&\quad - \nu a(\phi_j^n) (\psi_{j-1}^{n+1} - 2\psi_j^{n+1} + \psi_{j+1}^{n+1}) + \psi_j^N \delta_{n,N} \\
&\quad - \left[\lambda q(n\Delta t) \psi_{-1}^{n+1} - \nu a(\phi_0^n) (\psi_{-1}^{n+1} - \psi_0^{n+1}) \right] \delta_{j,0} \\
&\quad + \left[\lambda \max(f'_{bk}(\phi_{M-1}^n), 0) + \nu a(\phi_{M-1}^n) \right] \psi_M^{n+1} \delta_{j,M-1}.
\end{aligned}$$

If we consider that $(\partial \mathcal{L}_\Delta)/(\partial \phi_j^n)$ should be zero we have the following adjoint scheme to (1.22)

$$\begin{aligned}
\frac{\psi_j^n - \psi_j^{n+1}}{\Delta t} + q(n\Delta t) \frac{\psi_{j-1}^{n+1} - \psi_j^{n+1}}{\Delta z} + \frac{FA_1}{\Delta z} &= \frac{a(\phi_j^n)(\psi_{j-1}^{n+1} - 2\psi_j^{n+1} + \psi_{j+1}^{n+1})}{(\Delta z)^2}, \\
\psi_j^N &= \frac{\partial J_\Delta}{\partial \phi_j^N}, \\
q(n\Delta t) \psi_{-1}^{n+1} - a(\phi_0^n) \frac{\psi_{-1}^{n+1} - \psi_0^{n+1}}{\Delta z} &= 0, \\
\left[\lambda \max(f'_{bk}(\phi_{M-1}^n), 0) + \nu a(\phi_{M-1}^n) \right] \psi_M^{n+1} &= 0,
\end{aligned}$$

where

$$FA_1 = \min(f'_{bk}(\phi_j^n), 0) \psi_{j-1}^{n+1} + |f'_{bk}(\phi_j^n)| \psi_j^{n+1} - \max(f'_{bk}(\phi_j^n), 0) \psi_{j+1}^{n+1}. \quad (1.24)$$

Thus, we obtain the discrete gradient

$$\nabla \hat{J}_\Delta(\mathbf{e}) = -\nabla_{\mathbf{e}} E_\Delta \Delta z$$

$$\begin{aligned}
&= - \sum_{n,j} \left\{ \Delta t \nabla_{\mathbf{e}} f_{bk}^{EO}(\phi_j^n, \phi_{j+1}^n, \mathbf{e})(\psi_j^{n+1} - \psi_{j+1}^{n+1}) \right. \\
&\quad \left. - \lambda \nabla_{\mathbf{e}} A(\phi_j^n, \mathbf{e})(\psi_{j-1}^{n+1} - 2\psi_j^{n+1} + \psi_{j+1}^{n+1}) \right\} \\
&- \sum_{n=0}^{N-1} \left\{ \Delta t \nabla_{\mathbf{e}} f_{bk}^{EO}(\phi_{M-1}^n, \phi_M^n, \mathbf{e}) \psi_M^{n+1} - \lambda \left[\nabla_{\mathbf{e}} A(\phi_M^n, \mathbf{e}) \psi_{M-1}^{n+1} \right. \right. \\
&\quad \left. \left. - \nabla_{\mathbf{e}} A(\phi_{M-1}^n, \mathbf{e}) \psi_M^{n+1} - \nabla_{\mathbf{e}} A(\phi_0^n, \mathbf{e})(\psi_{-1}^{n+1} - \psi_0^{n+1}) \right] \right\}.
\end{aligned}$$

If we use (1.23) instead of (1.22) we get the adjoint scheme

$$\begin{aligned}
\frac{\psi_j^n - \psi_j^{n+1}}{\Delta t} + \frac{q(n\Delta t)Fq}{\Delta z} + \frac{FA_2}{\Delta z} &= \frac{a(\phi_j^n)(\psi_{j-1}^{n+1} - 2\psi_j^{n+1} + \psi_{j+1}^{n+1})}{\Delta z^2} \\
\psi_j^N &= \frac{\partial J_\Delta}{\partial \phi_j^N} \\
q(n\Delta t) \frac{\partial \phi_0^L}{\partial \phi_0^n} \psi_{-1}^{n+1} - a(\phi_0^n) \frac{\psi_{-1}^{n+1} - \psi_0^{n+1}}{\Delta z} &= 0 \\
\lambda \max(f'_{bk}(\phi_{M-1}^R), 0) \frac{\partial \phi_{M-1}^R}{\partial \phi_0^n} \psi_{M-1}^{n+1} + \nu a(\phi_{M-1}^n) \psi_M^{n+1} &= 0,
\end{aligned}$$

where

$$\begin{aligned}
FA_2 &= \min(f'_{bk}(\phi_{j-1}^L), 0) \frac{\partial \phi_{j-1}^L}{\partial \phi_j^n} (\psi_{j-2}^{n+1} - \psi_{j-1}^{n+1}) \\
&\quad + \left\{ \max(f'_{bk}(\phi_{j-1}^R), 0) \frac{\partial \phi_{j-1}^R}{\partial \phi_j^n} + \min(f'_{bk}(\phi_j^L), 0) \frac{\partial \phi_j^L}{\partial \phi_j^n} \right\} (\psi_{j-1}^{n+1} - \psi_j^{n+1}) \\
&\quad + \left\{ \max(f'_{bk}(\phi_j^R), 0) \frac{\partial \phi_j^R}{\partial \phi_j^n} + \min(f'_{bk}(\phi_{j+1}^L), 0) \frac{\partial \phi_{j+1}^L}{\partial \phi_j^n} \right\} (\psi_j^{n+1} - \psi_{j+1}^{n+1}) \\
&\quad + \max(f'_{bk}(\phi_{j+1}^R), 0) \frac{\partial \phi_{j+1}^R}{\partial \phi_j^n} (\psi_{j+1}^{n+1} - \psi_{j+2}^{n+1}), \\
Fq &= \left(\frac{\partial \phi_{j-1}^L}{\partial \phi_j^n} - \frac{\partial \phi_{j-1}^R}{\partial \phi_j^n} \right) (\psi_{j-2}^{n+1} - \psi_{j-1}^{n+1}) + \left(\frac{\partial \phi_j^L}{\partial \phi_j^n} - \frac{\partial \phi_j^R}{\partial \phi_j^n} \right) (\psi_{j-1}^{n+1} - \psi_j^{n+1}) \\
&\quad + \left(\frac{\partial \phi_{j+1}^L}{\partial \phi_j^n} - \frac{\partial \phi_{j+1}^R}{\partial \phi_j^n} \right) (\psi_j^{n+1} - \psi_{j+1}^{n+1}).
\end{aligned}$$

In this case the gradient is given by

$$\begin{aligned}
\nabla \hat{J}_\Delta(\mathbf{e}) = & - \sum_{n,j} \left\{ \Delta t \nabla_{\mathbf{e}} f_{bk}^{EO}(\phi_j^R, \phi_{j+1}^L, \mathbf{e})(\psi_j^{n+1} - \psi_{j+1}^{n+1}) \right. \\
& \quad \left. - \lambda \nabla_{\mathbf{e}} A(\phi_j^n, \mathbf{e})(\psi_{j+1}^{n+1} - 2\psi_j^{n+1} + \psi_{j+1}^{n+1}) \right\} \\
& - \sum_{n=0}^{N-1} \left\{ \Delta t \nabla_{\mathbf{e}} f_{bk}^{EO}(\phi_{M-1}^R, \phi_M^L, \mathbf{e}) \psi_{M-1}^{n+1} - \lambda \left[\nabla_{\mathbf{e}} A(\phi_M^n, \mathbf{e}) \psi_{M-1}^{n+1} \right. \right. \\
& \quad \left. \left. - \nabla_{\mathbf{e}} A(\phi_{M-1}^n, \mathbf{e}) \psi_M^{n+1} - \nabla_{\mathbf{e}} A(\phi_0^n, \mathbf{e}) (\psi_0^{n+1} - \psi_{-1}^{n+1}) \right] \right\}.
\end{aligned}$$

If the CFL condition

$$\lambda \max_{t \in [0, T]} |q(t)| + \lambda \max_{\phi \in [0, \phi_{max}]} |f'_{bk}(\phi)| + 2\nu \max_{\phi \in [\phi_c, \phi_{max}]} |a(\phi)| \leq 1,$$

is satisfied, then the first- and second-order schemes are stable. For another CFL condition see [19].

1.5.2 Implicit and semi-implicit schemes

The direct problem can be discretized by the semi-implicit scheme

$$\begin{aligned}
\frac{\phi_j^{n+1} - \phi_j^n}{\Delta t} + q(n\Delta t) \frac{\phi_{j+1}^n - \phi_j^n}{\Delta z} + \frac{f_{bk}^{EO}(\phi_j^n, \phi_{j+1}^n) - f_{bk}^{EO}(\phi_{j-1}^n, \phi_j^n)}{\Delta z} \\
= \frac{A(\phi_{j+1}^{n+1}) - 2A(\phi_j^{n+1}) + A(\phi_{j-1}^{n+1})}{(\Delta z)^2}, \quad j = 1, \dots, M-1 \quad (1.25)
\end{aligned}$$

or by the implicit scheme

$$\begin{aligned}
\frac{\phi_j^{n+1} - \phi_j^n}{\Delta t} + q((n+1)\Delta t) \frac{\phi_{j+1}^{n+1} - \phi_j^{n+1}}{\Delta z} + \frac{f_{bk}^{EO}(\phi_j^{n+1}, \phi_{j+1}^{n+1}) - f_{bk}^{EO}(\phi_{j-1}^{n+1}, \phi_j^{n+1})}{\Delta z} \\
= \frac{A(\phi_{j+1}^{n+1}) - 2A(\phi_j^{n+1}) + A(\phi_{j-1}^{n+1})}{(\Delta z)^2}, \quad j = 1, \dots, M-1 \quad (1.26)
\end{aligned}$$

In both cases the boundary condition at $z = L$ is discretized by $\phi_M^{n+1} = \phi_2((n+1)\Delta t)$ while the boundary condition at $z = 0$ is calculated by the interior scheme with the following approximation:

$$\left(f_{bk}(u) - \frac{\partial A(u)}{\partial z} \right) \Big|_{z=0} \approx f_{bk}^{EO}(\phi_{-1}^{n+1}, \phi_0^{n+1}) - \frac{A(\phi_0^{n+1}) - A(\phi_{-1}^{n+1})}{\Delta z} = 0.$$

Applying the same strategy as in the previous section to calculate the gradient, we have the adjoint state

$$\begin{aligned} \frac{\psi_j^n - \psi_j^{n+1}}{\Delta t} + q(n\Delta t) \frac{\psi_{j-1}^n - \psi_j^n}{\Delta t} + \frac{FA_1}{\Delta z} &= a(\phi_j^n) \frac{\psi_{j-1}^n - 2\psi_j^n + \psi_{j+1}^n}{(\Delta z)^2} \\ (1 + 2\nu a(\phi_j^N))\psi_j^N &= \frac{\partial J_\Delta}{\partial \phi_j^N} \\ q(n\Delta t)\psi_{-1}^{n+1} - a(\phi_0^n) \frac{\psi_{-1}^n - \psi_0^n}{\Delta z} &= 0 \\ \lambda \max\{f'_{bk}(\phi_{M-1}^n), 0\} \psi_M^{n+1} + \nu a(\phi_{M-1}^n) \psi_M^n &= 0 \end{aligned}$$

for (1.25) and

$$\begin{aligned} \frac{\psi_j^n - \psi_j^{n+1}}{\Delta t} + q(n\Delta t) \frac{\psi_{j-1}^n - \psi_j^n}{\Delta t} + \frac{FAI}{\Delta z} &= a(\phi_j^n) \frac{\psi_{j-1}^n - 2\psi_j^n + \psi_{j+1}^n}{(\Delta z)^2} \\ \psi_j^N &= \frac{\partial J_\Delta}{\partial \phi_j^{N-1}} \\ q(n\Delta t)\psi_{-1}^n - a(\phi_0^n) \frac{\psi_{-1}^n - \psi_0^n}{\Delta z} &= 0 \\ \{\lambda \max\{f'_{bk}(\phi_{M-1}^n), 0\} + \nu a(\phi_{M-1}^n)\} \psi_M^n &= 0 \end{aligned}$$

for (1.26), where FA_1 is given by (1.24) and

$$FAI = \min(f'_{bk}(\phi_j^n), 0) \psi_{j-1}^n + |f'_{bk}(\phi_j^n)| \psi_j^n - \max(f'_{bk}(\phi_j^n), 0) \psi_{j+1}^n.$$

The gradient in the semi-implicit case is given by

$$\begin{aligned} \nabla \hat{J}_\Delta(\mathbf{e}) &= -\nabla_{\mathbf{e}} E_\Delta \Delta z \\ &= -\sum_{n,j} \left\{ \Delta t \nabla_{\mathbf{e}} f_{bk}^{EO}(\phi_j^n, \phi_{j+1}^n, \mathbf{e})(\psi_j^{n+1} - \psi_{j+1}^{n+1}) \right. \\ &\quad \left. - \lambda \nabla_{\mathbf{e}} A(\phi_j^n, \mathbf{e})(\psi_{j-1}^n - 2\psi_j^n + \psi_{j+1}^n) \right\} \\ &\quad - \sum_{n=0}^{N-1} \left\{ \Delta t \nabla_{\mathbf{e}} f_{bk}^{EO}(\phi_{M-1}^n, \phi_M^n, \mathbf{e}) \psi_M^{n+1} - \lambda \left[\nabla_{\mathbf{e}} A(\phi_M^n, \mathbf{e}) \psi_{M-1}^n \right. \right. \\ &\quad \left. \left. - \nabla_{\mathbf{e}} A(\phi_{M-1}^n, \mathbf{e}) \psi_M^n - \nabla_{\mathbf{e}} A(\phi_0^n, \mathbf{e})(\psi_{-1}^n - \psi_0^n) \right] \right\}, \end{aligned}$$

whereas in the implicit case is

$$\begin{aligned} \nabla \hat{J}_\Delta(\mathbf{e}) &= -\nabla_{\mathbf{e}} E_\Delta \Delta z \\ &= -\sum_{n,j} \left\{ \Delta t \nabla_{\mathbf{e}} f_{bk}^{EO}(\phi_j^n, \phi_{j+1}^n, \mathbf{e})(\psi_j^n - \psi_{j+1}^n) \right. \end{aligned}$$

$$\begin{aligned}
& -\lambda \nabla_{\mathbf{e}} A(\phi_j^n, \mathbf{e})(\psi_{j-1}^n - 2\psi_j^n + \psi_{j+1}^n) \Big\} \\
& - \sum_{n=0}^{N-1} \left\{ \Delta t \nabla_{\mathbf{e}} f_{bk}^{EO}(\phi_{M-1}^n, \phi_M^n, \mathbf{e}) \psi_M^n - \lambda \left[\nabla_{\mathbf{e}} A(\phi_M^n, \mathbf{e}) \psi_{M-1}^n \right. \right. \\
& \quad \left. \left. - \nabla_{\mathbf{e}} A(\phi_{M-1}^n, \mathbf{e}) \psi_M^n - \nabla_{\mathbf{e}} A(\phi_0^n, \mathbf{e}) (\psi_{-1}^n - \psi_0^n) \right] \right\}.
\end{aligned}$$

The CFL condition for the semi-implicit scheme is given by

$$\lambda \left(\max_{t \in [0, T]} |q(t)| + \max_{\phi \in [0, \phi_{max}]} |f'_{bk}(\phi)| \right) \leq 1.$$

Meanwhile, the implicit scheme is “CFL free”, see [19].

The unconditional stability of the implicit scheme is useful because it allows the choice of a coarse grid without losing the convergence to the entropy solution. Therefore, the implicit scheme is a good alternative to simulate numerically the physical problem with a few (reasonable) steps of time as opposed to the cumbersome number of steps of time in the explicit and semi-implicit schemes.

1.5.3 Numerical tests

In this section we show some numerical experiments for our parameter identification method described below.

Identification from analytic data.

In our numerical experiments, we always assume that the observation data have some observation errors. In this example, we consider a solution ϕ_1 using exact data, and a noisy data given by $(1 + \delta)\phi_1$, with $\delta = 0.01$. Let us consider f_{bk} and σ_e as in (1.13) and (1.14) with the value of the parameters as given by table 1.1. Then the function $\phi_1(z, t) = z^2 + (t/30000)^2$ is the solution of the following IBVP:

$$\begin{aligned}
\frac{\partial \phi}{\partial t} + \frac{\partial}{\partial z}(f_{bk}(\phi)) &= \frac{\partial^2 A(\phi)}{\partial^2 z} + g_1(z, t) \quad (z, t) \in Q_T \\
\phi(z, 0) &= 0.05 \quad z \in [0, 1] \\
\phi(1, t) &= \phi_1(1, t) \quad t \in [0, T] \\
f_{bk}(\phi) - \frac{\partial A(\phi)}{\partial z} \Big|_{z=0} &= b_1(t) \quad t \in [0, T]
\end{aligned}$$

where the source term g_1 is defined by

$$g_1(z, t) = \frac{\partial \phi_1}{\partial t} + \frac{\partial}{\partial z}(f_{bk}(\phi_1)) - \frac{\partial^2 A(\phi_1)}{\partial^2 z}$$

C	α	β	ϕ_c
15.6	5.0	6.0	0.1

Table 1.1: Parameters for direct simulation [34].

ϕ^{obs}	J	C	α	β	ϕ_c
ϕ_1	8.245e-6	16.09203	5.50039	6.50071	0.10053
$1.01\phi_1$	4.553e-6	16.08856	5.50039	6.50069	0.10280

Table 1.2: Identified parameters with the second order EO explicit scheme.

and the boundary condition b_1 is given by

$$b_1(t) = f_{bk}(\phi_1(0, t)) - \frac{\partial A(\phi_1(0, t))}{\partial z}.$$

The identified parameters given in table 1.2 are obtained with the second-order explicit EO scheme ($\theta = 1.0$), the observed data $\phi_1(z, 12000)$ and the noisy form $1.01\phi_1(z, 12000)$. The grid parameters are $M = 200$ and $CFL = 0.98$. The initial guess corresponds to $C = 16.1, \alpha = 5.5, \beta = 6.5$ and $\phi_c = 0.2$. See Figures 1.1 and 1.2.

Figures 1.1 and 1.2 show the errors between the identified flux and the identified diffusion with respect to the exact data ϕ_1 and the noisy data $(1 + \delta)\phi_2$. From both figures and for this example, we can observe that the diffusion identification will be more sensitive to the measurement errors than the flux identification. Thus, we expect a better identification for the coefficient C , and larger errors for the coefficients of the diffusion α and β .

Identification from simulated data.

We present here a validation of the above Lagrangian method as well as a comparison between the four numerical schemes developed. Since we did not have access to real experimental data, the idea consists in using as an observation the result of a direct simulation given by Concha [34], which is very close to experimental data results. All tests are developed for batch sedimentation velocity $q = 0$, with an initially homogeneous suspension of concentration, namely $\phi_0 = 0.05$. The column is assumed to be closed, that is $\phi_2 = 0$. The physical constants considered are $u_\infty = -1.7200 \times 10^{-4}$, $\phi_{max} = 0.7$, $\Delta\rho = 1500$ and $g = 9.81$.

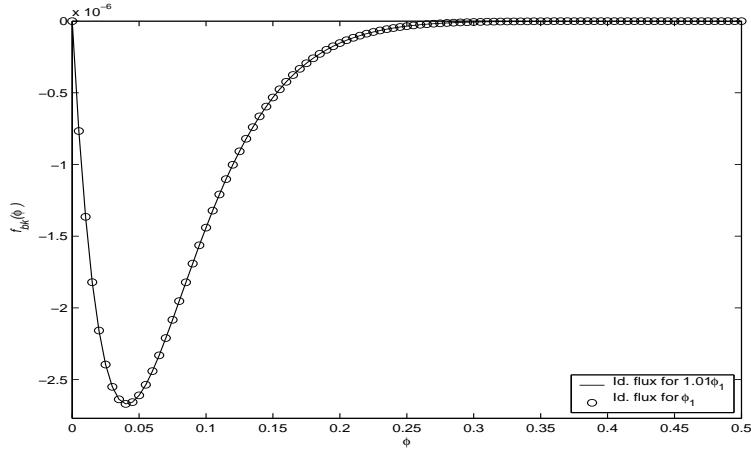


Figure 1.1: The identified flux for the observations ϕ_1 and $1.01\phi_1$, with the second order EO scheme at $T = 12000$ with $\Delta x = 0.005$, $CFL = 0.98$ and $\theta = 1$.

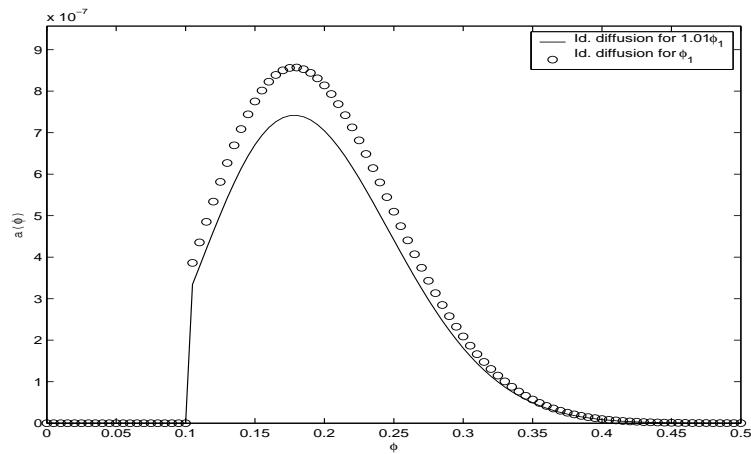


Figure 1.2: The identified diffusion for observations ϕ_1 and $1.01\phi_1$, with the second order EO scheme at $T = 12000$ with $\Delta x = 0.005$, $CFL = 0.98$ and $\theta = 1$.

M	EO1	EO2	EOS	EOI
10	15.34225	15.78555	17.21670	14.00079
50	15.62241	15.95830	15.50107	15.25519
100	15.60076	15.77236	15.70925	15.39562
200	15.58900	15.66536	15.64600	15.47951

Table 1.3: Identification of the flux: values of the parameter C .

We summarize in table 1.1 the parameters used for this simulation, which we want to recover in the inverse problem.

Three tests are performed, the first one is concerned only with the flux identification, parameter C . The second test identifies the diffusion together with the flux, that is parameters C , α and β , the critical concentration ϕ_c being fixed. Finally, we perform complete identification on the four parameters. We start with the following values: $\mathbf{e} = (16.5, 5.0, 6.0, 0.1)$, $\mathbf{e} = (16.5, 5.5, 6.5, 0.1)$ and $\mathbf{e} = (16.5, 5.5, 6.5, 0.2)$ for the test 1, 2 and 3 respectively. Several other initial points were considered with very similar and close results. The identification problem is solved at $T = 12144$ s in the three cases, and a simulation at $T = 30000$ s with the identified parameters is proposed.

For each test, we compare the four numerical schemes presented above, namely the first order scheme **EO1**, the second order scheme **EO2**, the semi-implicit scheme **EOS** and the fully implicit scheme **EOI**. For the explicit and semi-implicit schemes we employ $CFL = 0.5$ and in the case of explicit second order scheme we consider $\theta = 1$. For the implicit scheme we take $\Delta t = \Delta x$. Four different meshes were used, with $M = 10, 50, 100, 200$ steps.

In Figures 1.3 and 1.4 and table 1.3 we show numerical results for the first test. Figures 1.5 and 1.6 and table 1.4 shows the results of the second test and the results for the third test are given in figures 1.7 and 1.8 and table 1.5.

All schemes give satisfactory results, and it should be emphasized that the value of the critical concentration ϕ_c is correctly recovered. The explicit scheme is of course the simplest to implement, but turns out to be the worst in terms of computational time. Indeed, the stability restriction requires such a high number of time steps that the benefit of the computational simplicity is lost (here $\Delta t \approx \Delta x^2$). In the semi-explicit and fully implicit schemes, Gauss-Seidel and Newton methods are needed to solve linear and nonlinear systems, but this is compensated by the less restrictive CFL condition. For the semi-explicit scheme the restriction becomes the same as in the hyperbolic case, that is $\Delta t \approx \Delta x$, and finally the implicit scheme is CFL free.

M	Scheme	C	α	β
10	EO1	15.28574	5.48153	6.46587
	EO2	15.73031	5.48855	6.47817
	EOS	17.02054	5.53600	6.38967
	EOI	13.92483	5.47128	6.46003
50	EO1	15.59482	5.37159	6.07586
	EO2	15.90027	5.43575	6.25493
	EOS	15.72207	5.47796	6.46396
	EOI	15.16180	5.45272	6.43251
100	EO1	15.54214	5.30271	6.12221
	EO2	15.69580	5.42393	6.34055
	EOS	15.57037	5.44297	6.40494
	EOI	15.31271	5.37628	6.30522
200	EO1	15.59575	5.13254	5.86232
	EO2	15.66016	5.30389	6.02822
	EOS	15.58839	5.27715	6.13142
	EOI	15.48579	5.12949	5.88291

Table 1.4: Identification of the flux and diffusion except the critical concentration ϕ_c .

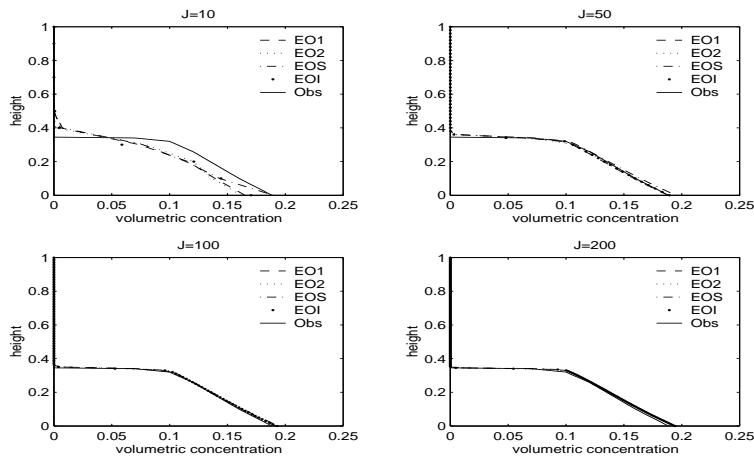


Figure 1.3: Numerical and observed concentration curves at $T = 12144$ for test 1.

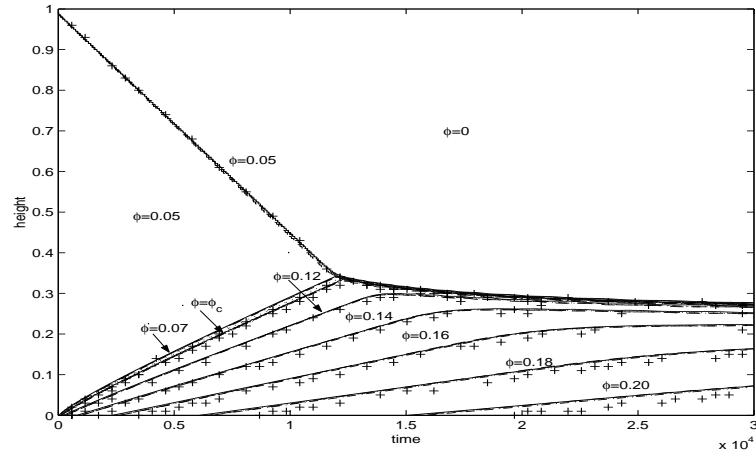


Figure 1.4: Isoconcentrations curves for test 1 with $M = 200$, $T = 30000$. We denote by \dots the results for EO1, by \cdots for EO2, by $-$ for EOS, by $--$ for EOI and by $++$ for the observation data.

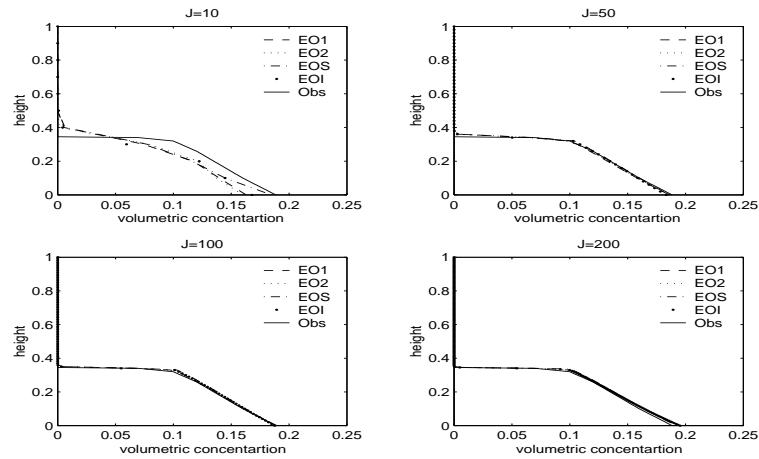


Figure 1.5: Numerical and observed concentration curves at $T = 12144$ for test 2.

