# Optimization on Gas Networks under Stochastic Boundary Conditions

Optimierung auf Gasnetzwerken
unter stochastischen Randdaten

Der Naturwissenschaftlichen Fakultät
der Friedrich-Alexander-Universität
Erlangen-Nürnberg
zur Erlangung des Doktorgrades

Dr. rer. nat.

vorgelegt von

**David Wintergerst**

aus Bamberg

Als Dissertation genehmigt

von der Naturwissenschaftlichen Fakultät

der Friedrich-Alexander-Universität Erlangen-Nürnberg

## Zusammenfassung

Wir betrachten die durch die isothermalen Eulergleichungen modellierte Gasdynamik auf Netzwerken. Die Abnahme der Kunden an den Ausgängen des Netzwerks stellt dabei eine Massenfluss-Randbedingung dar. Die wahren Abnahmen sind jedoch nicht exakt bekannt und werden mithilfe von Wahrscheinlichkeitsverteilungen abgebildet.

Zunächst werden die stationären Lösungen der quasilinearen isothermalen Eulergleichungen auf einem Rohr hergeleitet. Dabei wird ein Realgasmodell zugrunde gelegt. Es werden wichtige Monotonieeigenschaften der Lösung gezeigt, die für die Erweiterung der Lösung auf Netzwerke von entscheidender Bedeutung sind. Die Existenztheorie wird jedoch in einem allgemeineren Rahmen für Flussprobleme auf Netzwerken entwickelt. Damit ist es möglich unterschiedliche Kopplungsbedingungen an den Knoten sowie verschiedene – nicht auf Gas eingeschränkte – Modelle sowie aktive Elemente entlang der Kanten zu behandeln.

Für die Optimierung werden die unsicheren Randdaten mithilfe von Wahrscheinlichkeitsrestriktionen abgebildet. Diese garantieren, das System mit einer Mindestwahrscheinlichkeit innerhalb von technischen Druckschranken zu halten. Die Berechnung der Wahrscheinlichkeit erfolgt mittels einer Kombination aus Quasi-Monte-Carlo-Verfahren und der sphärisch radialen Zerlegung von gaußschen Zufallsvektoren. Für die gradientenbasierte Optimierung wird ein Resultat zur Gradientendarstellung auf den Fall von konvexen Mengen zulässiger Realisierungen erweitert. Es wird gezeigt, dass die benötigten Voraussetzungen für den Fall von Baumnetzwerken erfüllt sind. Zur numerischen Umsetzung wird ein Multilevelverfahren vorgeschlagen, das mithilfe der Lösung für niedriger Stichprobenzahlen einen Warmstart für das aufwendigere Modell berechnet. Dies führt zu einer deutlichen Verringerung der Rechenzeit.

Im transienten Fall wird die lineare Wellengleichung auf einer Kante betrachtet. Dabei sind die Anfangs- und die Randdaten an einem Intervallende unsicher und durch stochastische Prozesse gegeben. Das andere Intervallende wird durch ein Neumannfeedbackgesetz gesteuert. Die stochastischen Prozesse werden mithilfe des Karhunen-Loève-Theorems approximiert. Dadurch lassen sich die Anfangsranddaten mit endlich vielen Zufallsvariablen darstellen. Als Stabilitätsmaß wird dabei die Wahrscheinlichkeit in der $L^\infty$-Norm unterhalb einer gegebenen Schranke zu bleiben betrachtet. Das System lässt sich durch die Anpassung des Feedbackparameters stabilisieren.

## Abstract

We consider the gas dynamics on networks modeled by the isothermal Euler equations. The demands of the customers at the exits of the network represent a mass flow boundary condition. The true demands are not exactly known and are modeled by probability distributions.

First, the stationary solutions of the quasilinear isothermal Euler equations are derived on a single pipe. For this purpose, a real gas model is used. Important monotonicity properties of the solution, which are of crucial importance for the extension to networks, are shown. The existence theory is developed in a more general framework for flow problems on networks. This makes it possible to treat different coupling conditions at the nodes and to use various models that are not necessarily restricted to gas, or active elements along the edges.

For the optimization, the uncertain boundary conditions are represented by means of chance constraints. These guarantee that the system is kept within technical pressure limits with a certain probability. The calculation of the probability is realized by a combination of quasi-Monte Carlo methods and the spherical radial decomposition of Gaussian random vectors. For the gradient-based optimization, a result for a gradient representation is extended to include the case of convex sets of feasible realizations. It is shown that, for the case of tree networks, the required assumptions are fulfilled. For the numerical implementation, a multilevel method is proposed that uses the solution for low sample numbers as a warm start for the more complex model. This leads to a significant reduction in computational time.

In the transient case, the linear wave equation is considered on a single edge. The initial and boundary data at the end of an interval are uncertain and given by stochastic processes. The other end of the interval is controlled by a Neumann feedback law. The stochastic processes are approximated using the Karhunen-Loève theorem. This enables the representation of the initial-boundary data by a finite number of random variables. As a stability measure, the probability to stay below a given threshold in the $L^\infty$-norm, is considered. The system can be stabilized by adjusting the feedback parameter.

# Contents

# Preface

Before we dive into the depth of mathematical theory, there are a few people to mention without whom this dissertation would have taken longer to write, would contain more mistakes or would have been a dreadful experience to write due to the lack of joyful diversions. First and foremost I would like to thank my supervisor Martin Gugat for his insightful and helpful comments. Thanks to him, I had always a clear direction what to research next, which resulted in several joint publications; he is called the *paper machine* for a reason. Furthermore, I would like to thank Rüdiger Schultz and his group for offering their expertise during my visits in Essen, where we especially discussed the stochastic components of this thesis. Moreover, I would like to thank my family for their support and for proof reading this thesis. The practical joke "You should *optimize* this!" remained hilarious even after hearing it for the 100th time. I would also like to thank my (ex-)office mates Alex Keimer and Fabian Rüffler for the fruitful discussions on mathematical and—more often—non-mathematical topics. Furthermore, I would like to thank Martin Gugat, Günter Leugering, Alexander Martin, Martin Schmidt, and Mathias Sirvent for the great and effectively organized collaboration on the common publications. I would like to thank Johannes Semmler for his advise concerning Latex, Matlab, tikz and matters of aesthetics. I would like to thank Denis Aßmann and Robert Burlescu for their insightful comments on derivatives and linearizations. Last but not least I would like to thank the table soccer team consisting of Tobias Klein, Martin Knossalla, Tobias Kufner, Björn Liljegren-Sailer, Fabian Rüffler, Johannes Semmler, Mathias Sirvent, Michele Spinola, and Christoph Strohmeyer for the competitive environment and high level play in a thriving sport, which got me through long days of work.

by the German research foundation DFG[3]. The following articles are the cores of Chapter 1, 2 and 3 respectively:

- Martin Gugat, Rüdiger Schultz, and David Wintergerst. "Networks of pipelines for gas with nonconstant compressibility factor: stationary states". In: *Computational and Applied Mathematics* 37.2 (2018), pp. 1066–1097,

- David Wintergerst. "Application of Chance Constrained Optimization to Gas Networks". Friedrich-Alexander-Universität Erlangen-Nürnberg. Preprint. https://opus4.kobv.de/opus4-trr154/frontdoor/index/index/docId/158. 2017,

- Dennis Adelhütte, Denis Aßmann, Tatjana Gonzàlez Grandòn, Martin Gugat, Holger Heitsch, Renè Henrion, Frauke Liers, Sabrina Nitsche, Rüdiger Schultz, Michael Stingl, and David Wintergerst. "Joint model of probabilistic-robust (probust) constraints with application to gas network optimization". University of Duisburg-Essen, Weierstrass Institute Berlin, Humboldt University Berlin, Friedrich-Alexander Universität Erlangen-Nürnberg. Preprint. https://opus4.kobv.de/opus4-trr154/frontdoor/index/index/docId/215. 2018 (In Revision).

Further articles that were written in the past years, which have had a smaller contribution to this thesis, but empowered new ideas, broadened my horizon and introduced me to different fields of mathematics are

- Martin Gugat, Günter Leugering, Alexander Martin, Martin Schmidt, Mathias Sirvent, and David Wintergerst. "Towards simulation based mixed-integer optimization with differential equations". In: *Networks* 72.1 (2018), pp. 60–83,

- Martin Gugat, Günter Leugering, Alexander Martin, Martin Schmidt, Mathias Sirvent, and David Wintergerst. "MIP-based instantaneous control of mixed-integer PDE-constrained gas transport problems". In: *Computational Optimization and Applications* 70.1 (2018), pp. 267–294,

- Martin Gugat and David Wintergerst. "Transient Flow in Gas Networks: Traveling Waves". In: *International Journal of Applied Mathematics and Computer Science* 28.2 (2018), pp. 341–348.

---

[3]Deutsche Forschungsgemeinschaft

# Introduction

The turnaround in energy politics and the new focus on renewable energies led to an increased significance of gas as a buffer for power supply. Gas is tradable, quickly available and storable, which ensures its essential role over the next decades. It provides a reliable, clean, eco- and climate-friendly energy production.

The gas is transported through pipelines from the suppliers to the customers. Due to friction, the gas pressure drops along the pipe making it necessary to raise it again at compressor stations. This leads to a multitude of new mathematical challenges. The gas flow is modeled by a quasilinear system of hyperbolic balance laws—the isothermal Euler equations. The use of coupling condition at the junctions allows the extension to networks of pipelines. The demands of the customers translate to mass flow boundary conditions. The goal is to control the compressor stations and other active elements in a way that respects technical pressure bounds. This is a typical problem of mathematical optimization. In practice, however, the exact demands of the customers are not known, but statistical data is available. Hence, the resulting problem is an optimization problem with hyperbolic partial differential equations on networks under stochastic boundary conditions. In its entirety, the problem is completely out of reach for today's methods. The approach in this thesis is to break down the problem. First, we consider the existence and uniqueness of stationary states on networks. Then, the stochastic components are included and stationary stochastic optimization problems are solved. Finally, we consider a simplified transient model under stochastic initial-boundary conditions.

In Chapter 1, we discuss the stationary states and obtain explicit solutions on a single pipe. Instead of just focusing on the gas dynamics for the extension to networks, we derive a powerful existence-and-uniqueness result for general flow problems on networks. This does not only allow the use of varying node coupling conditions, but also different gas models, the inclusion of active elements and the transfer of the result to water networks or arterial blood flow. The proof is constructive. Hence, the stationary states on the network can be directly calculated, which is demonstrated by various examples.

In Chapter 2, we consider chance constrained optimization problems. An overview of Monte Carlo and quasi-Monte Carlo methods to evaluate the probability function is presented. For gaussian random variables a parametrization of the integral describing the probability—the spherical radial decomposition—provides a smaller variance for the sampling. We extend a known result for the gradient representation of this parametrization to the case, where the set of feasible realizations is convex. We then proceed to prove that the assumptions for the gradient representation are met if we consider tree shaped gas networks. This allows the application of gradient based optimization methods. A multilevel algorithm that generates a warm start with a lower sample number is proposed and leads to a significant improvement in computation time.

In Chapter 3, the linear wave equation is considered under stochastic inital-boundary conditions. The initial data and the boundary data on one end of the space interval is given by a stochastic process, while the other end of the space interval is governed by a feedback law. The Karhunen-Loève theorem allows to approximate stochastic processes by a finite number of random variables. We consider the probability that the maximum of the system state stays under a prescribed bound. This gives a measure for the stability of the system. The feedback parameter can be adapted to increase the stability of the system.

# 1. Stationary Solutions of the Isothermal Euler Equations

*Everything flows and nothing abides,*
*everything gives way and nothing stays fixed.*

(Heraclitus)

In this chapter, we discuss the gas flow through pipeline networks in an equilibrium. The gas dynamic along a single pipeline as shown in Figure 1.1 is modeled by the isothermal Euler equations.



Figure 1.1.: Gas pipeline next to the B-145, province of El Loa, Chile, Photo by Diego Delso, delso.photo, License CC-BY-SA

We derive analytic stationary solutions and discuss the properties of the pressure solution. It turns out that the mathematical model shares the physically expected characteristics: The pressure decreases in flow direction due to the friction in the pipe, a higher initial pressure leads to a higher pressure at the end of the pipe and a higher flow rate leads to a higher pressure drop.

The monotonicity properties of the pressure are the basis for the existence theory on networks. Mathematically, the pipelines correspond to arcs in a graph and the junctions correspond to nodes in a graph. By using a very general setting of arc coupling functions that couple the state values along the arcs, and node coupling functions that couple the state and flow values of adjacent arcs, the tackled problem class exceeds mere gas transport. In the gas transport case, however, the existence theory allows the inclusion of different node coupling conditions, like the Bernoulli invariant instead of

the usual pressure continuity condition and the inclusion of active elements like resistors or turbo compressors.

## Literature Survey

In Gugat, Hante, et al. 2015, stationary states of the ideal gas model were analyzed on specific networks. For real gas on a single pipe, the stationary states were discussed in Schmidt, Steinbach, and Willert 2014. In Gugat, Schultz, and Wintergerst 2018 a concrete representation of the stationary states was derived and the existence and uniqueness on arbitrary graphs was shown.

A survey of flow problems on networks modeled by hyperbolic balance laws is provided in Bressan, Čanić, et al. 2014. The well-posedeness of general networked systems of balance laws was studied in Gugat, Herty, et al. 2012. The articles Colombo, Guerra, et al. 2009 and Gugat and Herty 2011 examined the optimal control of gas networks and networked systems of canals. Conservation laws without friction were considered in Garavello and Piccoli 2009 for the case of traffic flow and in Colombo and Marcellini 2010 for the $p$-system describing gas flow. For the pressure law considered in this thesis, Bakhvalov's condition holds and hence the $p$-system for barotropic gas has a global entropy weak solution with uniformly bounded total variation; see Bressan, Chen, et al. 2015 and Bakhvalov 1970. However, it is pointed out in Bressan, Chen, et al. 2015 that general pressure laws can lead to a solution with arbitrary large total variation.

In Reigstad 2014, numerical models for isothermal junction flow without friction were investigated and alternative coupling conditions were proposed. In Reigstad et al. 2015, the conservation of energy over junctions was examined for different coupling conditions. Entropy-preserving coupling conditions for the temperature dependent Euler equations were considered in Lang and Mindt 2018.

The Weymouth equation, which results from a simplification of the stationary isothermal Euler equations, is used for evaluation of gas network capacities in Koch et al. 2015. The existence and uniqueness of gas networks governed by the Weymouth equation was shown in Ríos-Mercado et al. 2002, where the simpler structure ensuing from the Weymouth equation allows an elegant argument using strongly monotonic operators. Models, structures and algorithms for stationary gas flows are discussed in the thesis Stangl 2014.

An overview over optimization on gas networks is given in Rios-Mercado and Borraz-Sanchez 2015. Mixed Integer models for stationary gas network

optimization were considered in Martin, Möller, and Moritz 2006. In Gugat, Leugering, Martin, et al. 2018b, the stationary pressure loss equations on each pipe were treated as "black box"-nonlinearity leading to a linear mixed integer formulation, where the unknown nonlinearity is adaptively refined by solving a nonlinear optimization problem.

## 1.1. Stationary Solutions on a Single Pipe

We consider the isothermal Euler equations without slope using the notation explained in Table 1.1

$$\begin{cases} \partial_t \rho + \partial_x q = 0, \\ \partial_t q + \partial_x \left( p + \dfrac{q^2}{\rho} \right) = -\dfrac{\theta}{2} \dfrac{q|q|}{\rho}. \end{cases} \tag{Iso}$$

The first equation models the conservation mass, while the second equation describes the balance of momentum. When the systems reaches an equilibrium, the time derivatives vanish and we obtain the stationary isothermal Euler equations

$$\begin{cases} \partial_x q = 0, \\ \partial_x \left( p + \dfrac{q^2}{\rho} \right) = -\dfrac{\theta}{2} \dfrac{q|q|}{\rho}. \end{cases} \tag{IsoStat}$$

The relation between density and pressure is described by the state equation for real gas

$$p = RT z(p) \rho. \tag{StateEq}$$

For the compressibility factor $z(p)$, we use a model by the *American Gas Association* described in Starling, Savidge, Association, et al. 1992 and evaluated in De Almeida, Velásquez, and Barbieri 2014. The model is linearly decreasing in the pressure

$$z(p) = 1 + \alpha p, \tag{AGA}$$

with the strictly negative (AGA)-constant $\alpha$. For ideal gas, the compressibility factor would be constant.

Table 1.1.: Notation

| Variable | Symbol | Unit | Range |
|---|---|---|---|
| pressure | $p$ | Pa | $\mathbb{R}_{>0}$ |
| density | $\rho$ | $\mathrm{kg\,m^{-3}}$ | $\mathbb{R}_{>0}$ |
| flow | $q$ | $\mathrm{kg\,s^{-1}\,m^{-2}}$ | $\mathbb{R}$ |
| friction coefficient | $\theta$ | $\mathrm{m^{-1}}$ | $\mathbb{R}_{\geq 0}$ |
| length | $L$ | m | $\mathbb{R}_{>0}$ |
| compressibility factor | $z(p)$ | 1 | $(0,1)$ |
| (AGA)-constant | $\alpha$ | $\mathrm{Pa^{-1}}$ | $(-1,0)$ |
| specific gas constant | $R$ | $\mathrm{J\,kg^{-1}\,K^{-1}}$ | $\mathbb{R}_{\geq 0}$ |
| temperature | $T$ | K | $\mathbb{R}_{\geq 0}$ |
| gas velocity | $v$ | $\mathrm{m\,s^{-1}}$ | $\mathbb{R}$ |
| speed of sound | $c$ | $\mathrm{m\,s^{-1}}$ | $\mathbb{R}_{>0}$ |

## 1.1.1. Real Gas

First, we consider the case of real gas.

*Assumption* 1. Throughout this section we assume for all $x \in [0, L]$

**subsonic flow,** that is

$$\frac{q^2 RT}{p^2} < 1,$$

**positive compressibility factor,** that is

$$z(p) > 0.$$

*Remark* 1 (Subsonic Flow). The subsonic flow condition states that the *squared Mach number* $\eta := (v/c)^2$ shall be smaller than one. It is defined as quotient of the *gas velocity* $v := q/\rho$ and the *speed of sound* $c := \sqrt{\partial p / \partial \rho}$. By (StateEq) and (AGA) we see that

$$\frac{1}{c^2} = \frac{\partial \rho}{\partial p} = \frac{1}{RT} \partial_p \left[ \frac{p}{z(p)} \right] = \frac{pz'(p) - z(p)}{z(p)^2} = \frac{1}{RT z(p)^2}.$$

For the squared velocity, we obtain

$$v^2 = (RT)^2 \frac{q^2 z(p)^2}{p^2}.$$

In a nutshell, this yields

$$\eta = \frac{v^2}{c^2} = RT\frac{q^2}{p^2}.$$

The stationary states of (Iso) with nonconstant compressibility factor are derived in Gugat, Schultz, and Wintergerst 2018. Solutions for a constant compressibility factor are considered in Gugat, Hante, et al. 2015.