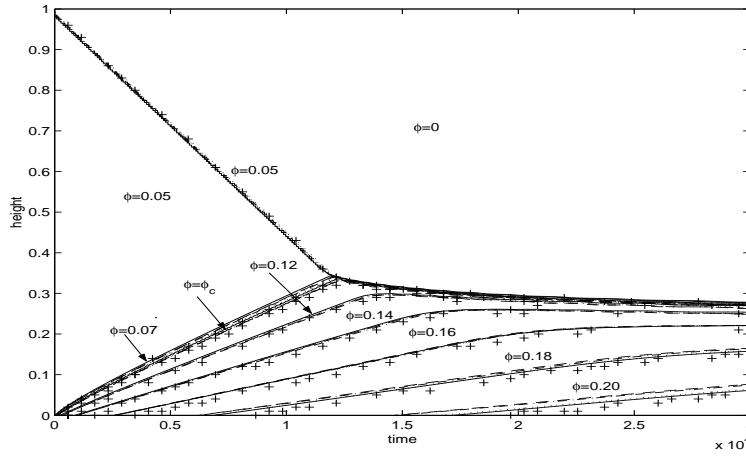


Figure 1.6: Isoconcentrations curves for test 2 with $M = 200$, $T = 30000$. We denote by $\dots - \dots$ the results for EO1, by $\cdots \cdots$ for EO2, by $-$ for EOS, by $--$ for EOI and by $++$ for the observation data.

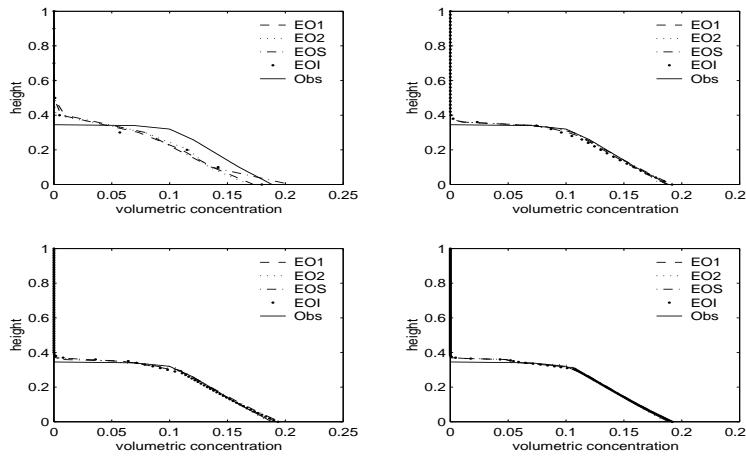


Figure 1.7: Numerical and observed concentration curves at $T = 12144$ for test 3.

M	Scheme	C	α	β	ϕ_c
10	EO1	15.55524	5.50000	6.50000	0.20000
	EO2	15.97224	5.50000	6.50000	0.20000
	EOS	17.67741	5.49970	6.49966	0.32618
	EOI	14.29822	5.50000	6.50000	0.20000
50	EO1	15.82567	5.50008	6.50008	0.10650
	EO2	16.04526	5.50008	6.50008	0.10196
	EOS	16.05873	5.50019	6.50023	0.10593
	EOI	15.83821	5.50004	6.50003	0.11159
100	EO1	16.06191	5.50020	6.50023	0.10849
	EO2	16.07651	5.50023	6.50026	0.10225
	EOS	16.08611	5.50028	6.50039	0.10886
	EOI	16.09204	5.50019	6.50021	0.11058
200	EO1	16.08896	5.50026	6.50035	0.10551
	EO2	16.10939	5.50028	6.50037	0.10183
	EOS	16.09543	5.50031	6.50047	0.10635
	EOI	16.09969	5.50026	6.50034	0.10711

Table 1.5: Identification of the flux and diffusion.

1.6 Parameters identification results using experimental data as observation.

In this section, which is not considered in [37], we present the identification from an experimental data profile given by Pamela Garrido (see [54]). The observation data of volumetric solid-concentration is the conversion of conductivity measurements performed by electronic sensors at different fixed heights of the settling column and for a finite number of time values.

The physical parameters in the experiment are: $L = 0.477$ m, $T = 15873$ s, $\Delta\rho = 1500$ kg/m³, $u_\infty = -1.40 \times 10^{-4}$ m/s and the initial concentration is $\phi_0 = 0.16$.

The flux density function is usually approximated by performing several settling experiments with different initial conditions. The parameters of solid effective stress may be obtained from the momentum equation at equilibrium. This experimental determination of the coefficients of the model is detailed in [33]. The resulting parameters usually have an error, but can be used as initial guess for the gradient method.

In this section, we consider σ_e given by

$$\sigma_e(\phi) = \begin{cases} \text{Cte.} & \text{for } \phi \leq \phi_c, \\ \alpha \exp(\beta\phi) & \text{for } \phi > \phi_c, \end{cases}$$

instead of (1.14).

The Figure 1.9 shows the identified profile using the EOI scheme with $M = 150$ and $N = 500$. The set of identified parameters are

$$C = 15.000357, \quad \alpha = 10.074600, \quad \beta = 5.500340, \quad \phi_c = 0.25001.$$

The initial guess considered is: $C = 10.81, \alpha = 43.77, \beta = 7.22$ and $\phi_c = 0.2$.

We observe that the algorithm gets a set of parameters such that the numerical solution of the phenomenological model predicts the behavior of flocculated suspensions reasonably well. However, the agreement of observed and identified profiles are not of the similar quality of those calculated in the numerical tests (see Section 1.5.3). This leads to the problem of convergence and sensitivity study, which remains to be done in a future work.

Finally, we should comment that this section is a part of the validation study of the numerical method of parameter identification using experimental measurements which is in progress (see [36]).

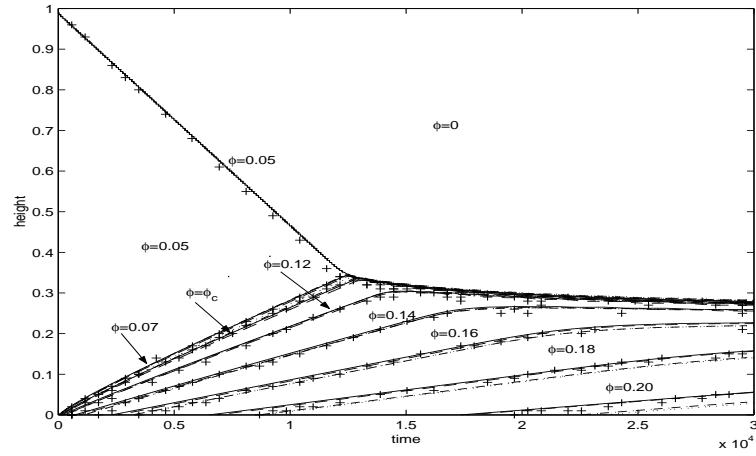


Figure 1.8: Isoconcentrations curves for test 2 with $M = 200$, $T = 30000$. We denote by \dots the results for EO1, by \cdots for EO2, by $-$ for EOS, by $--$ for EOI and by $++$ for the observation data.

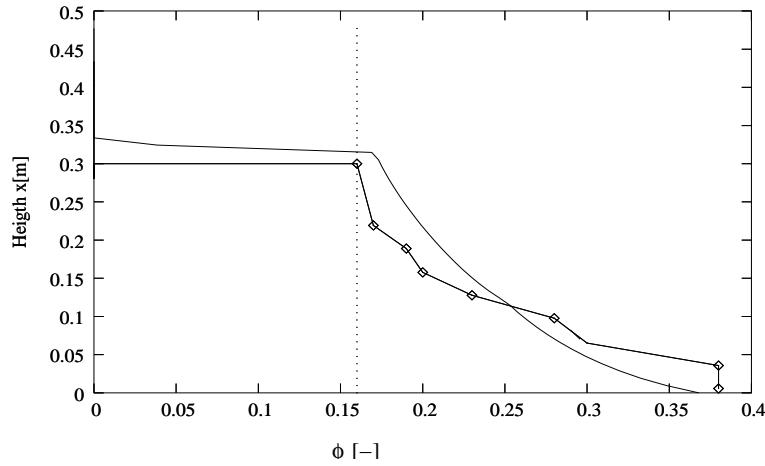


Figure 1.9: Identified and experimental concentration profiles at $T = 15873$ s. The continuous line correspond to the identified profile, the line with \square to the linear interpolation of experimental observed data and the doted line to the initial concentration in the column.

Chapter 2

Numerical identification of parameters for a strongly degenerate convection-diffusion problem modelling centrifugation of flocculated suspensions

This chapter presents the identification of parameters in the flux and diffusion functions for a quasilinear strongly degenerate parabolic equation which models the centrifugation of flocculated suspensions. We consider both a rotating tube and a basket centrifuge at a given angular velocity, and assume that the radius (i.e., the distance to the center of rotation) is the only spatial coordinate. The identification problem is formulated as the problem of minimization of a suitable cost function. Its formal gradient is derived by means of an adjoint equation, which is a backward linear degenerate parabolic equation with discontinuous coefficients. For the numerical approach, the direct problem is discretized by the Engquist-Osher scheme. The discrete Lagrangian formulation provides an associated discrete adjoint state. The conjugate gradient method permits to find numerically the physical parameters. In particular, it allows to identify the critical concentration value at which the model equation changes from second-order parabolic to first-order hyperbolic type. Physically, this critical value is the concentration value at which the solid particles begin to touch each other and determines the change of parabolic to hyperbolic behaviour in the model equation. The new feature as compared to the previous treatment of gravity settling in a column of sedimentation is the dependence of the flux function on the space variable.

2.1 Introduction

There is a large list of authors who have proposed analytical and numerical methods for inverse problems in evolution partial differential equations. For example, parameter identification methods for parabolic PDE can be found in [29, 30, 57, 70, 71] and the references cited in these papers. Recently, in [92] Yamamoto and Zou reconstructed the radiative coefficient and the initial data for a linear parabolic equation using a piecewise linear finite element method for the discretization. Analytical and numerical parameter identification for hyperbolic equations can be found in [61, 62], where it is confirmed that these problems are highly ill-posed, which implies non-uniqueness in most situations. In this chapter, we consider a nonlinear second-order parabolic equation which degenerates to first-order hyperbolic type, where the location of type change is unknown a priori and therefore part of the solution of the problem. For the discretization we consider a finite volume scheme in conservative form with an Engquist-Osher approximation for the numerical flux [11, 19, 42]. For a similar case of gravity settling in a column, such a technique, based on numerical parameter identification, was presented recently by Coronel, James and Sepúlveda [37]. It is the purpose of the present chapter to demonstrate that an analogous treatment also applies to centrifugal separations. In particular, the numerical technique obtained herein could also be applied to extract these model functions from measurements obtained from the newly developed laboratory centrifuge (the so-called “LumiFuge”) described in [49].

We adopt a spatially one-dimensional model for centrifugation of flocculated suspensions that is described in detail in [11], and which is a special case of a spatially multi-dimensional mathematical framework for these mixtures provided by [26]. For a (brief) mathematical analysis of the model, we refer to [20], while an extension to polydisperse flocculated suspensions is studied in [4]. Figure 2.1 shows the two configurations considered: (a) a tube and (b) a basket centrifuge, both rotating at a given angular velocity ω . To distinguish between these cases, we introduce a parameter σ taking the values $\sigma = 0$ and $\sigma = 1$ in the former and latter case, respectively. The unique spatial coordinate is the radius r , which varies between an inner radius $R_0 > 0$ and an outer radius $R > R_0$, corresponding to the suspension meniscus and the outer wall, respectively.

The resulting mathematical model is an initial-boundary value problem for the following strongly degenerate quasilinear parabolic partial differential equation for the solids concentration ϕ as a function of radius r and time t (see [11] for its detailed

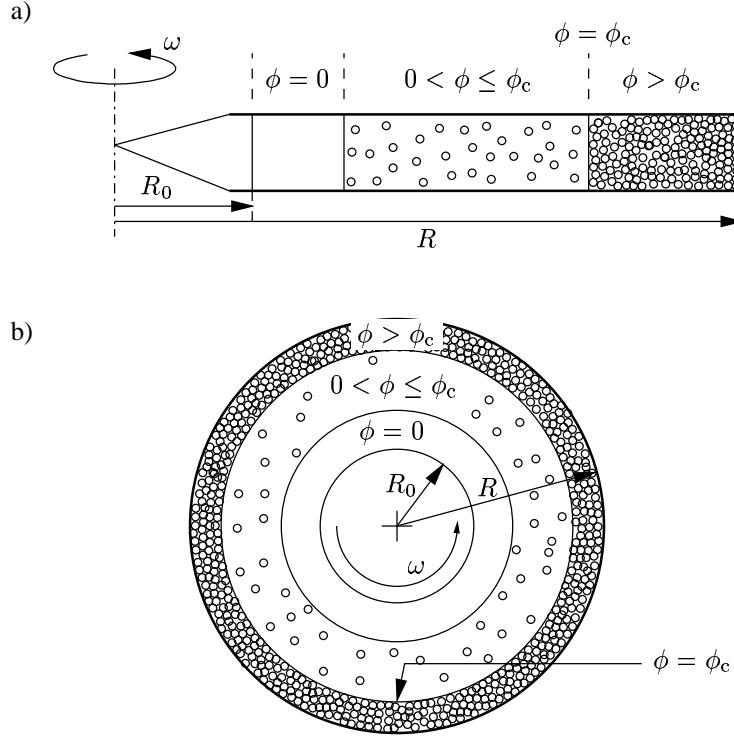


Figure 2.1: (a) Rotating tube with constant cross-section ($\sigma = 0$), (b) rotating axisymmetric cylinder ($\sigma = 1$). The concentration zones are the clear liquid ($\phi = 0$), the hindered settling zone ($0 < \phi \leq \phi_c$) and the compression zone $\phi > \phi_c$.

derivation):

$$\partial_t \phi + \frac{1}{r^\sigma} \partial_r \left(-\frac{\omega^2}{g} r^{1+\sigma} f_{\text{bk}}(\phi) \right) = \frac{1}{r^\sigma} \partial_r (r^\sigma \partial_r A(\phi)), \quad (r, t) \in Q_T, \quad (2.1)$$

where $Q_T := (R_0, R) \times (0, T)$ and g is the acceleration of gravity. The functions $f_{\text{bk}}(\phi)$ and $A(\phi)$ are the Kynch batch flux density function and the integrated diffusion coefficient, respectively, which account for hindered settling and sediment compressibility, respectively. Generically, we assume that $f_{\text{bk}}(\phi)$ is a Lipschitz continuous function satisfying $f_{\text{bk}}(\phi) = 0$ for $\phi \leq 0$ and $\phi \geq \phi_{\max}$, where ϕ_{\max} is the maximum solids concentration, and $f_{\text{bk}}(\phi) < 0$ for $0 < \phi < \phi_{\max}$. The function $A(\phi)$ is given by

$$A(\phi) = \int_0^\phi a(s) ds, \quad a(\phi) := -\frac{f_{\text{bk}}(\phi) \sigma'_e(\phi)}{\Delta \varrho g \phi},$$

where $\Delta \varrho$ is the solid-fluid density difference, σ_e is the effective solid stress function, and σ'_e is its derivative. The effective solid stress is assumed to be zero as long as

the solid flocs are in hindered settling and not in contact, which occurs wherever ϕ does not exceed a critical concentration ϕ_c , and to be a strictly increasing function of ϕ for $\phi > \phi_c$, i.e., we have

$$\sigma_e(\phi) \begin{cases} = 0 & \text{for } \phi \leq \phi_c, \\ > 0 & \text{for } \phi > \phi_c, \end{cases} \quad \sigma'_e(\phi) \begin{cases} = 0 & \text{for } \phi \leq \phi_c, \\ > 0 & \text{for } \phi > \phi_c. \end{cases}$$

Combining the assumptions on f_{bk} and on σ_e , we see that $a(\phi) = 0$ for $\phi \leq \phi_c$ and $\phi \geq \phi_{\max}$ and $a(\phi) > 0$ for $\phi_c < \phi < \phi_{\max}$. Thus, (2.1) is a first-order hyperbolic partial differential equation for $\phi \leq \phi_c$ and $\phi \geq \phi_{\max}$ and a second-order parabolic partial differential equation for $\phi_c < \phi < \phi_{\max}$. Since the degeneracy to hyperbolic type takes place on an interval of solution values of positive length, (2.1) is called *strongly degenerate parabolic*.

In this work, we limit the discussion to two common parametric forms of the model functions f_{bk} and σ_e . We assume that according to the common formulas by Michaels and Bolger [78] (where $0 < \phi_m \leq \phi_{\max}$) and Richardson and Zaki [80] (where $\phi_m = 1$), f_{bk} is given by

$$f_{bk}(\phi) = \begin{cases} u_\infty \phi (1 - \phi/\phi_m)^C & \text{for } 0 < \phi < \phi_{\max}, \\ 0 & \text{for } \phi \leq 0 \text{ and } \phi \geq \phi_{\max}, \end{cases} \quad u_\infty < 0, C \geq 1, \quad (2.2)$$

while the function σ_e is defined by the power-law function [83]

$$\sigma_e(\phi) = \begin{cases} 0 & \text{for } \phi \leq \phi_c, \\ \sigma_0 ((\phi/\phi_c)^k - 1) & \text{for } \phi > \phi_c, \end{cases} \quad \sigma_0 > 0, k \geq 1. \quad (2.3)$$

We limit the treatment furthermore to the case $\phi_m = 1$ corresponding to the Richardson and Zaki [80] flux density function. Observe that by the use of (2.3), the function $\sigma'_e(\phi)$ and therefore also the diffusion coefficient $a(\phi)$ do not only vanish on $[0, \phi_c]$, but are in general even allowed to be discontinuous at $\phi = \phi_c$. This case is, however, covered by the well-posedness analysis of the direct problem. Solutions of (2.1) are in general discontinuous and have to be defined as weak solutions with an additional selection criterion or entropy condition. The complete model for the centrifugation of a suspension of an initial concentration $\phi_0 = \phi_0(r)$ is given by (2.1) together with the initial condition

$$\phi(r, 0) = \phi_0(r), \quad r \in [R_0, R], \quad (2.4)$$

where we assume $\phi_0(r) \in [0, \phi_{\max}]$ for all $r \in [R_0, R]$, and the kinematic boundary conditions

$$\left(\frac{\omega^2 r_b}{g} f_{bk}(\phi) + \partial_r A(\phi) \right) (r_b, t) = 0, \quad t > 0, r_b \in \{R_0, R\}, \quad (2.5)$$

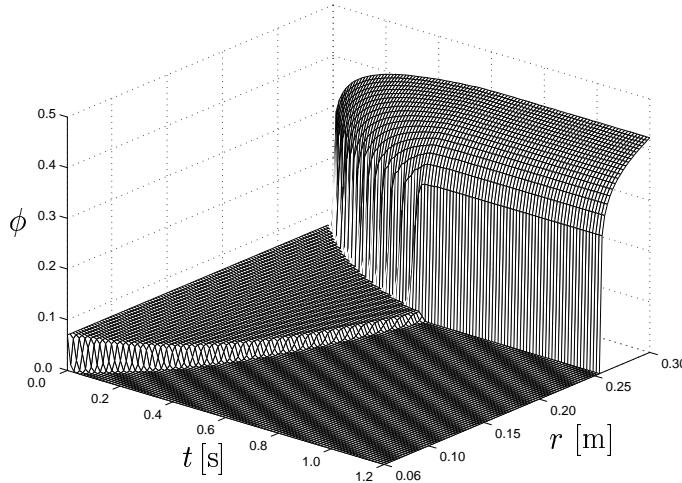


Figure 2.2: Numerical simulation of batch centrifugation of a flocculated suspension of initial concentration $\phi_0 = 0.07$ in a rotating tube (case $\sigma = 0$).

which express that the flux through $r = R_0$ and $r = R$ is zero.

Before proceeding with the discussion, we include here a numerical example to illustrate some basic features of the model. To this end, we use the parameters of Bürger and Concha [11], which shall also be employed for the numerical examples of parameter identification in this chapter, and consider a rotating tube ($\sigma = 0$). The initial concentration is chosen homogeneously as $\phi_0 = 0.07$ on the domain $r \in [R_0, R] = [0.06 \text{ m}, 0.3 \text{ m}]$; the flux function is chosen in accordance with (2.2), where $\phi_m = 1$, $u_\infty = 0.0001 \text{ m/s}$, $C = 5$ and $\phi_{\max} = 0.66$, and the angular velocity ω is such that $R\omega^2 = 10000 \text{ g}$. Additionally, we consider the power law function (2.3) with $\sigma_0 = 5.7 \text{ Pa}$, $k = 9$ and $\phi_c = 0.1$ for the effective solid stress ; and finally, the density $\Delta\varrho = 1660 \text{ kg/m}^3$ and the usual gravitational acceleration $g = 9.81 \text{ m/s}^2$. Figure 2.2 shows the numerical solution for this problem calculated for $0 \leq t \leq T = 1.2 \text{ s}$ and using an explicit version of the second-order numerical scheme of Section 2.7 with $[R_0, R]$ being subdivided into 200 intervals. Figure 2.2 displays some typical features of the solutions produced by the centrifugation model. There are two moving interfaces, the suspension-supernate interface moving towards the outer radius, and the suspension-sediment interface, at which the concentration exceeds the critical concentration ϕ_c , rising from the outer wall. The former is a curved shock that merges with the sediment-suspension interface after roughly 0.6 s. Furthermore, observe that the concentration of the suspension is not constant but decreases as a function of time. Finally, note that the system quickly attains a steady-state sediment profile.

Although it has been possible in recent years to embed this and related sedimentation-consolidation models into a well-posedness framework for strongly degenerate parabolic equations as well as to design numerical discretizations for their solution [16, 19], the main disadvantage of the model, namely the necessity to determine the functions $f_{bk}(\phi)$ and $\sigma_e(\phi)$ by experimentation, has persisted. In some cases [13, 14, 50] it has been possible to determine suitable model functions for published settling experiments by combining published information on the material properties with a trial-and-error procedure. For this procedure to produce reliable results, independent concentration and pore pressure measurements are necessary. The latter are usually obtained by standpipes or transducers, while the concentration is measured by X-ray or γ -ray equipment. This considerable experimental apparatus can be reduced by means of mathematical techniques that permit to obtain these model functions merely from the concentration data, and additional measurement techniques such as computerized axial tomography or small-scale light extinction sensors are can be applied when the pore pressure is not recorded. This chapter presents such a technique.

We restrict ourselves to the parametric forms given by (2.2) and (2.3) and assume that u_∞ is known, such that the problem of determining suitable model functions $f_{bk}(\phi)$ and $\sigma_e(\phi)$ from observations reduces to that of identifying the parameter vector

$$\mathbf{e} = (C, \phi_c, \sigma_0, k)^T \in \mathbb{R}^4.$$

The method consists in minimizing an appropriate cost function which indicates the difference between the solution of the direct problem and the experimentally measured observations. The physical parameters as solutions of the inverse problem are found by minimizing the cost function by a gradient method that in turn relies on the solution of the adjoint equation. The existence of solutions for the inverse problem is a consequence of the continuous dependence of the entropy solutions on the flux and the diffusion for a degenerate parabolic equation (see Theorem 4.1 below, and [32, 37, 46]).

The new aspect here is the dependence of the flux function on the space variable due to the varying body force and the rotating frame of reference. Moreover, in the application observations all over in the physical domain (not only at a single time or location, but at least on a grid) are permitted, and thus need to be considered in the inverse problem.

An aspect that deserves some discussion is the one-dimensionality (and the possibility of extension to several space dimensions) of the model and the parameter identification method. It should be emphasized that the reduction of the sedimentation-

consolidation model to one radial space dimension represents a strong simplification that is acceptable under several restrictions only. In particular, we presume that the angular velocity ω is large enough such that the centrifugal body force is dominant and the gravitational can be neglected, and on the other hand, and on the other hand we assume that ω is not so excessively large that Coriolis effects would become important. Physical phenomena that otherwise appear, and that are therefore excluded by the one-dimensional model, include the effects of gravity settling and the sedimentation of particles onto the backward wall (in the sense of rotation) of a tube centrifuge or a basket centrifuge with compartment walls. We refer to the review paper by Schaflinger [82] for an overview on these and other effects. The limitations of such one-dimensional models, which were first introduced by Anestis and Schneider [1, 2] (see also [77]), are clearly discussed by Ungarish [85]. We assume here the viewpoint that the conditions allowing for the above-mentioned simplification are satisfied. This view is supported by a series of recently published centrifugation experiments [49, 75], which exhibit good agreement with the predictions of one-dimensional models. It should also be mentioned that most measurement principles that are utilized in practice to obtain the raw data for parameter identification, including light extinction measurements [49, 74, 75] and capillary suctioning of samples at fixed radial positions [41], are implicitly based on the assumption that during the centrifugation process, the solids concentration varies in a spatially one-dimensional manner only.

For a more accurate description of the centrifugation process, multi-dimensional versions of the sedimentation-consolidation model utilized here are available [4]. However, passing to several space dimensions does not just simply mean that a multi-dimensional version of the scalar convection-diffusion equation (2.1) has to be solved; rather, we also have to solve equations of motion (for example, a variant of the Navier-Stokes system for incompressible flow) that are strongly coupled to the convection-diffusion equation. These equations of motion yield the volume average velocity of the mixture and the excess pore pressure; in one dimension, the former quantity is determined by boundary conditions (in particular, the mixture velocity vanishes in a closed settling tube or centrifuge), while the latter may be calculated a posteriori from the concentration distribution. Thus, an extension of the model used herein to several space dimensions is feasible at the cost of a significant increase of complexity.

The formal calculus that is employed here to derive a numerical scheme is very general, and may be applied to a large class of problems, including systems of equations (as in [60, 61]) and problems with multi-dimensional spatial domain [57, 70, 71, 92], and could probably also be applied to the more complex multi-

dimensional version of the sedimentation-consolidation model. Finally, in the examples shown in this paper we limit ourselves to ‘one-dimensional’ observation data (concentration profiles given either as a function of radius at a fixed time or as a function of time at a fixed radial position); however, the description of the algorithm also admits ‘two-dimensional’ observation data, i.e. a dense spatial and temporal grid of concentration measurements.

The remainder of this chapter is organized as follows. In Section 2.2, the appropriate entropy solution concept for the direct problem is stated and a known existence and uniqueness result is reviewed. In Section 2.3 the parameter identification problem is formulated as an optimization problem. The existence of a solution to the inverse problem is a direct consequence of a continuous dependence result for entropy solutions of the direct problem, which is stated in Section 2.4. The numerical optimization scheme for the solution of the parameter identification problem, which mimics the formal exact calculus of Section 2.3, is developed in Section 2.5. Calculations of the appropriate weak and discrete weak formulations are collected in an Appendix (Section A.1). A delicate ingredient of the numerical scheme are derivatives of the numerical flux, for which explicit expressions are provided in Section 2.6. The performance of the numerical parameter identification scheme is demonstrated in Section 2.7 by three numerical examples.

2.2 Entropy solutions of the direct problem

2.2.1 General form of the equations

For the analysis of the direct initial-boundary value problem (2.1), (2.4), (2.5), it is useful to study equations of the general form

$$\partial_t \phi + \partial_r f(\phi, r) = \partial_r^2 A(\phi) + g(\phi, r), \quad (2.6)$$

i.e., equations in conservative form with a source term, which includes (2.1) if we choose

$$f(\phi, r) := -\frac{\omega^2 r}{g} f_{\text{bk}}(\phi) - \frac{\sigma}{r} A(\phi), \quad g(\phi, r) := \sigma \left[\frac{\omega^2}{g} f_{\text{bk}}(\phi) + \frac{A(\phi)}{r^2} \right]. \quad (2.7)$$

For the analysis it is useful to recall that $f(\phi, r)$ and $g(\phi, r)$ can be written as

$$f(\phi, r) = k_1(r) f^1(\phi) + k_2(r) f^2(\phi), \quad (2.8)$$

$$k_1(r) := -\omega^2 r / g, \quad k_2(r) := -\sigma / r, \quad f^1(\phi) := f_{\text{bk}}(\phi), \quad f^2(\phi) := A(\phi), \quad (2.9)$$

$$g(\phi, r) = g^1(\phi) + k_3(r) g^2(\phi), \\ k_3(r) := \sigma / r^2, \quad g^1(\phi) := \frac{\sigma \omega^2}{g} f_{\text{bk}}(\phi), \quad g^2(\phi) := A(\phi). \quad (2.10)$$

The boundary conditions (2.5) then take the form

$$(f(\phi, r_b) - \partial_r A(\phi))(r_b, t) = -\frac{\sigma}{r_b} A(\phi(r_b, t)), \quad t > 0, \quad r_b \in \{R_0, R\}. \quad (2.11)$$

2.2.2 Entropy solutions of the direct problem

It is well known that solutions of the initial-boundary value problem (2.1), (2.4), (2.5) (or equivalently, of the initial-boundary value problem (2.4), (2.6), (2.11)) develop discontinuities due to both the nonlinearity of the flux and the degenerate diffusion term, and have to be characterized as weak solutions. To ensure uniqueness, weak solutions have to be defined as entropy solutions.