**Theorem 1** (Stationary Solutions for Real Gas).
*Let a pressure $p_u > 0$ with $z(p_u) > 0$ and a subsonic flow $q_a \in \mathbb{R}$ at $x = 0$ be given. Then, the system* (IsoStat) *under the boundary conditions*

$$q(0) = q_a,$$
$$p(0) = p_u,$$

*has the unique solution*

$$q(x) = q_a, \qquad\qquad\qquad \text{for all } x \in [0, L],$$
$$p(x) = F^{-1}\left( F(p_u, q_a) - RT q_a |q_a| \int_0^x \tfrac{\theta(s)}{2}\, \mathrm{d}s \right), \qquad \text{for all } x \in [0, L],$$

*where the function $F : (0, |\alpha|^{-1}) \times \mathbb{R} \to \mathbb{R}$ is defined as*

$$F(p, q) := \frac{p}{\alpha} - \frac{1}{\alpha^2}\ln(z(p)) + q^2 RT \ln\left( \frac{z(p)}{p} \right)$$

*and the inverse $F^{-1} : \mathbb{R} \to (0, |\alpha|^{-1})$ is to be understood with respect to $p$ for fixed flow $q$.*

*Proof.* The first equation in (IsoStat) states that the flow is constant on the interval $[0, L]$. The state equation (StateEq) allows the substitution of the density in (IsoStat) by $1/\rho = RT z(p)/p$. Inserting this in the second equation of (IsoStat) yields

$$\partial_x\left( p + q^2 RT\frac{z(p)}{p} \right) = -\frac{\theta}{2} q|q| RT \frac{z(p)}{p}.$$

Carrying out the differentiation results in

$$\left( 1 - RT\frac{q^2}{p^2} \right)\partial_x p = -\frac{\theta}{2} q|q| RT \frac{z(p)}{p}. \tag{1.1}$$

The subsonic flow condition $RT\frac{q^2}{p^2} < 1$ prevents the space derivative of the pressure from blowing up. Defining

$$f(p,q) := \frac{z(p)}{p}\left(1 - RT\frac{q^2}{p^2}\right)^{-1} \tag{1.2}$$

and

$$g(x;q) := -\frac{\theta(x)}{2}RTq|q| \tag{1.3}$$

allows us to write Equation (1.1) in the standard form of a seperable ordinary differential equation

$$\partial_x p = f(p,q)\,g(x;q). \tag{1.4}$$

To solve the differential equation explicitly, it is mandatory to find the primitive (w.r.t. $p$) of

$$\frac{1}{f(p,q)} = \frac{p}{z(p)} - RTq^2\frac{1}{z(p)p} = \frac{1}{\alpha} - \frac{1}{\alpha z(p)} - RTq^2\left(\frac{1}{p} - \frac{\alpha}{z(p)}\right). \tag{1.5}$$

As the terms on the right hand side have elementary antiderivatives, it is given by

$$F(p,q) := \int \frac{1}{f(y,q)}\,\mathrm{d}y = \frac{p}{\alpha} - \frac{1}{\alpha^2}\ln(z(p)) - RTq^2[\ln(p) - \ln(z(p))]$$

$$= \frac{p}{\alpha} - \frac{1}{\alpha^2}\ln(z(p)) + RTq^2\ln\left(\frac{z(p)}{p}\right). \tag{1.6}$$

The arguments of the logarithmic functions are positive, wherefore it is not necessary to take absolute values. The separation of variables theorem (see e.g. Burkard 2014, Section Separable Differential Equations) states that a nonzero pressure solution of (1.4) is given by solving the implicit equation

$$F(p,q) - F(p_u,q) = \int_0^x g(s;q)\,\mathrm{d}s. \tag{1.7}$$

The only step missing to obtain the pressure solution is to justify the use of the inverse function $F^{-1} : \mathbb{R} \to (0,|\alpha|^{-1})$. Lemma 1 below shows that $F$ is strictly increasing in the first argument. Hence, the inverse $F^{-1}$ exists. Applying it to (1.7) and using $q(x) = q_a$ for all $x \in [0,L]$ yields

$$p(x) = F^{-1}\left(F(p_u,q_a) - RTq_a|q_a|\int_0^x \frac{\theta(s)}{2}\,\mathrm{d}s\right). \qquad\qquad \square$$

**Lemma 1** (Strict Monotonicity of $F$)**.**
*The function* $F : (0, |\alpha|^{-1}) \times \mathbb{R} \to \mathbb{R}$ *is* strictly increasing *in the first argument.*

*Proof.* We argue using the first derivative w.r.t the first argument. By the definition (1.6) of $F$, the fundamental theorem of calculus and Equation (1.5), it is given by

$$\partial_p F(p, q) = \frac{1}{f(p, q)} = \frac{p}{z(p)} - RTq^2 \frac{1}{z(p)p} = \frac{p}{z(p)} \left( 1 - RT \frac{q^2}{p^2} \right). \quad (1.8)$$

By the assumptions of subsonic flow $\frac{q^2 RT}{p^2} < 1$ and positive compressibility factor $z(p) > 0$ it follows that $\partial_p F(p, q) > 0$. $\qquad\square$

The pressure and the corresponding Mach number along a pipe are displayed in Figure 1.2. The dependency of the initial pressure value is shown in Figure 1.3. One can observe that for a Mach number close to one the derivative tends to minus infinity. At $\eta = 1$ the solution breaks down due to the blowup in the space derivative.

If the friction factor $\theta$ is a constant function, the condition of subsonic flow at every point $x \in [0, L]$ can be rewritten in terms of a critical length up to which the solution is defined. Just for Lemma 2 we drop Assumption 1.

**Lemma 2** (Critical Length)**.**
*Let the friction factor be constant, i.e.* $\theta(x) = \theta_{const} > 0$ *for all* $x \in [0, L]$*. Let an initial flow* $q_a \neq 0$ *and an initial pressure* $p_u \in \left( |q_a| \sqrt{RT}, |\alpha|^{-1} \right)$ *be given. Define the critical length (note that* $g(x; q_a) \neq 0$ *is constant)*

$$L_c := |g(x; q_a)|^{-1} \big[ F(p_u, q_a) - F\big( |q_a| \sqrt{RT}, q_a \big) \big]. \quad (1.9)$$

*Then, the unique solution* $p$ *of* (IsoStat) *as defined in Theorem 1 exists for all* $x \in [0, L_c)$*.*

*Proof.* It needs to be shown that $p(x)$ stays in the interval

$$\left( |q_a| \sqrt{RT}, |\alpha|^{-1} \right) \quad \text{for} \quad x \in [0, L_c)$$

or in other words that $x \in [0, L_c)$ implies Assumption 1.

For positive flow $q_a > 0$, the inequality

$$x < |g(x; q_a)|^{-1} \big[ F(p_u, q_a) - F\big( |q_a| \sqrt{RT}, q_a \big) \big],$$

which holds by assumption, implies

$$F(p_u, q_a) + g(x; q_a)\, x > F\big(|q_a|\sqrt{RT}, q_a\big), \qquad (1.10)$$

because $g(x; q_a)$ is negative for $q_a > 0$. Lemma 1 states that the function $F$ is strictly increasing in the first argument. Taking the inverse of $F$ with respect to the first argument on both sides of the Inequality (1.10) and using the representation of Theorem 1 yields

$$p(x) > |q_a|\sqrt{RT}, \ \text{ for all } x \in [0, L_c)$$

for positive flow $q_a$.

For negative flow $q_a < 0$ the pressure strictly increases in $x$, which can be quickly seen by Equation (1.1) and is also stated in Lemma 4. Subsequently, for $q_a < 0$,

$$p(x) > p_u > |q_a|\sqrt{RT}$$

holds for all $x > 0$.

It remains to show that $x \in [0, L_c)$ implies $p(x) < |\alpha|^{-1}$. For positive flow $q_a > 0$, the function $p$ strictly decreases in $x$ due to Equation (1.1) or Lemma 4. Therefore, for $q_a > 0$,

$$p(x) < p_u < |\alpha|^{-1}$$

holds for all $x > 0$.

The case $q_a < 0$ is left. Suppose there is a point $\bar{x}$ such that $p(\bar{x}) = |\alpha|^{-1}$. Then, by Equation (1.4), we have $\partial_x p = 0$. The function $f(p, q)$ defined in (1.2) is continuously differentiable with respect to $p$ for all $p > |q_a|\sqrt{RT}$. Thus it is locally Lipschitz continuous and, by the Picard Lindelöf Theorem, the differential equation (1.4) has a unique solution. However, we see that $f(|\alpha|^{-1}, q_a) = 0$ and therefore $p_x(\bar{x}) = 0$, which implies that the constant solution $p(x) = |\alpha|^{-1}$ solves (1.4). This contradicts the uniqueness of the solution, because for $x$ close to 0, the function value $p(x)$ is close to $p_u < |\alpha|^{-1}$ by the continuity of $p$. Consequently, for $q_a < 0$,

$$p(x) < |\alpha|^{-1}$$

has to hold for all $x \in [0, L_c)$. This concludes the proof. $\qquad \square$

The representation in Theorem 1 is quite concrete but it is not yet clear, how to evaluate the function $F^{-1}$. In the following, we proceed to show the convexity of $F$, which allows the direct application of Newton's method. Therefore, we will see that $F^{-1}$ can be evaluated numerically with high accuracy using only a few Newton steps.

**Lemma 3** (Strict Convexity of $F$)**.**
*The function $F : (0, |\alpha|^{-1}) \times \mathbb{R} \to \mathbb{R}$ is strictly convex with respect to the first argument.*

*Proof.* The first derivative of $F$ is given by Equation (1.8). Differentiating once more yields

$$
\begin{aligned}
\partial_{pp}^2 F(p,q) &= \frac{1}{z(p)^2}\left(1 - RT\frac{q^2}{p^2}\right) + \frac{p}{z(p)}\frac{2RTq^2}{p^3} \\
&= \frac{1}{z(p)^2}\left(1 - RT\frac{q^2}{p^2}\right) + \frac{1}{z(p)}\frac{2RTq^2}{p^2}.
\end{aligned}
\tag{1.11}
$$

The two assumptions $\frac{q^2 RT}{p^2} < 1$ and $z(p) > 0$ imply $\partial_{pp}^2 F(p,q) > 0$. $\square$

*Remark 2* (Numerical Evaluation of $F^{-1}$). Consider $z \in \operatorname{im} F$ and $q \in \mathbb{R}$ and define
$$
\varphi(y) = F(y,q) - z.
$$
The evaluation of $F^{-1}(z)$ is equivalent to finding a root of $\varphi$. The function $F$ is strictly increasing and strictly convex in $y$. Subsequently the function $\varphi$ has the same properties. By the subsonic flow condition and the positive compressibility factor, we consider $y$ on the interval $(|q|\sqrt{RT}, |\alpha|^{-1})$ and therefore $z > F(|q|\sqrt{RT}, q)$. This implies

$$
\varphi(|q|\sqrt{RT}) < 0
$$

and by the asymptotic properties of $F(\,\cdot\,, q)$, we obtain

$$
\lim_{y \nearrow |\alpha|^{-1}} \varphi(y) = \infty.
$$

Consequently, Bolzano's intermediate value theorem implies the existence of a unique root $y^* \in (|q|\sqrt{RT}, |\alpha|^{-1})$.

Because the function $\varphi$ is strictly increasing and strictly convex, the Newton iteration defined by

$$
y_{n+1} = y_n - \frac{\varphi(y_n)}{\varphi'(y_n)}
$$

generates a monotonically decreasing sequence converging to $y^*$ for each starting point $y_0 \in (y^*, |\alpha|^{-1})$; see Ortega and Rheinboldt 1970, Theorem 13.3.7.

Figure 1.2.: The pressure $p(\,\cdot\,, q_a, p_u)$ and the corresponding Mach number $\sqrt{\eta}$ along the pipe for $\theta = 0.014261$, $p_u = 60 \cdot 10^5\,\mathrm{Pa}$, $q_a = 453.1495\,\mathrm{kg\,m^{-2}\,s^{-1}}$.

Figure 1.3.: The pressure $p(L, q_a, \cdot)$ at $x = L$ and the corresponding Mach number $\sqrt{\eta}$ as a function of the initial pressure for $\theta = 0.014261$, $L = 50\,\mathrm{km}$, $q_a = 453.1495\,\mathrm{kg\,m^{-2}\,s^{-1}}$

## 1.1.2. Ideal Gas

The case of ideal gas, i.e., the compressibility factor $z(p) = z_m$ is constant, was treated in Gugat, Hante, et al. 2015. In the subsonic case, that is,

$$\eta := \frac{q^2 c^2}{p^2} < 1,$$

the solution of (IsoStat) under the boundary conditions

$$q(0) = q_a,$$
$$p(0) = p_u,$$

is given by

$$p(L) = c\,|q_a|\sqrt{-W_{-1}\Big(-\exp\big(-C + \text{sign}(q_a)\int_0^L \theta_a(s)\,\mathrm{d}s\big)\Big)},$$

where the constant $C$ is defined as

$$C := \frac{1}{\eta_u} + \ln(\eta_u), \quad \eta_u := \frac{q_a^2 c^2}{p_u^2}.$$

# 1.2. Properties of the Pressure Solution for Real Gas

In this section, we investigate important properties of the stationary solutions derived in Theorem 1. The monotonicity of the solution with respect to the initial values will turn out to be essential for the existence proof on networks of multiple pipes. The concavity of the pressure with respect to the initial pressure proves to be important in the context of optimization; see Gugat, Leugering, Martin, et al. 2018b. Assumption 1 is still valid.

### Monotonicity Properties of the Pressure Solution

**Lemma 4** (Monotonicity of the Pressure Solution in the Space Variable). *For $p_u > 0$ and $q_a \neq 0$, the pressure solution is* strictly decreasing *in space for positive flow and* strictly increasing *in space for negative flow, i.e.,*

$$\text{sign}(\partial_x p) = -\text{sign}(q_a).$$

*Proof.* The space derivative of $p$ is given by Equation (1.4). For subsonic flow and positive compressibility factor, the term $f(p, q)$ defined in (1.2) is positive. The definition of $g$ in (1.3) implies $\mathrm{sign}(g(x; q)) = -\mathrm{sign}(q)$ and the result follows. $\qquad\square$

Next, we determine the monotonicity with respect to the initial values. The expected physical behavior occurs: A higher ingoing pressure leads to a higher pressure at the outlet and a higher flow leads to a higher pressure drop.

**Lemma 5** (Monotonicity of the Pressure Solution in the Initial Values). *For $x > 0$, the pressure solution $p$ is* strictly increasing *as a function of the initial pressure*

$$\partial_{p_u} p > 0 \tag{1.12}$$

*and* strictly decreasing *as a function of the flow for $q_a \neq 0$, i.e.*

$$\partial_{q_a} p < 0. \tag{1.13}$$

*Proof.* Implicit differentiation of Equation (1.7) yields

$$\partial_p F(p, q) \partial_{p_u} p - \partial_p F(p_u, q) = 0$$

and therefore we obtain

$$\partial_{p_u} p = \frac{\partial_p F(p_u, q)}{\partial_p F(p, q)}. \tag{1.14}$$

By Lemma 1, both the numerator and the denominator are positive, which approves the first claim (1.12).

To show (1.13), we start by calculating

$$\partial_q F(p, q) = 2RTq \ln\left(\frac{z(p)}{p}\right) \tag{1.15}$$

and

$$\partial_q g(x; q) = -\theta(x) RT |q|. \tag{1.16}$$

Implicit differentation of (1.7) with respect to $q_a$ leads to

$$\partial_p F(p, q) \partial_{q_a} p + \partial_q F(p, q) - \partial_q F(p_u, q) = \int_0^x \partial_q g(s; q) \, \mathrm{d}s, \tag{1.17}$$

which is equivalent to

$$\partial_{q_a} p = \frac{\partial_q F(p_u, q) - \partial_q F(p, q) + \int_0^x \partial_q g(s; q) \, \mathrm{d}s}{\partial_p F(p, q)}. \tag{1.18}$$

By Equation (1.16), the term $\int_0^x \partial_q g(s; q)\, \mathrm{d}s$ is negative for $x > 0$ and by Lemma 1 the denominator $\partial_p F(p, q)$ is positive. We inspect the remaining terms:

$$\partial_q F(p_u, q) - \partial_q F(p, q) = 2RTq\left[\ln\left(\frac{z(p_u)}{p_u}\right) - \ln\left(\frac{z(p)}{p}\right)\right]$$
$$= 2RTq\left[\ln\left(\frac{z(p_u)}{z(p)}\right) + \ln\left(\frac{p}{p_u}\right)\right].$$

By Lemma 4, the pressure $p$ is decreasing in $x$ for positive flow and increasing in $x$ for negative flow. Because the compressibility factor is decreasing in $p$, this implies

$$\operatorname{sign}\left(\ln\left(\tfrac{p}{p_u}\right)\right) = -\operatorname{sign}(q) \quad \text{and} \quad \operatorname{sign}\left(\ln\left(\tfrac{z(p_u)}{z(p)}\right)\right) = -\operatorname{sign}(q).$$

To sum things up: We showed for the terms in (1.18) that

$$\partial_q F(p_u, q) - \partial_q F(p, q) < 0, \quad \int_0^x \partial_q g(s; q)\, \mathrm{d}s < 0 \quad \text{and} \quad \partial_p F(p, q) > 0.$$

This asserts the second claim (1.13). □

The concavity of the pressure solution w.r.t. the initial value for positive flow was derived in Gugat, Leugering, Martin, et al. 2018b, Theorem 3 to obtain the correctness of a decomposition method for stationary mixed integer gas network optimization problems. For the existence theory on networks only the first order properties are relevant.

### Second Order Properties of the Pressure Solution

**Lemma 6** (Second Order Dependency on the Initial Pressure)**.**
*For $x > 0$ and nonzero flow $q_a$, the pressure solution $p$ is* strictly concave *for positive flow and* strictly convex *for negative flow as a function of the initial pressure $p_u$, i.e.,*

$$\operatorname{sign}\left(\partial^2_{p_u p_u} p\right) = -\operatorname{sign}(q_a).$$

*Proof.* The first derivative (1.14) in terms of $f$ as defined in (1.2) can be written as

$$\partial_{p_u} p = \frac{f(p, q)}{f(p_u, q)}. \tag{1.19}$$

Figure 1.4.: The function $h$ on the domain $(0, 1)$.

Differentiating once more yields

$$\partial^2_{p_u p_u} p = \frac{f(p_u, q) \partial_p f(p, q) \partial_{p_u} p - f(p, q) \partial_p f(p_u, q)}{f(p_u, q)^2}$$

and inserting (1.19) leads to

$$f(p_u, q)^2 \, \partial^2_{p_u p_u} p = f(p, q)[\partial_p f(p, q) - \partial_p f(p_u, q)]. \qquad (1.20)$$

As $f(p_u, q) > 0$ holds under Assumption 1, we have to check the sign of the bracketed term. The product rule applied to (1.2) with use of the notation $\eta = RT \frac{q^2}{p^2}$ leads to

$$\partial_p f(p, q) = \frac{z(p)}{p} \frac{\partial_p \eta}{(1 - \eta)^2} - \frac{1}{p^2} (1 - \eta)^{-1}$$

$$= -(RTq^2)^{-1} \left[ 2z(p) \frac{\eta^2}{(1 - \eta)^2} + \frac{\eta}{1 - \eta} \right]. \qquad (1.21)$$

The function

$$h(\eta) := \frac{\eta}{1 - \eta} \qquad (1.22)$$

is strictly increasing in the squared Mach number $\eta$, for $\eta \in (0, 1)$ and therefore $h$ and $h^2$ are strictly decreasing as a function of the pressure. For a quick view see Figure 1.4.

The compressibility factor $z$ is also strictly decreasing in the pressure. Furthermore, $h, h^2$ and $z$ are positive valued. Hence, the bracketed term

Figure 1.5.: The pressure $p(L, \cdot, p_u)$ at $x = L$ as a function of the initial flow and the corresponding Mach number $\sqrt{\eta}$ for $\theta = 0.014261$, $L = 50 \, \text{km}, \;\; p_u = 40 \cdot 10^5 \, \text{Pa}$

in (1.21) is strictly decreasing in $p$ and as a consequence of this, $\partial_p f(\,\cdot\,, q)$ is a strictly increasing function. Because of Lemma 4, it is true that $p < p_u$ for $q_a > 0$ and $p_u < p$ for $q_a < 0$. Therefore,

$$\text{sign}[\partial_p f(p, q) - \partial_p f(p_u, q)] = -\,\text{sign}(q)$$

holds, which, by Equation (1.20), leads us to the conclusion that

$$\text{sign}(\partial^2_{p_u p_u} p) = -\,\text{sign}(q_a).$$

This had to be shown. □

The dependence of the pressure on the initial flow for fixed initial pressure and fixed length is displayed in Figure 1.5.

For the second order dependency on the flow we start with the more abstract setting of an parameter dependent ordinary differential equation $u_x(x) = s(x; u, z), u(0) = u_0$ with a parameter $z \in \mathbb{R}$. We show sufficient conditions only depending on derivatives of the source term $s$ under which the solution is strictly concave as a function of the parameter $k$. Then, we proceed to show that the result is applicable to the system (IsoStat). To obtain the result, a version of the Gronwall-Bellman inequality (see Teschl 2012, Lemma 2.7) is needed.

**Lemma 7** (Gronwall-Bellman Inequality)**.**
*Consider the real interval $[0, L]$. Let $\alpha$ be a real-valued function on $[0, L]$ and*

*let $\beta$ be nonnegative and continuous. Furthermore, let $u$ be a continuous function and let the negative part of $\alpha$ be integrable on every compact subinterval of $[0, L]$. Suppose $u$ satisfies*

$$u(x) \leq \alpha(x) + \int_0^x \beta(s)u(s)\,\mathrm{d}s, \quad for\ x \in [0, L].$$

*Then,*

$$u(x) \leq \alpha(x) + \int_0^x \alpha(s)\beta(s) \exp\left(\int_s^x \beta(r)\,\mathrm{d}r\right)\mathrm{d}s,\ for\ x \in [0, L]$$

*holds. If $\alpha$ is an increasing function on $[0, L]$, then*

$$u(x) \leq \alpha(x) \exp\left(\int_0^x \beta(s)\,\mathrm{d}s\right), \quad for\ x \in [0, L].$$

Note that only the sign of $\beta$ needs to be nonnegative. The functions $\alpha$ and $u$ do not have a sign restriction. Therefore an inequality for the second derivative with $\alpha < 0$ can be used to show the concavity of $u$.

**Lemma 8** (Sufficient Conditions for Concavity).
*Consider the parametric initial value problem*

$$\begin{cases} u_x(x, z) = s(x; u(x, z), z) & x \in (0, L], \\ u(0, z) = u_0, \end{cases} \qquad \text{(param-IVP)}$$

*for a source term $s$ that is continuous in $x$ and two times continuously differentiable in $(u, z)$, a parameter $z \in \mathbb{R}$ and an initial value $u_0 \in \mathbb{R}$. Suppose that $u$ is two times continuously differentiable with respect to $z$ and that*

*a) $\partial_u s(x; u(x, z), z) \geq 0$ for all $x \in [0, L]$*

*b) $\begin{pmatrix} \partial_z u(x, z) & 1 \end{pmatrix} \nabla^2 s(x; u(x, z), z) \begin{pmatrix} \partial_z u(x, z) \\ 1 \end{pmatrix} < 0$ for all $x \in (0, L]$,*

*where $\nabla^2 s$ denotes the Hessian of $s$ with respect to $u$ and $z$. Then*

$$\partial_{zz}^2 u(x, z) < 0 \quad for\ all\ x \in (0, L].$$

*Proof.* The function $s$ is two times continuously differentiable. This implies that $s$ is Lipschitz continuous with respect to $u$ on every compact interval. Because $u(\,\cdot\,, z)$ is a continuous function, the Weierstraß theorem implies

that the first argument of the source term can be considered on the compact interval

$$u(x, z) \in [\min_{y \in [0,L]} u(y, z), \max_{y \in [0,L]} u(y, z)].$$

Therefore, by the Picard-Lindelöff theorem (Teschl 2012, Theorem 2.2), the initial value problem (param-IVP) has a unique solution $u$.

We integrate the differential equation to obtain

$$u(x, z) - u(0, z) = \int_0^x s(r; u(r, z), z) \, dr.$$

The differentiation with respect to $z$ yields

$$\partial_z u(x, z) = \int_0^x \nabla s(r; u(r, z), z)^T \begin{pmatrix} \partial_z u(r, z) \\ 1 \end{pmatrix} dr, \qquad (1.23)$$

where $\nabla s(r; u, z) = \begin{pmatrix} \partial_u s(r; u, z) & \partial_z s(r; u, z) \end{pmatrix}^T$ denotes the gradient of $s$ with respect to $u$ and $z$. Differentiating once more leads to

$$\partial_{zz}^2 u(x, z) = \int_0^x \begin{pmatrix} \partial_z u(r, z) & 1 \end{pmatrix} \nabla^2 s(r; u(r, z), z) \begin{pmatrix} \partial_z u(r, z) \\ 1 \end{pmatrix}$$
$$+ \partial_u s(r; u(r, z), z) \partial_{zz}^2 u(r, z) \, dr. \qquad (1.24)$$

Define

$$\alpha(z) := L \max_{x \in [0,L]} \left\{ \begin{pmatrix} \partial_z u(x, z) & 1 \end{pmatrix} \nabla^2 s(x; u(x, z), z) \begin{pmatrix} \partial_z u(x, z) \\ 1 \end{pmatrix} \right\}.$$

and

$$\beta(r, z) := \partial_u s(r; u(r, z), z).$$

The maximum in the definition of $\alpha$ is attained, because $u(\cdot, z), \partial_z u(\cdot, z)$ and $\nabla^2 s(\cdot, z)$ are continuous functions by assumption and the interval $[0, L]$ is compact. Therefore, Equation (1.24) can be estimated by

$$\partial_{zz}^2 u(x, z) \le \alpha(z) + \int_0^x \beta(r, z) \partial_{zz}^2 u(r, z) \, dr \quad \text{for } x \in [0, L] \qquad (1.25)$$

By assumption a), we know that $\beta(r, z) \ge 0$ and by assumption b), we know that $\alpha(z) < 0$. We are now in the setting of Lemma 7 for the function $\partial_{zz}^2 u(\cdot, z)$, which gives us the desired estimate

$$\partial_{zz}^2 u(x, z) \le \alpha(z) \exp\left( \int_0^x \beta(r, z) \, dr \right) < 0 \quad \text{for } x \in (0, L].$$

$\square$

*Remark* 3 (Discussion of Assumption b) in Lemma 8). Assumption b) in Lemma 8 is weaker than the assumption of strict concavity of the source term $s$. Strict concavity of the source term would be equivalent to

$$y^T \nabla^2 s(x; u, z) y < 0 \quad \text{for all } y \in \mathbb{R}^2,$$

whereas we only require the condition for the direction $y = \begin{pmatrix} \partial_z u(x, z) & 1 \end{pmatrix}^T$.

Now, treating the flow $q_a$ as the parameter in Lemma 8, we obtain the concavity of the stationary pressure solution $p$ with respect to a positive initial flow $q_a$.

**Lemma 9** (Concavity of the Pressure with Respect to the Initial Flow). *For $x > 0$ and a positive flow $q_a > 0$, the pressure solution $p$ is* strictly concave *as a function of the initial flow $q_a$, i.e.,*

$$\partial^2_{q_a q_a} p < 0.$$

*Proof.* We show that the requirements of Lemma 8 are fulfilled for $u = p$, $u_0 = p_u$ and $z = q$. Using the definitions (1.2), (1.3) and (1.21) and omitting arguments for the sake of presentation, the source term is given by

$$s(x; p, q) := f(p, q)g(x; q) = -\frac{\theta(x)}{2} RTq|q|\frac{z(p)}{p}\left(1 - RT\frac{q^2}{p^2}\right)^{-1}$$
$$= -\frac{\theta(x)}{2}\operatorname{sign}(q)pz(p)h(\eta).$$

For the derivative with respect to $p$, we get, using (1.21) with the notation (1.22),

$$\partial_p s(x; p, q) = g(x; q)\partial_p f(p, q)$$
$$= -(RTq^2)^{-1}g(x; q)\left[2z(p)h(\eta)^2 + h(\eta)\right]$$
$$= \operatorname{sign}(q)\frac{\theta(x)}{2}\left[2z(p)h(\eta)^2 + h(\eta)\right]. \tag{1.26}$$

The sign of this term depends only on the sign of $q$, i.e.,

$$\operatorname{sign}(\partial_p s(x, p, q)) = \operatorname{sign}(q). \tag{1.27}$$

For $q \geq 0$ this is condition a) of Lemma 8. Differentiating the source term with respect to $q$ gives us

$$\partial_q s(x; p, q) = -\theta(x)\frac{z(p)}{p}\frac{RT|q|}{(1 - \eta)^2}. \tag{1.28}$$

For the second derivatives, we obtain

$$\partial_{pp}^2 s(x; p, q) = \mathrm{sign}(q) \frac{\theta(x)}{2} \left[ 2\alpha h(\eta)^2 + [4z(p)h(\eta) + 1]\frac{\partial_p \eta}{(1 - \eta)^2} \right] \quad (1.29\mathrm{a})$$

$$\partial_{qq}^2 s(x; p, q) = -\mathrm{sign}(q)\, \theta(x) RT \frac{z(p)}{p} \frac{1 + 3\eta}{1 - \eta} \quad (1.29\mathrm{b})$$

$$\partial_p \partial_q s(x; p, q) = \frac{\theta(x)}{|q|}(1 - \eta)^{-1} \left[ 4z(p)h(\eta)^2 + h(\eta) \right]. \quad (1.29\mathrm{c})$$

Condition b) of Lemma 8 reads

$$[\partial_{pp}^2 s(x; p, q)\partial_{q_a} p + 2\partial_p \partial_q s(x; p, q)]\partial_{q_a} p + \partial_{qq}^2 s(x; p, q) < 0. \quad (1.30)$$

By Lemma 5, we know that $\partial_{q_a} p < 0$. To prove that (1.30) indeed holds, we proceed to show

$$[\partial_{pp}^2 s(x; p, q)\partial_{q_a} p + 2\partial_p \partial_q s(x; p, q)] > 0 \quad \text{and} \quad \partial_{qq}^2 s(x; p, q) < 0. \quad (1.31)$$

It is easy to see from (1.29b) that the second inequality holds under Assumption 1 for $q > 0$ . Looking at the bracketed term in (1.29a), we notice due to $\alpha < 0$ and $\partial_p \eta = -2\eta/p < 0$ that

$$\mathrm{sign}(\partial_{pp}^2 s(x; p, q)) = -\mathrm{sign}(q)$$

and hence $\partial_{pp}^2 s(x; p, q) < 0$ for $q > 0$. By Equation (1.29c) it becomes obvious that

$$\partial_p \partial_q s(x; p, q) > 0.$$

This leads us to the conclusion that (1.31) holds, which implies (1.30). Therefore, Lemma 8 is applicable and

$$\partial_{q_a q_a}^2 p < 0$$

follows. □

After showing the concavity of the pressure with respect to the initial flow, we turn to the concavity with respect to the space variable.

**Lemma 10** (Concavity of the Pressure with Respect to the Space Variable). *Consider a piecewise differentiable friction factor $\theta : [0, L] \to \mathbb{R}_{>0}$ and denote the points, where it is differentiable by $A \subset [0, L]$. Furthermore, let $|\theta'(x)/\theta(x)|$ be small in the sense that*

$$\left| \frac{\theta'(x)}{\theta(x)} \right| < |\partial_p f(p, q)\, g(x; q)| \text{ for } x \in A, \quad (1.32)$$

*with the functions $f$ and $g$ as defined in (1.2) and (1.3).*

*Then, for $x \in A \setminus \{0\}$ and $q_a \neq 0$, the pressure solution $p$ fulfills*

$$\partial_{xx}^2 p < 0.$$

*If $A = [0, L]$ this is equivalent to $p$ being* strictly concave *on $(0, L)$.*

*Proof.* Differentiation of Equation (1.4) with use of the product rule, the chain rule and the fact that $\partial_x q = 0$, leads to the equation

$$\partial_{xx}^2 p = \partial_p f(p, q) \, g(x; q) \, \partial_x p + f(p, q) \, \partial_x g(x; q).$$

Inserting (1.4) and $\partial_x g(x; q) = g(x; q) \frac{\theta'(x)}{/} \theta(x)$ yields

$$\partial_{xx}^2 p = \left[ \partial_p f(p, q) \, g(x; q) + \frac{\theta'(x)}{\theta(x)} \right] f(p, q) g(x; q). \tag{1.33}$$

By (1.2) and (1.3) one can observe that

$$\mathrm{sign}(f(p, q)g(x; q)) = -\mathrm{sign}(q). \tag{1.34}$$

The expression

$$\partial_p f(p, q) \, g(x; q) + \frac{\theta'(x)}{\theta(x)}$$

has the sign of $\partial_p f(p, q) \, g(x; q)$ by the smallness assumption (1.32).

We proceed to determine the sign of $\partial_p f(p, q) \, g(x; q)$. The derivative $f_p(p, q)$ was calculated in (1.21) and the notation (1.22) together with the use of the squared Mach number $\eta = RT \frac{q^2}{p^2}$ allows to write it as

$$\partial_p f(p, q) = -(RTq^2)^{-1} \left[ 2z(p)h(\eta)^2 + h(\eta) \right].$$

Since the function $h$ has positive values for $\eta \in (0, 1)$, we see that

$$\partial_p f(p, q) < 0.$$

Furthermore, the equation

$$\mathrm{sign}(g(x; q)) = -\mathrm{sign}(q)$$

holds. Hence,

$$\mathrm{sign}[\partial_p f(p, q)g(x; q)] = \mathrm{sign}(q). \tag{1.35}$$

Combining (1.34), (1.35), (1.32) and (1.33) yields the result. $\qquad \square$

*Remark* 4 (Assumption for the Friction Factor). The assumption (1.32) for the friction factor in Lemma 10 is certainly fulfilled for a constant friction factor $\theta(x) = \theta_{\mathrm{const}} > 0$ for all $x \in [0, L]$. In this case, the set of points $A$, where the function is differentiable, is the whole interval $[0, L]$.

## 1.3. Existence and Uniqueness for Flow Problems on Finite Graphs

The next goal is to show the existence and uniqueness of the stationary isothermal Euler equations coupled by Kirchhoff type conditions on arbitrary finite graphs. We can show this for a broader class of problems on graphs that contains the gas flow problem as a special case. The results for gas networks were presented in Gugat, Schultz, and Wintergerst 2018, Section 4.4. Consider the following problem:

Let a connected finite graph $\mathcal{G} = (\mathcal{V}, \mathcal{A})$ with nodes $\mathcal{V}$ and arcs $\mathcal{A}$ and be given. The node arc incidence matrix $A \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{A}|}$ contains the information, which nodes are connected by which arcs. We take the freedom to use the nodes $u$ and arcs $a$ as indices. The convention

$$A_{ua} = \begin{cases} -1, & \text{if } a = (u, v) \text{ for a } v \in \mathcal{V} \\ +1, & \text{if } a = (v, u) \text{ for a } v \in \mathcal{V} \\ 0, & \text{if the arc } a \text{ is not connected to the node } u \end{cases}$$

is used. In a node $u \in \mathcal{V}$ denote the set of ingoing arcs, the set of outgoing arcs and the set of arcs incident to $u$ by

$$\delta_{\text{in}}(u) := \{a \in \mathcal{A} \mid A_{ua} = 1\} \tag{1.36a}$$

$$\delta_{\text{out}}(u) := \{a \in \mathcal{A} \mid A_{ua} = -1\} \tag{1.36b}$$

$$\delta(u) := \delta_{\text{in}} \cup \delta_{\text{out}}. \tag{1.36c}$$

If one assigns a flow value $q_a$ to each arc $a$ and names the vector of all flow values through the arcs $q \in \mathbb{R}^{|\mathcal{A}|}$, the Kirchhoff coupling condition can be written as

$$Aq = q^{\text{out}}.$$

It was formulated for electric current in Kirchhoff 1847 and models the conservation of mass in each node.

Each arc $a = (u, v) \in \mathcal{A}$ corresponds to an interval $[0, L_a]$. The state values $y_a(0)$ and $y_a(L_a)$ at both ends are assigned to each arc $a \in \mathcal{A}$. On each arc, continuous functions

$$f_a : [0, L_a] \times \mathbb{R} \times \mathbb{R} \to \mathbb{R}, \quad (x, y, q) \mapsto f_a(x; y, q)$$

that couple the state values at both ends of an arc are given. The functions fulfill $f_a(0, y, q) = y$. The coupling on an arc $a = (u, v)$ reads

$$y_a(L_a) = f_a(L_a; y_a(0), q_a).$$

In each node $u \in \mathcal{V}$, consider continuous functions

$$h_u : \mathbb{R} \times \mathbb{R} \to \mathbb{R}, \quad (y, q) \mapsto h_u(y, q)$$

that couple state values of all pairs of arcs $a, b$ incident to $u$ via the set of equations

$$h_u(y_a(x_a(u)), q_a) = h_u(y_u(x_b(u)), q_b),$$

where

$$x_a(u) = \begin{cases} 0, & \text{if } A_{ua} = -1 \\ L_a, & \text{if } A_{ua} = +1. \end{cases}$$

Consider the space of balanced outflows

$$\mathcal{Q} := \{q^{\text{out}} \in \mathbb{R}^{|\mathcal{V}|} \mid \sum\nolimits_{u \in \mathcal{V}} q_u^{\text{out}} = 0\} \qquad (1.37)$$

and choose a vector

$$q^{\text{out}} \in \mathcal{Q}.$$

Furthermore, prescribe a value $y_r \in \mathbb{R}$ in a boundary node $r$, that is a node, which is only connected to one other node by an arc $a = (r, u)$ (the direction of the arc is without loss of generality). This means

$$y_a(0) = y_r. \qquad \text{(Initial State)}$$

The question of interest is, whether the system

$$\begin{align} Aq &= q^{\text{out}}, & \text{(Kirchhoff)} \\ y_a(L_a) &= f_a(L_a; y_a(0), q_a), \quad \forall a = (u, v) \in \mathcal{A}, & \text{(Arc Coupling)} \\ h_u(y_a(x_a(u)), q_a) &= h_u(y_b(x_b(u)), q_b), & \text{(Node Coupling)} \\ & \qquad \forall u \in \mathcal{V} \text{ and all } a, b \in \mathcal{A} \text{ incident to } u, \end{align}$$

has a unique solution.