Definition 2.2.1 *A function $\phi \in L^\infty(Q_T) \cap BV(Q_T)$ is an entropy solution of the initial-boundary value problem (2.4), (2.6), (2.11) if the following conditions are satisfied:*

1. *The integrated diffusion coefficient has the regularity $\partial_r A(\phi) \in L^2(Q_T)$.*

2. *The boundary condition (2.11) holds in the following sense:*

$$\gamma(r_b, t) \left(f(\phi, r_b) - \partial_r A(\phi) + \frac{\sigma}{r_b} A(\phi(r_b, t)) \right) = 0, \quad t > 0, \quad r_b \in \{R_0, R\} \quad (2.12)$$

where $\gamma(\cdot, t)$ is the trace operator.

3. *The initial condition (2.4) holds in the following sense:*

$$\lim_{t \downarrow 0} \phi(r, t) = \phi_0(r) \quad \text{for almost all } r \in (R_0, R).$$

4. *The following entropy inequality holds:*

$$\begin{aligned} \forall \varphi \in C_0^\infty(Q_T), \quad \varphi \geq 0, \quad \forall k \in \mathbb{R}: \int \int_{Q_T} & \left\{ |\phi - k| \partial_t \varphi + \operatorname{sgn}(\phi - k) \right. \\ & \times \left[(f(\phi, r) - f(k, r) - \partial_r A(\phi)) \partial_r \varphi + (f_r(k, r) - g(\phi, r)) \varphi \right] \right\} dt dr \geq 0. \end{aligned}$$

The proof of existence of an entropy solution of the direct problem, following the standard method of vanishing viscosity, is outlined in [20]. That paper also presents a sketch of the uniqueness proof, which in particular relies on results by Carrillo [31] that permit applying Kružkov's "doubling of the variables" technique [72] to strongly degenerate parabolic equations. Both existence and uniqueness proofs are slight modifications of the detailed treatments given in [15]. These results allow us to state the following theorem.

Theorem 2.2.1 *The initial-boundary value problem (2.4), (2.6), (2.11) has a unique entropy solution.*

2.3 Identification as optimization

2.3.1 The inverse problem as an optimization problem with PDE constraint

The parameters of the constitutive functions, which are collected in the common parameter vector \mathbf{e} , depend on the material properties of the suspension considered. The observation data $\hat{\phi}(r, t)$ are assumed to be piecewise constant on rectangles of size $\Delta r \times \Delta t$, and are thus given on a structured grid with

$$(r, t) \in \hat{Q} := \{r_1, \dots, r_j\} \times \{t_1, \dots, t_{\hat{N}}\} \subset \overline{Q}_T := [R_0, R] \times [0, T].$$

The aim is to determine the parameter vector \mathbf{e} for which the solution of the model problem, $\phi(r, t)$, approximates best the observed data $\hat{\phi}(r, t)$ (in a sense yet to be described). That solution $\phi = \phi(\mathbf{e})$ depends on the chosen parameters since the constitutive functions $f = f(\mathbf{e})$ and $A = A(\mathbf{e})$ do. This universal dependence of both the solution and the constitutive functions on the parameters will be suppressed for notational convenience.

The parameter identification problem can be written as a constrained optimization problem, where the constraint is given by the direct initial-boundary value problem (2.1), (2.4), (2.5) in its appropriate weak formulation, see Definition 2.2.1 below. Thus, the optimization problem can be written as

$$\text{minimize } \mathcal{J}(\phi) \text{ under the constraint } \phi = \phi(\mathbf{e}),$$

where the ‘cost function’ $\mathcal{J} = \mathcal{J}(\phi)$ measures the quality of approximation. That cost function depends on the parameter vector \mathbf{e} mediated by the model solution. A natural choice is the L^2 distance between the observed data $\hat{\phi}$ and the solution $\phi = \phi(\mathbf{e})$ of the model function, which gives rise to the cost function

$$\mathcal{J}(\phi(\mathbf{e})) := \frac{1}{2} \int_{\hat{Q}} (\phi(r, t) - \hat{\phi}(r, t))^2 dt dr.$$

Since first-order equations generally have discontinuous solutions, the governing equation (2.1) as constraint on $\phi = \phi(\mathbf{e})$ is replaced by its weak form

$$\begin{aligned} E(\phi, p; \mathbf{e}) := & - \int \int_{Q_T} (\phi \partial_t p + f(\phi, r) \partial_r p + A(\phi) \partial_r^2 p + g(\phi, r) p) dt dr \\ & + \int_{R_0}^R \phi p \Big|_{t=0}^T dr + \int_0^T A(\phi) \left(\partial_r p - \sigma \frac{p}{r} \right) \Big|_{r=R_0}^R dt = 0, \end{aligned}$$

where p is a test function. Summarizing, we have formulated the parameter identification problem, where a parametrization of the model equations for a given observation is sought, as an optimization problem, where the deviation of the model solution (which has to satisfy the model equations as constraint) from the observations is minimized with respect to the set of all parameters.

2.3.2 Lagrangian formulation

In classical optimization, it is a common technique to reformulate the optimization problem by adding (or subtracting) the constraint to the cost function. Thus, we consider the following Lagrangian for the problem “minimize $\mathcal{J}(\phi(\mathbf{e}))$ with respect to \mathbf{e} ”

$$\mathcal{L}(\phi, p; \mathbf{e}) := \mathcal{J}(\phi) - E(\phi, p; \mathbf{e}).$$

The test function p appears here as a generalized Lagrange multiplier related to the constraint $\phi = \phi(\mathbf{e})$. Furthermore, since $E(\phi(\mathbf{e}), p; \mathbf{e}) = 0$, we have that

$$\mathcal{L}(\phi(\mathbf{e}), p; \mathbf{e}) = \mathcal{J}(\phi(\mathbf{e})).$$

In the current application, the cost function is not parametrized by the parameters but only depends on the parameters via the solution of the constraining PDE. Therefore, the cost function cannot simply be differentiated with respect to the parameters. However, optimization algorithms for nonlinear equations (as the conjugate gradient or the Newton method) rely on the total derivative of the cost function, which can here be rewritten and specified with the help of the Lagrangian formulation as

$$\begin{aligned} \frac{d\mathcal{J}(\phi(\mathbf{e}))}{d\mathbf{e}} &= \frac{d\mathcal{L}(\phi(\mathbf{e}), p; \mathbf{e})}{d\mathbf{e}} + \frac{dE(\phi(\mathbf{e}), p; \mathbf{e})}{d\mathbf{e}} \\ &= \left\langle \partial_\phi \mathcal{L}(\phi(\mathbf{e}), p; \mathbf{e}), \frac{d\phi(\mathbf{e})}{d\mathbf{e}} \right\rangle + \frac{d\mathcal{L}(\phi(\mathbf{e}), p; \mathbf{e})}{d\mathbf{e}}, \end{aligned}$$

where $dE/d\mathbf{e}$ vanishes since $\phi(\mathbf{e})$ is considered to remain on the manifold of solutions to the weak formulation. This formal calculation of the total derivative of the cost function splits the problem of finding the total derivative of the cost function up into two parts, corresponding to the two terms in the last sum.

1. The gradient $\nabla_{\mathbf{e}}\phi(\mathbf{e})$ (and therefore $d\phi(\mathbf{e})/d\mathbf{e}$) cannot be calculated, since the solution $\phi(\mathbf{e})$ is not an explicit function of the parameters. This problem can be circumvented if we require that $\partial_\phi \mathcal{L}$ vanishes. The requirement $\partial_\phi \mathcal{L} = 0$ leads to adjoint equations that restrict the test function p . That idea has been introduced and exploited in previous works by James, Sepúlveda, and co-workers [56, 60, 61, 62].

2. Now, given a test function which lets the term $\partial_\phi \mathcal{L} d\phi(\mathbf{e})/d\mathbf{e}$ vanish, the calculation of the total derivative of the cost function reduces to the calculation of the gradient of the Lagrangian with respect to the parameter vector.

2.3.3 Adjoint state

The adjoint state is given by the requirement that $\partial_\phi \mathcal{L} = 0$, which ensures that the term $\partial_\phi \mathcal{L} d\phi(\mathbf{e})/d\mathbf{e}$ vanishes. The conditions on the test function p are obtained after the following straightforward derivation of the derivative of \mathcal{L} taken in the direction of $\delta\phi$:

$$\begin{aligned}\langle \partial_\phi \mathcal{L}(\phi, p; \mathbf{e}), \delta\phi \rangle &= \langle \partial_\phi \mathcal{J}(\phi) - \partial_\phi E(\phi, p; \mathbf{e}), \delta\phi \rangle \\ &= \int \int_{Q_T} \delta\phi(\phi(r, t) - \hat{\phi}(r, t)) \delta_{(r,t) \in \hat{Q}} dt dr \\ &\quad + \int \int_{Q_T} \delta\phi(\partial_t p + \partial_\phi f(\phi, r) \partial_r p + \partial_\phi A(\phi) \partial_r^2 p + \partial_\phi g(\phi, r) p) dt dr \\ &\quad - \int_{R_0}^R \delta\phi p(r, T) dr + \int_0^T \delta\phi \partial_\phi A(\phi) \left(\partial_r p - \sigma \frac{p}{r} \right) \Big|_{r=R_0}^R dt.\end{aligned}$$

The test function p has to be determined in such a way that this quantity vanishes, which leads to the adjoint equation

$$\partial_t p + \partial_\phi f(\phi, r) \partial_r p + \partial_\phi A(\phi) \partial_r^2 p = -(\phi - \hat{\phi}) \delta_{(r,t) \in \hat{Q}} - \partial_\phi g(\phi, r) p \quad \text{for } (r, t) \in Q_T,$$

which is a conservation equation for the unknown function p that arises as a backward problem with end and boundary conditions

$$p(r, T) = 0 \quad \text{for } r \in [R_0, R] \quad \text{and} \quad \left(\partial_r p - \sigma \frac{p}{r_b} \right)(r_b, t) = 0 \quad \text{for } t < T, r_b \in \{R_0, R\}.$$

The adjoint equation is ill-posed since its solution is not unique; different initial settings could lead to the same prescribed end state.

2.3.4 Gradient of cost function

Under the condition that the test function p satisfies the adjoint equations and noting that the cost function $\mathcal{J}(\phi(\mathbf{e}))$ is not a function of the parameter vector \mathbf{e} (thus the gradient $\nabla_{\mathbf{e}} \mathcal{J}(\phi(\mathbf{e}))$ vanishes), the total derivative of the cost function is given by

$$\frac{d\mathcal{J}(\phi(\mathbf{e}))}{d\mathbf{e}} = \frac{\mathcal{L}(\phi(\mathbf{e}), p; \mathbf{e})}{\mathbf{e}}$$

$$\begin{aligned}
&= \frac{d\mathcal{J}(\phi(\mathbf{e}))}{d\mathbf{e}} - \frac{dE(\phi(\mathbf{e}), p; \mathbf{e})}{d\mathbf{e}} \\
&= -\frac{dE(\phi(\mathbf{e}), p; \mathbf{e})}{d\mathbf{e}} \\
&= \int \int_{Q_T} \left(\frac{df(\phi, r)}{d\mathbf{e}} \partial_r p + \frac{dA(\phi, r)}{d\mathbf{e}} \partial_r^2 p + \frac{dg(\phi, r)}{d\mathbf{e}} p \right) dt dr, \quad (2.13)
\end{aligned}$$

which can be used to employ any gradient algorithm in order to minimize the cost function.

2.4 Existence of a solution of the inverse problem

We now establish a sufficient condition for the existence of a solution of the inverse problem. The existence follows from the continuous dependence of the entropy solution of the direct problem with respect to the nonlinearities. The continuous dependence for an initial-value problem with spatially dependent flux was studied in [46, 68]. It is straightforward to extend these results, which in turn are based on works of Carrillo [31] and Cockburn and Gripenberg [32], to the present initial-boundary value problem. The difference to the analysis in [37] consists in slightly different boundary conditions and the presence of a source term in (2.6). We first state the following lemma, where we use the following approximation of the sign function:

$$\operatorname{sgn}_\varepsilon(x) = \begin{cases} \operatorname{sgn}(x) & \text{for } |x| > \varepsilon, \\ x/\varepsilon & \text{for } x \leq \varepsilon. \end{cases}$$

Lemma 2.4.1 *Assume that the function $A(\cdot)$ is smooth and satisfies $A'(s) > 0$. Then the following inequality holds for any $\varphi \in C_0^\infty(Q_T)$ with $\varphi \geq 0$ and $k \in \mathbb{R}$:*

$$\begin{aligned}
&\int \int_{Q_T} \left\{ |\phi - k| \partial_t \varphi + \operatorname{sgn}(\phi - k) (f(\phi, r) - f(k, r) - \partial_r A(\phi)) \partial_r \varphi - \operatorname{sgn}(\phi - k) \right. \\
&\quad \times \left. (\partial_r f(\phi, r) - g(\phi, r)) \varphi \right\} dt dr = \lim_{\varepsilon \downarrow 0} \iint_{Q_T} A'(\phi) (\partial_r \phi)^2 \operatorname{sgn}'_\varepsilon(\phi - k) \varphi dt dr. \quad (2.14)
\end{aligned}$$

Proof. As in [37], we define

$$\psi_\varepsilon(z) := -\operatorname{sgn}_\varepsilon(A^{-1}(z) - k), \quad A_{\psi_\varepsilon}(\phi) := \int_k^\phi \psi_\varepsilon(A(s)) ds.$$

In the proof of this lemma, let $\langle \cdot, \cdot \rangle$ denote the usual pairing between $H^{-1}(a, b)$ and $C_0^1(R_0, R)$. Then the “weak chain rule” (see [31, 68]) implies

$$-\int_0^T \langle \partial_t \phi, -\operatorname{sgn}_\varepsilon(\phi - k) \varphi \rangle dt = \iint_{Q_T} A_{\psi_\varepsilon} \partial_t \varphi dt dr. \quad (2.15)$$

On the other hand, from Definition 2.2.1 we obtain

$$\begin{aligned} - \int_0^T \langle \partial_t \phi, \operatorname{sgn}_\varepsilon(\phi - k) \varphi \rangle dt + \iint_{Q_T} \left\{ (f(\phi, r) - f(k, r) - \partial_t A(\phi)) \partial_r (\operatorname{sgn}_\varepsilon(\phi - k) \varphi) \right. \\ \left. - (\partial_r f(\phi, r) - g(\phi, r)) \operatorname{sgn}_\varepsilon(\phi - k) \varphi \right\} dt dr = 0. \end{aligned} \quad (2.16)$$

Inequality (2.14) follows by combining (2.15) and (2.16) and letting $\varepsilon \downarrow 0$. \square

Theorem 2.4.1 *Let u and v be the entropy solutions of the initial-boundary value problems*

$$\begin{aligned} \partial_t u + \partial_r f_1(u, r) &= \partial_r^2 A(u) + g_1(u, r), \quad (r, t) \in Q_T, \\ u(r, 0) &= u_0(r), \quad r \in (R_0, R), \\ (f_1(u, r_b) - \partial_r A(u))(r_b, t) &= -\frac{\sigma}{r_b} A(u(r_b, t)), \quad r_b \in \{R_0, R\}, t > 0 \end{aligned} \quad (2.17)$$

and

$$\begin{aligned} \partial_t v + \partial_r f_2(v, r) &= \partial_r^2 B(v) + g_2(v, r), \quad (r, t) \in Q_T, \\ v(r, 0) &= v_0(r), \quad r \in (R_0, R), \\ (f_2(v, r_b) - \partial_r B(v))(r_b, t) &= -\frac{\sigma}{r_b} B(v(r_b, t)), \quad r_b \in \{R_0, R\}, t > 0, \end{aligned} \quad (2.18)$$

respectively, where $f_i(u, r) = k_1(r)f_i^1(u) + k_2(r)f_i^2(u)$, $i = 1, 2$, $f_1^2(u) = A(u)$ and $f_2^2(u) = B(u)$, and $k_1(r)$ and $k_2(r)$ are specified in 2.9. Then there exist constants C_1 and C_2 such that the inequality

$$\begin{aligned} \|u(\cdot, t) - v(\cdot, t)\|_{L^1} &\leq \exp(\tilde{C}_3 t) \left\{ \|u_0 - v_0\|_{L^1} + t \left[C_1 \left(\|f_1^1 - f_2^1\|_{\text{Lip}} + \|f_1^2 - f_2^2\|_{\text{Lip}} \right) \right. \right. \\ &\quad \left. \left. + C_2 \left(\|g_1^1 - g_2^1\|_{L^\infty[0, \phi_{\max}]} + \|g_1^2 - g_2^2\|_{L^\infty} \right) \right] + C_D \sqrt{t} \|\sqrt{a} - \sqrt{b}\|_{L^\infty} \right\} \end{aligned} \quad (2.19)$$

holds for almost all $t \in [0, T]$, where $L^1 = L^1(R_0, R)$, $L^\infty = L^\infty[0, \phi_{\max}]$, $a(u) = A'(u)$, $b(u) = B'(u)$, and

$$\tilde{C}_3 := \|g_2^1\|_{\text{Lip}} + \|k_3\|_{L^\infty[R_0, R]} \|g_2^2\|_{\text{Lip}}.$$

Observe that the convective flux in our problem is given by the sum of two terms of the form “ $k(r)f(u)$ ”, as in [68] (see also [69]). However, in our application the functions $k_1(r)$ and $k_2(r)$ do not express a material property for which identification is sought, and therefore they are considered to be the same for two different solutions that are compared. Thus, the continuous dependence estimate (2.19) does not express continuity with respect to k_1 or k_2 .

Proof. The proof is a straightforward extension of the analysis by Evje, Karlsen and Risebro [46] including an application of Gronwall’s inequality. \square

The following corollary is a direct consequence of Theorem 2.4.1.

Corollary 2.4.1 *The mapping*

$$\tilde{\mathcal{J}} : [\text{Lip} \cap L_{[0,\phi_{\max}]}^\infty] \times L_{[0,\phi_{\max}]}^\infty \times [\text{Lip} \cap L_{[0,\phi_{\max}]}^\infty] =: \mathcal{M} \ni (f, A, g) \mapsto \mathcal{J} \in \mathbb{R}$$

is continuous. Further, if $(f, g, A) \in \mathcal{F}$, where \mathcal{F} is a compact subset of \mathcal{M} , then there exists at least one solution of the inverse problem.

2.5 Optimization scheme for identification

Since a simple discretization of the formal gradient (2.13) leads to an incorrect and unstable scheme, the formal exact calculus of Section 2.3 needs to be transferred to its discrete version. We introduce a standard rectangular grid on Q_T by choosing $J, N \in \mathbb{N}$ and setting $\Delta r := (R - R_0)/J$, $\Delta t := T/N$, $r_j := R_0 + j\Delta r$ and $t_n := n\Delta t$. The numerical scheme for the solution of the direct problem is written in conservative form as a marching formula for the interior points (“interior scheme”), for $j = 1, \dots, J-1$,

$$\phi_j^{n+1} = \phi_j^n - \lambda_j(F_{j+1/2}^n(\mathbf{e}) - F_{j-1/2}^n(\mathbf{e})) + \mu_j(\mathcal{A}_{j+1/2}^n(\mathbf{e}) - \mathcal{A}_{j-1/2}^n(\mathbf{e})), \quad (2.20)$$

where $\lambda_j = \mu_j := \Delta t/(r_j^\sigma \Delta r)$, supplemented by the initial condition

$$\phi_j^0 = \phi_j^{\text{init}}, \quad j = 0, \dots, J \quad (2.21)$$

and the following discrete versions of the boundary conditions (2.5):

$$\lambda_0 F_{-1/2}^n(\mathbf{e}) - \mu_0 \mathcal{A}_{-1/2}^n(\mathbf{e}) = 0, \quad (2.22)$$

$$\lambda_J F_{J+1/2}^n(\mathbf{e}) - \mu_J \mathcal{A}_{J+1/2}^n(\mathbf{e}) = 0. \quad (2.23)$$

Inserting (2.22) and (2.23) into the formula for the interior scheme, (2.20), we obtain update formulae for the boundary solution values ϕ_0^n and ϕ_J^n , respectively (“boundary scheme”).

The numerical flux

$$F_{j+1/2}^n = F_{j+1/2}^n(\phi_{j-K+1}^n, \dots, \phi_{j+K}^n, r_{j+1/2})$$

and the numerical diffusion term

$$\mathcal{A}_{j+1/2}^n = \mathcal{A}_{j+1/2}^n(\phi_{j-\bar{K}+1}^n, \dots, \phi_{j+\bar{K}}^n, r_{j+1/2})$$

are specified in the next section.

The discrete versions of the unknown ϕ and the test function p are denoted by ϕ_Δ and p_Δ , and ϕ_j^n and p_j^n are the constant values of ϕ_Δ and p_Δ at (r_j, t_n) , $(j, n) \in Q_\Delta$, respectively, where $Q_\Delta := (0, \dots, J-1) \times (0, \dots, N-1)$. The calculus for the discrete formulation is analogous to the formal continuous one and thus will also be presented in an analogous structuring. Whereas the formal calculus has focused on the formal motivation, the discrete calculus is focused on an efficient scheme as result.

2.5.1 Discrete optimization with PDE as constraint

The discrete minimization problem is stated as

$$\text{minimize } \mathcal{J}_\Delta(\phi_\Delta(\mathbf{e})) \text{ with respect to } \mathbf{e},$$

where the discrete cost function is

$$\mathcal{J}_\Delta(\phi_\Delta(\mathbf{e})) := \frac{\Delta r \Delta t}{2} \sum_{(j,n) \in \hat{Q}_\Delta} (\phi_j^n(\mathbf{e}) - \hat{\phi}_j^n)^2,$$

where we assume that the identification points in \hat{Q} are actually grid points and $\hat{Q}_\Delta \subset Q_\Delta$ is the index set associated with \hat{Q} . The constraint to the optimization problem,

$$E_\Delta(\phi_\Delta(\mathbf{e}), p_\Delta; \mathbf{e}) = 0,$$

serves as a numerical scheme for the computation of p_Δ and is defined by the expression

$$\begin{aligned} E_\Delta(\phi_\Delta, p_\Delta; \mathbf{e}) := & \sum_{(j,n) \in Q_\Delta} \left\{ \phi_j^n(p_j^n - p_j^{n+1}) + F_{j+1/2}^n(\mathbf{e})(\lambda_j p_j^{n+1} - \lambda_{j+1} p_{j+1}^{n+1}) \right. \\ & \left. - \mathcal{A}_{j+1/2}^n(\mathbf{e})(\mu_j p_j^{n+1} - \mu_{j+1} p_{j+1}^{n+1}) \right\} + \sum_{j=0}^{J-1} (\phi_j^N p_j^N - \phi_j^0 p_j^0), \end{aligned} \quad (2.24)$$

which is derived by multiplying the scheme with p_j^{n+1} and summing over j and n such that in the final form, the sums are taken over differences of the test function. This imitates the continuous weak form, as is detailed in the Appendix (see Section A.1).

2.5.2 Discrete Lagrangian formulation

The discrete Lagrangian formulation

$$\mathcal{L}_\Delta(\phi_\Delta, p_\Delta; \mathbf{e}) := \frac{1}{\Delta t \Delta r} \mathcal{J}_\Delta(\phi_\Delta) - E_\Delta(\phi_\Delta, p_\Delta; \mathbf{e}) \quad (2.25)$$

is again used to allow an explicit expression for the total derivative of the cost function

$$\begin{aligned} \frac{d\mathcal{J}_\Delta(\phi_\Delta)}{d\mathbf{e}} &= \Delta t \Delta r \left[\frac{d\mathcal{L}_\Delta(\phi_\Delta, p_\Delta; \mathbf{e})}{d\mathbf{e}} + \frac{dE_\Delta(\phi_\Delta, p_\Delta; \mathbf{e})}{d\mathbf{e}} \right] \\ &= \Delta t \Delta r \left\langle \partial_{\phi_\Delta} \mathcal{L}_\Delta(\phi_\Delta, p_\Delta; \mathbf{e}), \frac{d\phi}{d\mathbf{e}} \right\rangle + \Delta t \Delta r \frac{d\mathcal{L}_\Delta(\phi_\Delta, p_\Delta; \mathbf{e})}{d\mathbf{e}}, \end{aligned}$$

which again splits the problem up into two parts: If, firstly, the adjoint equation prescribes the test function p_Δ such that $\partial_{\phi_\Delta} \mathcal{L} = 0$ and the corresponding term vanishes, then, secondly, the gradient with respect to the parameters of the Lagrangian gives a descent direction of the parameter vector for the algorithm.

2.5.3 Discrete adjoint state

From

$$\begin{aligned}\partial_{\phi_j^n} \mathcal{L}_\Delta &= \partial_{\phi_j^n} \mathcal{J}_\Delta(\phi_\Delta) - \partial_{\phi_j^n} E_\Delta(\phi_\Delta, p_\Delta; \mathbf{e}) \\ &= p_j^n - p_j^{n+1} + p_j^N \delta_{nN} + \sum_{k=-K}^{K-1} \partial_{\phi_j^n} F_{j+k+1/2}^n(\mathbf{e})(\lambda_{j+k} p_{j+k}^{n+1} - \lambda_{j+k+1} p_{j+k+1}^{n+1}) \\ &\quad - \sum_{\ell=-\bar{K}}^{\bar{K}-1} \partial_{\phi_j^n} \mathcal{A}_{j+\ell+1/2}^n(\mathbf{e})(\mu_{j+\ell} p_{j+\ell}^{n+1} - \mu_{j+\ell+1} p_{j+\ell+1}^{n+1}) + (\phi_j^n(\mathbf{e}) - \hat{\phi}_j^n) \delta_{(j,n) \in \hat{Q}_\Delta}\end{aligned}$$

we see that the adjoint scheme for the discrete test function p_j^n is given by

$$\begin{aligned}p_j^n &= p_j^{n+1} - \sum_{k=-K}^{K-1} \partial_{\phi_j^n} F_{j+k+1/2}^n(\mathbf{e})(\lambda_{j+k} p_{j+k}^{n+1} - \lambda_{j+k+1} p_{j+k+1}^{n+1}) \\ &\quad + \sum_{\ell=-\bar{K}}^{\bar{K}-1} \partial_{\phi_j^n} \mathcal{A}_{j+\ell+1/2}^n(\mathbf{e})(\mu_{j+\ell} p_{j+\ell}^{n+1} - \mu_{j+\ell+1} p_{j+\ell+1}^{n+1}) - (\phi_j^n(\mathbf{e}) - \hat{\phi}_j^n) \delta_{(j,n) \in \hat{Q}_\Delta} \\ &\quad \text{for } j = 0, 1, \dots, J \text{ and } n = N-1, N-2, \dots, 0\end{aligned}\tag{2.26}$$

with the end condition $p_j^N = 0$ for $j \in \{0, \dots, \max(K, \bar{K})\} \cup \{J - \max(K, \bar{K}) + 1, \dots, J\}$, and we consider the conventional notation $F_{k+1/2}^n = \mathcal{A}_{\ell+1/2}^n = 0$ for $\ell, k \leq -1$ and $\ell, k \geq J$.

2.5.4 Discrete gradient of cost function

The discrete gradient of the cost function

$$\nabla_{\mathbf{e}} \mathcal{J}_\Delta(\mathbf{e}) = \Delta r \Delta t \nabla_{\mathbf{e}} \mathcal{L}_\Delta(\phi_\Delta(\mathbf{e}), p_\Delta; \mathbf{e}) = -\Delta r \Delta t \nabla_{\mathbf{e}} E_\Delta(\phi_\Delta(\mathbf{e}), p_\Delta; \mathbf{e})$$

where

$$\nabla_{\mathbf{e}} E_\Delta = \sum_{(j,n) \in Q_\Delta} \nabla_{\mathbf{e}} F_{j+1/2}^n(\mathbf{e})(\lambda_j p_j^{n+1} - \lambda_{j+1} p_{j+1}^{n+1}) - \nabla_{\mathbf{e}} \mathcal{A}_{j+1/2}^n(\mathbf{e})(\mu_j p_j^{n+1} + \mu_{j+1} p_{j+1}^{n+1})$$

finally gives the steepest descent direction. Summarizing, the discrete calculus provides a precise procedure for the treatment of the optimization problem. Note that the calculus performed up to now is independent of the numerical flux function. Now, it only remains to specify how the derivatives of the numerical flux (as desired in (2.27)) are obtained, i.e. “explicit” expressions for the gradients $\nabla_{\mathbf{e}} F_{j+1/2}^n$ and $\nabla_{\mathbf{e}} \mathcal{A}_{j+1/2}^n$ and, more important, for the partial derivatives $\partial_{\phi_j^n} F_{j+k+1/2}^n$ for $k = -K, \dots, K-1$ and $\partial_{\phi_j^n} \mathcal{A}_{j+\ell+1/2}^n$ for $\ell = -\bar{K}, \dots, \bar{K}-1$ need to be derived. These quantities are required in the adjoint scheme 2.26.

2.6 Derivatives of numerical fluxes

The numerical scheme (2.20) is based on the non-conservative form (2.1) of the governing equation. Thus, the numerical fluxes $F_{j+1/2}^n$ approximate the physical flux $f(\phi, r) = -\omega^2 r^{1+\sigma} f_{\text{bk}}(\phi)/g$, and $\mathcal{A}_{j+1/2}^n$ is an approximation of

$$\hat{A} := r^\sigma \partial_r A(\phi).$$

Here, we consider the forward finite difference approximation of $\partial_r A$, which gives

$$\mathcal{A}_{j+1/2}^n := \frac{A(\phi_{j+1}^n) - A(\phi_j^n)}{\Delta r} r_{j+1/2}^\sigma.$$

2.6.1 Engquist-Osher numerical flux

We employ the numerical flux function corresponding to the Engquist-Osher generalized upwind scheme [11, 19, 42] defined by

$$F^{\text{EO}}(u, v, r) := f(0) + \int_0^u \max\{\partial_s f(s, r), 0\} ds + \int_0^v \min\{\partial_s f(s, r), 0\} ds.$$

In the present application, where the dependence of the numerical flux on the position r is of multiplicative type, and the function $f(\cdot, r)$ has only one single maximum, denoted u_m , the integrals in this definition can be easily evaluated, and leads to the explicit formula

$$F^{\text{EO}}(u, v, r) = \begin{cases} f(u, r) & \text{for } u \leq u_m, \quad v \leq u_m, \\ f(u, r) + f(v, r) - f(u_m, r) & \text{for } u \leq u_m, \quad v > u_m, \\ f(u_m, r) & \text{for } u > u_m, \quad v \leq u_m, \\ f(v, r) & \text{for } u > u_m, \quad v > u_m, \end{cases}$$

which is used to the evaluation of $F_{j+1/2}^n$ as will be specified below in (2.27).

2.6.2 Differentiation with respect to the parameters

In view of (2.27), and (2.27), the problem of calculating the gradient of the numerical flux $F_{j+1/2}^n$ with respect to the parameters is shifted to the calculation of the gradient of the flux f . In addition, from (2.27) the calculation of $\nabla_{\mathbf{e}} \mathcal{A}_{j+1/2}^n$ is given in terms of $\nabla_{\mathbf{e}} \hat{A}$. In our application, the terms introducing the dependence on r do not depend on the parameters, i.e. we do not deal with shape optimization. Thus, $\nabla_{\mathbf{e}} F_{j+1/2}^n$ is calculated in terms of $\nabla_{\mathbf{e}} f(\phi, r) = r^{1+\sigma} \nabla_{\mathbf{e}} f_{\text{ck}}(\phi)$, where $f_{\text{ck}}(\phi) = -\omega^2 f_{\text{bk}}(\phi)/g$ and $\nabla_{\mathbf{e}} \mathcal{A}_{j+1/2}^n$ is calculated in terms of $\nabla_{\mathbf{e}} A$ by

$$\nabla_{\mathbf{e}} \mathcal{A}_{j+1/2}^n := \frac{r_{j+1/2}^\sigma}{\Delta r} (\nabla_{\mathbf{e}} A(\phi_{j+1}^n) - \nabla_{\mathbf{e}} A(\phi_j^n)).$$

2.6.3 Differentiation with respect to the unknown

We note that $\bar{K} = 1$ in (2.27), which implies the following local derivatives of $\mathcal{A}_{j+1/2}^n$:

$$\partial_{\phi_j^n} \mathcal{A}_{j+1/2}^n = -\frac{a(\phi_j^n) r_{j+1/2}^\sigma}{\Delta r} \quad \text{and} \quad \partial_{\phi_{j+1}^n} \mathcal{A}_{j+1/2}^n = \frac{a(\phi_{j+1}^n) r_{j+1/2}^\sigma}{\Delta r}.$$

The local derivatives of the numerical fluxes with respect to the solution ϕ_j^n are, for the first-order scheme with $F_{j+1/2}^n = F^{\text{EO}}(\phi_j^n, \phi_{j+1}^n, r_{j+1/2}^n)$ (i.e., $K = 1$) and as a consequence of (2.27), given by

$$\partial_{\phi_j^n} F_{j+1/2} = \max\{\partial_{\phi_j^n} f(\phi_j^n, r_{j+1/2}), 0\}, \quad \partial_{\phi_{j+1}^n} F_{j+1/2} = \min\{\partial_{\phi_{j+1}^n} f(\phi_{j+1}^n, r_{j+1/2}), 0\}.$$

On the basis of the numerical flux, a second-order scheme can be constructed with a linear reconstruction of the unknown by the slopes

$$s_j^n = \text{MM}\left(\theta \frac{\phi_j^n - \phi_{j-1}^n}{\Delta r}, \frac{\phi_{j+1}^n - \phi_{j-1}^n}{2\Delta r}, \theta \frac{\phi_{j+1}^n - \phi_j^n}{\Delta r}\right),$$

where $\theta \in [0, 2]$ is a parameter and the standard minmod function

$$\text{MM}(a, b, c) := \begin{cases} \min\{a, b, c\} & \text{if } a, b, c > 0, \\ \max\{a, b, c\} & \text{if } a, b, c \leq 0, \\ 0 & \text{otherwise} \end{cases}$$

is used to ensure the TVD property. In order to facilitate the latter calculus, we introduce the local slopes

$$s_j^- := \theta \frac{\phi_j^n - \phi_{j-1}^n}{\Delta r}, \quad s_j^0 := \frac{\phi_{j+1}^n - \phi_{j-1}^n}{2\Delta r}, \quad s_j^+ := \theta \frac{\phi_{j+1}^n - \phi_j^n}{\Delta r}$$

and the indicator functions

$$\chi_j^* := \begin{cases} 1 & \text{if } s_j^* = \max\{s_j^-, s_j^0, s_j^+\} > 0 \text{ or } s_j^* = \min\{s_j^-, s_j^0, s_j^+\} \leq 0, \\ 0 & \text{otherwise,} \end{cases}$$

which select the active slope such that

$$s_j^n = \text{MM}(s_j^-, s_j^0, s_j^+) = \chi_j^- s_j^- + \chi_j^0 s_j^0 + \chi_j^+ s_j^+.$$

In terms of the reconstructions of the solution at the cell boundaries

$$\phi_j^R = \phi_j^n + \frac{\Delta r}{2} s_j^n, \quad \phi_j^L = \phi_j^n - \frac{\Delta r}{2} s_j^n,$$

the numerical flux for the second-order scheme used in (2.20) is based on the evaluation

$$F_{j+1/2}^n = F^{\text{EO}}(\phi_j^{\text{R}}, \phi_{j+1}^{\text{L}}, r_{j+1/2}).$$

From 2.27, we see that, in this case $K = 2$. An application of the chain rule shows that the local derivatives of the second order numerical flux are given by

$$\begin{aligned}\partial_{\phi_{j-1}^n} F_{j+1/2}^n &= \partial_{\phi_j^{\text{R}}} F_{j+1/2}^n \partial_{\phi_{j-1}^n} \phi_j^{\text{R}}, \\ \partial_{\phi_j^n} F_{j+1/2}^n &= \partial_{\phi_j^{\text{R}}} F_{j+1/2}^n \partial_{\phi_j^n} \phi_j^{\text{R}} + \partial_{\phi_{j+1}^{\text{L}}} F_{j+1/2}^n \partial_{\phi_j^n} \phi_{j+1}^{\text{L}}, \\ \partial_{\phi_{j+1}^n} F_{j+1/2}^n &= \partial_{\phi_j^{\text{R}}} F_{j+1/2}^n \partial_{\phi_{j+1}^n} \phi_j^{\text{R}} + \partial_{\phi_{j+1}^{\text{L}}} F_{j+1/2}^n \partial_{\phi_{j+1}^n} \phi_{j+1}^{\text{L}}, \\ \partial_{\phi_{j+2}^n} F_{j+1/2}^n &= \partial_{\phi_{j+1}^{\text{L}}} F_{j+1/2}^n \partial_{\phi_{j+2}^n} \phi_{j+1}^{\text{L}},\end{aligned}$$

where

$$\partial_{\phi_j^{\text{R}}} F_{j+1/2}^n = \max\{\partial_{\phi_j^{\text{R}}} f(\phi_j^{\text{R}}, r_{j+1/2}), 0\} \text{ and } \partial_{\phi_{j+1}^{\text{L}}} F_{j+1/2}^n = \max\{\partial_{\phi_{j+1}^{\text{L}}} f(\phi_{j+1}^{\text{L}}, r_{j+1/2}), 0\}$$

and where the derivatives of ϕ_j^{R} and ϕ_{j+1}^{L} are evaluated from

$$\begin{aligned}\partial_{\phi_{j-1}^n} \phi_j^{\text{R}} &= \frac{\Delta r}{2} \partial_{\phi_{j-1}^n} s_j^n, \quad \partial_{\phi_j^n} \phi_j^{\text{R}} = 1 + \frac{\Delta r}{2} \partial_{\phi_j^n} s_j^n, \quad \partial_{\phi_{j+1}^n} \phi_j^{\text{R}} = \frac{\Delta r}{2} \partial_{\phi_{j+1}^n} s_j^n, \\ \partial_{\phi_{j-1}^n} \phi_j^{\text{L}} &= -\frac{\Delta r}{2} \partial_{\phi_{j-1}^n} s_j^n, \quad \partial_{\phi_j^n} \phi_j^{\text{L}} = 1 - \frac{\Delta r}{2} \partial_{\phi_j^n} s_j^n, \quad \partial_{\phi_{j+1}^n} \phi_j^{\text{L}} = -\frac{\Delta r}{2} \partial_{\phi_{j+1}^n} s_j^n,\end{aligned}$$

which are a consequence of (2.27). The slope derivatives are in turn given by

$$\partial_{\phi_{j-1}^n} s_j^n = -\frac{\theta}{\Delta r} \chi_j^- - \frac{1}{2\Delta r} \chi_j^0, \quad \partial_{\phi_j^n} s_j^n = \frac{\theta}{\Delta r} (\chi_j^- - \chi_j^+), \quad \partial_{\phi_{j+1}^n} s_j^n = \frac{1}{2\Delta r} \chi_j^0 + \frac{\theta}{\Delta r} \chi_j^+.$$

2.7 Numerical examples

In these numerical examples, we utilize the parameters introduced in Section 2.1, and consider two kinds of observation data: a profile of concentration at $t = T$ as a function of r , and a solution profile at the fixed position $r = R$ as a function of time. Moreover, we consider the case of a rotating tube ($\sigma = 0$) only. These observations are generated by a numerical simulation of the direct problem with the explicit second-order Engquist-Osher-scheme and with the discretization parameters $J = 200$ and N such that the following CFL condition (cf. [11]) holds:

$$\frac{R\omega^2}{g} \max_{\phi \in [0, \phi_{\max}]} |f'_{\text{bk}}(\phi)| \frac{\Delta t}{\Delta r} + 2 \max_{\phi \in [0, \phi_{\max}]} a(\phi) \frac{\Delta t}{(\Delta r)^2} < 1. \quad (2.27)$$

The stability requirement imposed by (2.27) to the explicit scheme implies that we need extremely small values of $\Delta t (\approx (\Delta x)^2)$, which considerably increases computational time. This disadvantage is removed by considering a fully implicit scheme which is unconditionally stable. Thus, for the identification results presented we consider the following first-order implicit discretization of (2.13):

$$\begin{aligned}\phi_j^{n+1} &= \phi_j^n - \lambda_j(F_{j+1/2}^{n+1}(\mathbf{e}) - F_{j-1/2}^{n+1}(\mathbf{e})) + \mu_j(\mathcal{A}_{j+1/2}^{n+1}(\mathbf{e}) - \mathcal{A}_{j-1/2}^{n+1}(\mathbf{e})), \\ \phi_j^0 &= \phi_j^{\text{init}}, \quad \lambda_0 F_{-1/2}^{n+1}(\mathbf{e}) - \mu_0 \mathcal{A}_{-1/2}^{n+1}(\mathbf{e}) = \lambda_J F_{J+1/2}^{n+1}(\mathbf{e}) - \mu_J \mathcal{A}_{J+1/2}^{n+1}(\mathbf{e}) = 0.\end{aligned}\quad (2.28)$$

Since the scheme converges to an entropy solution only for $\Delta x \rightarrow 0$ and $\Delta t \rightarrow 0$ [8] (even though the ratio does not need to be fixed), we find it convenient to link Δx and Δt by setting $N = J$. Choosing Δt smaller or larger will increase or decrease the temporal accuracy, which is critical since our method is first order in time only. More importantly, however, for our identification problem, J and N are also the numbers of data points of the “observed” spatial and temporal concentration profiles used for parameter identification. To be able to compare the numerical results of both choices of observed data it seemed useful to us to set $J = N$.

The weak formulation $E_\Delta = E_\Delta(\phi_\Delta(\mathbf{e}), p_\Delta; \mathbf{e})$ is given by

$E_\Delta = E_\Delta(\phi_\Delta(\mathbf{e}), p_\Delta; \mathbf{e})$ is given by

$$\begin{aligned}E_\Delta &= \sum_{(j,n) \in Q_\Delta} \left\{ \phi_j^n(p_j^n - p_j^{n+1}) + F_{j+1/2}^n(\lambda_j p_j^n - \lambda_{j+1} p_{j+1}^n) - \mathcal{A}_{j+1/2}^n(\mu_j p_j^n - \mu_{j+1} p_{j+1}^n) \right\} \\ &\quad + \sum_{j=0}^J \left\{ [\phi_j^N + \lambda_j(F_{j+1/2}^N - F_{j-1/2}^N) - \mu_j(\mathcal{A}_{j+1/2}^N - \mathcal{A}_{j-1/2}^N)] p_j^N \right. \\ &\quad \left. - [\phi_j^0 + \lambda_j(F_{j+1/2}^0 - F_{j-1/2}^0) - \mu_j(\mathcal{A}_{j+1/2}^0 - \mathcal{A}_{j-1/2}^0)] p_j^0 \right\}.\end{aligned}\quad (2.29)$$

The adjoint scheme and the gradients which are given below are obtained with the methodology developed in Sections 2.5 and 2.6.

We use the conjugate gradient method in the modified form of Polak and Ribière to minimize the objective function $\mathcal{J}(\phi(\mathbf{e}))$, starting with the initial vector $\mathbf{e} = (5.5, 6.5, 9.5, 0.08)$. In order to solve the linear minimization step with the conjugate gradient algorithm we employ Wolfe’s linear search algorithm as described in [58].

2.7.1 Example 1: Profile of concentration at $t = T$ as observation

In this example we consider the radial profiles at fixed times $\hat{\phi}(r, T)$ with $T \in \{0.1 \text{ s}, 0.3 \text{ s}, 1.2 \text{ s}\}$ as observation data, such that the cost function is given by

$$\mathcal{J}(\phi) = \frac{1}{2} \int_{R_0}^R (\phi(r, T) - \hat{\phi}(r, T))^2 dr.$$

T	J	C	σ_0	k	ϕ_c	L^2 error	rate
0.1	IG	5.500000	6.500000	9.500000	0.080000	1.437×10^{-2}	—
	100	5.500019	6.499972	9.499787	0.103858	3.074×10^{-3}	—
	150	5.499859	6.499976	9.499794	0.106275	2.294×10^{-3}	0.722
	200	5.500190	6.499977	9.499799	0.107336	2.188×10^{-3}	0.164
0.3	IG	5.500000	6.500000	9.500000	0.080000	2.567×10^{-2}	—
	100	5.500066	6.499970	9.499746	0.109268	2.800×10^{-3}	—
	150	5.499963	6.499959	9.499700	0.108991	2.766×10^{-3}	0.030
	200	5.500081	6.499966	9.499729	0.108849	2.766×10^{-3}	0.000
1.2	IG	5.500000	6.500000	9.500000	0.080000	3.933×10^{-2}	—
	100	5.500434	6.499983	9.499801	0.110288	1.427×10^{-3}	—
	150	5.500173	6.499975	9.499766	0.109467	7.190×10^{-4}	1.691
	200	5.500364	6.499981	9.499813	0.109645	6.116×10^{-4}	0.562

Table 2.1: Example 1: Initially guessed (IG) and numerically identified parameters for the observation profiles at $T = 0.1$, $T = 0.3$ and $T = 1.2$ along with the L^2 errors and estimated convergence rates (where applicable).

Let $P^n := (p_0^n, \dots, p_J^n)$, and denote, for the implicit calculation of the adjoint scheme, by \mathbf{A}_n the $J \times J$ tridiagonal matrix with entries

$$\begin{aligned} a_{j,j-1}^n &= \lambda_{j-1} \partial_{\phi_j^n} F_{j-1/2}^n - \mu_{j-1} \partial_{\phi_j^n} \mathcal{A}_{j-1/2}^n, \quad j = 2, \dots, J, \\ a_{j,j}^n &= 1 + \lambda_j (\partial_{\phi_j^n} F_{j+1/2}^n - \partial_{\phi_j^n} F_{j-1/2}^n) - \mu_j (\partial_{\phi_j^n} \mathcal{A}_{j+1/2}^n - \partial_{\phi_j^n} \mathcal{A}_{j-1/2}^n), \\ &\qquad\qquad\qquad j = 1, \dots, J-1, \\ a_{j,j+1}^n &= -\lambda_{j+1} \partial_{\phi_j^n} F_{j+1/2}^n + \mu_{j+1} \partial_{\phi_j^n} \mathcal{A}_{j+1/2}^n, \quad j = 1, \dots, J-1, \\ a_{0,0}^n &= 1 + \lambda_0 \partial_{\phi_0^n} F_{1/2}^n - \mu_0 \partial_{\phi_0^n} \mathcal{A}_{1/2}^n, \\ a_{J,J}^n &= 1 - \lambda_J \partial_{\phi_J^n} F_{J-1/2}^n + \mu_J \partial_{\phi_J^n} \mathcal{A}_{J+1/2}^n. \end{aligned}$$

Then P^0 is the solution of the linear implicit adjoint scheme

$$\mathbf{A}_n P^n = P^{n+1}, \quad \text{for } n = N-1, \dots, 0,$$

with the end condition

$$p_j^N = \frac{\phi_j^N - \hat{\phi}_j^N}{a_{j,j}^N} \quad \text{for } j \in \{0, 1, \dots, J\}.$$

The gradient of the discrete cost function is

$$\nabla_{\mathbf{e}} \mathcal{J}_\Delta(\mathbf{e}) = -\Delta r \nabla_{\mathbf{e}} E_\Delta(\phi_\Delta(\mathbf{e}), p_\Delta; \mathbf{e}),$$

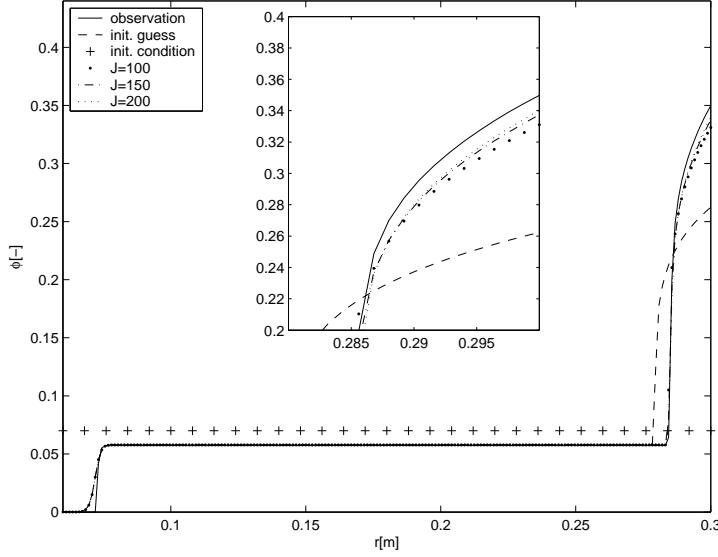


Figure 2.3: Example 1: Comparison of the observed and identified profiles at $T = 0.1$.

where the gradient of the discrete weak form with respect to the parameters is evaluated from

$$\begin{aligned} \nabla_{\mathbf{e}} E_{\Delta} &= \sum_{(j,n) \in Q_{\Delta}} \left\{ \nabla_{\mathbf{e}} F_{j+1/2}^n (\lambda_j p_j^n - \lambda_{j+1} p_{j+1}^n) - \nabla_{\mathbf{e}} \mathcal{A}_{j+1/2}^n (\mu_j p_j^n - \mu_{j+1} p_{j+1}^n) \right\} \\ &\quad + \sum_{j=0}^M \left\{ [\lambda_j (\nabla_{\mathbf{e}} F_{j+1/2}^N - \nabla_{\mathbf{e}} F_{j-1/2}^N) - \mu_j (\nabla_{\mathbf{e}} \mathcal{A}_{j+1/2}^N - \nabla_{\mathbf{e}} \mathcal{A}_{j-1/2}^N)] p_j^N \right. \\ &\quad \left. - [\lambda_j (\nabla_{\mathbf{e}} \mathcal{A}_{j+1/2}^0 - \nabla_{\mathbf{e}} \mathcal{A}_{j-1/2}^0) - \mu_j (\nabla_{\mathbf{e}} \mathcal{A}_{j+1/2}^0 - \nabla_{\mathbf{e}} \mathcal{A}_{j-1/2}^0)] p_j^0 \right\}. \end{aligned}$$

The identified parameters are shown in Table 2.1 and the profiles for the three different observation times are presented in Figures 2.3, 2.4 and 2.5. These figures present the results with several step sizes of resolution and thus show the convergence of the numerical identification scheme when the accuracy increases.

2.7.2 Example 2: Profile of concentration at $z = \bar{R} \in [R_0, R]$ as observation

In the second example, a profile $\phi(\bar{R}, t)$ with $t \in [0, T]$ and $\bar{R} = R = 0.3$ m is considered, leading to the cost function

$$\mathcal{J}(\phi) = \frac{1}{2} \int_0^T (\phi(R, t) - \hat{\phi}(R, t))^2 dt. \quad (2.30)$$

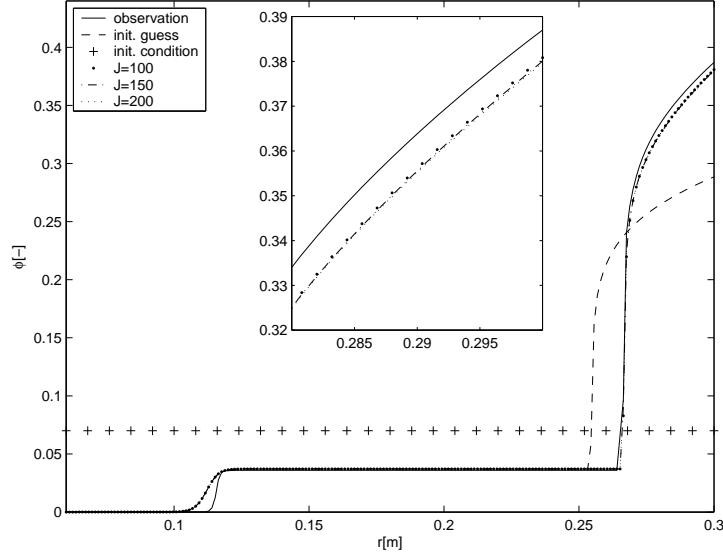


Figure 2.4: Example 1: Comparison of the observed and identified profiles at $T = 0.3$.

J	C	σ_0	k	ϕ_c	L^2 error	rate
IG	5.500000	6.500000	9.500000	0.080000	1.138×10^{-1}	—
100	5.500259	6.499978	9.499804	0.112781	6.127×10^{-3}	—
150	5.500328	6.499982	9.499830	0.111601	7.331×10^{-3}	-0.443
200	5.500368	6.499986	9.499840	0.111226	3.553×10^{-3}	2.518

Table 2.2: Example 2: Initially guessed (IG) and numerically identified parameters for the observation profile given at $\hat{R} = R = 0.3$ along with the L^2 errors and estimated convergence rates (where applicable).

In this case, the adjoint scheme is given by

$$\mathbf{A}_n P^n = P^{n+1} + \mathbf{c} \quad \text{for } n = N-1, \dots, 0,$$

where $\mathbf{c} = (0, \dots, 0, \phi_J^n - \hat{\phi}_J^n)^T$ with the end condition $p_J^N = 0$. The gradient is calculated from

$$\nabla_{\mathbf{e}} \mathcal{J}_{\Delta}(\mathbf{e}) = -\Delta t \nabla_{\mathbf{e}} E_{\Delta}(\phi_{\Delta}(\mathbf{e}), p_{\Delta}; \mathbf{e}),$$

where $\nabla_{\mathbf{e}} E_{\Delta}$ is given by (2.30). The numerical results are shown in Table 2.2 and the profiles are given in Figure 2.6.

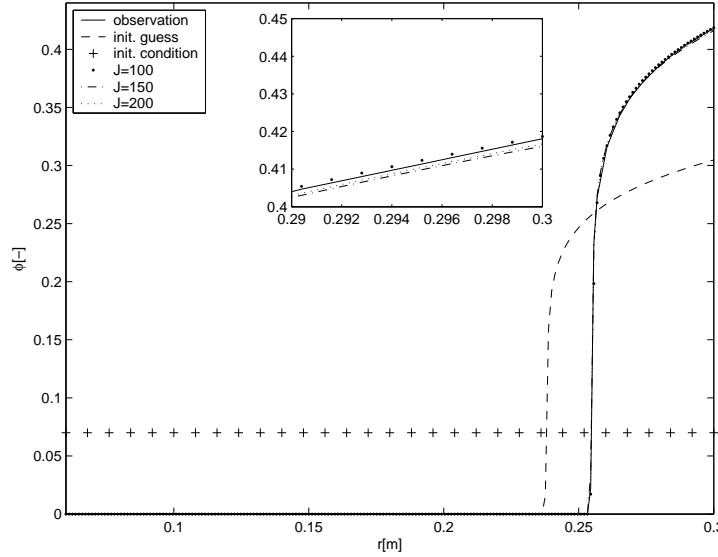


Figure 2.5: Example 1: Comparison of the observed and identified profiles at $T = 1.2$.

J	C	σ_0	k	ϕ_c
100	5.500140	6.499936	9.499538	0.109727
200	5.500043	6.499962	9.499705	0.110450

Table 2.3: Example 2: The identified parameters for the observation profile at $\hat{R} = 0.286$.

2.7.3 Example 3: Profile of concentration with $\sigma = 1$ at $r = \bar{R} \in [R_0, R]$ as observation

In this example, we consider a cylindrical centrifuge ($\sigma = 1$) and assume that a profile $\phi(\bar{R}, t)$ with $t \in [0, 1.2]$ and $\bar{R} = 0.286$ m is observed, that is, concentrations are measured at a fixed (radial) location as a function of time. This leads to the cost function of the form (2.30).

In this case, the adjoint scheme is given by $\mathbf{A}_n P^n = P^{n+1} + \mathbf{c}$ for $n = N-1, \dots, 0$ with the end condition $p_j^N = 0$, and where the column vector $\mathbf{c} = (c_1, \dots, c_J)^T$ is given by

$$c_j = \begin{cases} \phi_j^n - \hat{\phi}_j^n & \text{if } \bar{R} \in [r_{j-1/2}, r_{j+1/2}], \\ 0 & \text{otherwise,} \end{cases} \quad j = 1, \dots, J.$$

The gradient is calculated from $\nabla_{\mathbf{e}} \mathcal{J}_{\Delta}(\mathbf{e}) = -\Delta t \nabla_{\mathbf{e}} E_{\Delta} \phi_{\Delta}(\mathbf{e}), p_{\Delta}; \mathbf{e}$, where $\nabla_{\mathbf{e}} E_{\Delta}$ is

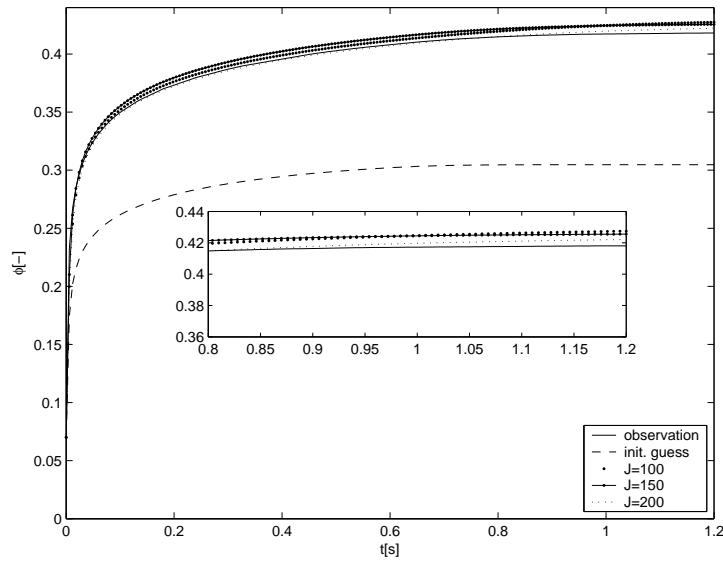


Figure 2.6: Example 2: Comparison of the observed and identified profiles at the boundary $r = R = 0.3$ with temporal resolution $\Delta t = 0.006$.