*Remark* 5 (Nontopological Arcs). Since the arc direction only determines the sign of the flow, the condition (Arc Coupling) is always to be understood in both directions, especially

$$y_a(0) = f_a(L_a; y_a(L_a), -q_a) \quad \forall a = (u, v) \in \mathcal{A}$$

should also hold.

The Kirchhoff condition is illustrated in Figure 1.6. The arc coupling condition is depicted in Figure 1.7 and the node coupling condition can be seen in Figure 1.8.

First, we derive important properties of the flow system described by the Kirchhoff conditions, while the states $y_a$ on $a \in \mathcal{A}$ are ignored for the moment.

Figure 1.6.: An example network with the flow variables on the edges, the outflow values at the nodes and the Kirchhoff coupling conditions at each node.

Figure 1.7.: The arc coupling condition illustrated on two connected arcs $a = (r, u)$ and $b = (u, v)$.

$$h_u(y_a(L_a), q_a) = h_u(y_b(0), q_b) = h_u(y_c(0), q_c)$$



Figure 1.8.: An example network with the state variables on edge $a$ and $e$ and the node coupling conditions in the nodes $u, v$ and $w$.

**Lemma 11** (Solvability of the Flow System)**.**
*Let $A \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{A}|}$ be the node arc incidence matrix of a finite connected graph $\mathcal{G} = (\mathcal{V}, \mathcal{A})$. The linear system*

$$Aq = q^{out}$$

*has a solution if and only if*

$$q^{out} \in \mathcal{Q}.$$

*Proof.* We need to show that $\mathcal{Q} = \text{im}(A)$. Clearly, by the definition of $\mathcal{Q}$; see Equation (1.37), we have

$$\mathcal{Q} = \text{span}\big(\{\mathbf{1}_{|\mathcal{V}|}\}\big)^{\perp},$$

where $\mathbf{1}_{|\mathcal{V}|} \in \mathbb{R}^{|\mathcal{V}|}$ is the vector of only ones. The null space of $A^T$ is given by

$$\ker(A^T) = \text{span}\big(\{\mathbf{1}_{|\mathcal{V}|}\}\big). \tag{1.38}$$

This can be seen as follows: Each arc $a = (u, v) \in \mathcal{A}$ has exactly two incident nodes in $\mathcal{V}$. Consequently, each row in $A$ has exactly two nonzero

entries; one being $-1$ and the other one being $1$. Hence, a vector $q^{\text{ker}}$ with

$$q_u^{\text{ker}} = q_v^{\text{ker}} = \mu \in \mathbb{R}, \tag{1.39}$$

would ensure that the $a$-th row of the system $A^T q^{\text{ker}} = 0$ is fulfilled. However, as the graph is connected, at least one of the nodes $u$ and $v$ is adjacent to another node $w$. Without loss of generality, we assume that $v$ and $w$ are connected by the arc $b = (v, w)$. Then, by the $b$-th row of $Aq^{\text{ker}} = 0$, it follows that $q_w^{\text{ker}} = \mu$. Proceeding inductively leads to the conclusion that $q^{\text{ker}} = \mu \mathbf{1}_{|\mathcal{V}|}$. As $\mu$ in Equation (1.39) was arbitrary and the construction allows no vectors not having the form $q^{\text{ker}} = \mu \mathbf{1}_{|\mathcal{V}|}$, this shows (1.38). Putting things together, yields

$$\mathcal{Q} = \text{span}\big(\{\mathbf{1}_{|\mathcal{V}|}\}\big)^\perp = \ker(A^T)^\perp = \text{im}(A), \tag{1.40}$$

which completes the proof. $\qquad\square$

For the following results, we need the definition of a separator. The direction of arcs is only relevant for the flow direction but not for the topology of the graph. This means if the nodes $u$ and $v$ are incident to $a = (u, v)$, then the node $u$ is not just connected to $v$ but the node $v$ is also connected to $u$.

**Definition 1** (Separator)**.** Let a connected finite graph $\mathcal{G} = (\mathcal{V}, \mathcal{A})$ with nodes $\mathcal{V}$ and edges $\mathcal{A}$ be given. A set of edges $\mathcal{A}_{\mathcal{S}} \subsetneq \mathcal{A}$ such that the removal of $\mathcal{A}_{\mathcal{S}}$ from $\mathcal{A}$ separates the graph $\mathcal{G}$ into two distinct connected subgraphs

$$\mathcal{G}_l = (\mathcal{V}_l, \mathcal{A}_l) \quad \text{and} \quad \mathcal{G}_r = (\mathcal{V}_r, \mathcal{A}_r)$$

in a way that

$$\mathcal{V} = \mathcal{V}_l \cup \mathcal{V}_r \quad \text{and} \quad \mathcal{A} = \mathcal{A}_l \cup \mathcal{A}_{\mathcal{S}} \cup \mathcal{A}_r$$

is called a *separator*.

The nodes that are incident to an edge $a \in \mathcal{A}_{\mathcal{S}}$ are denoted by $\mathcal{V}_{\mathcal{S}}$ and the graph $\mathcal{G}_{\mathcal{S}} = (\mathcal{V}_{\mathcal{S}}, \mathcal{A}_{\mathcal{S}})$ is called *separating graph*.

The incidence matrices of $\mathcal{G}_l, \mathcal{G}_r$ and $\mathcal{G}_{\mathcal{S}}$ are denoted by $A^l, A^r$ and $A^{\mathcal{S}}$.

Figure 1.9 shows a separator for the graph used in Figures 1.6 and 1.8.

Figure 1.9.: The separator $\mathcal{A}_{\mathcal{S}} = \{b, c\}$ (dashed) seperates the graph $\mathcal{G} = (\{r, u, v, w, s\}, \{a, b, c, d, e\})$ into the two disjoint subgraphs $\mathcal{G}_l = (\{r, v\}, \{a\})$ and $\mathcal{G}_r = (\{v, w, s\}, \{d, e\})$.

**Lemma 12** (Flow Behavior on Subgraphs).
*Let a finite connected graph $\mathcal{G} = (\mathcal{V}, \mathcal{A})$ and two distinct connected subgraphs $\mathcal{G}_l = (\mathcal{V}_l, \mathcal{A}_l)$ and $\mathcal{G}_r = (\mathcal{V}_r, \mathcal{A}_r)$ separated by the separating graph $\mathcal{G}_{\mathcal{S}} = (\mathcal{V}_{\mathcal{S}}, \mathcal{A}_{\mathcal{S}})$ be given. Let $A$ be the incidence matrix of $G$ and let the outflow vector $q^{out} \in \mathcal{Q}$ be given. Consider one specific solution $q \in \mathbb{R}^{|\mathcal{A}|}$ of the system*

$$Aq = q^{out}.$$

*Denote the restrictions $q_{\mathcal{S}}^{out} := (q_u^{out})_{u \in \mathcal{V}_{\mathcal{S}}}$ and $q_{\mathcal{S}} := (q_a)_{a \in \mathcal{A}_{\mathcal{S}}}$, define*

$$\tilde{q}^{out} := q_{\mathcal{S}}^{out} - A^{\mathcal{S}} q_{\mathcal{S}}, \quad \tilde{q}^{out} \in \mathbb{R}^{|\mathcal{V}_{\mathcal{S}}|}$$

*and set the new outflows to*

$$f_u := \begin{cases} \tilde{q}_u^{out} & \text{if } u \in \mathcal{V}_{\mathcal{S}}, \\ q_u^{out} & \text{if } u \notin \mathcal{V}_{\mathcal{S}}. \end{cases}$$

*Then, the restrictions $q^l := (q_u)_{u \in \mathcal{V}_l}$ and $q^r := (q_u)_{u \in \mathcal{V}_r}$ of the flow solution are solutions of the two linear systems*

$$A^l q^l = f^l \quad \text{and} \quad A^r q^r = f^r,$$

*with the restrictions $f^l := (f_u)_{u \in \mathcal{V}_l}$ and $f^r := (f_u)_{u \in \mathcal{V}_r}$ as right hand sides. Furthermore, the restrictions $f^l$ and $f^r$ of the new outflow vector sum to zero on the subgraphs $\mathcal{G}_l$ and $\mathcal{G}_r$, i.e.,*

$$\sum_{u \in \mathcal{V}_l} f_u^l = \sum_{v \in \mathcal{V}_r} f_v^r = 0.$$

*Proof.* Define the node arc incidence matrices with respect to the graphs $(\mathcal{V}, \mathcal{A}_l), (\mathcal{V}, \mathcal{A}_{\mathcal{S}})$ and $(\mathcal{V}, \mathcal{A}_r)$ as $B^l, B^{\mathcal{S}}$ and $B^r$ (these matrices have additional rows containing zeros compared to $A^l, A^{\mathcal{S}}$ and $A^r$). With appropriate ordering of the arcs, the vector $q$ and the matrix $A$, can be written as

$$q = \begin{pmatrix} q^l \\ q_{\mathcal{S}} \\ q^r \end{pmatrix} \qquad A = \begin{bmatrix} B^l & B^{\mathcal{S}} & B^r \end{bmatrix}.$$

This leads to the linear system

$$Aq = B^l q^l + B^{\mathcal{S}} q_{\mathcal{S}} + B^r q^r = q^{\text{out}}.$$

The definitions of $f$ and $\tilde{q}^{\text{out}}$ imply $f = q^{\text{out}} - B^{\mathcal{S}} q_{\mathcal{S}}$ and thus

$$B^l q^l + B^r q^r = f.$$

The rows of $B^l$ that correspond to nodes in $\mathcal{V}_r$ and the rows of $B^r$ that correspond to nodes in $\mathcal{V}_l$ contain only zeros. Hence, the first claim of the Lemma

$$A^l q^l = f^l \quad \text{and} \quad A^r q^r = f^r$$

holds.

For the second part observe that the nodes incident to each separating arc $a = (u, v) \in \mathcal{A}_{\mathcal{S}}$ fulfill

$$\tilde{q}_u^{\text{out}} + \tilde{q}_v^{\text{out}} = q_u^{\text{out}} - A_{ua}^{\mathcal{S}} q_a + q_v^{\text{out}} - A_{va}^{\mathcal{S}} q_a = q_u^{\text{out}} + q_v^{\text{out}}, \qquad (1.41)$$

because $A_{ua}^{\mathcal{S}}$ and $A_{va}^{\mathcal{S}}$ have opposite signs. Due to $q^{\text{out}} \in \mathcal{Q}$, Equation (1.41) and the definition of $f$, we have

$$0 = \sum_{u \in \mathcal{V}} q_u^{\text{out}} = \sum_{u \in \mathcal{V}_{\mathcal{S}}} \tilde{q}_u^{\text{out}} + \sum_{v \notin \mathcal{V}_{\mathcal{S}}} q_v^{\text{out}} = \sum_{u \in \mathcal{V}_l} f_u^l + \sum_{v \in \mathcal{V}_r} f_v^r. \qquad (1.42)$$

Suppose $\sum_{u \in \mathcal{V}_l} f_u^l \neq 0$. Then, Lemma 11 applied to the graph $\mathcal{G}_l = (\mathcal{V}_l, \mathcal{A}_l)$ with incidence matrix $A^l$ states that the flow system $A^l \tilde{q} = f^l$ is not solvable.

This is a contradiction, as we have already shown that $q^l$ solves $A^l q^l = f^l$. Therefore,

$$\sum_{u \in \mathcal{V}_l} f^l = 0$$

and subsequently, by Equation (1.42), we obtain

$$\sum_{v \in \mathcal{V}_r} f^r = 0. \qquad \qquad \square$$

Figure 1.10.: The new outflows $f^l$ of Example 1. The restrictions of the specific solution $q$ (numbers at the edges) are solutions of the flow problem on the subgraphs $\mathcal{G}_l$ and $\mathcal{G}_r$ with right hand sides $f^l$ and $f^r$.

**Example 1.** Consider the graphs

$$\mathcal{G} = (\{r, u, v, w, s\}, \{a, b, c, d, e\}), \qquad \mathcal{G}_l = (\{r, v\}, \{a\}),$$
$$\mathcal{G}_\mathcal{S} = (\{u, v, w\}, \{b, c\}) \qquad \text{and} \qquad \mathcal{G}_r = (\{v, w, s\}, \{d, e\})$$

shown in Figure 1.9. The incidence matrices $A, A^l, A^{\mathcal{S}}$ and $A^r$ (with node and arc indices for orientation) are given by

| $A$ | $a$ | $b$ | $c$ | $d$ | $e$ |
|---|---|---|---|---|---|
| $r$ | $-1$ | $0$ | $0$ | $0$ | $0$ |
| $u$ | $+1$ | $-1$ | $-1$ | $0$ | $0$ |
| $v$ | $0$ | $+1$ | $0$ | $-1$ | $0$ |
| $w$ | $0$ | $0$ | $+1$ | $+1$ | $-1$ |
| $s$ | $0$ | $0$ | $0$ | $0$ | $+1$ |

| $A^l$ | $a$ |
|---|---|
| $r$ | $-1$ |
| $u$ | $+1$ |

| $A^{\mathcal{S}}$ | $b$ | $c$ |
|---|---|---|
| $u$ | $-1$ | $-1$ |
| $v$ | $+1$ | $0$ |
| $w$ | $0$ | $+1$ |

| $A^r$ | $d$ | $e$ |
|---|---|---|
| $v$ | $-1$ | $0$ |
| $w$ | $+1$ | $-1$ |
| $s$ | $0$ | $+1.$ |

Consider the outflow vector

$$(q^{\text{out}})^T = (-1,\ 0,\ 0,\ 0,\ 1)^T$$

and the specific solution

$$q^T = (1,\ 0.5,\ 0.5,\ 0.5,\ 1)^T$$

of the system $Aq = q^{\text{out}}$. The new outflows for the nodes of the separating graph are

$$\tilde{q}_u^{\text{out}} = q_u^{\text{out}} - A_{ub}^{\mathcal{S}} q_b - A_{uc}^{\mathcal{S}} q_c = 1,$$
$$\tilde{q}_v^{\text{out}} = q_v^{\text{out}} - A_{vb}^{\mathcal{S}} q_b \qquad\quad = -0.5,$$
$$\tilde{q}_w^{\text{out}} = q_w^{\text{out}} - A_{wc}^{\mathcal{S}} q_c \qquad\quad = -0.5$$

and thus

$$(f)^T = (q_r^{\text{out}}, \tilde{q}_u^{\text{out}}, \tilde{q}_v^{\text{out}}, \tilde{q}_w^{\text{out}}, q_s^{\text{out}})^T = (-1,\ 1,\ -0.5,\ -0.5,\ 1)^T.$$

The flow values are shown in Figure 1.10. Lemma 12 states that $q^l = 1$ solves $A^l q^l = f^l$ and $(q^r)^T = (0.5,\ 1)^T$ solves $A^r q^r = f^r$, with the right hand sides

$$(f^l)^T = (-1,\ 1)^T \quad \text{and} \quad (f^r)^T = (-0.5,\ -0.5,\ 1)^T.$$

Let us verify this:

$$A^l q^l = \begin{bmatrix} -1 \\ +1 \end{bmatrix} 1 = \begin{pmatrix} -1 \\ 1 \end{pmatrix} = f^l \quad \text{and}$$

$$A^r q^r = \begin{bmatrix} -1 & 0 \\ +1 & -1 \\ 0 & +1 \end{bmatrix} \begin{pmatrix} 0.5 \\ 1 \end{pmatrix} = \begin{pmatrix} -0.5 \\ -0.5 \\ 1 \end{pmatrix} = f^r.$$

We recognize that the first part of Lemma 12 indeed holds. Furthermore,

$$f_r^l + f_u^l = f_v^r + f_w^r + f_s^r = 0$$

validates the second part of the Lemma.

Lemma 12 leads to two results that are essential for the proof of Theorem 2. The first one describes a monotonicity property of the network: If the outflow one node increases and the outflow in another node decreases, there is also a path between the two nodes with increased flow. This path may change for different prescribed values yet Lemma 13 ensures its existence. The second result states that sufficiently high inflow in one node and equal outflow in another node imply the existence of a path of positive flow between both nodes.

**Lemma 13** (Path of Increasing Flow)**.**
*On the graph $\mathcal{G} = (\mathcal{V}, \mathcal{A})$, consider two distinct nodes $u, v \in \mathcal{V}$. Let the outflow vector $q^{out}$ be given and a positive scalar $\lambda > 0$ be given. Define a new outflow vector by*

$$q_\lambda^{out} = q^{out} - \lambda e_u + \lambda e_v,$$

*where $e_u, e_v$ are the corresponding unit vectors. This means that the ingoing flow in $u$ and the outgoing flow in $v$ increases. Let $q$ be a specific solution of $Aq = q^{out}$ and let $q^\lambda$ be a specific solution of $Aq^\lambda = q_\lambda^{out}$.*
*Then, there exists a path $\mathcal{P} \subset \mathcal{A}$ between $u$ and $v$ (without loss of generality we assume it is directed from $u$ to $v$; see Remark 6) such that flow increases on each arc $a$ of this path in the sense that*

$$q_a^\lambda \geq q_a \quad \forall a \in \mathcal{P}.$$

*Additionally, the inequality*

$$q_a^\lambda > q_a$$

*holds for at least one $a \in \mathcal{P}$. The choice of the path $\mathcal{P}$ may depend on the value of $\lambda$.*

*Proof.* Suppose no such path exists. Then, there is a separator

$$\mathcal{G}_\mathcal{S} = (\mathcal{V}_\mathcal{S}, \mathcal{A}_\mathcal{S})$$

that separates the graph $\mathcal{G}$ into the two distinct subgraphs

$$\mathcal{G}_l = (\mathcal{V}_l, \mathcal{A}_l) \quad \text{and} \quad \mathcal{G}_r = (\mathcal{V}_r, \mathcal{A}_r)$$

such that the claim for the flow values does not hold on $\mathcal{A}_\mathcal{S}$. Without loss of generality, we assume all arcs $a \in \mathcal{A}_\mathcal{S}$ point from a node in $\mathcal{V}_l$ to a node in $\mathcal{V}_r$. Furthermore, let $u \in \mathcal{V}_l$ and $v \in \mathcal{V}_r$. The assumption for the proof by contradiction then reads

$$q_a^\lambda < q_a \quad \forall a \in \mathcal{A}_\mathcal{S}. \tag{1.43}$$

Define $f \in \mathbb{R}^{|\mathcal{V}|}$ and $f^\lambda \in \mathbb{R}^{|\mathcal{V}|}$ analogously to Lemma 12, where the definition of $f^\lambda$ uses $q_\lambda^{\text{out}}$ instead of $q^{\text{out}}$. By Lemma 12, we have

$$0 = \sum_{w \in \mathcal{V}_l} f_w = \sum_{w \in \mathcal{V}_l \cap \mathcal{V}_\mathcal{S}} f_w + \sum_{w \in \mathcal{V}_l \backslash \mathcal{V}_\mathcal{S}} f_w = \sum_{w \in \mathcal{V}_l \cap \mathcal{V}_\mathcal{S}} \tilde{q}_w^{\text{out}} + \sum_{w \in \mathcal{V}_l \backslash \mathcal{V}_\mathcal{S}} q_w^{\text{out}}$$