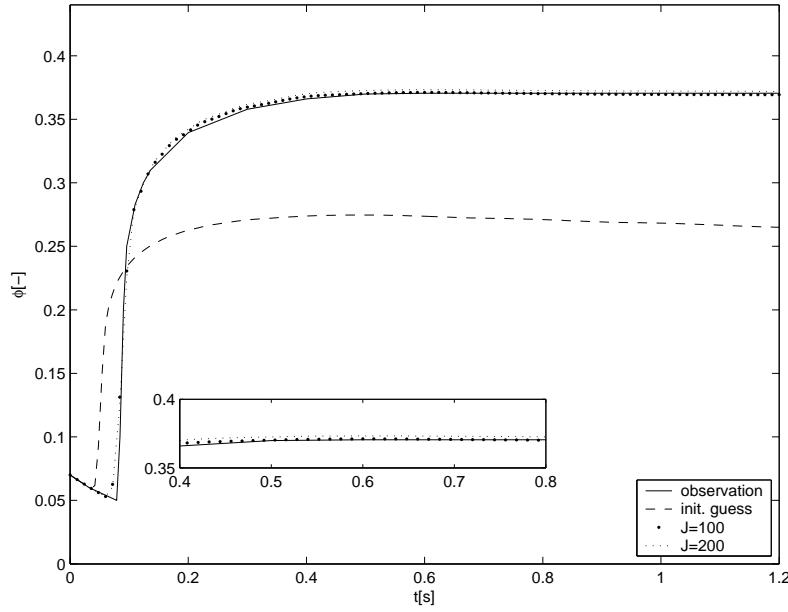


Figure 2.7: Example 3: Comparison of the observed and identified profiles at $r = \bar{R} = 0.286$ with temporal resolution $\Delta t = 0.012$.

$\hat{\phi}$	\mathcal{J}	C	σ_0	k	ϕ_c
ϕ_1	1.0667784e-06	5.494702	6.499970	9.500670	0.104333
$1.01\phi_1$	1.0992580e-06	5.499093	6.499970	9.499876	0.104148

Table 2.4: Example 3: The identified parameters from analytic observation and noisy data.

given by (2.30).

The numerically identified parameters, are shown in Table 2.3 and the profiles are given in Figure 2.6.

2.7.4 Example 4: Analytic data as observation

The test performed before does not consider measurement errors. As the inverse problem is ill-posed, small perturbations in the observation data can have large effects on the considered solutions. In order to test the sensitivity of the method with respect to perturbations of the observation data we consider the analytic data $\hat{\phi}(r, t) = \phi_1(r, t)$ defined by

$$\phi_1(r, t) = \frac{\phi_{\max}}{2} \left[1 + \cos(t) \sin\left(\frac{4\pi}{R - R_0}(r - R_0) - \frac{\pi}{3}\right) \right].$$

This function satisfies the equation

$$\partial_t \phi + \partial_r \left(-\frac{\omega^2}{g} r f_{\text{bk}}(\phi) \right) = \partial_r^2 A(\phi) + s(r, t), \quad (r, t) \in Q_T,$$

with the initial condition $\phi(r, 0) = \phi_1(r, 0)$ and the boundary conditions

$$\left(-\frac{\omega^2}{g} f_{\text{bk}}(\phi) r_b + \partial_r A(\phi) \right)(r_b, t) = c(r_b, t), \quad t > 0, \quad r_b \in \{R_0, R\},$$

where

$$\begin{aligned} s(r, t) &= \partial_t \phi_1 + \partial_r \left(-\frac{\omega^2 r}{g} f_{\text{bk}}(\phi_1) \right) - \partial_r^2 A(\phi_1), \\ c(r, t) &= -\frac{\omega^2 r}{g} f_{\text{bk}}(\phi_1) + \partial_r A(\phi_1). \end{aligned}$$

We consider the perturbation $(1 + \delta)\phi_1$ with $\delta = 0.01$ and we obtain the results shown in Table 2.4. From these results we see that the scheme is stable with respect to noisy data.

2.7.5 Concluding remarks

We now comment on the quality of the parameter identification in these examples, which more precisely present problems of parameter recognition. One would expect that the scheme accurately reproduces the parameters that have actually been used for the simulation. However, we observe in Tables 2.1 and 2.2 that the identification moves away very little from the initially guessed values of C , σ_0 and k , while there are considerable changes in the ϕ_c component, namely from 0.08 to values around 0.11, whereas the value that is actually used is 0.10. Of course, this is no flaw of the method, since Figures 2.3–2.6 and the generally observed decrease of the L^2 error, measured between the numerical solution and the observation data in each case, clearly illustrate that the profiles calculated using the parameters identified by our numerical method approximate well the observed data. Rather, these examples seem to alert to the ill-posedness of the problem, which means that the solution to the identification problem in general fails to be unique, and therefore the method may converge to a solution that is not the intended one. Moreover, it appears that the solution profiles of the direct problem depend very critically on the choice of ϕ_c , and to a lesser degree on the other parameters. This leads to the problem of determining the *sensitivities* of the solution with respect to each of the parameters involved, which remains to be done in future work.

Finally, we recall that the method used for the solution of the direct problem is formally first-order in time and second order in space. This characterization is rigorously valid for smooth solutions, and near discontinuities our scheme falls back to first order. Since we are interested in solutions that include both traveling discontinuities (as the suspension-clear liquid interface) as well as steady-state discontinuities (as the sediment-clear liquid interface), we should expect overall convergence rates to be lower than one. This is consistent with the recorded error histories for the direct problem in [19]. It is interesting to note that in our Example 1, the errors and convergence rates are most favourable for the identification time $T = 1.2$ s, when the solution has attained steady state (see Figure 2.2).

The choice of the scheme has in part been motivated by its amenability to mathematical (convergence) analysis [8, 19, 47]. The basic first-order in space and time scheme analyzed in [8, 47] was extended in [19] to second order in space by MUSCL (*Monotonic Upstream-Centered Scheme for Conservation Laws*) extrapolation, and its implicit variant is used here for the solution of the direct problem as well as for the construction of a scheme of the adjoint problem. Though numerical results of the extrapolated scheme have turned out satisfactory throughout, a rigorous convergence analysis is still lacking. To attain a full second-order scheme, it would be necessary to consider also second-order accuracy in time. This can be done by re-

placing the forward Euler time discretization by a linear multi-step method or a Runge-Kutta type discretization; general references to these methods are [59, 76]. Alternatively, the MUSCL idea can also be applied to attain formal second-order accuracy in time. The resulting so-called MUSCL-Hancock schemes are outlined, for example, in [84]. Since a suitable discretization of the adjoint problem, and therefore of the parameter identification problem, can be obtained in a straightforward manner from an accurate scheme for the direct problem, second-order in both space and time schemes for the parameter identification problem appear to be feasible by modifying the scheme (2.28) by Runge-Kutta time stepping or a MUSCL-Hancock extrapolation.

Chapter 3

Convergence of an upwind scheme for an initial-boundary value problem of a strongly degenerate parabolic equation modelling sedimentation-consolidation processes

We prove the convergence of an explicit monotone finite difference scheme approximating an initial-boundary value problem for a spatially one-dimensional quasi-linear strongly degenerate parabolic equation, which is supplied with two different inhomogeneous flux-type boundary conditions. This problem arises in the modeling of sedimentation-consolidation process. We formulate the definition of entropy solution of the model in the sense of Kružkov and prove the convergence of the numerical scheme to the unique BV entropy solution of the problem, up to satisfaction of one of the boundary conditions.

3.1 Introduction

In this paper we prove convergence of an upwind difference scheme for an initial-boundary value problem for scalar strongly degenerate parabolic equation. Our interest in this problem, and some of the assumptions entering the analysis, arise from a model of batch or continuous sedimentation-consolidation processes of in-

dustrial particulate suspensions [7, 15, 26, 28]. However, under slight modifications, the initial-boundary value problem studied herein can also be understood as a spatially finite model of two-phase flow in porous media [43] or traffic flow with driver reaction [17, 79].

The initial-boundary value problem (IBVP) is given on the rectangular domain $Q_T := I \times \mathcal{T}$, $I := (0, 1)$, $\mathcal{T} := (0, T)$ by

$$\partial_t \phi + \partial_x(f(\phi, t)) = \partial_x^2 A(\phi), \quad (x, t) \in Q_T, \quad (3.1)$$

$$\phi(x, 0) = \phi_0(x), \quad x \in I, \quad (3.2)$$

$$(f(\phi, t) - \partial_x A(\phi))(1, t) = \Psi(t), \quad t \in \mathcal{T}, \quad (3.3)$$

$$(b(\phi) - \partial_x A(\phi))(0, t) = 0, \quad t \in \mathcal{T}, \quad (3.4)$$

where t is time, x is the spatial coordinate, T is the final time, ϕ is the unknown function, $f(\phi, t) = q(t)\phi + b(\phi)$ is the total flux, where $q(t)$ is a control function and $b(\phi)$ is a material-dependent flux density function, and A is the integrated diffusion function, i.e.

$$A(\phi) = \int_0^\phi a(s) ds,$$

where $a(\phi) \geq 0$ is a diffusion function that is allowed to vanish on entire intervals of positive length.

In the framework of the sedimentation-consolidation model, the coordinate x increases vertically, and ϕ is the sought solids volume fraction. The IBVP (3.1)-(3.4) describes a one-dimensional ideal continuous thickener for the settling of industrial suspensions; see [7, 15, 26, 28] and the references cited therein for details. In some of these works, a Dirichlet boundary condition is considered instead of (3.3). However, our flux boundary condition (3.3) describes the same physical process in a simpler way: at the top boundary $x = 1$, a feed source is located, through which fresh suspension is fed into the unit at the feed rate $\Psi(t)$ (see [19] for details). At the bottom ($x = 0$), the total solids flux $f(\phi, t) - \partial_x A(\phi)$ is reduced to its convective part $q(t)\phi(0, t)$, which leads to the boundary condition (3.4). A subcase included here is that of batch settling of a suspension in a closed column, which corresponds to setting $q \equiv 0$.

The basic assumptions on the coefficient functions $q(t)$, $b(\phi)$ and $A(\phi)$ and on the initial and boundary data arising from the sedimentation-consolidation model are the following. The function $q(t)$ is the non-positive volume-averaged velocity of the suspension, which can be controlled externally. The function $b(\phi)$ is a continuous,

piecewise smooth function satisfying

$$b(\phi) = \begin{cases} = 0 & \text{for } \phi \leq 0 \text{ and } \phi \geq \phi_{\max}, \\ < 0 & \text{for } \phi \in (0, \phi_{\max}), \end{cases} \quad (3.5)$$

where $\phi_{\max} \in (0, 1]$ is the maximum concentration value. This function $b(\phi)$ is sometimes called Kynch batch flux density function [28, 73] or hindered settling function, and models the concentration-dependent hindrance of the settling of a solid particle due to the presence of other particles. The degenerating diffusion function $a(\phi)$ models the sediment compressibility, and is assumed to satisfy

$$a(\phi) = \begin{cases} = 0 & \text{for } \phi \leq \phi_c \text{ and } \phi \geq \phi_{\max}, \\ > 0 & \text{for } \phi_c < \phi < \phi_{\max}, \end{cases} \quad 0 \leq \phi_c \leq \phi_{\max}, \quad (3.6)$$

where ϕ_c is a critical concentration or gel point at which the solid particles get into contact with each other. The functions $a(\phi)$ and $b(\phi)$ reflect the specific material properties of the suspension being considered. Since (3.1) reduces to a first-order conservation law on the interval of positive length $[0, \phi_c]$, this equation is called strongly degenerate. Finally, the initial function specified by (3.2) is a piecewise continuous function with $0 \leq \phi_0(x) \leq \phi_{\max}$, while the feed flux $\Psi(t)$ satisfies $\Psi(t) \leq 0$ and $\Psi(t)$ must be larger than the minimum of $f(\cdot; t)$.

To put this paper in the proper perspective, we mention that an analysis of the IBVP (3.1)-(3.4) is given in [15]. Since solutions of equation (3.1) are discontinuous in general, they need to be defined as weak solutions along with an entropy condition to select the physically relevant weak solution. In [15] the existence of *BV* entropy weak solutions to (3.1)-(3.4) in the sense of Kružkov [72] and Vol'pert and Hudjaev [88, 89] is shown via the vanishing viscosity method, while their uniqueness is shown by a technique introduced by Carrillo [31]. On the other hand, Evje and Karlsen [47] show that explicit monotone finite difference schemes, which were first introduced by Crandall and Majda [38] for conservation laws, converge to *BV* entropy solutions for initial-value problems of equation (3.1) (in the slightly simpler case that the flux does not depend on t). These results are extended to implicit schemes in [45], and to several space dimensions in [67]. The extension of these schemes to initial-boundary value problems with flux-type boundary conditions is utilized for the simulation of spatially one-dimensional sedimentation-consolidation processes defined by the IBVP (3.1)-(3.4) in a number of papers including [7, 16, 21, 51, 52]. We also refer to these articles for numerical examples illustrating the scheme and physical model used in this paper.

In [19], the analyses of [15] and [47] are summarized, and a detailed error study of the monotone scheme presented herein (as well as of a MUSCL extrapolation

to formal second-order spatial accuracy) is presented. However, convergence of the scheme is not proved in [19]. The present contribution supplies this convergence analysis, which in part relies on [47, 67] but includes some new proofs required by the presence of boundary conditions. Convergence of monotone schemes towards an entropy solution has also been proved for conservation laws and strongly degenerate convection-diffusion equations with discontinuous flux [22, 18, 21, 65, 66]. Such equations arise, for example, if the sedimentation-consolidation model studied herein is extended to so-called clarifier-thickener units. In addition, the convergence proof of monotone finite difference schemes is important for the justification of the numerical parameter identification scheme for the sedimentation-consolidation model advanced in [37].

The remainder of this paper is organized as follows. In Section 3.2, we recall the definition of an entropy solution of (3.1)-(3.4) and the characterization of the traces for this kind of degenerate parabolic equations, and we describe the numerical schemes. In Section 3.3, we derive BV and L^∞ estimates for the numerical solution of the finite difference scheme, and the Lipschitz continuity estimates of the discrete integrated diffusion function. Finally, in Section 3.4, we show that the schemes satisfy a cell entropy inequality which permits to prove that the discrete solutions converge to a limit that satisfies the entropy condition. We also show that the limit satisfies the initial condition, and one of the boundary conditions. Combining the results of Section 3.3 and Section 3.4, we obtain that the scheme converges to a solution that satisfies all ingredients of the definition of entropy weak solutions except for the boundary condition at $x = 0$. Available numerical results, however, indicate that also this boundary condition is properly approximated, and for the special case $q = 0$, we prove in Section 3.4 that the limit does satisfy the boundary condition at $x = 0$, and therefore is the unique entropy weak solution of (3.1)-(3.4).

3.2 Preliminaries

3.2.1 Definition of entropy solution

We consider a concept of entropy solutions that is similar to the one introduced for “Problem B” in [15].

Definition 3.2.1 *A measurable function $\phi = \phi(x, t)$ is an entropy solution of the initial-boundary value problem (3.1)-(3.4) if the following conditions are satisfied:*

$$(S1) \quad \phi \in L^\infty(Q_T) \cap BV(Q_T).$$

$$(S2) \quad A(\phi) \in C^{1,1/2}(Q_T).$$

(S3) For all test functions $\varphi \in C_0^\infty(Q_T)$, $\varphi \geq 0$, and any $k \in \mathbb{R}$, the following entropy inequality holds:

$$\iint_{Q_T} \left\{ |\phi - k| \partial_t \varphi + \operatorname{sgn}(\phi - k) [(f(\phi, t) - f(k, t)) - \partial_x A(\phi)] \partial_x \varphi \right\} dx dt \geq 0. \quad (3.7)$$

(S4) The boundary condition at $x = 0$ is satisfied in the following sense:

$$\gamma_0(b(\phi) - \partial_x A(\phi)) = 0 \quad \text{for almost all } t \in \mathcal{T}, \quad (3.8)$$

where $\gamma_0 v$ denotes the trace of v with respect to $x \downarrow 0$.

(S5) The boundary condition at $x = 1$ is satisfied in the following sense:

$$\gamma_1(f(\phi, t) - \partial_x A(\phi)) = \Psi(t) \quad \text{for almost all } t \in \mathcal{T}, \quad (3.9)$$

where $\gamma_1 v$ denotes the trace of v with respect to $x \uparrow 1$.

(S6) The initial condition is satisfied in the following sense:

$$\lim_{t \rightarrow 0} \phi(x, t) = \phi_0(x) \quad \text{for almost all } x \in I. \quad (3.10)$$

The traces $\gamma_0 \phi = (\gamma \phi)(0, t)$ and $\gamma_1 \phi = (\gamma \phi)(1, t)$ are well defined by the result given in [91].

Note that in [15], it was first assumed that $\partial_x A(u) \in L^2(Q_T)$, and then $A(u) \in C^{1,1/2}(Q_T)$ was proved as an additional regularity property. Here, we impose this regularity property a priori. The existence, uniqueness and stability of the entropy solutions for the initial-boundary value problem (3.1)-(3.4) was proved in [15] under the following hypotheses:

- (H1) The function b is continuous, piecewise differentiable with $\|b'\|_\infty \leq \infty$, and satisfies (3.5).
- (H2) The function q is locally Lipschitz continuous such that $q(t) \leq 0$ for almost all $t \in \mathcal{T}$, $TV_{\mathcal{T}}(q) < \infty$ and $TV_{\mathcal{T}}(q') < \infty$.
- (H3) The function a is bounded, piecewise continuous with $\operatorname{supp} a \subset \operatorname{supp} b$, and satisfies (3.6).
- (H4) As a consequence of (H3), A is a monotonically non-increasing Lipschitz continuous function. In particular, $\operatorname{sgn}(k_1 - k_2)(A(k_1) - A(k_2)) = |A(k_1) - A(k_2)|$ for all $k_1, k_2 \in \mathbb{R}$.

(H5) We assume that $\text{TV}_{\mathcal{T}}(\Psi) < \infty$ and

$$\min_{\phi \in [0, \phi_{\max}]} f(\phi, t) \leq \Psi(t) \leq 0, \quad \Psi(t) \geq f(\phi_{\max}, t), \quad t \in \overline{\mathcal{T}}.$$

(H6') The initial datum ϕ_0 belongs to the set

$$\begin{aligned} \mathcal{B}' := \{&\phi \in BV(I) : \phi(x) \in [0, \phi_{\max}] \forall x \in \overline{I} \\ &\wedge \exists M_0 > 0 : \text{TV}_I(\partial_x A_\varepsilon(\phi)) < M_0 \text{ uniformly in } \varepsilon\}, \end{aligned}$$

where

$$A_\varepsilon(\phi) := \int_0^\phi a_\varepsilon(s) ds = \int_0^\phi ((a + \varepsilon) * \omega_\varepsilon)(s) ds, \quad \varepsilon > 0,$$

where ω_ε is a standard C^∞ mollifier with support in $(-\varepsilon, \varepsilon)$. This condition is in particular satisfied if ϕ_0 is a constant.

The hypotheses (H1)–(H5) and (H6') are the mathematical statements of the modelling assumptions stated in Section 3.1 in order to get well-posedness of the IVP (3.1)–(3.4). Condition (H6'), stated in [15], refers to the viscous approximation of (3.1)–(3.4), and is necessary to make a uniform estimate of the spatial variation of the time derivative of the solutions of the regularized, strictly parabolic problem possible. For the analysis of the difference schemes, we need to replace (H6') by a condition that involves the spatial discretization. To this end, let $J \in \mathbb{N}$ denote the number of space steps, $\Delta x := 1/J$, $x_{-1/2} := 0$, $x_{J+1/2} := 1$, $x_{j+1/2} := (j + 1/2)\Delta x$ for $j = 0, \dots, J - 1$, $I_j := [x_{j-1/2}, x_{j+1/2})$, and

$$\phi_j^0 := \frac{1}{\Delta x} \int_{I_j} \phi_0(x) dx, \quad j = 0, \dots, J. \quad (3.11)$$

The new condition can be stated as follows.

(H6) The initial datum belongs to the set

$$\begin{aligned} \mathcal{B} := \Big\{&\phi_0 \in BV(I) : \phi_0(x) \in [0, \phi_{\max}] \forall x \in \overline{I} \wedge \exists M_0 > 0 : \\ &\sum_{j=1}^{J-1} |A(\phi_{j+1}^0) - 2A(\phi_j^0) + A(\phi_{j-1}^0)| \leq M_0 \Delta x \text{ uniformly in } \Delta x\Big\}. \end{aligned}$$

This condition is satisfied, for example, if the initial datum is constant. Moreover, in [15] the uniform bounds for the spatial total variation were derived under the condition that either $\Psi = 0$ or there exist positive constants ξ and M_g such that $\xi a(\phi) - (q(t) + b'(\phi)) \geq M_g$. This condition, which is due to Wu [90], is not required in this work. Instead, we impose the following condition.

- (H7) There exists a finite number \mathcal{N} of piecewise continuous functions $\varphi_1(t) \leq \varphi_2(t) \leq \dots \leq \varphi_{\mathcal{N}}(t)$ such that

$$f(u, t) = \Psi(t) \implies u \in \{\varphi_1(t), \dots, \varphi_{\mathcal{N}}(t)\} \quad \forall t \in \mathcal{T},$$

i.e., the functions $\varphi_1(t), \dots, \varphi_{\mathcal{N}}(t)$ denote the time-dependent intersections of the functions $f_1(u; t) := f(u, t)$ and $f_2(u; t) = \Psi(t)$. Moreover, we assume that there exists a constant C_1 such that $\text{TV}_{\mathcal{T}}(\varphi_j) \leq C_1$ for $j = 1, \dots, \mathcal{N}$.

The hypotheses (H7) is satisfied under very general assumptions, for example when b is a polynomial function of u and Ψ is a smoothly varying function of t . Note that (H7) is an assumption on the relationship between the model function $b(u)$ and the control function $\Psi(t)$, and does not involve any evaluation of the unknown solution itself.

3.2.2 The difference schemes

In addition to the notation introduced in the previous section, let $N \in \mathbb{N}$ be the number of time steps, $\Delta t = T/N$, and $I^n := [t^n, t^{n+1})$, where $t^n = n\Delta t$ for $n = 0, \dots, N$. We denote by ϕ_j^n the numerical solution at (x_j, t_n) , assume that the values for $n = 0$ are given by (3.11), and consider a numerical approximation to the solutions of the IBVP (3.1)-(3.4) obtained by a 3-point explicit finite difference scheme. The scheme is defined by an ‘interior’ formula, which is consistent to the governing equation (3.1), and two ‘boundary’ schemes that are produced by inserting the discrete versions of the boundary conditions (3.4) and (3.3) into the interior formula, respectively. The scheme is defined by

$$\phi_0^{n+1} = \phi_0^n - \lambda g_{1/2}^n + \mu d_{1/2}^n + \lambda q(t^n) \phi_0^n, \quad (3.12)$$

$$\phi_j^{n+1} = \phi_j^n - \lambda(g_{j+1/2}^n - g_{j-1/2}^n) + \mu(d_{j+1/2}^n - d_{j-1/2}^n), \quad j = 1, \dots, J-1, \quad (3.13)$$

$$\phi_J^{n+1} = \phi_J^n + \lambda g_{J-1/2}^n - \mu d_{J-1/2}^n - \lambda \Psi(t^n). \quad (3.14)$$

The scheme (3.12)-(3.14) is defined for $n = 0, \dots, N-1$, and we let $\lambda := \Delta t/\Delta x$, $\mu := \Delta t/(\Delta x)^2$, $d_{j+1/2}^n := A(\phi_{j+1}^n) - A(\phi_j^n)$ and $g_{j+1/2}^n := g(\phi_j^n, \phi_{j+1}^n, t^n)$, where $g(u, v, t)$ is the numerical flux. We assume that $g : \mathbb{R}^3 \rightarrow \mathbb{R}$ is a locally Lipschitz continuous function such that:

- (G1) The function g restricted to $[0, \phi_{\max}] \times [0, \phi_{\max}] \times [0, T]$ is nondecreasing in its first argument and nonincreasing in its second.
- (G2) The function g is consistent with f in the standard sense, i.e., $g(\phi, \phi, t) = f(\phi, t)$ for all $(\phi, t) \in [0, \phi_{\max}] \times [0, T]$.

An example of g is the upwind numerical flux considered in [19], where $g(u, v, t) := q(t)v + b^{\text{EO}}(u, v)$, and b^{EO} is the Engquist-Osher [42] numerical flux given by

$$b^{\text{EO}}(u, v) := b(0) + \int_0^u \max\{b'(s), 0\} ds + \int_0^v \min\{b'(s), 0\} ds.$$

3.2.3 Compactness criterion

We carry out the pass to the limit in the approximate solutions ϕ_Δ by appealing to the embedding of $L^\infty(Q_T) \cap BV(Q_T)$ in $L^1(Q_T)$ (see [44]) and using the following well-known L^1 compactness criterion.

Lemma 3.2.1 *Let $\{z_h\}_{h>0}$ be a sequence of functions defined on Q_T . Assume that there exist constants C_0, \dots, C_3 which may depend on T , but not on h , such that*

$$\begin{aligned} \|z_h\|_{L^\infty(Q_T)} &\leq C_0, \quad \|z_h\|_{L^1(Q_T)} \leq C_1, \\ \|z_h(\cdot + y, \cdot) - z_h(\cdot, \cdot)\|_{L^1(Q_T)} &\leq C_2(|y| + h) \quad \text{for all } y \in \mathbb{R} \text{ as } h \downarrow 0, \\ \|z_h(\cdot, \cdot + \tau) - z_h(\cdot, \cdot)\|_{L^1(Q_T)} &\leq C_3(\tau + h) \quad \text{for all } \tau > 0 \text{ as } h \downarrow 0. \end{aligned}$$

Then $\{z_h\}_{h>0}$ is compact in the strong topology of $L^1_{\text{loc}}(Q_T)$. Moreover, any limit point of $\{z_h\}_{h>0}$ belongs $L^1(Q_T) \cap L^\infty(Q_T) \cap C([0, T], L^1(I))$.

3.2.4 Mollifiers and related functions

Let $\omega \in C_0^\infty(\mathbb{R})$ be a function that satisfies $\omega \geq 0$, $\text{supp } \omega \subset (-1, 1)$ and $\|\omega\|_{L^1(\mathbb{R})} = 1$. A standard mollifier with support in $(-h, h)$, $h > 0$, is then defined by $\omega_h(x) := \omega(x/h)/h$. For sufficiently small $h > 0$ we define the functions

$$\rho_h(x) := \int_{-\infty}^x \omega_h(\xi) d\xi, \quad \mu_h(x) := 1 - \rho_h(x - 2h), \quad \nu_h(x) := \rho_h(x - (1 - 2h)),$$

which satisfy $|\mu'_h(x)|, |\nu'_h(x)| \leq C_\mu/h$ and $|\mu''_h(x)|, |\nu''_h(x)| \leq C_\mu/h^2$ with a constant C_μ that is independent of h , and which have the following property (see [15]).

Lemma 3.2.2 *Let $v \in L^1(\mathcal{T}; L^\infty(I))$. If the traces $\gamma_0 v := (\gamma v)(0, t)$ and $\gamma_1 v := (\gamma v)(1, t)$ exist a.e. in \mathcal{T} , then we have for $\varphi \in C^\infty(Q_T)$*

$$\lim_{h \downarrow 0} \iint_{Q_T} \partial_x(\varphi(x, t)(\mu_h(x) + \nu_h(x))) v(x, t) dt dx = \int_0^T (\varphi(1, t) \gamma_1 v - \varphi(0, t) \gamma_0 v) dt.$$

3.3 Estimates on the approximate solutions

The aim of this section is to prove the stability and some regularity properties of the approximate solution values $\{\phi_j^n\} = \{\phi_j^n : 0 \leq j \leq J, 0 \leq n \leq N\}$. We mainly follow the classical work of Crandall and Majda [38] and its generalizations given by Karlsen et al. in [45, 47, 67].

3.3.1 L^∞ stability

Lemma 3.3.1 *If the CFL condition*

$$2\lambda \max_{t \in [0, T]} \|\partial_\phi f(., t)\|_\infty + 2\mu \|a\|_\infty \leq 1 \quad (3.15)$$

holds, then the approximate solution $\{\phi_j^n\}$ of the IBVP (3.1)-(3.4) obtained by the explicit scheme (3.12)-(3.14) is L^∞ stable. More precisely,

$$\phi_j^n \in [0, \phi_{\max}] \quad \text{for } 0 \leq j \leq J, 1 \leq n \leq N. \quad (3.16)$$

Proof. The proof is by induction on n . Let us first introduce the following notation. For $j = 0, \dots, J-1$, we set

$$\begin{aligned} B_{j+1/2}^n &:= \frac{g(\phi_j^n, \phi_{j+1}^n, t^n) - g(\phi_j^n, \phi_j^n, t^n)}{\phi_{j+1}^n - \phi_j^n}, \\ C_{j+1/2}^n &:= \frac{g(\phi_j^n, \phi_{j+1}^n, t^n) - g(\phi_{j+1}^n, \phi_{j+1}^n, t^n)}{\phi_j^n - \phi_{j+1}^n}, \quad D_{j+1/2}^n := \frac{d_{j+1/2}^n}{\phi_{j+1}^n - \phi_j^n} \end{aligned}$$

if $\phi_{j+1}^n \neq \phi_j^n$ and $B_{j+1/2}^n = C_{j+1/2}^n = D_{j+1/2}^n = 0$ otherwise. We can then rewrite the interior explicit scheme as

$$\phi_j^{n+1} = \tilde{B}_j^n \phi_{j-1}^n + \tilde{C}_j^n \phi_j^n + \tilde{D}_j^n \phi_{j+1}^n \quad \text{for } j = 1, \dots, J-1, \quad (3.17)$$

where we define for $j = 1, \dots, J-1$

$$\begin{aligned} \tilde{B}_j^n &:= \lambda C_{j-1/2}^n + \mu D_{j-1/2}^n, \\ \tilde{C}_j^n &:= 1 + \lambda(B_{j+1/2}^n - C_{j-1/2}^n) - \mu(D_{j+1/2}^n + D_{j-1/2}^n), \\ \tilde{D}_j^n &:= -\lambda B_{j+1/2}^n + \mu D_{j+1/2}^n. \end{aligned}$$

The monotonicity property (G1) of the numerical flux, the hypothesis (H4) and the CFL condition (3.15) imply that \tilde{B}_j^n , \tilde{C}_j^n and \tilde{D}_j^n are non-negative. Furthermore, since the coefficients in (3.17) sum up to one for each j , we conclude that if $\phi_j^n \in [0, \phi_{\max}]$ for $j = 0, \dots, J$, then $\phi_j^{n+1} \in [0, \phi_{\max}]$ for $j = 1, \dots, J-1$. It remains

to deal with the boundary schemes. A straightforward computation shows that the boundary scheme for $x = 0$ is equivalent to

$$\begin{aligned}\phi_0^{n+1} &= \left(1 - \lambda \frac{g(\phi_0^n, \phi_1^n, t^n) - g(\phi_0^n, \phi_0^n, t^n)}{\phi_0^n - \phi_1^n} - \mu D_{1/2}^n\right) \phi_0^n \\ &\quad + \left(\lambda \frac{g(\phi_0^n, \phi_1^n, t^n) - g(\phi_0^n, \phi_0^n, t^n)}{\phi_0^n - \phi_1^n} + \mu D_{1/2}^n\right) \phi_1^n - \lambda b(\phi_0^n).\end{aligned}\quad (3.18)$$

In view of the CFL condition, $D_{1/2}^n \geq 0$ and since g is monotonically decreasing in its second argument, we see that the coefficients of ϕ_0^n and ϕ_1^n in (3.18) are nonnegative, and obviously sum up to one. Since $\lambda b(\phi_0^n) \leq 0$, we deduce from (3.18) that if $\phi_0^n \geq 0$ and $\phi_1^n \geq 0$, then $\phi_0^{n+1} \geq 0$. Applying a similar argument to the boundary scheme (3.14), rewritten as

$$\begin{aligned}\phi_J^{n+1} &= \left(\lambda \frac{g(\phi_{J-1}^n, \phi_J^n, t^n) - g(0, \phi_J^n, t^n)}{\phi_{J-1}^n} + \mu D_{J-1/2}^n\right) \phi_{J-1}^n \\ &\quad + \left(1 + \lambda \frac{g(0, \phi_J^n, t^n) - g(0, 0, t^n)}{\phi_J^n} - \mu D_{J-1/2}^n\right) \phi_J^n - \lambda \Psi(t^n),\end{aligned}\quad (3.19)$$

we see that also $\phi_J^{n+1} \geq 0$ if $\phi_{J-1}^n \geq 0$ and $\phi_J^n \geq 0$. Furthermore, to bound the boundary solution values from above, we rewrite the boundary scheme at $x = 0$ as

$$\begin{aligned}\phi_0^{n+1} &= \left(1 - \lambda \frac{g(\phi_0^n, \phi_1^n, t^n) - g(\phi_{\max}, \phi_1^n, t^n)}{\phi_0^n - \phi_{\max}} + \lambda q(t^n) - \mu D_{1/2}^n\right) \phi_0^n + \mu D_{1/2}^n \phi_1^n \\ &\quad + \left(\lambda \frac{g(\phi_0^n, \phi_1^n, t^n) - g(\phi_{\max}, \phi_1^n, t^n)}{\phi_0^n - \phi_{\max}} - \lambda q(t^n)\right) \phi_{\max} + P_0, \\ P_0 &:= \lambda(g(\phi_{\max}, \phi_{\max}, t^n) - g(\phi_{\max}, \phi_1^n, t^n)).\end{aligned}$$

Again, we can argue that the coefficients of ϕ_0^n , ϕ_1^n and ϕ_{\max} are positive and sum up to one. Moreover, $P_0 \leq 0$ since g is monotonically decreasing with respect to its second argument. We conclude that if $\phi_0^n \leq \phi_{\max}$ and $\phi_1^n \leq \phi_{\max}$, then $\phi_0^{n+1} \leq \phi_{\max}$. Finally, we rewrite the boundary scheme at $x = 1$ as

$$\begin{aligned}\phi_J^{n+1} &= \mu D_{J-1/2}^n \phi_{J-1}^n + \left(1 + \lambda \frac{g(\phi_{\max}, \phi_J^n, t^n) - g(\phi_{\max}, \phi_{\max}, t^n)}{\phi_J^n - \phi_{\max}} - \mu D_{J-1/2}^n\right) \phi_J^n \\ &\quad + \left(-\lambda \frac{g(\phi_{\max}, \phi_J^n, t^n) - g(\phi_{\max}, \phi_{\max}, t^n)}{\phi_J^n - \phi_{\max}}\right) \phi_{\max} + P_J, \\ P_J &:= \lambda(g(\phi_{J-1}^n, \phi_J^n, t^n) - g(\phi_{\max}, \phi_J^n, t^n)) + \lambda(f(\phi_{\max}, t^n) - \Psi(t^n)).\end{aligned}$$

The first and the second term in P_J are non-positive due to the monotonicity of g with respect to the first argument and due to assumption (H5), respectively. Thus, we have $P_J \leq 0$, and using the same arguments as for the boundary $x = 0$, we

conclude that $\phi_J^n \leq \phi_{\max}$ if $\phi_{J-1}^n \leq \phi_{\max}$ and $\phi_J^n \leq \phi_{\max}$. We now have shown that $\phi_j^n \in [0, \phi_{\max}]$ for $j = 0, \dots, J$ implies $\phi_j^{n+1} \in [0, \phi_{\max}]$ for $j = 0, \dots, J$. Thus, the explicit scheme (3.12)-(3.14) satisfies the discrete maximum principle. \square

3.3.2 BV estimates

Lemma 3.3.2 *Under the assumptions of Lemma 3.3.1, the numerical solution ϕ_Δ of the IVP (3.1)-(3.4) produced by the explicit scheme (3.12)-(3.14) satisfies the inequality*

$$\sum_{j=0}^{J-1} |\phi_{j+1}^n - \phi_j^n| \leq \sum_{j=0}^{J-1} |\phi_{j+1}^0 - \phi_j^0| + C_4,$$

where C_4 is a constant independent of Δ .

Proof. We define $w_{j+1/2}^n := \phi_{j+1}^n - \phi_j^n$ for $j = 0, \dots, J-1$. From (3.17) we see that

$$w_{j+1/2}^{n+1} = \bar{A}_{j+1/2}^n w_{j-1/2}^n + \bar{B}_{j+1/2}^n w_{j+1/2}^n + \bar{C}_{j+1/2}^n w_{j+3/2}^n, \quad j = 1, \dots, J-2, \quad (3.20)$$

where we define

$$\bar{A}_{j+1/2}^n := \lambda C_{j-1/2}^n + \mu D_{j-1/2}^n, \quad j = 1, \dots, J, \quad (3.21)$$

$$\bar{B}_{j+1/2}^n := 1 - \lambda(C_{j+1/2}^n - B_{j+1/2}^n) - 2\mu D_{j+1/2}^n, \quad j = 0, \dots, J-1, \quad (3.22)$$

$$\bar{C}_{j+1/2}^n := \mu D_{j+3/2}^n - \lambda B_{j+3/2}^n, \quad j = -1, \dots, J-2. \quad (3.23)$$

From the boundary scheme (3.12) we obtain

$$w_{1/2}^{n+1} = \lambda b(\phi_0^n) + \bar{B}_{1/2}^n w_{1/2}^n + \bar{C}_{1/2}^n w_{3/2}^n, \quad (3.24)$$

while the boundary scheme (3.14) leads to

$$w_{J-1/2}^{n+1} = \bar{A}_{J-1/2}^n w_{J-3/2}^n + \bar{B}_{J-1/2}^n w_{J-1/2}^n + \lambda(f(\phi_J, t^n) - \Psi(t^n)). \quad (3.25)$$

For our proof, we need two auxiliary solution values ϕ_{-1}^n and $\phi_{J+1/2}^n$. We define ϕ_{-1}^n to be the largest solution $u = \phi_{-1}^n$ of the equation

$$\Upsilon(u; \phi_0^n) := -b^{\text{EO}}(u, \phi_0^n) - \mu(A(u) - A(\phi_0^n)) = 0. \quad (3.26)$$

Note that ϕ_{-1}^n is well defined, since Υ is a monotonically decreasing, continuous function of u with $\Upsilon(\phi_0^n; \phi_0^n) = -b(\phi_0^n) \geq 0$ and

$$\begin{aligned} \Upsilon(\phi_{\max}; \phi_0^n) &= -b^{\text{EO}}(\phi_{\max}, \phi_0^n) - \mu(A(\phi_{\max}) - A(\phi_0^n)) \\ &\leq -b^{\text{EO}}(\phi_{\max}, \phi_{\max}) - \mu(A(\phi_{\max}) - A(\phi_0^n)) \\ &= \mu(A(\phi_0^n) - A(\phi_{\max})) \leq 0. \end{aligned}$$

This discussion implies in particular that $\phi_0^n \leq \phi_{-1}^n$. Moreover, we define ϕ_{J+1}^n to be the smallest solution $v = \phi_{J+1}^n$ of the equation

$$\Theta(v; \phi_J^n) := \mu(A(v) - A(\phi_J^n)) - \lambda g(\phi_J^n, v, t^n) = -\lambda \Psi(t^n). \quad (3.27)$$

Again, we note that ϕ_{J+1}^n is well defined, since Θ is a monotonically increasing function of v with

$$\begin{aligned} \Theta(0; \phi_J^n) &= -\mu A(\phi_J^n) - \lambda g(\phi_J^n, 0, t) \\ &\leq -\mu A(\phi_J^n) - \lambda g(0, 0, t) = -\mu A(\phi_J^n) \leq 0 \leq -\lambda \Psi(t^n), \\ \Theta(\phi_{\max}; \phi_J^n) &= \mu(A(\phi_{\max}) - A(\phi_J^n)) - \lambda g(\phi_J^n, \phi_{\max}, t^n) \geq -\lambda g(\phi_J^n, \phi_{\max}, t^n) \\ &\geq -\lambda g(\phi_{\max}, \phi_{\max}, t^n) = -\lambda f(\phi_{\max}, t^n) \geq -\lambda \Psi(t^n), \end{aligned}$$

where the last inequality is a consequence of (H5). We now set

$$C_{-1/2}^n := \frac{g(\phi_{-1}^n, \phi_0^n, t^n) - g(\phi_0^n, \phi_0^n, t^n)}{\phi_{-1}^n - \phi_0^n}, \quad D_{-1/2}^n := \frac{A(\phi_{-1}^n) - A(\phi_0^n)}{\phi_{-1}^n - \phi_0^n}$$

if $\phi_{-1}^n \neq \phi_0^n$ and $C_{-1/2}^n = D_{-1/2}^n = 0$ otherwise, and extend the definition (3.21) to $j = 0$ by setting $\bar{A}_{1/2}^n := \mu C_{-1/2}^n + \mu D_{-1/2}^n$. Likewise, we set

$$B_{J+1/2}^n := \frac{g(\phi_J^n, \phi_{J+1}^n, t^n) - g(\phi_J^n, \phi_J^n, t^n)}{\phi_{J+1}^n - \phi_J^n}, \quad D_{J+1/2}^n := \frac{A(\phi_{J+1}^n) - A(\phi_J^n)}{\phi_{J+1}^n - \phi_J^n}$$

if $\phi_{J+1}^n \neq \phi_J^n$ and $B_{J+1/2}^n = D_{J+1/2}^n = 0$ otherwise, and extend (3.23) to $j = J - 1$ by setting $\bar{C}_{J-1/2}^n := \mu D_{J+1/2}^n - \lambda B_{J+1/2}^n$. Defining $w_{-1/2}^n := \phi_0^n - \phi_{-1}^n$ and $w_{J+1/2}^n := \phi_{J+1}^n - \phi_J^n$, we can now rewrite (3.24) and (3.25) as

$$w_{1/2}^{n+1} = \bar{A}_{1/2}^n w_{-1/2}^n + \bar{B}_{1/2}^n w_{1/2}^n + \bar{C}_{1/2}^n w_{3/2}^n, \quad (3.28)$$

$$w_{J-1/2}^{n+1} = \bar{A}_{J-1/2}^n w_{J-3/2}^n + \bar{B}_{J-1/2}^n w_{J-1/2}^n + \bar{C}_{J-1/2}^n w_{J+1/2}^n. \quad (3.29)$$

Thus, (3.20) is valid for $j = 0, \dots, J - 1$. Moreover, we have

$$\bar{A}_{j+1/2}^n, \bar{B}_{j+1/2}^n, \bar{C}_{j+1/2}^n \geq 0, \quad \bar{A}_{j+3/2}^n + \bar{B}_{j+1/2}^n + \bar{C}_{j-1/2}^n = 1, \quad j = 0, \dots, J - 1 \quad (3.30)$$

We now obtain from (3.20), (3.28) and (3.29)

$$\begin{aligned} \sum_{j=0}^{J-1} |w_{j+1/2}^{n+1}| &\leq \bar{A}_{1/2}^n \operatorname{sgn}(w_{-1/2}^n) w_{-1/2}^n + (\bar{B}_{1/2}^n + \bar{A}_{3/2}^n) |w_{1/2}^n| \\ &\quad + (\bar{B}_{J-1/2}^n + \bar{C}_{J-3/2}^n) |w_{J-1/2}^n| + \bar{C}_{J-1/2}^n \operatorname{sgn}(w_{J+1/2}^n) w_{J+1/2}^n \\ &\quad + \sum_{j=1}^{J-2} (\bar{A}_{j+3/2}^n + \bar{B}_{j+1/2}^n + \bar{C}_{j-1/2}^n) |w_{j+1/2}^n|. \end{aligned} \quad (3.31)$$

Taking into account that the boundary condition (3.12) can be written as

$$\phi_0^{n+1} = \phi_0^n + (\mu D_{1/2}^n - \lambda B_{1/2}^n)(\phi_1^n - \phi_0^n) - \lambda b(\phi_0^n) = \phi_0^n + \bar{C}_{-1/2} w_{1/2}^n - \lambda b(\phi_0^n),$$

we get

$$\bar{A}_{1/2} w_{-1/2}^n = \lambda b(\phi_0^n) = \phi_0^n - \phi_0^{n+1} + \bar{C}_{-1/2} w_{1/2}^n. \quad (3.32)$$

Similarly, rewriting the boundary condition (3.14) as

$$\begin{aligned} \phi_J^{n+1} &= \phi_J^n - (\lambda C_{J-1/2}^n + \mu D_{J-1/2}^n)(\phi_J^n - \phi_{J-1}^n) + \lambda(f(\phi_J^n, t^n) - \Psi(t^n)) \\ &= \phi_J^n - \bar{A}_{J+1/2} w_{J-1/2}^n + \lambda(f(\phi_J^n, t^n) - \Psi(t^n)) \end{aligned}$$

yields

$$\bar{C}_{J-1/2}^n w_{J+1/2}^n = \lambda(f(\phi_J^n, t^n) - \Psi(t^n)) = \phi_J^{n+1} - \phi_J^n + \bar{A}_{J+1/2}^n w_{J-1/2}^n. \quad (3.33)$$

Our construction implies $\operatorname{sgn}(w_{-1/2}^n) \in \{-1, 0\}$. If $\operatorname{sgn}(w_{-1/2}^n) = -1$, then (3.32) implies

$$\operatorname{sgn}(w_{-1/2}^n) \bar{A}_{1/2} w_{-1/2}^n = \phi_0^{n+1} - \phi_0^n - \bar{C}_{-1/2} w_{1/2}^n; \quad (3.34)$$

if $\operatorname{sgn}(w_{-1/2}^n) = w_{-1/2}^n = 0$, then (3.32) implies that the right-hand side of (3.34) vanishes. Consequently, inserting (3.32) and (3.34) into (3.31), we obtain

$$\begin{aligned} \sum_{j=0}^{J-1} |w_{j+1/2}^{n+1}| &\leq \phi_0^{n+1} - \phi_0^n + \sum_{j=0}^{J-1} (\bar{A}_{j+3/2}^n + \bar{B}_{j+1/2}^n + \bar{C}_{j-1/2}^n) |w_{j+1/2}^n| \\ &\quad + \operatorname{sgn}(w_{J+1/2}^n) (\phi_J^{n+1} - \phi_J^n). \end{aligned}$$

In light of (3.30), this inequality implies

$$\sum_{j=0}^{J-1} |w_{j+1/2}^{n+1}| \leq \phi_0^{n+1} - \phi_0^n + \sum_{j=0}^{J-1} |w_{j+1/2}^n| + \operatorname{sgn}(w_{J+1/2}^n) (\phi_J^{n+1} - \phi_J^n). \quad (3.35)$$

Let $M \in \{1, \dots, N\}$. Clearly, (3.35) leads to

$$\sum_{j=0}^{J-1} |w_{j+1/2}^M| \leq \phi_0^M - \phi_0^0 + \sum_{j=0}^{J-1} |w_{j+1/2}^n| + S, \quad S := \sum_{n=0}^{M-1} \operatorname{sgn}(w_{J+1/2}^n) (\phi_J^{n+1} - \phi_J^n). \quad (3.36)$$

It remains to estimate S . To this end, we recall that $\bar{C}_{J-1/2}^n \geq 0$, and $\bar{C}_{J-1/2}^n = 0$ if and only if $w_{J+1/2}^n = \phi_{J+1}^n - \phi_J^n = 0$. Thus, we obtain from (3.33)

$$\operatorname{sgn}(w_{J+1/2}^n) = \operatorname{sgn}(f(\phi_J^n, t^n) - \Psi(t^n)),$$

and therefore

$$S = \sum_{n=0}^{M-1} \operatorname{sgn}(f(\phi_J^n, t^n) - \Psi(t^n))(\phi_J^{n+1} - \phi_J^n). \quad (3.37)$$

We cannot simply estimate all summands in (3.37) by $|\phi_J^{n+1} - \phi_J^n|$, since we do not have control of the temporal variation of the discrete solution values at a fixed location. However, we may still estimate S . To this end, let m be the smallest integer such there are numbers $n_0 := -1, n_1, n_2, \dots, n_{m-1}, n_m := M - 1$ with

$$n_0 + 1 = 0 \leq n_1, n_1 + 1 \leq n_2, n_2 + 1 \leq n_3, \dots, n_{m-1} + 1 \leq n_m$$

such that

$$\begin{aligned} \operatorname{sgn}(f(\phi_J^n, t^n) - \Psi(t^n)) &\geq 0 \text{ or } \operatorname{sgn}(f(\phi_J^n, t^n) - \Psi(t^n)) \leq 0 \\ \text{for } n = n_l + 1, \dots, n_{l+1}, l = 0, \dots, m-1. \end{aligned} \quad (3.38)$$

Defining $\tau_l := t^{n_l+1}$ for $l = 0, \dots, m-1$ and $\sigma_l := \operatorname{sgn}(f(\phi_J^{n_l+1}, \tau_l) - \Psi(\tau_l))$ for $l = 0, \dots, m-1$, we can rewrite S as

$$S = \sum_{l=0}^{m-1} \sigma_l \sum_{k=n_l+1}^{n_{l+1}} (\phi_J^{k+1} - \phi_J^k) = \sum_{l=0}^{m-1} \sigma_l (\phi_J^{n_{l+1}+1} - \phi_J^{n_l+1}).$$

Note that $\sigma_l = 0$ cannot occur, and that $\sigma_l \neq \sigma_{l+1}$ for $l = 0, \dots, m-2$, since otherwise in both cases we could merge two neighbouring index segments where (3.38) holds to one, in contradiction to the assumption that m is minimal. Consequently, we have $\sigma_l = (-1)^l$ or $\sigma_l = -(-1)^l$, and thus

$$S = \pm \sum_{l=0}^{m-1} (-1)^l (\phi_J^{n_{l+1}+1} - \phi_J^{n_l+1}), \quad (3.39)$$

where the sign multiplying the sum depends on σ_0 . To bound (3.39), we assume without loss of generality that $\sigma_0 = 1$, and rewrite the sum as

$$\begin{aligned} S &= -\phi_J^{n_0+1} + 2\phi_J^{n_1+1} - 2\phi_J^{n_2+1} \\ &\quad + 2\phi_J^{n_3+1} - \dots + (-1)^{m-1} \cdot 2\phi_J^{n_m+1} - (-1)^{m-1} \phi_J^{n_m+1}. \end{aligned}$$

Now recall that due to (H5), $f(0, t) = 0 \geq \Psi(t)$ and $f(\phi_{\max}, t) \leq \Psi(t)$ for all $t \in \mathcal{T}$. Consequently, the assumption $\sigma_0 = 1$, i.e., $f(\phi_J^{n_0+1}, \tau_0) \geq \Psi(\tau_0)$, implies that there exists an intersection point $\varphi_{k(0)}(\tau_0)$ such that $\phi_J^{n_0+1} \leq \varphi_{k(0)}(\tau_0)$, where we assume that $k(0) \in \{1, \dots, \mathcal{N}\}$ is minimal. Similarly, since $\sigma_1 = -1$, we have $f(\phi_J^{n_1+1}, \tau_1) \leq$

$\Psi(\tau_1)$, such that $\phi_J^{n_1+1} \leq \varphi_{k(1)}(\tau_1)$, where we assume that $k(1) \in \{1, \dots, \mathcal{N}\}$ is maximal, and so on. Proceeding in this way, we obtain

$$\begin{aligned} S &\leq -\varphi_{k(0)}(\tau_0) + 2\varphi_{k(1)}(\tau_1) - 2\varphi_{k(2)}(\tau_2) + 2\varphi_{k(3)}(\tau_3) \\ &\quad - \dots + (-1)^{m-1} \cdot 2\varphi_{k(m-1)}(\tau_{m-1}) - (-1)^{m-1} \varphi_{k(m)}(\tau_m) \\ &= [\varphi_{k(1)}(\tau_1) - \varphi_{k(0)}(\tau_0)] - [\varphi_{k(2)}(\tau_2) - \varphi_{k(1)}(\tau_1)] + [\varphi_{k(3)}(\tau_3) - \varphi_{k(2)}(\tau_2)] \\ &\quad + \dots + (-1)^{m-1} [\varphi_{k(m)}(\tau_m) - \varphi_{k(m-1)}(\tau_{m-1})]. \end{aligned} \quad (3.40)$$

Consider now the first difference in squared brackets. Since $\varphi_i(t) \leq \varphi_j(t)$ for $i \leq j$, $1 \leq i, j \leq \mathcal{N}$, we can estimate this expression from above by

$$\varphi_{k(1)}(\tau_1) - \varphi_{k(0)}(\tau_0) \leq \varphi_{p_1}(\tau_1) - \varphi_{p_1}(\tau_0), \quad p_1 := \max\{k(0), k(1)\}.$$

Similarly, we can estimate the second difference in squared brackets from below by

$$\varphi_{k(2)}(\tau_2) - \varphi_{k(1)}(\tau_1) \geq \varphi_{p_2}(\tau_2) - \varphi_{p_2}(\tau_1), \quad p_2 := \min\{k(1), k(2)\}.$$

Proceeding in analogous way for all differences in (3.40) by setting

$$p_l := \begin{cases} \max\{k(l), k(l-1)\} & \text{if } l \text{ is odd,} \\ \min\{k(l), k(l-1)\} & \text{if } l \text{ is even,} \end{cases}$$

we obtain

$$|S| \leq \sum_{l=1}^m |\varphi_{p_l}(\tau_l) - \varphi_{p_l}(\tau_{l-1})| \leq \sum_{p=1}^{\mathcal{N}} \sum_{l=1}^m |\varphi_p(\tau_l) - \varphi_p(\tau_{l-1})| \leq \sum_{p=1}^{\mathcal{N}} \text{TV}_{[0, t_M]}(\varphi_p). \quad (3.41)$$

By our assumption (H7), the last sum in (3.41) is bounded independently of the discretization. Combining (3.36) and (3.41), we finally obtain

$$\begin{aligned} \sum_{j=0}^{J-1} |\phi_{j+1}^M - \phi_j^M| &\leq \sum_{j=0}^{J-1} |\phi_{j+1}^0 - \phi_j^0| + \phi_0^M - \phi_0^0 + \sum_{p=1}^{\mathcal{N}} \text{TV}_{[0, t_M]}(\varphi_p) \\ &\leq \phi_{\max} + \text{TV}(\phi_{\Delta}^0) + \sum_{p=1}^{\mathcal{N}} \text{TV}_{[0, t_M]}(\varphi_p), \end{aligned} \quad (3.42)$$

which concludes the proof. \square

Lemma 3.3.3 *Under the CFL condition (see (3.15)) the approximate solution $\{\phi_j^n\}$ of the IVP (3.1)-(3.4) obtained by the scheme (3.12)-(3.14), satisfies*

$$\Delta x \sum_{j=0}^J |\phi_j^{n+1} - \phi_j^n| \leq C_5 \Delta t \quad \text{for } n = 0, \dots, N-1, \quad (3.43)$$

where C_5 is a constant independent of Δ .