(by the definition of $\tilde{q}^{\text{out}}$)

$$= \sum_{w \in \mathcal{V}_l \cap \mathcal{V}_\mathcal{S}} \Big[ q_w^{\text{out}} - \sum_{a \in \delta(w)} A_{wa} q_a \Big] + \sum_{w \in \mathcal{V}_l \backslash \mathcal{V}_\mathcal{S}} q_w^{\text{out}}$$

$$= q_u^{\text{out}} + \sum_{\substack{w \in \mathcal{V}_l \cap \mathcal{V}_\mathcal{S} \\ w \neq u}} \Big[ q_w^{\text{out}} - \sum_{a \in \delta(w)} A_{wa} q_a \Big] + \sum_{\substack{w \in \mathcal{V}_l \backslash \mathcal{V}_\mathcal{S} \\ w \neq u}} q_w^{\text{out}}$$

(as $A_{wa} = -1$ and the Inequality (1.43) states $q_a^\lambda < q_a$ for $a \in \mathcal{A}_\mathcal{S}$)

$$> q_u^{\text{out}} + \sum_{\substack{w \in \mathcal{V}_l \cap \mathcal{V}_\mathcal{S} \\ w \neq u}} \Big[ q_w^{\text{out}} + \sum_{a \in \delta(w)} q_a^\lambda \Big] + \sum_{\substack{w \in \mathcal{V}_l \backslash \mathcal{V}_\mathcal{S} \\ w \neq u}} q_w^{\text{out}}$$

(remember $(q_\lambda^{\text{out}})_u = q_u^{\text{out}} - \lambda$ and the definition of $f^\lambda$)

$$> (q_\lambda^{\text{out}})_u + \sum_{\substack{w \in \mathcal{V}_l \cap \mathcal{V}_\mathcal{S} \\ w \neq u}} f_w^\lambda + \sum_{\substack{w \in \mathcal{V}_l \backslash \mathcal{V}_\mathcal{S} \\ w \neq u}} q_w^{\text{out}} = \sum_{w \in \mathcal{V}_l} f^\lambda = 0,$$

where the last equality is due to Lemma 12 applied to the outflow $q_\lambda^{\text{out}}$. This is a contradiction and hence proves that there exists a path $\mathcal{P}$ from $u$ to $v$ such that

$$q_a^\lambda \geq q_a \quad \forall a \in \mathcal{P}. \tag{1.44}$$

The second claim is clear if either $u$ or $v$ are simple nodes, i.e., $|\delta(u)| = 1$ or $|\delta(v)| = 1$. Then, the one incident edge $a$ has the flow $q_a^\lambda = q_a + \lambda > q_a$. In the case, where $|\delta(u)| > 1$ and $|\delta(v)| > 1$, we proceed inductively. If Inequality (1.44) does not hold strictly for any $a \in \mathcal{P}$, this implies

$$q_a^\lambda = q_a \quad \forall a \in \mathcal{P}.$$

Denote the nodes $w$ on the path $\mathcal{P}$ with $|\delta(w)| = 2$ by $\mathcal{V}_{\mathcal{P}}$. Define

$$n := \min\{|\delta(u)|, |\delta(v)|\}, \quad \mathcal{G}^n := \mathcal{G} \quad \text{and} \quad q_n^{\text{out}} := q^{\text{out}}.$$

The removal of the edges in $\mathcal{P}$ and the nodes in $\mathcal{V}_{\mathcal{P}}$ leads to the graph

$$\mathcal{G}^{n-1} = (\mathcal{V}^{n-1}, \mathcal{A}^{n-1}) := (\mathcal{V} \setminus \mathcal{V}_{\mathcal{P}}, \mathcal{A} \setminus \mathcal{P})$$

with incidence matrix $A^{n-1}$. The restriction $q^{n-1} = (q_a)_{a \in \mathcal{A}^{n-1}}$ solves the system

$$A^{n-1}q^{n-1} = q_{n-1}^{\text{out}},$$

where

$$(q_{n-1}^{\text{out}})_w := (q_n^{\text{out}})_w - \sum_{a \in \delta(w) \cap \mathcal{P}} A_{wa} q_a$$

and, with analogous notation, $q_\lambda^{n-1}$ solves

$$A^{n-1}q_\lambda^{n-1} = q_{\lambda,n-1}^{\text{out}}.$$

Reapplying the result above for the graph $\mathcal{G}^{n-1}$ yields that there is a path $\mathcal{P}^{n-1}$ with

$$q_a^\lambda \geq q_a \quad \forall a \in \mathcal{P}^{n-1}.$$

Assume again that

$$q_a^\lambda = q_a \quad \forall a \in \mathcal{P}^{n-1}$$

and remove the path and the nodes $\mathcal{V}_{\mathcal{P}}^{n-1}$, defined similarly as above, from the graph. After at most $n-1$ steps, we arrive at a graph $\mathcal{G}^k$ and a path $\mathcal{P}^k$ with

$$q_a^\lambda \geq q_a \quad \forall a \in \mathcal{P}^k$$

such that one arc $a = (w, s)$ is a separator (if this is possible with a single arc, it is also called a *bridge* in graph theory) separating $\mathcal{G}_l^k$ from $\mathcal{G}_r^k$ such that $\{u, w\} \subset \mathcal{G}_l^k$ and $\{v, s\} \subset \mathcal{G}_r^k$. By Lemma 12, this implies that

$$q_a^\lambda = (q_{\lambda,k}^{\text{out}})_w = -(q_{\lambda,k}^{\text{out}})_u - \sum_{\substack{t \in \mathcal{V}_l^k \\ t \neq u}} (q_{\lambda,k}^{\text{out}})_t$$

(because $(q_{\lambda,k}^{\text{out}})_u = q_k^{\text{out}} - \lambda$)

$$= \lambda - (q_k^{\text{out}})_u - \sum_{\substack{t \in \mathcal{V}_l^k \\ t \neq u}} (q_k^{\text{out}})_t = q_a + \lambda.$$

This shows the second claim of the Lemma. $\qquad\square$

*Remark* 6 (No Loss of Generality). If one allows arcs $a = (w, s)$ in $\mathcal{P}$ with $A_{sa} = -1$, then the result of Lemma 13 reads

$$A_{sa} q_a^\lambda \geq A_{sa} q_a \quad \forall a = (w, s) \in \mathcal{P}.$$

In the proof, instead of (1.43), the inequality

$$A_{sa} q_a^\lambda < A_{sa} q_a \quad \forall a = (w, s) \in \mathcal{A}_{\mathcal{S}},$$

is used. The adjustments to the subsequent steps are straightforward.

The result from Lemma 13 can be slightly generalized in the sense that we can add an additional reference node. The proof is analogous with a bit more notation due to the additional node.

**Corollary 1** (Path of Increasing Flow to a Reference Node)**.**
*Let the assumptions of Lemma 13 hold and consider an additional reference node $r \in \mathcal{V} \setminus \{u, v\}$.*

*Then there exists a path $\mathcal{P} \subset \mathcal{A}$ between $u$ and $r$ (without loss of generality we assume it is directed from $u$ to $r$) such that flow increases on each arc $a$ of this path in the sense that*

$$q_a^\lambda \geq q_a \quad \forall a \in \mathcal{P}.$$

*Additionally, the inequality*

$$q_a^\lambda > q_a$$

*holds for at least one $a \in \mathcal{P}$.*

*Proof.* The proof of Lemma 13 can be directly adapted by separating $u$ from $r$ instead of $u$ from $v$. $\qquad\square$

A result similar to Lemma 13 is developed in Lemma 14. It states that if a high outflow in one node and a high inflow in another node is prescribed, then there must be a path of unidirectional flow between both nodes.

**Lemma 14** (Path of Unidirectional Flow)**.**
*On the graph $\mathcal{G} = (\mathcal{V}, \mathcal{A})$, consider two distinct nodes $u, v \in \mathcal{V}$. Assume that the outflow values $q_w^{out}$ are prescribed in all nodes $w \in \mathcal{V} \setminus \{u, v\}$ except in $u$ and $v$. Set*

$$q_u^{out} := -\sum_{w \in \mathcal{V} \setminus \{u, v\}} |q_w^{out}| \quad and \quad q_v^{out} := \sum_{w \in \mathcal{V} \setminus \{u, v\}} |q_w^{out}|$$

*and let $q$ be a specific solution of $Aq = q^{out}$.*

*Then, there exists a path $\mathcal{P}$ between $u$ and $v$ (without loss of generality we assume the arcs on the path are directed from $u$ to $v$) such that*

$$q_a \geq 0 \quad \forall a \in \mathcal{P}.$$

*The inequality is strict for at least one $a \in \mathcal{P}$ if not all outflows are zero.*

*Proof.* The idea of the proof is similar to the proof of Lemma 13. We assume the contrary of the claim. This means, there is a separator $\mathcal{G}_\mathcal{S} = (\mathcal{V}_\mathcal{S}, \mathcal{A}_\mathcal{S})$ that separates the graph $\mathcal{G}$ into the two distinct subgraphs $\mathcal{G}_l = (\mathcal{V}_l, \mathcal{A}_l)$ and $\mathcal{G}_r = (\mathcal{V}_r, \mathcal{A}_r)$ such that

$$q_a < 0 \quad \forall a \in \mathcal{A}_\mathcal{S}.$$

Here, we assumed that all arcs in $\mathcal{A}_\mathcal{S}$ are directed from a node in $w \in \mathcal{V}_l \cap \mathcal{V}_\mathcal{S}$ to a node $s \in \mathcal{V}_r \cap \mathcal{V}_\mathcal{S}$ and that $u$ lies in $\mathcal{V}_l$ and $v$ in $\mathcal{V}_r$. Define $f^l$ and $f^r$ as in Lemma 12. By Lemma 12, we have

$$0 = \sum_{w \in \mathcal{V}_l} f_w^l = \sum_{w \in \mathcal{V}_l \cap \mathcal{V}_\mathcal{S}} f_w + \sum_{w \in \mathcal{V}_l \backslash \mathcal{V}_\mathcal{S}} f_w = \sum_{w \in \mathcal{V}_l \cap \mathcal{V}_\mathcal{S}} \tilde{q}_w^{\mathrm{out}} + \sum_{w \in \mathcal{V}_l \backslash \mathcal{V}_\mathcal{S}} q_w^{\mathrm{out}}$$

(by the definition of $\tilde{q}^{\mathrm{out}}$)

$$= \sum_{w \in \mathcal{V}_l \cap \mathcal{V}_\mathcal{S}} \left[ q_w^{\mathrm{out}} - \sum_{a \in \delta(w)} A_{wa} q_a \right] + \sum_{w \in \mathcal{V}_l \backslash \mathcal{V}_\mathcal{S}} q_w^{\mathrm{out}}$$

(as $A_{wa} = -1$ and $q_a < 0$ for $a \in \mathcal{A}_\mathcal{S}$)

$$< \sum_{w \in \mathcal{V}_l \cap \mathcal{V}_\mathcal{S}} q_w^{\mathrm{out}} + \sum_{w \in \mathcal{V}_l \backslash \mathcal{V}_\mathcal{S}} q_w^{\mathrm{out}} = q_u^{\mathrm{out}} + \sum_{\substack{w \in \mathcal{V}_l \cap \mathcal{V}_\mathcal{S} \\ w \neq u}} q_w^{\mathrm{out}} + \sum_{\substack{w \in \mathcal{V}_l \backslash \mathcal{V}_\mathcal{S} \\ w \neq u}} q_w^{\mathrm{out}}$$

(by the definition of $q_u^{\mathrm{out}}$)

$$= -\sum_{w \in \mathcal{V} \backslash \{u,v\}} |q_w^{\mathrm{out}}| + \sum_{\substack{w \in \mathcal{V}_l \cap \mathcal{V}_\mathcal{S} \\ w \neq u}} q_w^{\mathrm{out}} + \sum_{\substack{w \in \mathcal{V}_l \backslash \mathcal{V}_\mathcal{S} \\ w \neq u}} q_w^{\mathrm{out}} \leq 0.$$

This is a contradiction and thus the first part of the Lemma is proven. The second claim is clear if either $u$ or $v$ are simple nodes. Else, a reduction similar to the second part of the proof of Lemma 13 can be applied to show that there is a edge on the path $\mathcal{P}$ with $q_a > 0$. $\qquad \square$

Another key component for proving Theorem 2 is the inversion of the node coupling conditions via the implicit function theorem. For the convenience of the reader, we state it in the appendix; see Appendix A.

**Lemma 15** (Implicit Function Theorem for the Coupling Conditions)**.**
*In the node $u$ with incident arcs $a \in \delta_{in}(u), b \in \delta_{out}(u)$ consider the coupling conditions*

$$h_u(y_a, q_a) = h_u(y_b, q_b).$$

*If the function $h_u$ is continuously differentiable and the function $h_u(\,\cdot\,, q)$ strictly increasing and surjective for each $q \in \mathbb{R}$, then there exists a uniquely determined continuously differentiable function*

$$Y_a : \mathbb{R}^3 \to \mathbb{R}, \quad (y_b, q_a, q_b) \mapsto Y_a(y_b, q_a, q_b)$$

*such that*

$$h_u(Y_a(y_b, q_a, q_b), q_a) = h_u(y_b, q_b). \tag{1.45}$$

*The partial derivatives of $Y_a$ are given by*

$$\partial_{q_a} Y_a(y_b, q_a, q_b) = -\frac{\partial_q h_u(Y_a(y_b, q_a, q_b), q_a)}{\partial_y h_u(Y_a(y_b, q_a, q_b), q_a)}, \tag{1.46a}$$

$$\partial_{q_b} Y_a(y_b, q_a, q_b) = +\frac{\partial_q h_u(y_b, q_b)}{\partial_y h_u(Y_a(y_b, q_a, q_b), q_a)}, \tag{1.46b}$$

$$\partial_{y_b} Y_a(y_b, q_a, q_b) = +\frac{\partial_y h_u(Y_a(y_b, q_a, q_b), q_a)}{\partial_y h_u(y_b, q_b)}. \tag{1.46c}$$

*Analogously, there exists a uniquely determined continuously differentiable function*

$$Y_b : \mathbb{R}^3 \to \mathbb{R}, \quad (y_a, q_a, q_b) \mapsto Y_b(y_a, q_a, q_b)$$

*such that*

$$h_u(y_a, q_a) = h_u(Y_b(y_a, q_a, q_b), q_b).$$

*The partial derivatives of $Y_b$ with respect to the flow variables are given by*

$$\partial_{q_a} Y_b(y_a, q_a, q_b) = +\frac{\partial_q h_u(y_a, q_a)}{\partial_y h_u(Y_b(y_a, q_a, q_b), q_b)}, \tag{1.47a}$$

$$\partial_{q_b} Y_b(y_a, q_a, q_b) = -\frac{\partial_q h_u(Y_b(y_a, q_a, q_b), q_b)}{\partial_y h_u(Y_b(y_a, q_a, q_b), q_b)}, \tag{1.47b}$$

$$\partial_{y_a} Y_b(y_a, q_a, q_b) = +\frac{\partial_y h_u(Y_b(y_a, q_a, q_b), q_b)}{\partial_y h_u(y_a, q_a)}. \tag{1.47c}$$

*Proof.* We show the claim for the function $Y_a$. The result for $Y_b$ follows analogously. The function $h_u(\,\cdot\,, q)$ is invertible by assumption. Therefore, for each triple $(y_b, q_a, q_b)$, there exists a uniquely determined state $y_a$ such that

$$h_u(y_a, q_a) = h_u(y_b, q_b).$$

The implicit function theorem (see Appendix A) guarantees the existence of neighborhoods $U(y_a)$ and $U(q_a, y_b, q_b)$ and a uniquely determined continuously differentiable mapping

$$Y_a : U(q_a, y_b, q_b) \to U(y_a), \quad (y_b, q_a, q_b) \mapsto Y_a(y_b, q_a, q_b)$$

such that (1.45) and (1.46) hold for all $(q_a, y_b, q_b) \in U(q_a, y_b, q_b)$. Since the point $(q_a, y_b, q_b) \in \mathbb{R}^3$ can be chosen arbitrarily and the function $Y_a$ is uniquely determined as a consequence of the implicit function theorem, this yields that $Y_a$ can be defined globally, i.e.,

$$Y_a : \mathbb{R}^3 \to \mathbb{R}, \ (y_b, q_a, q_b) \mapsto Y_a(y_b, q_a, q_b).$$

The function $Y_a$ is continuously differentiable, because the denominator in Equation (1.46) is not changing sign as

$$\partial_y h_u(Y_a(y_b, q_a, q_b), q_a) > 0 \quad \text{and} \quad \partial_y h_u(y_b, q_b) > 0$$

hold by assumption. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

With the representation for the derivatives of the state derived from solving the node coupling conditions it is possible to ensure that the states are monotonously decreasing in the flow. The result of the following Lemma is a key ingredient for the proof of Theorem 2. We reference to the succeeding theorem's assumptions within the lemma.