Proof. We define

$$\hat{B}_j^{n+1/2} := \frac{g(\phi_{j-1}^n, \phi_j^{n+1}, t^{n+1}) - g(\phi_{j-1}^n, \phi_j^n, t^{n+1})}{\phi_j^{n+1} - \phi_j^n}$$

if $\phi_j^{n+1} \neq \phi_j^n$ and $\hat{B}_j^{n+1/2} := 0$ otherwise for $j = 1, \dots, J$,

$$\hat{C}_j^{n+1/2} := \frac{g(\phi_j^{n+1}, \phi_{j+1}^{n+1}, t^{n+1}) - g(\phi_j^n, \phi_{j+1}^{n+1}, t^{n+1})}{\phi_j^{n+1} - \phi_j^n}$$

if $\phi_j^{n+1} \neq \phi_j^n$ and $\hat{C}_j^{n+1/2} := 0$ otherwise for $j = 0, \dots, J-1$, and

$$\hat{D}_j^{n+1/2} := \frac{A(\phi_j^{n+1}) - A(\phi_j^n)}{\phi_j^{n+1} - \phi_j^n}$$

if $\phi_j^{n+1} \neq \phi_j^n$ and $\hat{D}_j^{n+1/2} := 0$ otherwise for $j = 0, \dots, J$. Thus, we can write the numerical scheme (3.12)-(3.14) in terms of $u_j^{n+1/2} := \phi_j^{n+1} - \phi_j^n$ as follows:

$$u_0^{n+3/2} = (1 - \lambda \hat{C}_0^{n+1/2} - \mu \hat{D}_0^{n+1/2} + \lambda q(t^{n+1})) u_0^{n+1/2} \quad (3.44)$$

$$+ (-\lambda \hat{B}_1^{n+1/2} + \mu \hat{D}_1^{n+1/2}) u_1^{n+1/2} + \lambda (q(t^{n+1}) - q(t^n)) (\phi_0^n - \phi_1^n),$$

$$u_j^{n+3/2} = (\lambda \hat{C}_{j-1}^{n+1/2} + \mu \hat{D}_{j-1}^{n+1/2}) u_{j+1}^{n+1/2} \quad (3.45)$$

$$+ (1 - \lambda \hat{C}_j^{n+1/2} + \lambda \hat{B}_j^{n+1/2} - 2\mu \hat{D}_j^{n+1/2}) u_j^{n+1/2}$$

$$+ (-\lambda \hat{B}_{j+1}^{n+1/2} + \mu \hat{D}_{j+1}^{n+1/2}) u_{j+1}^{n+1/2}, \quad j = 1, \dots, J-1,$$

$$u_J^{n+3/2} = (\lambda \hat{C}_{J-1}^{n+1/2} + \mu \hat{D}_{J-1}^{n+1/2}) u_{J-1}^{n+1/2} + (1 + \lambda \hat{B}_J^{n+1/2} - \mu \hat{D}_J^{n+1/2}) u_J^{n+1/2} \quad (3.46)$$

$$+ \lambda (q(t^{n+1}) - q(t^n)) \phi_J^n - \lambda (\Psi(t^{n+1}) - \Psi(t^n)).$$

Using the CFL condition (3.15) and a summation by parts, we obtain the inequality

$$\begin{aligned} \sum_{j=0}^J |u_j^{n+3/2}| &\leqslant \sum_{j=0}^J |u_j^{n+1/2}| + \lambda |q(t^{n+1}) - q(t^n)| |\phi_0^n - \phi_1^n| \\ &\quad + \lambda |q(t^{n+1}) - q(t^n)| \phi_J^n + \lambda |\Psi(t^{n+1}) - \Psi(t^n)|. \end{aligned} \quad (3.47)$$

We can now proceed inductively to obtain

$$\begin{aligned} \sum_{j=0}^J |u_j^{n+3/2}| &\leqslant \sum_{j=0}^J |u_j^{1/2}| \\ &\quad + \lambda \sum_{\nu=0}^n \left\{ |q(t^{\nu+1}) - q(t^\nu)| (|\phi_0^\nu - \phi_1^\nu| + \phi_J^\nu) + |\Psi(t^{\nu+1}) - \Psi(t^\nu)| \right\} \\ &\leqslant \sum_{j=0}^J |u_j^{1/2}| + 2\lambda \phi_{\max} \text{TV}_{[0,t^n]}(q) + \lambda \text{TV}_{[0,t^n]}(\Psi). \end{aligned}$$

To conclude the proof, we need to show that the sum in the last right-hand side is bounded. However,

$$\begin{aligned}
\Delta x \sum_{j=0}^J |u_j^{1/2}| &= \Delta x \sum_{j=0}^J |\phi_j^1 - \phi_j^0| \\
&= \Delta x |\lambda g_{1/2}^0 - \mu d_{1/2}^0 - \lambda q(0) \phi_0^0| \\
&\quad + \Delta x \sum_{j=1}^{J-1} |\lambda(g_{j+1/2}^0 - g_{j-1/2}^0) - \mu(d_{j+1/2}^0 - d_{j-1/2}^0)| \\
&\quad + \Delta x |\lambda g_{J-1/2}^0 - \mu d_{J-1/2}^0 - \lambda \Psi(t^0)| \\
&\leq \Delta t \left| \frac{g(\phi_0^0, \phi_1^0, 0) - g(0, \phi_1^0, 0)}{\phi_0^0} \right| |\phi_0^0| + \Delta t |q(t^0)| |\phi_1^0 - \phi_0^0| \\
&\quad + \Delta t \sum_{j=1}^{J-1} |B_{j+1/2}^0| |\phi_{j+1}^0 - \phi_j^0| + \Delta t \sum_{j=1}^{J-1} |C_{j-1/2}^0| |\phi_j^0 - \phi_{j-1}^0| \\
&\quad + \lambda \sum_{j=1}^{J-1} |A(\phi_{j+1}^0) - 2A(\phi_j^0) + A(\phi_{j-1}^0)| \\
&\quad + \Delta t \left| \frac{g(\phi_{J-1}^0, \phi_J^0, 0) - g(\phi_{J-1}^0, 0, 0)}{\phi_J^0} \right| |\phi_J^0| + |\Psi(t^0)|.
\end{aligned}$$

The last right-hand side is uniformly bounded due the property (G1) and the assumptions (H5) and (H6). \square

3.3.3 Global estimates on $A(\phi_\Delta)$

In this section we present the space and time translate estimates for $A(\phi_\Delta)$, where ϕ_Δ is an approximate solution of the IBVP (3.1)-(3.4). The proof of the following Lemma 3.3.4 is a slight modification of a proof given in [67], which includes the boundary conditions, while the proof of Lemma 3.3.5 does not appeal to boundary conditions, and is therefore a sub-case of the derivation leading to Theorem 4.1 in [67].

Lemma 3.3.4 *If the CFL condition (3.15) holds, then there exists a constant C_6 independent of Δ such that*

$$|A(\phi_j^n) - A(\phi_i^n)| \leq C_6 |j - i| \Delta x \quad \text{for } 0 \leq i, j \leq J$$

for ϕ_j^n obtained by the explicit scheme (3.12)-(3.14).

Proof. We choose $j \in \{0, \dots, J-1\}$ and $i = j+1$; then the quantity to be estimated is $|d_{j+1/2}^n|$. For $j = J-2$, the statement follows immediately from the following inequality, which is a consequence of the boundary scheme (3.14):

$$\frac{1}{\Delta x} |d_{J-1/2}^n| \leq \frac{1}{\lambda} |\phi_J^{n+1} - \phi_J^n| + |g_{J-1/2}^n| + \|\Psi\|_\infty.$$

For $j = 0, \dots, J-2$, we may use Lemma 3.3.3 to obtain

$$\begin{aligned} \frac{1}{\Delta x} |d_{j+1/2}^n| - |g_{j+1/2}^n| &\leq \left| g_{j+1/2}^n - \frac{1}{\Delta x} d_{j+1/2}^n \right| \\ &= \left| \sum_{m=1}^j \Delta_+ \left(g_{m-1/2}^n - \frac{1}{\Delta x} d_{m-1/2}^n \right) + g_{1/2}^n - \frac{1}{\Delta x} d_{1/2}^n \right| \\ &\leq \frac{1}{\lambda} \sum_{m=0}^j |\phi_m^{n+1} - \phi_m^n| + \phi_{\max} \|q\|_\infty \leq C_5 + \phi_{\max} \|q\|_\infty, \end{aligned}$$

where we use the standard notation $\Delta_+ V_j^n := V_{j+1}^n - V_j^n$, $\Delta_- V_j^n := V_j^n - V_{j-1}^n$. The statement of Lemma 3.3.4 follows with $C_6 = C_5 + \phi_{\max} \|q\|_\infty + \|g\|_\infty$. The proof for the semi-implicit is analogous. \square

An immediate consequence of Lemma 3.3.4, we obtain that there exists a constant C_7 independent of Δ such that

$$\lambda \sum_{j=0}^{J-1} \sum_{n=0}^{N-1} (A(\phi_{j+1}^n) - A(\phi_j^n))^2 \leq C_7. \quad (3.48)$$

Lemma 3.3.5 *Under the assumptions of Lemma 3.3.4, there exists a constant C_8 independent of Δ such that*

$$|A(\phi_j^n) - A(\phi_j^m)| \leq C_8 \sqrt{|n-m|\Delta t}, \quad 0 \leq n, m \leq N, \quad j = 0, \dots, J. \quad (3.49)$$

The proof of Lemma 3.3.5 is analogous to the proof of Lemma 4.3 in [65], which in turn is an adaptation of a technique introduced in [47], and therefore is omitted.

3.4 Convergence Analysis

In the sequel, we denote by ϕ_Δ (where $\Delta = (\Delta x, \Delta t)$) the interpolant of degree one associated with the data points $\{\phi_j^n\}$, see [45, 47]. Note that ϕ_Δ is continuous everywhere and differentiable almost everywhere. From Lemmas 3.3.1–3.3.3 we deduce that there is a constant $C_9 = C_9(T)$ such that

$$\|\phi_\Delta\|_{L^\infty(Q_T)} + \text{TV}_{Q_T}(\phi_\Delta) \leq C_9, \quad (3.50)$$

while Lemmas 3.3.4 and 3.3.5 imply that there is a constant C_{10} such that

$$\begin{aligned} & |A(\phi_\Delta(y, \tau)) - A(\phi_\Delta(x, t))| \\ & \leq C_{10} \left(|x - y| + \sqrt{|t - \tau|} + \Delta x + \sqrt{\Delta t} \right) \quad \forall (x, t), (y, \tau) \in Q_T. \end{aligned} \quad (3.51)$$

Consequently, in view of the embedding of $L^\infty(Q_T) \cap BV(Q_T)$ in $L^1(Q_T)$ [44], there exists a sequence $\{\Delta_i\}_{i \in \mathbb{N}}$ with $\Delta_i \rightarrow 0$ for $i \rightarrow \infty$ and a function $\phi \in L^\infty(Q_T) \cap BV(Q_T)$ such that $\phi_\Delta \rightarrow \phi$ a.e. on Q_T . Furthermore, the Arzelà-Ascoli theorem implies that $A(\phi_\Delta) \rightarrow A(\phi)$ uniformly on Q_T , and we have that $A(\phi) \in C^{1,1/2}(Q_T)$. It remains to show that ϕ satisfies the entropy condition (S3), and we need to examine the boundary conditions.

3.4.1 The discrete entropy inequality

Lemma 3.4.1 *For the numerical approximation ϕ_Δ of the IBVP (3.1)-(3.4) obtained by the explicit scheme (3.12)-(3.14), under the CFL condition (3.15) and for all $k \in \mathbb{R}$, the discrete entropy inequalities*

$$|\phi_0^{n+1} - k| - |\phi_0^n - k| + \lambda G_{1/2}^n - \mu \mathcal{A}_{1/2}^n - \lambda q(t^n)|\phi_0^n - k| \leq 0, \quad (3.52)$$

$$\begin{aligned} & |\phi_j^{n+1} - k| - |\phi_j^n - k| + \lambda(G_{j+1/2}^n - G_{j-1/2}^n) \\ & \quad - \mu(\mathcal{A}_{j+1/2}^n - \mathcal{A}_{j-1/2}^n) \leq 0, \quad j = 1, \dots, J-1, \end{aligned} \quad (3.53)$$

$$|\phi_J^{n+1} - k| - |\phi_J^n - k| - \lambda G_{J-1/2}^n + \mu \mathcal{A}_{J-1/2}^n \leq 0 \quad (3.54)$$

hold, with the discrete entropy flux and diffusion functions

$$\begin{aligned} G_{j+1/2}^n &:= G_{j+1/2}^n(\phi_j^n, \phi_{j+1}^n, t^n, k) \\ &:= g(\phi_j^n \top k, \phi_{j+1}^n \top k, t^n) - g(\phi_j^n \perp k, \phi_{j+1}^n \perp k, t^n), \\ \mathcal{A}_{j+1/2}^n &:= \mathcal{A}_{j+1/2}^n(\phi_j^n, \phi_{j+1}^n, k) \\ &:= A(\phi_{j+1}^n \top k) - A(\phi_j^n \top k) - A(\phi_{j+1}^n \perp k) + A(\phi_j^n \perp k), \end{aligned}$$

where we use the standard notation $a \top b := \max\{a, b\}$ and $a \perp b := \min\{a, b\}$.

Proof. We begin by introducing the functions $\mathcal{H}_0, \dots, \mathcal{H}_J$ such that the explicit scheme (3.12)-(3.14) can be written in the form $\phi_j^n = \mathcal{H}_j(\phi_0^n, \phi_1^n, \dots, \phi_J^n, t^n)$ for $j = 0, \dots, J$. In the sequel, we write out only the relevant arguments of $\mathcal{H}_0, \dots, \mathcal{H}_J$. Thus, the functions $\mathcal{H}_0(\phi_0^n, \phi_1^n, t^n)$, $\mathcal{H}_j(\phi_{j-1}^n, \phi_j^n, \phi_{j+1}^n, t^n)$ for $j = 1, \dots, J-1$ and $\mathcal{H}_J(\phi_{J-1}^n, \phi_J^n, t^n)$ are defined by the right-hand sides of (3.12), (3.13), and (3.14), respectively. Clearly, these functions satisfy

$$\begin{aligned} \mathcal{H}_0(k, k, t^n) &= (1 + \lambda q(t^n))k, \quad \mathcal{H}_j(k, k, k, t^n) = 0, \quad j = 1, \dots, J, \\ \mathcal{H}_J(k, k, t^n) &= k - \lambda \Psi(t^n). \end{aligned}$$

From the CFL condition (3.15), (G1) and (H4) we can see that the functions $\mathcal{H}_0, \dots, \mathcal{H}_J$ are increasing functions of their ϕ -arguments. This implies

$$\begin{aligned}\phi_0^{n+1} &\leq \mathcal{H}_0(\phi_0^n \top k, \phi_1^n \top k, t^n), \\ \phi_j^{n+1} &\leq \mathcal{H}_j(\phi_{j-1}^n \top k, \phi_j^n \top k, \phi_{j+1}^n \top k, t^n), \quad j = 1, \dots, J, \\ \phi_J^{n+1} &\leq \mathcal{H}_J(\phi_{J-1}^n \top k, \phi_J^n \top k, t^n), \\ k &\leq \mathcal{H}_0(\phi_0^n \top k, \phi_1^n \top k, t^n) - \lambda q(t^n)k, \\ k &\leq \mathcal{H}_j(\phi_{j-1}^n \top k, \phi_j^n \top k, \phi_{j+1}^n \top k, t^n), \quad j = 1, \dots, J, \\ k &\leq \mathcal{H}_J(\phi_{J-1}^n \top k, \phi_J^n \top k, t^n) + \lambda \Psi(t^n).\end{aligned}$$

These inequalities imply

$$\phi_0^{n+1} \top k \leq \mathcal{H}_0(\phi_0^n \top k, \phi_1^n \top k, t^n) - \lambda q(t^n)k, \quad (3.55)$$

$$\phi_j^{n+1} \top k \leq \mathcal{H}_j(\phi_{j-1}^n \top k, \phi_j^n \top k, \phi_{j+1}^n \top k, t^n), \quad j = 1, \dots, J, \quad (3.56)$$

$$\phi_J^{n+1} \top k \leq \mathcal{H}_J(\phi_{J-1}^n \top k, \phi_J^n \top k, t^n) + \lambda \Psi(t^n). \quad (3.57)$$

We proceed analogously to the proof of

$$\phi_0^{n+1} \perp k \geq \mathcal{H}_0(\phi_0^n \perp k, \phi_1^n \perp k, t^n) - \lambda q(t^n)k, \quad (3.58)$$

$$\phi_j^{n+1} \perp k \geq \mathcal{H}_j(\phi_{j-1}^n \perp k, \phi_j^n \perp k, \phi_{j+1}^n \perp k, t^n), \quad j = 1, \dots, J, \quad (3.59)$$

$$\phi_J^{n+1} \perp k \geq \mathcal{H}_J(\phi_{J-1}^n \perp k, \phi_J^n \perp k, t^n) + \lambda \Psi(t^n). \quad (3.60)$$

Subtracting (3.59) from (3.56), (3.58) from (3.55) and (3.60) from (3.57) and recalling that $|\phi_j^{n+1} - k| = \phi_j^{n+1} \top k - \phi_j^{n+1} \perp k$, we get the desired cell entropy inequalities (3.52), (3.53) and (3.54), respectively. \square

3.4.2 Satisfaction of the entropy inequality

We now show that the limit function ϕ satisfies the continuous entropy inequality stated in (S3). More precisely we have the following Lemma.

Lemma 3.4.2 *Let ϕ_Δ the numerical solution of the IBVP (3.1)-(3.4), $\phi \in L^\infty(Q_T) \cap BV(Q_T)$ such that $\phi_\Delta \rightarrow \phi$ in L^1_{loc} . If ϕ_Δ is obtained by the explicit scheme (3.12)-(3.14) under the CFL condition (3.15) then ϕ satisfies the entropy inequality (S3).*

Proof. Let ϕ_Δ denote the linearly interpolated approximation ϕ_Δ of the IBVP (3.1)-(3.4) obtained by (3.12)-(3.14) and let $\varphi \in C_0^\infty(Q_T)$ be a nonnegative function. We multiply the discrete entropy inequality (3.53) by $\varphi(x, t^n)$, integrate over $I_j =$

$[x_{j-1/2}, x_{j+1/2})$, and sum over $(j, n) \in \{0, \dots, J\} \times \{0, \dots, N - 1\}$. This yields an inequality $E_1^\Delta + E_2^\Delta + E_3^\Delta + E_4^\Delta \leq 0$, where

$$\begin{aligned} E_1^\Delta &:= \sum_{n=0}^{N-1} \sum_{j=0}^J (|\phi_j^{n+1} - k| - |\phi_j^n - k|) \int_{I_j} \varphi(x, t^n) dx, \\ E_2^\Delta &:= \sum_{n=0}^{N-1} \sum_{j=1}^{J-1} \lambda (G_{j+1/2}^n - G_{j-1/2}^n) \int_{I_j} \varphi(x, t^n) dx \\ &\quad + \lambda \sum_{n=0}^{N-1} \left\{ G_{1/2}^n \int_{I_0} \varphi(x, t^n) dx - G_{J-1/2}^n \int_{I_J} \varphi(x, t^n) dx \right\}, \\ E_3^\Delta &:= - \sum_{n=0}^{N-1} \sum_{j=1}^{J-1} \mu (\mathcal{A}_{j+1/2}^n - \mathcal{A}_{j-1/2}^n) \int_{I_j} \varphi(x, t^n) dx \\ &\quad + \mu \sum_{n=0}^{N-1} \left\{ \mathcal{A}_{J-1/2}^n \int_{I_J} \varphi(x, t^n) dx - \mathcal{A}_{1/2}^n \int_{I_0} \varphi(x, t^n) dx \right\}, \\ E_4^\Delta &:= - \sum_{n=0}^{N-1} \lambda q(t^n) |\phi_0^n - k| \int_{I_0} \varphi(x, t^n) dx. \end{aligned}$$

Note that $E_4^\Delta \geq 0$ since $q(t) \leq 0$, which implies $E_1^\Delta + E_2^\Delta + E_3^\Delta \leq 0$. A summation by parts and using that φ has compact support yield

$$\begin{aligned} E_1^\Delta &= \sum_{j=0}^J \left\{ |\phi_j^N - k| \int_{I_j} \varphi(x, T) dx - |\phi_j^0 - k| \int_{I_j} \varphi(x, 0) dx \right. \\ &\quad \left. - \Delta t \sum_{n=0}^{N-2} |\phi_j^{n+1} - k| \int_{I_j} \frac{\varphi(x, t^{n+1}) - \varphi(x, t^n)}{\Delta t} dx \right\} \\ &= - \Delta t \sum_{j=0}^J \sum_{n=0}^{N-2} |\phi_j^{n+1} - k| \int_{I_j} \frac{\varphi(x, t^{n+1}) - \varphi(x, t^n)}{\Delta t} dx \end{aligned}$$

Another summation by parts, using the consistency of g with f and Lemma 3.3.2 imply

$$\begin{aligned} E_2^\Delta &= - \Delta t \sum_{n=0}^{N-1} \sum_{j=0}^{J-1} G_{j+1/2}^n \int_{I_j} \frac{\Delta_+ \varphi(x, t^n)}{\Delta x} dx \\ &= - \Delta t \sum_{n=0}^{N-1} \sum_{j=0}^{J-1} \operatorname{sgn}(\phi_j^n - k) (f(\phi_j, t^n) - f(k, t^n)) \int_{I_j} \frac{\Delta_+ \varphi(x, t^n)}{\Delta x} dx + \mathcal{O}(\Delta x), \end{aligned}$$

where $\Delta_- V(x, t^n) := V(x + \Delta x, t^n) - V(x, t^n)$. Finally, we get

$$\begin{aligned} E_3^\Delta &= -\Delta t \sum_{n=0}^{N-1} \sum_{j=1}^{J-1} |A(\phi_j^n) - A(k)| \int_{I_j} \frac{\Delta_+ \Delta_- \varphi(x, t^n)}{\Delta x^2} dx \\ &\quad + \lambda \sum_{n=0}^{N-1} \left\{ |A(\phi_J^n) - A(k)| \int_{I_{J-1}} \frac{\Delta_+ \varphi(x, t^n)}{\Delta x} dx \right. \\ &\quad \left. - |A(\phi_0) - A(k)| \int_{I_0} \frac{\Delta_+ \varphi(x, t^n)}{\Delta x} dx \right\} \\ &= -\Delta t \sum_{n=0}^{N-1} \sum_{j=1}^{J-1} |A(\phi_j^n) - A(k)| \int_{I_j} \frac{\Delta_+ \Delta_- \varphi(x, t^n)}{\Delta x^2} dx + \mathcal{O}(\Delta x). \end{aligned}$$

Since φ is smooth, we now may state the inequality $-E_1^\Delta - E_2^\Delta - E_3^\Delta \geq 0$ as

$$\begin{aligned} \iint_{Q_T} \left\{ |\phi_\Delta - k| \partial_t \varphi + \operatorname{sgn}(\phi_\Delta - k) (f(\phi_\Delta, t) - f(k, t)) + |A(\phi_\Delta) - A(k)| \partial_x^2 \varphi \right\} dx dt \\ \geq -C_{11}(\Delta t + \Delta x). \end{aligned}$$

Taking $\Delta \rightarrow 0$ and doing integration by parts, we see that the limit function ϕ satisfies (3.7). \square

3.4.3 Satisfaction of initial and boundary conditions

Lemma 3.4.3 *Let ϕ_Δ be the linearly interpolated numerical solution of the IBVP (3.1)-(3.4) and $\phi \in L^\infty(Q_T) \cap BV(Q_T)$ such that $\phi_\Delta \rightarrow \phi$ in L^1 . If ϕ_Δ is obtained by the explicit scheme (3.12)-(3.14) under the CFL condition (3.15) then ϕ satisfies the boundary condition (S5) and the initial condition (S6).*

Proof. Multiplying the boundary scheme (3.14) by $\int_{I_J} \varphi(x, t^n) dx$, where the test function φ is given by $\varphi(x, t) := \Phi(t) \nu_h(x)$ with $\Phi \in C_0^\infty(\mathcal{T})$, summing the result over $n = 0, \dots, N-1$, using summation by parts and (3.13) and (3.14), we get

$$\begin{aligned} 0 &= \sum_{n=0}^{N-1} \lambda \Psi(t^n) \Phi(t^n) \int_{I_J} \nu_h(x) dx \\ &\quad + \sum_{n=0}^{N-1} (\phi_J^{n+1} - \phi_J^n - \lambda g_{J-1/2}^n + \mu d_{J-1/2}^n) \Psi(t^n) \int_{I_J} \nu_h(x) dx \\ &= \Delta t \sum_{n=0}^{N-1} \sum_{j=0}^{J-1} \Psi(t^n) \Phi(t^n) \int_{I_j} \frac{\Delta_+ \nu_h(x)}{\Delta x} dx \end{aligned}$$

$$\begin{aligned}
& + \sum_{n=0}^{N-1} \left\{ (\phi_J^{n+1} - \phi_J^n) \int_{I_J} \nu_h(x) dx - (\phi_{J-1}^{n+1} - \phi_{J-1}^n) \int_{I_{J-1}} \nu_h(x) dx \right. \\
& \quad - (\lambda g_{J-1/2}^n - \mu d_{J-1/2}^n) \int_{I_J} \nu_h(x) dx \\
& \quad \left. + (\lambda g_{J-3/2}^n - \mu d_{J-3/2}^n) \int_{I_{J-1}} \nu_h(x) dx \right\} \Phi(t^n) \\
& + \sum_{n=0}^{N-1} \sum_{j=1}^{J-1} \left\{ (\phi_{j+1}^{n+1} - \phi_{j+1}^n) \int_{I_{j+1}} \nu_h(x) dx - (\phi_j^{n+1} - \phi_j^n) \int_{I_j} \nu_h(x) dx \right. \\
& \quad - (\lambda g_{j+1/2}^n - \mu d_{j+1/2}^n) \int_{I_{j+1}} \nu_h(x) dx \\
& \quad \left. + (\lambda g_{j-1/2}^n - \mu d_{j-1/2}^n) \int_{I_j} \nu_h(x) dx \right\} \Phi(t^n) \\
& = \Delta t \sum_{n=0}^{N-1} \sum_{j=0}^{J-1} \Psi(t^n) \Phi(t^n) \int_{I_j} \frac{\Delta_+ \nu_h(x)}{\Delta x} dx + \sum_{n=0}^{N-1} (\phi_J^{n+1} - \phi_J^n) \Phi(t^n) \int_{I_J} \nu_h(x) dx \\
& \quad - \Delta t \sum_{n=0}^{N-1} \left(g_{J-1/2}^n - \frac{d_{J-1/2}^n}{\Delta x} \right) \Phi(t^n) \int_{I_{J-1}} \frac{\Delta_+ \nu_h(x)}{\Delta x} dx \\
& \quad - \sum_{n=0}^{N-1} \sum_{j=0}^{J-2} \left\{ (\lambda g_{j+3/2}^n - \mu d_{j+3/2}^n) \int_{I_{j+1}} \nu_h(x) dx \right. \\
& \quad \left. - (\lambda g_{j+1/2}^n - \mu d_{j+1/2}^n) \int_{I_j} \nu_h(x) dx \right\} \Phi(t^n) \\
& =: \mathcal{S}_1 + \mathcal{S}_2 + \mathcal{S}_3 + \mathcal{S}_4.
\end{aligned}$$

Firstly, we have that

$$\mathcal{S}_1 = \iint_{Q_T} \Psi(t) \Phi(t) \nu'_h(x) dx + \mathcal{O}(\Delta x). \quad (3.61)$$

Assuming that Δx is chosen small enough such that $\nu_h = 0$ on I_0 , we get

$$\mathcal{S}_2 = \left\{ \phi_J^N \Phi(t^{N-1}) - \sum_{n=1}^{N-1} \phi_J^n (\Phi(t^n) - \Phi(t^{n-1})) \right\} \int_{I_J} \nu_h(x) dx, \quad (3.62)$$

which implies

$$|\mathcal{S}_2| \leq \phi_{\max} (\|\Phi\|_\infty + \text{TV}_T(\Phi)) \Delta x = C_{12} \Delta x, \quad (3.63)$$

where the constant C_{12} depends on the choice of Φ , but not on Δ . Moreover, we assume that $\Delta x < h$, which implies $\mathcal{S}_3 = 0$. Next, we get

$$\begin{aligned}\mathcal{S}_4 &= -\sum_{n=0}^{N-1} \sum_{j=0}^{J-2} \left\{ \Delta_+ (\lambda g_{j+1/2}^n - \mu d_{j+1/2}^n) \int_{I_{j+1}} \nu_h(x) dx \right. \\ &\quad \left. + (\lambda g_{j+1/2}^n - \mu d_{j+1/2}^n) \int_{I_j} \Delta_+ \nu_h(x) dx \right\} \Phi(t^n) \\ &= -\sum_{n=0}^{N-1} \sum_{j=1}^{J-1} (\phi_j^{n+1} - \phi_j^n) \Phi(t^n) \int_{I_j} \nu_h(x) dx \\ &\quad - \Delta t \sum_{n=0}^{N-1} \sum_{j=1}^{J-1} \left(g_{j-1/2}^n - \frac{d_{j-1/2}^n}{\Delta x} \right) \Phi(t^n) \int_{I_{j-1}} \frac{\Delta_+ \nu_h(x)}{\Delta x} dx =: \mathcal{S}_4^1 + \mathcal{S}_4^2.\end{aligned}$$

Note that

$$\begin{aligned}\mathcal{S}_4^1 &= -\sum_{j=1}^{J-1} \left(\sum_{n=0}^{N-1} (\phi_j^{n+1} - \phi_j^n) \Phi(t^n) \right) \int_{I_j} \nu_h(x) dx \\ &= -\sum_{j=1}^{J-1} \left(\phi_j^N \Phi(t^{N-1}) + \sum_{n=1}^{N-1} (\Phi(t^n) - \Phi(t^{n-1})) \phi_j^n \right) \int_{I_j} \nu_h(x) dx,\end{aligned}$$

which implies the estimate

$$|\mathcal{S}_4^1| \leq \phi_{\max}(\|\Phi\|_\infty + \text{TV}_T(\Phi)) \int_I \nu_h(x) dx \leq C_{13} h, \quad (3.64)$$

where the constant C_{13} does not depend on Δ or h . Finally, we see that

$$\begin{aligned}\mathcal{S}_4^2 &= -\Delta t \sum_{n=0}^{N-1} \sum_{j=1}^{J-1} \left(g_{j-1/2}^n - \frac{d_{j-1/2}^n}{\Delta x} \right) \Phi(t^n) \int_{I_{j-1}} \nu'_h(x) dx + \mathcal{O}(\Delta x) \\ &= -\Delta t \sum_{n=0}^{N-1} \left\{ \sum_{j=1}^{J-1} g_{j-1/2}^n \Phi(t^n) \int_{I_{j-1}} \nu'_h(x) dx + \frac{A(\phi_0^n)}{\Delta x} \int_{I_0} \nu'_h(x) dx \right. \\ &\quad \left. - \frac{A(\phi_{J-1}^n)}{\Delta x} \int_{I_{J-2}} \nu'_h(x) dx + \Delta x \sum_{j=1}^{J-2} A(\phi_j^n) \int_{I_{j-1}} \nu''_h(x) dx \right\} \Phi(t^n) + \mathcal{O}(\Delta x).\end{aligned}$$

We assume that Δx is sufficiently small, such that $\nu'_h = 0$ on I_0 and I_{J-2} . This assumption implies

$$\mathcal{S}_4^2 = -\Delta t \sum_{n=0}^{N-1} \left\{ \sum_{j=1}^{J-1} g_{j-1/2}^n \int_{I_{j-1}} \nu'_h(x) dx + \sum_{j=1}^{J-2} A(\phi_j^n) \int_{I_{j-1}} \nu''_h(x) dx \right\} \Phi(t^n)$$

$$\begin{aligned}
& + \mathcal{O}(\Delta x) \\
= & -\Delta t \sum_{n=0}^{N-1} \left\{ \sum_{j=1}^{J-1} \int_{I_{j-1}} \left(g_{j-1/2}^n \nu'_h(x) + A(\phi_j^n) \nu''_h(x) \right) \Phi(t^n) dx \right\} \\
& + \mathcal{O}(\Delta x(1 + h^{-2})) \\
= & - \iint_{Q_T} (f(\phi_\Delta, t) \nu'_h(x) + A(\phi_\Delta) \nu''_h(x)) \Phi(t) dx dt + \mathcal{O}(\Delta x(1 + h^{-2})).
\end{aligned}$$

Using an integration by parts, we get

$$\mathcal{S}_4^2 = - \iint_{Q_T} (f(\phi_\Delta, t) - \partial_x A(\phi_\Delta)) \Phi(t) \nu'_h(x) dx dt + \mathcal{O}(\Delta x(1 + h^{-2})). \quad (3.65)$$

Combining (3.61), (3.63), $\mathcal{S}_3 = 0$, (3.64) and (3.65), we obtain

$$\iint_{Q_T} (\Psi(t) - f(\phi_\Delta, t) + \partial_x A(\phi_\Delta)) \Phi(t) \nu'_h(x) dx dt = \mathcal{O}(\Delta x(1 + h^{-2})). \quad (3.66)$$

Letting jointly $\Delta \rightarrow 0$ and $h \rightarrow 0$ in such a way that $\Delta x/h^2 \rightarrow 0$ (for example, by choosing $\Delta x = h^3$), and considering that $f(\phi_\Delta, t) - \partial_x A(\phi_\Delta)$ as well as its spatial total variation is uniformly bounded, we may pass to the limit in (3.66) to obtain

$$\int_0^T \gamma_1 (\Psi(t) - f(\phi_\Delta, t) + \partial_x A(\phi_\Delta)) \Phi(t) dt = 0, \quad (3.67)$$

which implies that boundary condition (S5) is satisfied.

From Lemma 3.3.3 we deduce that

$$\int_I |\phi_\Delta(x, \Delta t) - \phi_\Delta(x, 0)| dx \leq C \Delta t.$$

Taking $\Delta t \rightarrow 0$ we see that the initial condition (S6) is satisfied. \square

For the general case $q \not\equiv 0$, we have not been able to prove that the limit ϕ satisfies the boundary condition (S4). The basic difficulty becomes apparent when one attempts to repeat the proof of Lemma 3.4.3 for the boundary $x = 0$, starting from the discrete boundary scheme (3.12), multiplying that scheme by $\int_{I_0} \varphi(x, t^n) dx$, where $\varphi(x, t) = \Phi(t) \mu_h(x)$, and summing the result over $n = 0, \dots, N-1$. This procedure will lead to the necessity to estimate terms like

$$\sum_{n=1}^N \sum_{j=1}^J q(t^n) \Phi(t^n) (\phi_{j+1}^n - \phi_j^n) \int_{I_j} \mu_h(x) dx.$$

The latter would be possible if we had more accurate information of the behaviour of the discrete solution in an $\mathcal{O}(h)$ strip $[0, C \cdot h] \times \mathcal{T}$. By our present analysis,

we cannot exclude that a strongly oscillatory (in time) boundary layer forms near $x = 0$, although our numerical experiments [19, 51, 52] illustrate that this does not happen. A related question is whether the functions $\phi_\Delta(0, \cdot)$ converge to a meaningful function of t as $\Delta \rightarrow 0$.

Let us point out, however, that in the case $q \equiv 0$, which corresponds to the practically important case of batch settling in a closed column, we can straightforwardly prove that the boundary condition (S4) is satisfied by repeating the proof of Lemma 3.4.3 under the modifications given above.

Finally, the series of previously established lemmas form a proof of the following theorem.

Theorem 3.4.1 *Let us assume that (H1)–(H7) hold. Under the CFL condition (3.15) the interpolated approximate solution ϕ_Δ obtained by the explicit scheme (3.12)–(3.14) converge in the strong topology of $L^1(Q_T)$ for $\Delta \rightarrow 0$ to a function $\phi \in L^1(Q_T) \cap BV(Q_T)$, which has the properties (S1)–(S3), (S5) and (S6) stated in the definition of an entropy weak solution. In the special case $q \equiv 0$, also the boundary condition at $x = 0$, (S4), is satisfied, and the limit function is an entropy weak solution.*

Appendix A

A.1 Weak and discrete weak formulations

To derive the weak formulation of an initial-boundary value problem of the form

$$\partial_t u + \partial_r f(u, r) = \partial_r^2 A(u, r) + g(u, r) \quad (\text{A.1})$$

$$u(r, 0) = u_0(r), \quad (\text{A.2})$$

$$f(u(r, t)) - \partial_r A(u(r, t))|_{r \in \{R_0, R\}} = \Gamma(r, t)|_{r \in \{R_0, R\}}, \quad (\text{A.3})$$

Equation (A.1) is multiplied by a test function p and integrated over $Q_T = (R_0, R) \times (0, T)$ to obtain

$$\iint_{Q_T} (\partial_t u + \partial_r f(u, r) - \partial_r^2 A(u, r)) p \, dt \, dr = \iint_{Q_T} g(u, r) p \, dt \, dr.$$

Integrating by parts in combination with the initial and boundary conditions (A.2)-(A.3), the weak formulation $E(u(\mathbf{e}), p; \mathbf{e}) = 0$ is obtained with

$$\begin{aligned} E(u(\mathbf{e}), p; \mathbf{e}) &= - \iint_{Q_T} (u \partial_t p + f(u, r) \partial_r p + A(u, r) \partial_r^2 p + g(u, r) p) \, dt \, dr \\ &\quad + \int_{R_0}^R u p \Big|_{t=0}^T \, dr + \int_0^T f(u, r) p \Big|_{r=R_0}^R \, dt \\ &\quad - \int_0^T (\partial_r A(u, r) p - A(u, r) \partial_r p) \Big|_{r=R_0}^R \, dt \\ &= - \iint_{Q_T} (u \partial_t p + f(u, r) \partial_r p + A(u, r) \partial_r^2 p + g(u, r) p) \, dt \, dr \\ &\quad + \int_{R_0}^R u p \Big|_{t=0}^T \, dr + \int_0^T (\Gamma(r, t) p + A(u, r) \partial_r p) \Big|_{r=R_0}^R \, dt. \end{aligned}$$

To derive the discrete weak formulation, we start from the scheme written in conservation form as

$$u_j^{n+1} = u_j^n - \lambda_j(F_{j+1/2}^n - F_{j-1/2}^n) + \mu_j(\mathcal{A}_{j+1/2}^n - \mathcal{A}_{j-1/2}^n) + g_j^n \quad (\text{A.4})$$

with the initial condition $u_j^0 = u_0(r_j, 0)$ and the boundary conditions

$$\lambda_j F_{j-1/2}^n - \mu_j \mathcal{A}_{j-1/2}^n = \Gamma_j^n \quad \text{for } j \in \{0, J+1\}.$$

Moreover, in the notation we suppress the dependence of the numerical fluxes on the parameter vector \mathbf{e} . Multiplying (A.4) by p_j^{n+1} and summing the result over $(j, n) \in Q_\Delta$, we obtain

$$\sum_{(j,n) \in Q_\Delta} \left\{ (u_j^{n+1} - u_j^n) + \lambda_j(F_{j+1/2}^n - F_{j-1/2}^n) - \mu_j(\mathcal{A}_{j+1/2}^n - \mathcal{A}_{j-1/2}^n) - g_j^n \right\} p_j^{n+1} = 0.$$

In view of the “summation by parts” identities

$$\begin{aligned} \sum_{n=0}^{N-1} u_j^{n+1} p_j^{n+1} &= \sum_{n=0}^{N-1} u_j^n p_j^n + u_j^N p_j^N - u_j^0 p_j^0, \\ \sum_{j=0}^J \lambda_j F_{j-1/2}^n p_j^{n+1} &= \sum_{j=0}^J \lambda_{j+1} F_{j+1/2}^n p_{j+1}^{n+1} + \lambda_0 F_{-1/2}^n p_0^{n+1} - \lambda_{J+1} F_{J+1/2}^n p_{J+1}^{n+1}, \\ \sum_{j=0}^J \mu_j \mathcal{A}_{j-1/2}^n p_j^{n+1} &= \sum_{j=0}^J \mu_{j+1} \mathcal{A}_{j+1/2}^n p_{j+1}^{n+1} + \mu_0 \mathcal{A}_{-1/2}^n p_0^{n+1} - \mu_{J+1} \mathcal{A}_{J+1/2}^n p_{J+1}^{n+1}, \end{aligned}$$

the discrete weak formulation becomes $E_\Delta(u_\Delta(\mathbf{e}), p_\Delta; \mathbf{e}) = 0$, where

$$\begin{aligned} E_\Delta &:= \sum_{(j,n) \in Q_\Delta} u_j^n (p_j^n - p_j^{n+1}) + \sum_{j=0}^J (u_j^N p_j^N - u_j^0 p_j^0) \\ &\quad + \sum_{(j,n) \in Q_\Delta} F_{j+1/2}^n (\lambda_j p_j^{n+1} - \lambda_{j+1} p_{j+1}^{n+1}) + \sum_{n=0}^{N-1} \lambda_0 F_{-1/2}^n p_0^{n+1} \\ &\quad - \sum_{n=0}^{N-1} \lambda_{J+1} F_{J+1/2}^n p_{J+1}^{n+1} - \sum_{(j,n) \in Q_\Delta} \mathcal{A}_{j+1/2}^n (\mu_j p_j^{n+1} - \mu_{j+1} p_{j+1}^{n+1}) \\ &\quad - \sum_{n=0}^{N-1} \mu_0 \mathcal{A}_{-1/2}^n p_0^{n+1} + \sum_{n=0}^{N-1} \mu_{J+1} \mathcal{A}_{J+1/2}^n p_{J+1}^{n+1} - \sum_{(j,n) \in Q_\Delta} g_j^n p_j^{n+1} \\ &= \sum_{(j,n) \in Q_\Delta} \left\{ \phi_j^n (p_j^n - p_j^{n+1}) + F_{j+1/2}^n (\lambda_j p_j^{n+1} - \lambda_{j+1} p_{j+1}^{n+1}) \right\} \end{aligned}$$

$$\begin{aligned}
& -\mathcal{A}_{j+1/2}^n(\mu_j p_j^{n+1} - \mu_{j+1} p_{j+1}^{n+1}) - g_j^n p_j^{n+1} \Big\} \\
& + \sum_{j=0}^J (u_j^N p_j^N - u_j^0 p_j^0) + \sum_{n=0}^{N-1} (\Gamma_0^n - \Gamma_{J+1}^n) p_0^{n+1}.
\end{aligned}$$

Bibliography

- [1] G. Anestis. *Eine eindimensionale Theorie der Sedimentation in Absetzbehältern veränderlichen Querschnitts und in Zentrifugen*. PhD thesis, Technical University of Vienna, 1981.
- [2] G. Anestis and W. Schneider. Application of the theory of kinematic waves to the centrifugation of suspensions. *Ing.-Arch.*, 53:399–407, 1983.
- [3] R. Becker. *Espesamiento continuo, diseño y simulación de espesadores*. Msc. thesis, Universidad de Concepción, 1982.
- [4] S. Berres and R. Bürger. On gravity and centrifugal settling of polydisperse suspensions forming compressible sediments. *Int. J. Solids Structures*, 40:4965–4987, 2003.
- [5] S. Berres, R. Bürger, A. Coronel, and M. Sepúlveda. Numerical identification of parameters for a strongly degenerate convection-diffusion problem modelling centrifugation of flocculated suspensions. *Appl. Numer. Math.*, to appear.
- [6] S. Berres, R. Bürger, A. Coronel, and M. Sepúlveda. Numerical identification of parameters for flocculated suspension from concentration measurements during batch centrifugation. *Chem. Eng. J.*, to appear.
- [7] S. Berres, R. Bürger, K. H. Karlsen, and E. M. Tory. Strongly degenerate parabolic-hyperbolic systems modeling polydisperse sedimentation with compression. *SIAM J. Appl. Math.*, 64:41–80, 2003.
- [8] R. Bürger, A. Coronel, and M. Sepúlveda. Convergence of upwind schemes for an initial-boundary value problem of a strongly degenerate parabolic equation modelling sedimentation-consolidation processes. Technical Report 2004-09, Departamento de Ingniería Matemática, Universidad de Concepción, Chile, 2004. Submitted to Math. Comp.

- [9] R. Bürger, M. C. Bustos, and F. Concha. Settling velocities of particulate systems: 9. phenomenological theory of sedimentation processes: Numerical simulation of the transient behavior of flocculated suspensions in an ideal batch or continuous thickener. *Int. J. Mineral Process.*, 55:267–282, 1999.
- [10] R. Bürger and F. Concha. Mathematical model and numerical simulation of the settling of flocculated suspensions. *Int. J. Multiphase Flow*, 24:1005–1023, 1998.
- [11] R. Bürger and F. Concha. Settling velocities of particulate systems: 12. batch centrifugation of flocculated suspensions. *Int. J. Mineral Process.*, 63:115–145, 2001.
- [12] R. Bürger, F. Concha, K.-K. Fjelde, and K. H. Karlsen. Numerical simulation of the settling of polydisperse suspensions of spheres. *Powder Technol.*, 113:30–54, 2000.
- [13] R. Bürger, F. Concha, and F. M. Tiller. Applications of the phenomenological theory to several published experimental cases of sedimentation processes. *Chem. Eng. J.*, 80:105–117, 2000.
- [14] R. Bürger, J. J. R. Damasceno, and K. H. Karlsen. A mathematical model for batch and continuous thickening in vessels with varying cross section. *Int. J. Mineral Process.*, 73:183–208, 2004.
- [15] R. Bürger, S. Evje, and K. H. Karlsen. On strongly degenerate convection-diffusion problems modeling sedimentation-consolidation processes. *J. Math. Anal. Appl.*, 247:517–556, 2000.
- [16] R. Bürger, S. Evje, K. H. Karlsen, and K. A. Lie. Numerical methods for the simulation of the settling of flocculated suspensions. *Chem. Eng. J.*, 80:91–104, 2000.
- [17] R. Bürger and K. H. Karlsen. On a diffusively corrected kinematic-wave traffic flow model with changing road surface conditions. *Math. Models Methods Appl. Sci.*, 13:1767–1799, 2003.
- [18] R. Bürger, K. H. Karlsen, N. H. Risebro, and J. D. Towers. Well-posedness in BV_t and convergence of a difference scheme for continuous sedimentation in ideal clarifier-thickener units. *Numer. Math.*, 97(1):25–65, 2004.

- [19] R. Bürger and K. H. Karlsen. On some upwind difference schemes for the phenomenological sedimentation-consolidation model. *J. Engrg. Math.*, 41:145–166, 2001.
- [20] R. Bürger and K. H. Karlsen. A strongly degenerate convection-diffusion problem modeling centrifugation of flocculated suspensions. In *Hyperbolic Problems: Theory, Numerics, Applications, Vol. I, II (Magdeburg, 2000)*, volume 141 of *Internat. Ser. Numer. Math.*, 140, pages 207–216. Birkhäuser, Basel, 2001.
- [21] R. Bürger, K. H. Karlsen, and J. D. Towers. A model of continuous sedimentation of flocculated suspensions in clarifier-thickener units. *SIAM J. Appl. Math.*, to appear.
- [22] R. Bürger, K. H. Karlsen, N. H. Risebro, and J. D. Towers. Monotone difference approximations for the simulation of clarifier-thickener units. *Comput. Visual. Sci.* 6:67–74, 2004.
- [23] R. Bürger and W. L. Wendland. Entropy boundary and jump conditions in the theory of sedimentation with compression. *Math. Methods Appl. Sci.*, 21:865–882, 1998.
- [24] R. Bürger and W. L. Wendland. Existence, uniqueness, and stability of generalized solutions of an initial-boundary value problem for a degenerating quasi-linear parabolic equation. *J. Math. Anal. Appl.*, 218:207–239, 1998.
- [25] R. Bürger and W. L. Wendland. Sedimentation and suspension flows: historical perspective and some recent developments. *J. Engrg. Math.*, 41:101–116, 2001.
- [26] R. Bürger, W. L. Wendland, and F. Concha. Model equations for gravitational sedimentation-consolidation processes. *ZAMM Z. Angew. Math. Mech.*, 80:79–92, 2000.
- [27] F. Bouchut and F. James. One-dimensional transport equations with discontinuous coefficients. *Nonlinear Anal. TMA*, 32(7):891–933, 1998.
- [28] M. C. Bustos, F. Concha, R. Bürger, and E. M. Tory. *Sedimentation and thickening*, volume 8 of *Mathematical Modelling: Theory and Applications*. Kluwer Academic Publishers, Dordrecht, 1999.
- [29] J. R. Cannon. Determination of certain parameters in heat conduction problems. *J. Math. Anal. Appl.*, 8:188–201, 1964.

- [30] J. R. Cannon and D. Zachmann Parameter determination in parabolic partial differential equations from overspecified boundary data. *Int. J. Engng. Sci.*, 20:779–788, 1982.
- [31] J. Carrillo. Entropy solutions for nonlinear degenerate problems. *Arch. Rat. Mech. Anal.*, 147:269–361, 1999.
- [32] B. Cockburn and G. Gripenberg. Continuous dependence on the nonlinearities of solutions of degenerate parabolic equations. *J. Diff. Equ.*, 151:231–251, 1999.
- [33] F. Concha. *Manual de Filtración y Separación*. Centro de Tecnología Mineral, CETTEM, fconcha@udec.cl, Concepción, Chile, pp. 184, 2001.
- [34] F. Concha. Private communication. 2002.
- [35] F. Concha and R. Bürger. Thickening in the 20th century: a historical perspective. *Minerals & Metallurgical Process.*, 20(2):57–67, 2003.
- [36] A. Coronel, M. Sepúlveda and F. Concha. *In preparation*, 2004.
- [37] A. Coronel, F. James, and M. Sepúlveda. Numerical identification of parameters for a model of sedimentation processes. *Inverse Problems*, 19(4):951–972, 2003.
- [38] M. G. Crandall and A. Majda. Monotone difference approximations for scalar conservation laws. *Math. Comp.*, 34(149):1–21, 1980.
- [39] G. Chavent. Identification of distributed parameters: about the output least square method, its implementation, and identifiability. In *Procc. of 5th. IFAC Symp. Identification and System Parameter Estimations*, volume 1, pages 85–97. Oxford: Pergamon, 1979.
- [40] C. M. Dafermos. Generalized characteristics and the structure of solutions of hyperbolic conservation laws. *Indiana Univ. Math. J.*, 26(6):1097–1119, 1977.
- [41] J. G. Dueck, D. Y. Kilimnik, L. L. Min'kov, T. Neeße. Measurement of the rate of sedimentation in a platelike centrifuge. *J. Eng. Phys. Thermophys.*, 76:748–759, 2003.
- [42] B. Engquist and S. Osher. One-sided difference approximations for nonlinear conservation laws. *Math. Comp.*, 36:321–351, 1981.
- [43] M. S. Espedal and K. H. Karlsen. Numerical solution of reservoir flow models based on large time step operator splitting algorithms. In *Filtration in porous*

- media and industrial application (Cetraro, 1998)*, volume 1734 of *Lecture Notes in Math.*, pages 9–77. Springer, Berlin, 2000.
- [44] L. C. Evans and R. F. Gariepy. *Measure theory and fine properties of functions*. Studies in Advanced Mathematics. CRC Press, Boca Raton, FL, 1992.
 - [45] S. Evje and K. H. Karlsen. Degenerate convection-diffusion equations and implicit monotone difference schemes. In *Hyperbolic problems: Theory, Numerics, Applications, Vol. I,II (Zürich, 1998)*, volume I of *M. Fey and R. Jeltsch (Eds.)*, pages 285–294. Birkhäuser, Basel, 1999.
 - [46] S. Evje, K. H. Karlsen, and N. H. Risebro. A continuous dependence result for nonlinear degenerate parabolic equations with spatially dependent flux function. In *Hyperbolic problems: Theory, Numerics, Applications, Vol. I, II (Magdeburg, 2000)*, volume 141 of *Internat. Ser. Numer. Math.*, 140, pages 337–346. Birkhäuser, Basel, 2001.
 - [47] S. Evje and K. H. Karlsen. Monotone difference approximations of BV solutions to degenerate convection-diffusion equations. *SIAM J. Numer. Anal.*, 37(6):1838–1860 (electronic), 2000.
 - [48] R. Eymard, T. Gallouët, and R. Herbin. Finite volume methods. In *Handbook of numerical analysis, Vol. VII*, Handb. Numer. Anal., VII, pages 713–1020. North-Holland, Amsterdam, 2000.
 - [49] D. Frömer and D. Lerche. An experimental approach to the study of the sedimentation of dispersed particles in a centrifugal field. *Arch. Appl. Mech.*, 72:85–95, 2002.
 - [50] P. Garrido, R. Bürger, and F. Concha. Settling velocities of particulate systems: 11. comparison of the phenomenological sedimentation-consolidation model with published experimental results. *Int. J. Mineral Process.*, 60:213–227, 2000.
 - [51] P. Garrido, R. Burgos, F. Concha, and R. Bürger. Software for the design and simulation of gravity thickeners. *Minerals Eng.*, 16:85–92, 2003.
 - [52] P. Garrido, R. Burgos, F. Concha, and R. Bürger. Settling velocities of particulate systems: 13. Software for the batch and continuous sedimentation of flocculated suspensions. *Int. J. Mineral Process.*, 73:131–144, 2004.
 - [53] P. Garrido, F. Concha, R. Bürger. Settling velocities of particulate systems: 14. Unified model of sedimentation, centrifugation and filtration of flocculated suspensions. *Int. J. Mineral Process.*, 72:57–74, 2003

- [54] P. Garrido Private communication. 2004.
- [55] L. Gosse and F. James. Numerical approximations of one-dimesional linear conservation equations with discontinuous coefficients *Math. Comp.*, 69:987–1015, 2003.
- [56] G. Guiochon, F. James, and M. Sepúlveda. Numerical results for the flux identification in a system of conservation laws. In *Hyperbolic problems: theory, numerics, applications, Vol. I (Zürich, 1998)*, volume 129 of *Internat. Ser. Numer. Math.*, pages 423–432. Birkhäuser, Basel, 1999.
- [57] S. Gutman. Identification of discontinuous parameters in flow equations. *SIAM J. Control Optim.*, 28:1049–1060, 1990.
- [58] J.-B. Hiriart-Urruty and C. Lemaréchal. *Convex analysis and minimization algorithms. I*, volume 305 of *Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]*. Springer-Verlag, Berlin, 1993. Fundamentals.
- [59] W. Hundsdorfer, J. G. Verwer. *Numerical solution of time-dependent advection-diffusion-reaction equations*. Springer-Verlag, Berlin 2003.
- [60] F. James, M. Postel, and M. Sepúlveda. Numerical comparison between relaxation and nonlinear equilibrium models. Application to chemical engineering. *Phys. D*, 138:316–333, 2000.
- [61] F. James and M. Sepúlveda. Parameter identification for a model of chromatographic column. *Inverse Problems*, 10:1299–1314, 1994.
- [62] F. James and M. Sepúlveda. Convergence results for the flux identification in a scalar conservation law. *SIAM J. Control Optim.*, 37:869–891 (electronic), 1999.
- [63] F. James and M. Sepúlveda. Parameter identification for a hyperbolic equation modelling chromatography. In *Nonlinear hyperbolic problems: theoretical, applied, and computational aspects (Taormina, 1992)*, volume 43 of *Notes Numer. Fluid Mech.*, pages 347–353. Vieweg, Braunschweig, 1993.
- [64] F. James, M. Sepúlveda, F. Charton, I. Quiñones and G. Guiochon. Determination of binary competitive equilibrium isotherms from the individual chromatographic band profiles. *Chem. Engng. Sci.*, 54:1677–1696, 1999.

- [65] K. H. Karlsen, N. H. Risebro, and J. D. Towers. Upwind difference approximations for degenerate parabolic convection-diffusion equations with a discontinuous coefficient. *IMA J. Numer. Anal.*, 22(4):623–664, 2002.
- [66] K. H. Karlsen, N. H. Risebro, and J. D. Towers. L^1 stability for entropy solutions of nonlinear degenerate parabolic convection-diffusion equations with discontinuous coefficients. *Skr. K. Nor. Vidensk. Selsk.*, (3):1–49, 2003.
- [67] K. H. Karlsen and N. H. Risebro. Convergence of finite difference schemes for viscous and inviscid conservation laws with rough coefficients. *M2AN Math. Model. Numer. Anal.*, 35(2):239–269, 2001.
- [68] K. H. Karlsen and N. H. Risebro. On the uniqueness and stability of entropy solutions of nonlinear degenerate parabolic equations with rough coefficients. *Discrete Contin. Dyn. Syst.*, 9(5):1081–1104, 2003.
- [69] K. H. Karlsen and M. Ohlberger. A note on the uniqueness of entropy solutions of nonlinear degenerate parabolic equations. *J. Math. Anal. Appl.*, 275:439–458, 2002.
- [70] Y. L. Keung and Y. Zou. Numerical identification of parameters in parabolic systems. *Inverse Problems*, 14:1299–1314, 1998.
- [71] K. Kunisch and L. White. The parameter estimation problem for parabolic equations and discontinuous observation operators. *SIAM J. Control Optim.*, 23:900–927, 1985.
- [72] S. N. Kružkov. First order quasilinear equations in several independent space variables. *Math. USSR Sb.*, 10:217–243, 1970.
- [73] G. J. Kynch. A theory of sedimentation. *Trans. Faraday Soc.*, 48:166–176, 1952.
- [74] D. Lerche. Dispersion stability and particle characterization by sedimentation kinetics in a centrifugal field, *J. Disp. Sci. Technol.* 23:699–709, 2002.
- [75] D. Lerche and D. Frömer. Theoretical and experimental analysis of the sedimentation kinetics of concentrated red cell suspensions in a centrifugal field: Determination of the aggregation and deformation of rbc by flux density and viscosity functions. *Biorheology*, 38:249–262, 2001.
- [76] R. J. Le Veque. *Finite volume methods for hyperbolic problems*. Cambridge University Press, Cambridge, UK, 2002.

- [77] R. M. Lueptow and W. Hübler. Sedimentation of a suspension in a centrifugal field. *J. Biomech. Eng.*, 113:485–491, 1991.
- [78] A. S. Michaels and J. C. Bolger. Settling rates and sediment volumes of flocculated Kaolin suspensions. *Ind. Engrg. Chem. Fund.*, 1:24–33, 1962.
- [79] P. Nelson. Traveling-wave solutions of the diffusively corrected kinematic-wave model. *Math. Comp. Modelling*, 35:561–579, 2002.
- [80] J. F. Richardson and W. N. Zaki. Sedimentation and fluidization: Part I. *Trans. Instn. Chem. Engrs. (London)*, 32:35–53, 1954.
- [81] J. F. Richardson and W. N. Zaki. The sedimentation of uniform spheres under conditions of viscous flow. *Chem. Eng. Sci.*, 3:65–73, 1954.
- [82] U. Schaflinger. Centrifugal separation of a mixture. *Fluid Dyn. Res.* 6:213–249, 1990.
- [83] F. M. Tiller and W. F. Leu. Basic data fitting in filtration. *J. Chin. Inst. Chem. Engrs.*, 11:61–70, 1980.
- [84] E. F. Toro. *Riemann Solvers and Numerical Methods for Fluid Dynamics*. Springer-Verlag, Berlin, 1997.
- [85] M. Ungarish. On the separation of a suspension in a tube centrifuge. *Int. J. Multiphase Flow*, 27:1285–1291, 2001.
- [86] S. Ulbrich. A sensitivity and adjoint calculus for discontinuous solutions of hyperbolic conservation laws with source terms. *SIAM J. Control Optim.*, 41(3):740–797 (electronic), 2002.
- [87] S. Ulbrich. Adjoint-based derivative computations for the optimal control of discontinuous solutions of hyperbolic conservation laws. *Systems & Control Lett.*, 48(3-4):313–328, 2003.
- [88] A. I. Vol'pert. Spaces BV and quasilinear equations. *Mat. Sb. (N.S.)*, 73 (115):255–302, 1967.
- [89] A. I. Vol'pert and S. I. Hudjaev. The Cauchy problem for second order quasilinear degenerate parabolic equations. *Mat. Sb. (N.S.)*, 78 (120):374–396, 1969.
- [90] Z. Q. Wu. *A boundary value problem for quasilinear degenerate parabolic equations*, volume 2484 of *MRC Technical Summary Report*. University of Wisconsin Center for the Mathematical Sciences, Madison, WI, 1983.

- [91] Z. Q. Wu and J. Y. Wang. Some results on quasilinear degenerate parabolic equations of second order. In *Proceedings of the 1980 Beijing Symposium on Differential Geometry and Differential Equations, Vol. 1, 2, 3 (Beijing, 1980)*, pages 1593–1609, Beijing, 1982. Science Press.
- [92] M. Yamamoto and Y. Zou. Simultaneous reconstruction of the initial temperature and heat radiative coefficient. *Inverse Problems*, 17:1181–1202, 2001.