**Lemma 16** (Flow-Monotonicity of the State)**.**
*Consider a graph $\mathcal{G} = (\mathcal{V}, \mathcal{A})$ of two edges $\mathcal{A} = \{a, b\} = \{(r, u), (u, v)\}$ with the state $y_r$ prescribed in the root node $r$. Assume that the flow values $q_a(\lambda), q_b(\lambda)$ on the arcs are continuously differentiable, increasing functions of a parameter $\lambda$ and at least one of them strictly increases in $\lambda$. The state values are coupled by* (Arc Coupling) *and* (Node Coupling)*. Furthermore, assume that the node and arc coupling functions on the graph $\mathcal{G}$ fulfill Assumptions a)(i)–a)(ii) and b)(i)–b)(vi) of Theorem 2.*

*Then, the state $y_b(L_b)$ is a strictly decreasing function of $\lambda$. The same remains true if the prescribed state $y_r(\lambda)$ is a continuously differentiable, decreasing function in $\lambda$.*

*Proof.* We can calculate the state values iteratively using the function $Y_b$ of Lemma 15,

$$y_a(L_a) = f_a(y_r, q_a), \quad y_b(0) = Y_b(y_a(L_a), q_a, q_b), \quad y_b(L_b) = f_b(y_b(0), q_b).$$

The chain rule yields

$$\partial_\lambda y_b(L_b) = \partial_y f_b(y_b(0), q_b)\, \partial_\lambda y_b(0) + \partial_q f_b(y_b(0), q_b)\, \partial_\lambda q_b, \tag{1.48a}$$

$$\partial_\lambda y_b(0) = \partial_{y_a} Y_b(y_a(L_a), q_a, q_b)\, \partial_\lambda y_a(L_a) + \partial_{q_a} Y_b(y_a(L_a), q_a, q_b)\, \partial_\lambda q_a$$
$$+ \partial_{q_b} Y_b(y(L_a), q_a, q_b)\, \partial_\lambda q_b, \tag{1.48b}$$

$$\partial_\lambda y_a(L_a) = \partial_y f_a(y_r, q_a)\, \partial_\lambda y_r + \partial_q f_a(y_r, q_a)\, \partial_\lambda q_a. \tag{1.48c}$$

Denote the terms

$$T_y := \partial_y f_b(y_b(0), q_b)\, \partial_{y_a} Y_b(y_a(L_a), q_a, q_b)\, \partial_y f_a(y_r, q_a),$$
$$T_a := \partial_{y_a} Y_b(y_a(L_a), q_a, q_b)\, \partial_q f_a(y_r, q_a) + \partial_{q_a} Y_b(y_a(L_a), q_a, q_b),$$
$$T_b := \partial_y f_b(y_b(0), q_b)\partial_{q_b} Y_b(y(L_a), q_a, q_b) + \partial_q f_b(y_b(0), q_b).$$

Inserting (1.48c) in (1.48b) results in

$$\partial_\lambda y_b(0) = \partial_{y_a} Y_b(y_a(L_a), q_a, q_b)\big[\partial_y f_a(y_r, q_a)\, \partial_\lambda y_r + \partial_q f_a(y_r, q_a)\, \partial_\lambda q_a\big]$$
$$+ \partial_{q_a} Y_b(y_a(L_a), q_a, q_b)\, \partial_\lambda q_a + \partial_{q_b} Y_b(y(L_a), q_a, q_b)\, \partial_\lambda q_b$$
$$= \partial_{y_a} Y_b(y_a(L_a), q_a, q_b)\, \partial_y f_a(y_r, q_a)\, \partial_\lambda y_r$$
$$+ \partial_{q_b} Y_b(y(L_a), q_a, q_b)\, \partial_\lambda q_b + T_a\, \partial_\lambda q_a \tag{1.49}$$

Inserting (1.49) in (1.48a) yields

$$\partial_\lambda y_b(L_b) = \partial_y f_b(y_b(0), q_b)\big[\partial_{y_a} Y_b(y_a(L_a), q_a, q_b)\, \partial_y f_a(y_r, q_a)\, \partial_\lambda y_r$$
$$+ \partial_{q_b} Y_b(y(L_a), q_a, q_b)\, \partial_\lambda q_b + T_a\, \partial_\lambda q_a\big] + \partial_q f_b(y_b(0), q_b)\, \partial_\lambda q_b$$
$$= T_y\, \partial_\lambda y_r + \partial_y f_b(y_b(0), q_b)T_a\, \partial_\lambda q_a + T_b\, \partial_\lambda q_b. \tag{1.50}$$

We show that $T_y > 0, T_a < 0$ and $T_b < 0$ hold. The expansion of the first term leads to

$$T_y = \partial_y f_b(y_b(0), q_b)\, \partial_{y_a} Y_b(y_a(L_a), q_a, q_b)\, \partial_y f_a(y_r, q_a)$$

(by Lemma 15, Equation (1.47c))

$$= \frac{\partial_y h_u(y_b(0), q_b)}{\partial_y h_u(y_a(L_a), q_a)}\, \partial_y f_b(y_b(0), q_b)\, \partial_y f_a(y_r, q_a) > 0 \tag{1.51}$$

as $\partial_y h_u > 0$ by Assumption b)(ii) and $\partial_y f_a$ and $\partial_y f_b$ are positive by Assumption a)(i). The second term can be written as

$$T_a = \partial_{y_a} Y_b(y_a(L_a), q_a, q_b) \, \partial_q f_a(y_r, q_a) + \partial_{q_a} Y_b(y_a(L_a), q_a, q_b)$$

(using Equation (1.47a) and Equation (1.47c))

$$= \frac{\partial_y h_u(y_b(0), q_b)}{\partial_y h_u(y_a(L_a), q_a)} \partial_q f_a(y_r, q_a) + \frac{\partial_q h_u(y_a(L_a), q_a)}{\partial_y h_u(y_b(0), q_b)} < 0 \qquad (1.52)$$

due to Assumption b)(vi). The final term expands to

$$T_b = \partial_y f_b(y_b(0), q_b) \partial_{q_b} Y_b(y(L_a), q_a, q_b) + \partial_q f_b(y_b(0), q_b)$$

(by Equation (1.47b))

$$= \partial_q f_b(y_b(0), q_b) - \partial_y f_b(y_b(0), q_b) \frac{\partial_q h_u(y_b(0), q_b)}{\partial_q h_u(y_b(0), q_b)} < 0 \qquad (1.53)$$

as a consequence of Assumption b)(v). Using (1.51), (1.52) and (1.53), the assumptions for $\partial_\lambda y_r, \partial_\lambda q_a$ and $\partial_\lambda q_b$ and $\partial_y f_b > 0$ shows the claim

$$\partial_\lambda y_b(L_b) = T_y \partial_\lambda y_r + \partial_y f_b(y_b(0), q_b) T_a \partial_\lambda q_a + T_b \partial_\lambda q_b < 0.$$

$\square$

Now we have the necessary tools available to show the existence and uniqueness of solutions on arbitrary networks. For the proof the monotonicity properties of the arc coupling functions are of crucial importance. The assumptions on the node coupling functions ensure that the composition of the arc coupling functions and the function that resolves the node coupling conditions preserves those properties.

Figure 1.11.: Cutting an inner edge in a graph with *one circle* results in a *tree* graph



Figure 1.12.: Cutting an inner edge in a graph with *two circles* results in a graph with *one circle*

**Theorem 2** (Existence and Uniqueness).
*Let a outflow vector $q^{out} \in \mathcal{Q}$ (see Equation (1.37)) and a state value $y_r$ be given. Assume the following:*

a) *The arc coupling functions $f_a$, $a \in \mathcal{A}$ are continuously differentiable and*

    (i) *strictly monotonic increasing with respect to the state, i.e.,*

$$\partial_y f_a(x; y, q) > 0,$$

    (ii) *strictly monotonic decreasing with respect to the flow, i.e.,*

$$\partial_q f_a(x; y, q) < 0, \quad \text{for } x \in (0, L_a],$$

    *for $q \neq 0$,*

    (iii) *strictly monotonic decreasing in flow direction, i.e.,*

$$\text{sign}(\partial_x f_a(x; y, q)) = -\text{sign}(q), \quad \text{for } x \in (0, L_a].$$

b) *The node coupling functions $h_u$, $u \in \mathcal{V}$ are continuously differentiable and*

    (i) *invertible with respect to the first argument,*

    (ii) *strictly monotonic increasing with respect to the state, i.e.,*

$$\partial_y h_u(y, q) > 0,$$

    (iii) *fulfill*

$$\partial_q h_u(y, q) \geq 0 \quad \forall q \geq 0,$$

    (iv) *are symmetric in the second argument, i.e.,*

$$h_u(y, q) = h_u(y, -q).$$

    (v) *have partial derivatives that fulfill the bound*

$$\frac{\partial_q h_u(y, q)}{\partial_y h_u(y, q)} > \frac{\partial_q f_a(L_a(u); y, q)}{\partial_y f_a(L_a(u); y, q)},$$

    *for all $a \in \delta(u)$ if $q \neq 0$,*

(vi) *are flow-monotonicity preserving in the sense that for two arcs* $a \in \delta_{in}(u), b \in \delta_{out}(u)$ *with states* $y_a, y_b$ *and nonzero flows* $q_a, q_b$ *that fulfill*

$$h_u(y_a(L_a), q_a) = h_u(y_b(0), q_b)$$

*and*

$$y_a(L_a) = f_a(L_a; y_a(0), q_a),$$

*the inequality*

$$\frac{\partial_q h_u(y_a(L_a), q_a)}{\partial_y h_u(y_b(0), q_b)} < -\frac{\partial_y h_u(y_b(0), q_b)}{\partial_y h_u(y_a(L_a), q_a)} \partial_q f_a(L_a; y_a(0), q_a)$$

*holds,*

(vii) *are space-monotonicity preserving in the sense that for arcs* $a \in \delta_{in}(u)$ *and* $b \in \delta_{out}(u)$ *with corresponding flows* $q_a > q_b \geq 0$, *and states* $y_a(0)$ *and* $y_b(0)$ *that fulfill*

$$h_u\big(f_a(L_a; y_a(0), q_a), q_a\big) = h_u(y_b(0), \xi),$$

*the inequality*

$$\frac{\partial_q h_u(y_b(0), \xi)}{\partial_y h_u(y_b(0), \xi)} < \frac{y_a(0) - f_a(L_a; y_a(0), q_a)}{q_a - q_b},$$

*holds for all* $\xi \in [q_b, q_a]$.

*Furthermore, assume that on arcs* $a$ *that are in a fundamental cycle, Assumption b)(v) holds if* $L_a$ *is replaced by* $L_a/2$.

*Then, the system of the Equations*

(Kirchhoff), (Arc Coupling) *and* (Node Coupling)

*under the condition* (Initial State) *has a unique solution consisting of the flow vector* $q \in \mathbb{R}^{|\mathcal{A}|}$ *and the state vectors*

$$y(0) := (y_a(0))_{a \in \mathcal{A}}, \quad y(L_{\mathcal{A}}) := (y_a(L_a))_{a \in \mathcal{A}}$$

*with* $L_{\mathcal{A}} := (L_a)_{a \in \mathcal{A}}$.

*Remark* 7 (Regarding the Assumptions). Assumption b)(vii) can be guaranteed if the inequality

$$\frac{\partial_q h_u(y_b(0), \xi)}{\partial_y h_u(y_b(0), \xi)} < \frac{L_a |\partial_x f_a(x; y_a(0), q_a)|}{q_a - q_b}$$

is fulfilled for all $x \in [0, L_a]$. This follows directly from the mean value theorem, which states that there is a $x \in [0, L_a]$ such that

$$L_a \partial_x f_a(x; y_a(0), q_a) = f(L_a; y_a(0), q_a) - f(0; y_a(0), q_a).$$

Furthermore, if $f_a$ is concave in $x$, only the point $x = 0$ has to be considered.

Due to the strict monotonicity of $f_a$, the two Assumptions b)(v) and b)(vii) certainly hold if $h_u$ is constant in $q$, e.g., if the coupling is given by the continuity of states over the node, i.e., $h_u(y, q) = y$.

*Remark* 8 (Idea of the Proof of Theorem 2). The proof shows the existence and uniqueness of the solution by mathematical induction over the number of fundamental circles in the graph. The initial step of the proof starts with a tree graph, where a solution can easily be constructed: For prescribed boundary flows the linear flow system has a unique solution. By the strict monotonicity of the coupling function, the node coupling conditions are invertible with respect to the state and the state values can be obtained iteratively, starting from the root node, where the state value is prescribed. The induction assumption requires the existence of a unique solution on a graph that has $\dim(\ker(A)) = n$ fundamental circles. For the induction step, on a graph with $\dim(\ker(A)) = n + 1$ fundamental cycles, we cut one arc in a fundamental cycle to arrive at a graph with $\dim(\ker(A)) = n$. For each value of the unknown flow $\lambda$ on this arc, the resulting graph has a unique solution by the induction assumption. This is true, because after the removal of a sufficient number of edges, the graph is a tree. To fulfill the coupling conditions in the two nodes that were generated by the cut, we introduce an auxiliary function $H$, which is given by the difference of the state values in the two new nodes. The monotonicity properties allow us to show that $H$ has a unique root, which is the desired flow value on the cut edge and can be used to construct the solution on the whole graph. A graphical illustration of this idea is depicted in Figure 1.11 and Figure 1.12.

*Proof of Theorem 2.* We show the statement of the theorem by mathematical induction over the dimension of the kernel of the incidence matrix $\dim(\ker(A))$, which is the number of fundamental cycles in the graph.

*Initial Step*
Consider the case $\dim(\ker(A)) = 0$, that is, $\mathcal{G}$ is a tree. By Lemma 11, the linear flow system $Aq = q^{\text{out}}$ is solvable for a right hand side $q^{\text{out}} \in \mathcal{Q}$ and, because $\dim(\ker(A)) = 0$, the flow solution $q$ is unique. Now, knowing the flow value on every arc, the states $y_a(0)$ and $y_a(L_a)$ on the arc $(u, v)$ can be obtained as follows: Denote the path from $r$ to $v$ by $\mathcal{P}$ and assume without

loss of generality that it is directed from $r$ to $v$. Since the graph is a tree, the path is unique. Denote the predecessor of an arc $a \in \mathcal{P}$ by $\pi(a)$. The state at $x = 0$ on an arc $a = (u, v)$ can be retained from the state $y_{\pi(a)}(L_a)$ and the flow values $q_{\pi(a)}$ and $q_a$ by the mapping

$$Y_a : \mathbb{R}^3 \to \mathbb{R}, \quad (y_{\pi(a)}(L_{\pi(a)}), q_{\pi(a)}, q_a) \mapsto y_a(0)$$

that is introduced in Lemma 15. The assumptions required for Lemma 15 are fulfilled due to Assumption b). The value $y_a(L_a)$ can be directly calculated by

$$y_a(L_b) = f_a(y_a(0), q_a),$$

where $f_a$ is the arc coupling function on the arc $a$ (throughout this proof, we omit the argument $L_a$ ). The argument $y_{\pi(a)}(L_{\pi(a)})$ of the function $Y_a$ can be obtained by the functions $Y_{\pi(a)}$ and $f_{\pi(a)}$. This defines an recursion to calculate all state values on the path, because the state in the node $r$ is prescribed by

$$y_b(0) = y_r,$$

where $\{b\} = \delta(r)$; see (Initial State). Because the functions $f_a$ and $Y_a$ are continuous, the solution $(y(0), y(L_\mathcal{A}), q)$ is continuous with respect to the prescribed data $y_r$ and $q^{\mathrm{out}}$.

*Induction Assumption*
If $\dim(\ker(A)) \leq n$ for a fixed $n \geq 0$, then for prescribed $y_r \in \mathbb{R}$ and $q^{\mathrm{out}} \in \mathcal{Q}$, there exists a unique solution $(y(0), y(L_\mathcal{A}), q) \in \mathbb{R}^{|\mathcal{A}|} \times \mathbb{R}^{|\mathcal{A}|} \times \mathbb{R}^{|\mathcal{A}|}$ fulfilling the Equations (Initial State), (Kirchhoff), (Arc Coupling) and (Node Coupling). The state is continuously depending on the values of $y_r$ and $q^{\mathrm{out}}$.

*Induction Step*
To reduce the graph $\mathcal{G} = (\mathcal{V}, \mathcal{A})$ with $\dim(\ker(A)) = n + 1$ to a graph $\mathcal{G}^n$ with $\dim(\ker(A^n)) = n$, we pick one basis vector $q^{\mathrm{ker}}$ of $\ker(A)$. Now, we cut one arc $c = (v_\mathrm{l}, v_\mathrm{r})$ that corresponds to a nonzero entry of $q^{\mathrm{ker}}$ in the middle. We obtain two new arcs $c_\mathrm{l} = (v_\mathrm{l}, v_{\mathrm{cl}})$ and $c_\mathrm{r} = (v_{\mathrm{cr}}, v_\mathrm{r})$ with length $L_{c_\mathrm{l}} = L_{c_\mathrm{r}} = L_c/2$. Here, $v_{\mathrm{cl}}$ and $v_{\mathrm{cr}}$ are the artificial nodes generated by the cut. Set

$$\mathcal{V}^n := \mathcal{V} \cup \{v_{\mathrm{cl}}, v_{\mathrm{cr}}\}, \ \mathcal{A}^n := (\mathcal{A} \setminus \{c\}) \cup \{c_\mathrm{l}, c_\mathrm{r}\} \text{ and } \mathcal{G}^n := (\mathcal{V}^n, \mathcal{A}^n)$$

and denote the node-arc incidence matrix of $\mathcal{G}^n$ by $A^n$. The procedure is shown in Figure 1.12 and Figure 1.11.

A basis for $\ker(A^n)$ can be obtained by extending the basis vectors of $\ker(A)$ with zero entries for the components corresponding to $c_\mathrm{l}$ and $c_\mathrm{r}$.

This is not possible for vectors with a nonzero component at $c$. One basis vector, for which this is not possible, is the vector $q^{\text{ker}}$ that was chosen to generate the cut. As the cut does not generate new cycles, this yields

$$\dim(\ker(A^n)) \leq \dim(\ker(A)) - 1 = n.$$

The outflow vector $z$ for the graph $\mathcal{G}^n$ is known for all nodes, except in $v_{\text{cl}}$ and $v_{\text{cr}}$, by the values of $q^{\text{out}}$ on the graph $\mathcal{G}$. The unknown flow $\lambda \in \mathbb{R}$ corresponds to the flow on the arc $c$ on the original graph. For each flow $\lambda \in \mathbb{R}$ and

$$z_{v_{\text{cl}}} = -\lambda, \quad z_{v_{\text{cr}}} = \lambda, \quad z_u = q_u^{\text{out}} \quad \forall u \in \mathcal{V},$$

the induction assumption states that there is a unique solution consisting of $y^\lambda(0), y^\lambda(L_\mathcal{A})$ and $q^\lambda$ on the graph $\mathcal{G}^n$. Due to the continuity of $f_c$, the values $y_{c_{\text{l}}}^\lambda(L_{c_{\text{l}}})$ and $y_{c_{\text{r}}}^\lambda(0)$ should coincide, as they correspond to the value of $f_c(L_c/2; y_c(0), q_a)$ on the original graph.

It is equivalent to find a root of the function

$$H(\lambda) := y_{c_{\text{l}}}^\lambda(L_{c_{\text{l}}}) - y_{c_{\text{r}}}^\lambda(0).$$

First, we show that the function $H$ is strictly monotonic. Consider a flow $\mu > \lambda$. Corollary 1 shows the existence of a path $\mathcal{P}_l$ from $r$ to $v_{\text{cl}}$ and a path $\mathcal{P}_r$ from $r$ to $v_{\text{cr}}$ (once again, we assume that the direction of the arcs coincides with the direction of the path) such that the flow on the two paths fulfills

$$\begin{align}
q_a^\mu &\leq q_a^\lambda \quad \forall a \in \mathcal{P}_l \qquad \text{(Note the direction of } \mathcal{P}_l) \tag{1.54a} \\
q_a^\mu &\geq q_a^\lambda \quad \forall a \in \mathcal{P}_r \tag{1.54b}
\end{align}$$

and each inequality holds strictly for at least one $a \in \mathcal{P}_l$ and one $a \in \mathcal{P}_r$.

By using the function $Y_{c_{\text{l}}}$ as in Lemma 15, it is possible to express the value of $y_{c_{\text{l}}}(L_{c_{\text{l}}})$ by (again, we omit the argument $L_{c_{\text{l}}}$ of $f_{c_{\text{l}}}$)

$$\begin{align}
y_{c_{\text{l}}}(L_{c_{\text{l}}}) &= f_{c_{\text{l}}}(y_{c_{\text{l}}}(0), q_{c_{\text{l}}}) \tag{1.55} \\
y_{c_{\text{l}}}(0) &= Y_{c_{\text{l}}}\big(y_{\pi(c_{\text{l}})}(L_{\pi(c_{\text{l}})}), q_{\pi(c_{\text{l}})}, q_{c_{\text{l}}}\big), \tag{1.56}
\end{align}$$

for fixed $y_r$, where the state $y_{\pi(c_{\text{l}})}(L_{\pi(c_{\text{l}})})$ is determined by the equation

$$y_{\pi(c_{\text{l}})}(L_{\pi(c_{\text{l}})}) = f_{\pi(c_{\text{l}})}(y_{\pi(c_{\text{l}})}(0), q_{\pi(c_{\text{l}})}).$$

This can be expanded up to the node $y_r$, where the state is prescribed. The iterative application of Lemma 16, starting at the root node, where the

state is prescribed, shows that the states $y_{c_l}(L_{c_l})$ strictly decreases in $\lambda$, as (1.54a) holds strictly for at least one arc. Therefore,

$$y_{c_l}^{\lambda}(L_{c_l}) < y_{c_l}^{\mu}(L_{c_l}) \tag{1.57}$$

holds. On the path $\mathcal{P}_r$, the argumentation is analogous, but, compared to (1.54b), the inequality (1.54b) is reversed. Consequently,

$$y_{c_r}^{\lambda}(0) > y_{c_r}^{\mu}(0) \tag{1.58}$$

holds. The combination of (1.57) and (1.58) shows

$$H(\lambda) = y_{c_l}^{\lambda}(L_{c_l}) - y_{c_r}^{\lambda}(0) < y_{c_l}^{\mu}(L_{c_l}) - y_{c_r}^{\mu}(0) = H(\mu),$$

for $\lambda < \mu$. Therefore, the function $H$ is strictly increasing. Furthermore, it is continuous, because the states depend continuously on $q^{\mathrm{out}}$.

If we find a flow value $\lambda_-$ with $H(\lambda_-) < 0$ and a flow value $\lambda_+ > \lambda_-$ with $H(\lambda_+) > 0$, Bolzano's intermediate value theorem ensures that there is a unique solution $\lambda^* \in (\lambda_-, \lambda_+)$ fulfilling $H(\lambda^*) = 0$. Choose

$$\lambda_+ := \sum_{v \in \mathcal{V}} |q_v^{\mathrm{out}}| \quad \text{and} \quad z_{v_l} = -\lambda_+.$$

Then, by Lemma 14, there exists a path $\mathcal{P}$ from $v_{cl}$ to $v_{cr}$ such that

$$q_a \geq 0 \quad \forall a \in \mathcal{P}$$

with $q_a > 0$ for at least one $a \in \mathcal{P}$. Consider a node $u$ on the path with $a \in \delta_{\mathrm{in}}(u) \cap \mathcal{P}$ and $b \in \delta_{\mathrm{in}}(u) \cap \mathcal{P}$. Our goal is to show the decrease of the state values on the path. If the flow values on the incident arcs are equal, i.e., $q_a = q_b$, then the coupling condition $h_u(y_a(L_a), q_a) = h_u(y_b(0), q_b)$ implies $y_b(0) = y_a(L_a)$ and, by Assumption a)(iii), we have $y_a(L_a) < y_a(0)$ if $q_a > 0$ and $y_a(L_a) = y_a(0)$ if $q_a = 0$.

If $q_a < q_b$, then, we have, by (1.47b) and Assumption b)(iii)

$$y_b(0) = Y_b(y_a(L_a), q_a, q_b) \leq Y_b(y_a(L_a), q_a, q_a) = y_a(L_a),$$

where the last equation follows from the fact that $Y_b$ is the unique solution of $h_u(y_a, q_a) = h_u(Y_b(y_a(L_a), q_a, q_a), q_a)$.

If $q_a > q_b$, then, using the first order Taylor expansion around the point $(y_a(L_a), q_a, q_a)$, we obtain

$$y_b(0) = Y_b(y_a(L_a), q_a, q_b) = Y_b(y_a(L_a), q_a, q_a) + (q_b - q_a)\partial_{q_b} Y_b(y_a, q_a, \xi)$$

$\big($by $Y_b(y_a(L_a), q_a, q_a) = y_a(L_a)$ and Equation (1.47b)$\big)$

$$= y_a(L_a) - (q_b - q_a)\frac{\partial_q h_u(Y_b(y_a, q_a, \xi), \xi)}{\partial_y h_u(Y_b(y_a, q_a, \xi), \xi)},$$
$$= f_a(L_a; y_a(0), q_a) - (q_b - q_a)\frac{\partial_q h_u(Y_b(y_a, q_a, \xi), \xi)}{\partial_y h_u(Y_b(y_a, q_a, \xi), \xi)}$$
$$< y_a(0),$$

where the last inequality is due to Assumption b)(vii). We have shown that the states are decreasing along the path (using that $Y_b$ is strictly increasing in the first argument). By Assumption a)(iii) the decrease is strict on the arc $a$, where $q_a > 0$. This shows

$$y_{c_l}^{\lambda_+}(L_{c_l}) > y_{c_r}^{\lambda_+}(0)$$

and therefore $H(\lambda_+) > 0$. Using the same monotonicity arguments with reversed signs of the flows, we also see that the choice

$$\lambda_- := -\sum_{v \in \mathcal{V}} \left|q_v^{\text{out}}\right|$$

leads to $H(\lambda_-) < 0$. Therefore, by Bolzano's intermediate value theorem, there exists a unique $\lambda^* \in (\lambda_-, \lambda_+)$ such that $H(\lambda^*) = 0$. The induction assumption states that for prescribed $\lambda^*$ there is a unique solution $(y^{\lambda^*}(0), y^{\lambda^*}(L), q^{\lambda^*})$ on the graph $\mathcal{G}^n$. A solution on the original graph $\mathcal{G}$ can be easily reconstructed by setting

$$q_c := \lambda^*, \quad y_c(0) := y_{c_l}^{\lambda^*}(0), \quad y_c(L_c) := y_{c_r}^{\lambda^*}(L_{c_r})$$

and the remaining flow and state values to their corresponding parts on the graph $\mathcal{G}^n$. The solution $(y(0), y(L), q)$ is continuous with respect to the prescribed data, because the solution $(y^{\lambda^*}(0), y^{\lambda^*}(L), q^{\lambda^*})$ depends continuously on the prescribed data by the induction assumption. Since the function $H$ is even continuously differentiable with respect to $p_r$ and $q^{\text{out}}$ (as a composition of the continuously differentiable functions $Y_a$ and $f_a$) the implicit function theorem ensures that $\lambda^*$ is continuously depending on $p_r$ and $q^{\text{out}}$. $\qquad\square$

*Remark* 9 (The Solution's Independence of the Cut Choice). Each fundamental cycle corresponds to a basis vector of the node-arc incidence matrix. Therefore, the choice, which arc should be cut in the proof of Theorem 2,

can be interpreted as choosing a basis vector of the incidence matrix. The arcs in a fundamental cycle correspond to nonzero entries of a basis vector of the kernel of the incidence matrix. Hence, choosing different edges is equivalent to choosing a different basis for the kernel and does not interfere with the uniqueness of the solution. If a flow vector $q$ is retained in one basis, a change of the basis only changes its representation but not its value. Therefore, the solution is independent of which edges are being cut.

# 1.4. Existence and Uniqueness on Gas Networks

In the previous section, we developed a framework for general flow problems on graphs. In this section, we apply the results to gas networks. In this application the flow corresponds to the mass flow of the gas. The outflow can be thought of as the demand of customers withdrawing gas from the network. The relevant state values are the pressure values at both ends of a gas pipeline (the gas density instead of the gas pressure can be used with slight adjustments).

In Section 1.4.1, we focus on passive networks governed by the stationary isothermal Euler equations. Here, networks exclusively consisting of pipelines are considered under different coupling conditions.

In Section 1.4.2, we include active elements, namely resistors and turbo compressor, in the networks. A rigorous analysis using the ideal gas model with pressure continuity coupling is provided.

## 1.4.1. Passive Networks

Consider the gas flow through a network of pipelines. The pipelines are represented by the arcs $\mathcal{A}$ of a graph $\mathcal{G}$. The interconnections of multiple pipelines are represented by the nodes $\mathcal{V}$. We denote the pressure on a arc $a \in \mathcal{A}$ by $p_a$, the mass flow by $m_a$ and the mass flow divided by the cross sectional area of the pipe by $q_a = \frac{4m_a}{\pi D_a^2}$.

### Kirchhoff Conditions

The stationary gas dynamics along the pipes are described by the isothermal Euler equations (IsoStat). The conservation of mass in each node $u \in \mathcal{V}$

reads

$$\sum_{a \in \delta(u)} A_{ua} m_a = m_u^{\text{out}}$$

which is equivalent to the system

$$Am = m^{\text{out}}. \qquad\qquad\text{(KirchhoffM)}$$

## Arc Coupling Conditions

We consider the case of real and ideal gas.

### Real Gas

For real gas, Theorem 1 states that the pressure at both ends of a pipe is coupled via

$$p_a(L_a) = f_a(L_a; p_a(0), q_a), \text{ where}$$

$$f_a(L; p, q) := F^{-1}\left(F(p, q) - RTq|q| \int_0^L \frac{\theta(s)}{2} \, \mathrm{d}s\right), \qquad\text{(Real)}$$

$$F(p, q) := \frac{p}{\alpha} - \frac{1}{\alpha^2} \ln(z(p)) + q^2 RT \ln\left(\frac{z(p)}{p}\right).$$

The domain of $f_a$ is given by

$$\mathcal{D}(f_a) := \left\{(x, p, q) \in [0, L_a] \times \mathbb{R}_{>0} \times \mathbb{R} \mid p^2 > q^2 RT, p < |\alpha|^{-1}\right\}. \quad (1.59)$$

### Ideal Gas

For ideal gas, an analytic solution of the system (IsoStat) is presented in Gugat, Hante, et al. 2015, Lemma 3.5 and given by

$$p_a(L_a) = f_a(L_a; p_a(0), q_a), \text{ where}$$

$$f_a(L; p, q) := c \, |q| \sqrt{-W_{-1}\left(-\exp\left(-C + \text{sign}(q) \int_0^L \theta_a(s) \, \mathrm{d}s\right)\right)}, \quad\text{(Ideal)}$$

$$C := \frac{1}{\eta} + \ln(\eta), \quad \eta := \frac{q^2 c^2}{p^2}.$$

The Lambert-W-function is the inverse function of $x \mapsto x \exp(x)$. We use the branch $W_{-1}$ for arguments in $(-e^{-1}, 0)$; see Figure 1.13. For details on the Lambert-W-function see Veberic 2010 or Corless et al. 1996. Keep in mind that for ideal gas the sound speed $c = \sqrt{\partial p / \partial \rho}$ and the constant compressibility factor $z(p) = z_m$ are coupled via the state equation (StateEq). To be concrete, $c^2 = RT z_m$ holds. The domain of $f_a$ is given by

$$\mathcal{D}(f_a) := \left\{(x, p, q) \in [0, L_a] \times \mathbb{R}_{>0} \times \mathbb{R} \mid p^2 > q^2 c^2\right\}. \qquad (1.60)$$

Figure 1.13.: The two branches of the Lambert-W function. The branch $W_{-1}$ is used in the ideal gas case.

*Remark* 10 (Different Mass Flow Quantities). To apply the results of Section 1.3, the function $f_a$ should depend on $m_a$ instead of $q_a$. We use the more common quantity $q_a$, while keeping in mind that all relevant properties of $f_a$ carry over to a function $\tilde{f}_a : (L, p, m) \mapsto \tilde{f}_a(L; p, m) = f_a(L; p, q)$ as the flow argument is only scaled by a positive constant.

### Node Coupling Conditions

Different coupling function are in use; see Reigstad 2014, p. 74. For the case of ideal gas, we focus on the Bernoulli invariant. For real gas, we discuss the continuity of pressure.

### Pressure

$$h_u(p, q) = p, \tag{Pr}$$

### Momentum Flux

$$h_u(p, q) = \frac{p}{c^2}\left(1 + \frac{q^2 c^2}{p^2}\right) = \rho(1 + \eta), \tag{MF}$$

### Bernoulli Invariant

$$h_u(p, q) = \log\left(\frac{p}{c^2}\right) + \frac{1}{2}\frac{q^2 c^2}{p^2} = \log(\rho) + \frac{1}{2}\eta. \tag{BI}$$

Figure 1.14.: The different coupling functions over the pressure for the values $q = 400\,\mathrm{kg\,s^{-1}\,m^{-2}}, R = 448.66\,\mathrm{J\,kg^{-1}\,K^{-1}}, T = 290\,\mathrm{K}$ and $z(p) = 1$. To obtain similar order of magnitudes, the momentum flux and the Bernoulli invariant are scaled. The functions $c^2\rho(1+\eta)$ (scaled momentum flux) and $c^2\exp(\log(\rho)+\frac{1}{2}\eta)$ (scaled Bernoulli invariant) are depicted.

For the case of real gas, we consider $h_u$ on the domain

$$\mathcal{D}(h_u) := \left\{(p,q) \in \mathbb{R}_{>0} \times \mathbb{R} \mid p^2 > q^2RT, p < |\alpha|^{-1}\right\}. \qquad (1.61)$$

For the case of ideal gas, we consider $h_u$ on the domain

$$\mathcal{D}(h_u) := \left\{(p,q) \in \mathbb{R}_{>0} \times \mathbb{R} \mid p^2 > q^2c^2\right\}. \qquad (1.62)$$

In the following sections, we show that the functions $f_a$ and $h_u$ fulfill the assumptions of Theorem 2. The properties of the functions $f_a$ in the real gas case were discussed in Section 1.2.

**Existence and Uniqueness for Real Gas with Pressure Continuity**

First, we show the required assumptions for the arc coupling functions.

**Lemma 17** (Arc Coupling Assumptions)**.**
*The functions $f_a$ as defined in* (Real) *fulfill the Assumptions a) of Theorem 2.*

*Proof.* Lemma 5 shows Assumptions a)(i) and a)(ii). Lemma 4 shows Assumption a)(iii). □

Next, we show the required assumptions for the node coupling functions.

**Lemma 18** (Node Coupling Assumptions).
*Let $h_u$ be given by the pressure* (Pr). *Then Assumptions b) of Theorem 2 are fulfilled.*

*Proof.* We have $h_u(p, q) = p$. It is obvious that Assumptions b)(ii)–b)(iii) and b)(iv) hold. Assumptions b)(v) and b)(vi) are true, because

$$\partial_q h_u = 0, \ \partial_p h_u > 0 \quad \text{and} \quad \partial_q f_a < 0, \ \partial_p f_a > 0$$

hold. The remaining Assumption b)(vii) is true, because $\partial_q h = 0$ and $\partial_x f_a < 0$ for positive flow. Hence, the left hand side of the Equation in b)(vii) is zero and the right hand side is positive. This shows the properties for the choice $h_u(p, q) = p$. □

Now, we can apply Theorem 2 to show the existence and uniqueness of solutions on real gas networks.

**Theorem 3** (Existence and Uniqueness for Real Gas).
*Consider the arc coupling conditions* (Real), *the node coupling conditions* (Pr) *and the Kirchhoff conditions* (KirchhoffM). *Let the pressure $p_r > 0$ with $z(p_r) > 0$ be prescribed at the simple node $r \in \mathcal{V}$. Then there exists a constant $C(p_r, G) > 0$ such that for all $m^{out} \in \mathbb{R}^{|\mathcal{V}|}$ with*

$$\mathbf{1}^T m^{out} = 0, \qquad \|m^{out}\| < C(p_r, G) \quad and \quad \frac{q_r^{out2} RT}{p_r^2} < 1$$

*a* unique solution $(p(0), p(L), q) \in \mathbb{R}^{|\mathcal{A}|} \times \mathbb{R}^{|\mathcal{A}|} \times \mathbb{R}^{|\mathcal{A}|}$ *exists on the graph $G = (\mathcal{V}, \mathcal{A})$. This solution fulfills*

(KirchhoffM), (Arc Coupling) *and* (Node Coupling).

*Proof.* With the results from Lemma 17 and Lemma 18, we see that the assumptions of Theorem 2 hold on the domains $\mathcal{D}(f_a), a \in \mathcal{A}$ and $\mathcal{D}(h_u), u \in \mathcal{V}$ as defined in Equations (1.59) and (1.61). In the node $r$, the prescribed pressure and flow fulfill

$$p_r < |\alpha|^{-1}, \quad \text{and} \quad \frac{q_r^{out2} RT}{p_r^2} < 1$$

by assumption. The proof of Theorem 2 reveals that in each induction step, the pressure and flow values depend continuously on the boundary values. This implies that the values $p_a(0), p_a(L_a), q_a$ stay inside the domains of $f_a$ and $h_u$ for a sufficiently small prescribed flow $q^{\text{out}}$. Consequently, there exists a constant $C(p_r, G) > 0$ such that all obtained pressure-flow pairs are inside the corresponding domains $\mathcal{D}(f_a)$ and $\mathcal{D}(h_u)$ if $\|q^{\text{out}}\| < C(p_r, G)$ holds. The difference between the conditions (KirchhoffM) and (Kirchhoff) is discussed in Remark 10. □

Theorem 3 was first stated in Gugat, Schultz, and Wintergerst 2018, Theorem 1, where a direct proof for the specific setting was given. Here, we used the generalized version of the theorem (Theorem 2) for the proof.

**Existence and Uniqueness for Ideal Gas**

For the case of ideal, we follow the same steps: First, we show that the arc coupling assumptions are fulfilled.

**Lemma 19** (Arc Coupling Assumptions).
*The functions $f_a$ as defined in* (Ideal) *fulfill the Assumptions a) of Theorem 2.*

*Proof.* Assumption a)(i) and a)(ii) follow from Gugat, Hante, et al. 2015, Lemma 3.6. Assumption a)(iii) is implied by Equation (18) in Gugat, Hante, et al. 2015, Lemma 3.5. □

Next, the node coupling assumptions are shown.

**Lemma 20** (Node Coupling Assumptions for the Bernoulli Invariant).
*Let $h_u$ be given by the the Bernoulli Invariant* (BI) *on the domain $\mathcal{D}(h_u)$ defined in* (1.62). *For positive flow $q_a > 0$, let the sufficient decrease condition*

$$1 < \frac{p_a(L_a)^2}{p_b(0)^2} \frac{(1-\eta_b)^2}{(1-\eta_L)^2} \left( \frac{1}{\eta_a} - \frac{1}{\eta_L} \right), \tag{1.63}$$

*hold on two adjacent arcs $a, b$. Furthermore, assume that for negative flow $q < 0$, the sufficient increase condition*

$$f_a(p,q)^2 - p^2 > c^2 q^2 \quad \text{or equivalently} \quad \frac{1}{\eta_L} - \frac{1}{\eta} > 1, \tag{1.64}$$

*with $\eta_L = \frac{q^2 c^2}{f_a(p,q)^2}$ holds. Additionally, assume that the constant friction is sufficiently high in the sense that*

$$\frac{1}{2} \theta_a L_a > \frac{p_b(0)}{p_a(0)} \frac{h(\eta_b^{q_a})}{h(\eta_a)} \left( 1 - \frac{q_b}{q_a} \right), \tag{1.65}$$

*with* $\eta_b^{q_a} := \dfrac{q_a^2 c^2}{p_b(0)^2} < 1, \; \eta_a := \dfrac{q_a^2 c^2}{p_a(0)^2} < 1 \;\; and \;\; h(\eta) := \dfrac{\eta}{1 - \eta},$

*where* $a \in \delta_{in}(u), b \in \delta_{out}(u)$ *with the flows* $q_a > q_b \geq 0$ *and*

$$h_u(p_b(0), q_b) = h_u(f_a(L_a; p_a(0), q_a), q_a).$$

*Then, the Assumptions b) of Theorem 2 are fulfilled.*

*Remark* 11 (Assumptions of Lemma 20). Both of the Inequalities (1.64) and (1.65) can also be seen as a requirement for the length of the pipe to be sufficiently large. It is evident that Inequality (1.65) is easy to fulfill if the flows $q_a$ and $q_b$ are of similar size. For Inequality 1.63, observe that the fraction

$$\frac{p_a(L_a)^2}{p_b(0)^2} \frac{(1 - \eta_b)^2}{(1 - \eta_L)^2}$$

is close to one if the flows $q_a$ and $q_b$ are of similar size. The condition is fullfilled if the pressure decrease on the arc $a$ measured by $\frac{1}{\eta_a} - \frac{1}{\eta_b}$ is large enough.

*Proof.* The Bernoulli Invariant (BI) is given by

$$h_u(p, q) = \log\left(\frac{p}{c^2}\right) + \frac{1}{2}\frac{q^2 c^2}{p^2}.$$

Assumption b)(iv) can immediately be seen. We proceed to show assumption b)(i). Let $q_a \neq 0$ and $(q_b, p_b) \in \mathcal{D}(h_u)$ be given. We need to find $p_a$ such that

$$h_u(p_a, q_a) = h_u(p_b, q_b).$$

By definition (BI), we have, using the variables $\rho_a = \frac{p_a}{c^2}$ and $\eta_a = \frac{q_a^2 c^2}{p_a^2}$ ($\rho_b, \eta_b$ analogously),

$$\log(\rho_a) + \frac{1}{2}\eta_a = \log(\rho_b) + \frac{1}{2}\eta_b$$

or equivalently

$$\log\left(\frac{\rho_a}{\rho_b}\right) = \frac{1}{2}(\eta_b - \eta_a).$$

Hence,

$$\frac{\rho_a}{\rho_b} = \exp\left(\frac{1}{2}(\eta_b - \eta_a)\right).$$

Squaring the equation and multiplying by $\rho_b^2$ leads to

$$\rho_a^2 = \rho_b^2 \exp(\eta_b - \eta_a).$$

We multiply by $\frac{c^2}{q_a^2} \exp(\eta_a)$ to obtain

$$\eta_a^{-1} \exp(\eta_a) = \frac{c^2 \rho_b^2}{q_a^2} \exp(\eta_b).$$

Rearranging shows

$$(-\eta_a) \exp(-\eta_a) = -\frac{q_a^2}{c^2 \rho_b^2} \exp(-\eta_b).$$

Using the right branch of the Lambert-W function yields

$$\eta_a = -W_0\left(-\frac{q_a^2}{c^2 \rho_b^2} \exp(-\eta_b)\right),$$

which is the physically right solution as the range of $W_0$ for negative arguments is $[-1, 0]$. For $q_a \neq 0$, and $\frac{q_a^2}{c^2 \rho_b^2} < 1$, $\eta_b < 1$, the argument of $W_0$ is in $(-e^{-1}, 0)$ and therefore $\eta_a$ lies in $(0, 1)$. This means, we can explicitly calculate $p_a$ as

$$p_a = q_a c \sqrt{-\left[W_0\left(-\frac{q_a^2}{c^2 \rho_b^2} \exp(-\eta_b)\right)\right]^{-1}}$$

and Assumption b)(i) is shown.

We obtain the partial derivatives

$$\partial_p h_u(p, q) = \frac{1}{p}(1 - \eta) \quad \text{and} \quad \partial_q h_u(p, q) = \frac{\eta}{q}. \tag{1.66}$$

Clearly, $\eta < 1$ implies $\partial_p h_u(p, q) > 0$, which is Assumption b)(ii). Assumption b)(iii) is obvious for $p > 0$.

We proceed with Assumption b)(v). The first argument of $f_a$ is omitted. By Gugat, Hante, et al. 2015, Lemma 3.6, the partial derivatives of $f_a$ are given by

$$\partial_p f_a(p, q) = \frac{p}{f_a(p, q)} \frac{1 - \eta}{1 - \eta_L}, \quad \text{where} \quad \eta_L := \frac{q^2 c^2}{f_a(p, q)^2} \tag{1.67}$$

and

$$\partial_q f_a(p, q) = -\frac{f_a(p, q)}{q} \frac{\eta_L - \eta}{\eta(1 - \eta_L)}. \tag{1.68}$$

Assumption b)(v) holds for $q > 0$, because

$$\frac{\partial_q h_u(p, q)}{\partial_p h_u(p, q)} > 0 > \frac{\partial_q f_a(p, q)}{\partial_y f_a(p, q)}$$

holds. The case $q < 0$ remains. We obtain

$$\frac{\partial_q h_u(p, q)}{\partial_p h_u(p, q)} = \frac{p\eta}{q(1 - \eta)} \tag{1.69}$$

and

$$\frac{\partial_q f_a(p, q)}{\partial_p f_a(p, q)} = -\frac{f_a(p, q)^2(\eta_L - \eta)}{qp\eta(1 - \eta)}.$$

We have to show

$$-\frac{f_a(p, q)^2(\eta_L - \eta)}{qp\eta(1 - \eta)} < \frac{p\eta}{q(1 - \eta)}$$

or equivalently, by multiplication with the negative constant $\frac{p}{qc^2}(1 - \eta)$,

$$-\frac{f_a(p, q)^2(\eta_L - \eta)}{q^2 c^2 \eta} > \frac{p^2 \eta}{q^2 c^2}.$$

Inserting $\eta = \frac{q^2 c^2}{p^2}$ and $\eta_L = \frac{q^2 c^2}{f_a(p,q)^2}$, yields

$$\frac{(\eta - \eta_L)}{\eta_L \eta} > 1 \quad \text{or} \quad \frac{1}{\eta_L} - \frac{1}{\eta} > 1.$$

This is assumption (1.64) and shows that, since all steps were equivalent, the inequality

$$\frac{\partial_q h_u(p, q)}{\partial_p h_u(p, q)} > \frac{\partial_q f_a(p, q)}{\partial_p f_a(p, q)}$$

holds. Hence, Assumption b)(v) has been proven.

For Assumption b)(vii), we use the observations of Remark 7. We want to show

$$\frac{\partial_q h_u(p_b(0), \xi)}{\partial_p h_u(p_b(0), \xi)} < \frac{L_a |\partial_x f_a(0; p_a(0), q_a)|}{q_a - q_b}, \tag{1.70}$$

for $q_a > \xi > q_b \geq 0$. The absolute value of the space derivative of $f_a$ is given by

$$|\partial_x f_a(0; p_a(0), q_a)| = \frac{1}{2}\theta_a p_a(0)\frac{\eta_a}{1 - \eta_a}, \quad \text{where } \eta_a := \frac{q_a^2 c^2}{p_a(0)^2}; \tag{1.71}$$

see Gugat, Hante, et al. 2015, Equation (18). Inserting (1.69) and (1.71) in Inequality (1.70) yields

$$\frac{p_b(0)}{\xi}\frac{\eta_b}{1-\eta_b} < \frac{\frac{1}{2}\theta_a L_a p_a(0)}{q_a - q_b}\frac{\eta_a}{1-\eta_a}, \quad \text{where } \eta_b^\xi := \frac{\xi^2 c^2}{p_b(0)^2}. \tag{1.72}$$

This is equivalent to

$$\frac{1}{2}\theta_a L_a > \frac{p_b(0)}{p_a(0)}\frac{h(\eta_b)}{h(\eta_a)}\frac{q_a - q_b}{\xi}, \tag{1.73}$$

where $h(\eta) := \frac{\eta}{1-\eta}$ as in (1.22) is used. Using $q_a > \xi$ and the fact that $h$ is strictly increasing, we find a sufficient condition for (1.72) that does not depend on $\xi$:

$$\frac{1}{2}\theta_a L_a > \frac{p_b(0)}{p_a(0)}\frac{h(\eta_b^{q_a})}{h(\eta_a)}\left(1 - \frac{q_b}{q_a}\right) \tag{1.74}$$

Here, we used the definition $\eta_b^{q_a} := \frac{q_a^2 c^2}{p_b(0)^2}$. Inequality (1.74) is equal to assumption (1.65).

Last but not least, we discuss Assumption b)(vi). We notate

$$\eta_a = \frac{q_a^2 c^2}{p_a(0)^2}, \quad \eta_b = \frac{q_b^2 c^2}{p_b(0)^2}, \quad \eta_L = \frac{q_a^2 c^2}{p_a(L_a)^2}$$

and use $p_a(L_a) = f_a(L_a; p_a(0), q_a)$. We have to show

$$\frac{\partial_q h_u(p_a(L_a), q_a)}{\partial_p h_u(p_b(0), q_b))} < -\frac{\partial_p h_u(p_b(0), q_b)}{\partial_p h_u(p_a(L_a), q_a)}\partial_q f_a(L_a; p_a(0), q_a).$$

The inequality is always true for $q_a < 0$, therefore we consider $q_a > 0$. By Equation (1.66) and Equation (1.68) the inequality is equivalent to

$$\frac{p_b(0)}{q_a}\frac{\eta_L}{1-\eta_b} < \frac{p_a(L_a)}{p_b(0)}\frac{1-\eta_b}{1-\eta_L}\frac{p_a(L_a)}{q_a}\frac{\eta_L - \eta_a}{\eta_a(1-\eta_L)}.$$

This is equivalent to

$$1 < \frac{p_a(L_a)^2}{p_b(0)^2}\frac{(1-\eta_b)^2}{(1-\eta_L)^2}\frac{\eta_L - \eta_a}{\eta_L \eta_a}.$$

Hence, we have

$$1 < \frac{p_a(L_a)^2}{p_b(0)^2}\frac{(1-\eta_b)^2}{(1-\eta_L)^2}\left(\frac{1}{\eta_a} - \frac{1}{\eta_L}\right),$$

which is Assumption (1.63). $\qquad\square$

Finally, we can apply Theorem 2 once again and obtain the existence and uniqueness of real gas solutions on networks coupled by the Bernoulli invariant.

**Theorem 4** (Existence and Uniqueness for Ideal Gas)**.**
*Consider the arc coupling conditions* (Ideal)*, the node coupling conditions* (Pr) *or* (BI) *and the Kirchhoff conditions* (KirchhoffM)*. Let the pressure $p_r > 0$ be prescribed at the simple node $r \in \mathcal{V}$ and assume that the assumptions of Lemma 20 are fulfilled. Then there exists a constant $C(p_r, G) > 0$ such that for all massflows $m^{out} \in \mathbb{R}^{|\mathcal{V}|}$ with*

$$\mathbf{1}^T m^{out} = 0, \qquad \|m^{out}\| < C(p_r, G) \quad and \quad \frac{q_r^{out2} c^2}{p_r^2} < 1$$

*a unique solution $(p(0), p(L), q) \in \mathbb{R}^{|\mathcal{A}|} \times \mathbb{R}^{|\mathcal{A}|} \times \mathbb{R}^{|\mathcal{A}|}$ exists on the graph $G = (\mathcal{V}, \mathcal{A})$. This solution fulfills*

$$(\text{KirchhoffM}), (\text{Arc Coupling}) \ and \ (\text{Node Coupling}).$$

*Proof.* The proof is analogous to the proof of Theorem 3. With the results from Lemma 19 and Lemma 20, we see that the assumptions of Theorem 2 hold on the domains $\mathcal{D}(f_a), a \in \mathcal{A}$ and $\mathcal{D}(h_u), u \in \mathcal{V}$ as defined in Equations (1.60) and (1.62). In the node $r$, the prescribed pressure and flow fulfill $(q_r^{out})^2 c^2 < p_r^2$ by Assumption. The proof of Theorem 2 reveals that in each induction step, the pressure and flow values depend continuously on the boundary values. This implies that the values $p_a(0), p_a(L_a), q_a$ stay inside the domains of $f_a$ and $h_u$ for a sufficiently small prescribed flow $q^{out}$. Consequently, there exists a constant $C(p_r, G) > 0$ such that all obtained pressure-flow pairs are inside the corresponding domains $\mathcal{D}(f_a)$ and $\mathcal{D}(h_u)$ if $\|q^{out}\| < C(p_r, G)$ holds. The difference between the conditions (KirchhoffM) and (Kirchhoff) is discussed in Remark 10. $\quad\square$

## 1.4.2. Networks with Active Elements

In the following, we introduce active elements that are used to control the network. For the sake of simplicity, we restrict ourselves to the case of ideal gas with the pressure continuity as node coupling condition. We view all active elements as arcs that couple the pressure values via functions $f_a$. This is analogous to plain pipes with the only difference that the set of arcs decomposes into multiple subsets. We assume that all active elements conserve the mass flow. This means especially that compressors only use

external energy sources and do not use the transported gas as fuel. Active elements that only change the topology of the graph like valves are already included in our analysis up to this point. Since the pressure continuity coupling conditions fulfills Assumptions b) of Theorem 2 (see Lemma 18), we focus on showing the required assumptions for the functions $f_a$. The statement of the models is given in Koch et al. 2015, Chapter 2.

### Resistors

"Resistors are a surrogate modeling tool ...[for pressure loss induced by] flow diversion and turbulence in shaped pieces, measurement devices, curvature of the piping within compressor stations and pressure regulators, filter systems, reduced radii, and partially closed valves." (Koch et al. 2015, Section 2.3.2) The Darcy-Weisbach form of the pressure loss on an arc $a = (u, v) \in \mathcal{A}_{\mathrm{Re}} \subset \mathcal{A}$ is given by

$$f_a(x; p, q) := p - \frac{1}{2}\zeta \frac{q^2 c^2}{p}x, \quad x \in [0, 1], \tag{R}$$

where $\zeta > 0$ is a positive friction coefficient and the flow $q > 0$ is positive. The domain of $f_a$ is given by

$$\mathcal{D}(f_a) := \left\{(x, p, q) \in [0, 1] \times \mathbb{R}_{>0} \times \mathbb{R} \mid p^2 > q^2 c^2\right\}. \tag{1.75}$$

The space variable $x$ is artificially introduced with the goal to include the function $f_a$ in the framework of Theorem 2. The outgoing pressure and the ingoing pressure are therefore coupled via

$$p_a(1) = f(1; p_a(0), q_a).$$

We check the signs of the derivatives of $f_a$.

**Lemma 21** (Derivatives of the Resistor Function).
*Let $f_a$ be given by (R). Then, for $(x, p, q) \in \mathcal{D}(f_a)$ and $q > 0$, the following holds: The function $f_a$ is*

    *a)* strictly increasing *in the ingoing pressure, i.e.,*

$$\partial_p f_a(x; p, q) > 0,$$

    *b)* strictly decreasing *in the flow, i.e.,*

$$\partial_q f_a(x; p, q) < 0, \quad \text{for } x \in (0, 1]$$

c) strictly decreasing *in flow direction, i.e.,*

$$\partial_x f_a(x; p, q) < 0, \quad \text{for } x \in (0, 1].$$

*Proof.* The pressure derivative is given by

$$\partial_p f_a(x; p, q) = 1 + \frac{1}{2} \zeta \frac{q^2 c^2}{p^2} x > 0.$$

For the flow derivative, we obtain

$$\partial_q f_a(x; p, q) = -\zeta \frac{q c^2}{p} x < 0$$

and the space derivative is

$$\partial_x f_a(x; p, q) = -\frac{1}{2} \zeta \frac{q^2 c^2}{p} < 0. \qquad \square$$

### Turbo Compressors

Compressors and compressor station are included in the network to keep the pressure sufficiently high and counteract the pressure loss by friction. We focus on turbo compressors, which are described by a characteristic field for the specific change of adiabatic enthalpy $H_{\text{ad}}$ in dependence of the volume flow $Q = \frac{m}{\rho}$, where $m = qA$ is the mass flow through a pipe with cross-sectional area $A$. The specific change of the adiabatic enthalpy and can be thought of as the compression factor. It is given by

$$H_{\text{ad}}(p_u, p_v) := RTz_m \frac{\kappa}{\kappa - 1} \left[ \left( \frac{p_v}{p_u} \right)^{\frac{\kappa - 1}{\kappa}} - 1 \right], \qquad (1.76)$$

with the constant compressibility factor $z_m \in (0, 1]$ and the isentropic exponent $\kappa > 1$; see Koch et al. 2015, Section 2.3.5.1, Domschke et al. 2017, Section 7.4.7 and Cerbe 2008, Section 5.3.[1] The characteristic field of a turbo compressor is depicted in Figure 1.15. For a fixed compressor speed, the change of adiabatic enthalpy can be described by a continuously differentiable function (e.g. quadratic polynomials) $F : Q \mapsto F(Q)$ of the volume flow. Assume that the function $F$ is *positive* and *strictly monotonic decreasing* in the operational range $\mathcal{D}(F)$ (nonempty and open) of the

---

[1]Menon 2005, Section 4.4 gives values of $\kappa$ between 1.2 and 1.4, Domschke et al. 2017, Section 7.4.7 use $\kappa = 1.29$ and Koch et al. 2015, Section 2.3.5.1 use $\kappa = 1.296$.

Figure 1.15.: The characteristic field of a turbo compressor. The shaded area depicts the working range. The solid lines represent lines of equal compressor speed. The dashed lines represent lines of equal efficiency. The actual measure points are depicted by the solid dots.[2]

compressor. To obtain the outgoing pressure $p_a(1)$, where $a \in \mathcal{A}_{\mathrm{Co}} \subset \mathcal{A}$ is an arc with $L_a = 1$, one must solve the equation

$$H_{\mathrm{ad}} = F(Q)$$

with $Q = \frac{m_a}{\rho_a(0)}$. By Equation (1.76), we obtain

$$RT z_m \frac{\kappa}{\kappa - 1} \left[ \left( \frac{p_a(1)}{p_a(0)} \right)^{\frac{\kappa-1}{\kappa}} - 1 \right] = F(Q).$$

Hence,

$$\left( \frac{p_a(1)}{p_a(0)} \right)^{\frac{\kappa-1}{\kappa}} - 1 = (RT z_m)^{-1} \frac{\kappa - 1}{\kappa} F(Q),$$

which yields

$$\frac{p_a(1)}{p_a(0)} = \left[ (RT z_m)^{-1} \frac{\kappa - 1}{\kappa} F(Q) + 1 \right]^{\frac{\kappa}{\kappa-1}}.$$

[2]I am grateful to Martin Schmidt for his permission to use this picture.

Finally, we obtain

$$p_a(1) = \left[(RTz_m)^{-1}\frac{\kappa-1}{\kappa}F(Q)+1\right]^{\frac{\kappa}{\kappa-1}}p_a(0).$$

In our framework, once again utilizing the artificial space variable $x \in [0,1]$, this leads to the function

$$f_a(x;p,q) = p(1-x) + \left[(RTz_m)^{-1}\frac{\kappa-1}{\kappa}F(Q)+1\right]^{\frac{\kappa}{\kappa-1}}px, \qquad (C)$$

with $Q = \frac{Ac^2q}{p}$. It is defined on the domain

$$\mathcal{D}(f_a) := \left\{(x,p,q) \in [0,1] \times \mathbb{R}_{>0} \times \mathbb{R} \mid p^2 > q^2c^2, \ Q \in \mathcal{D}(F)\right\}. \qquad (1.77)$$

The signs of its derivatives are discussed in the following Lemma.

**Lemma 22** (Derivatives of the Compressor Function).
*Let $f_a$ be given by (C). Then, for $(x,p,q) \in \mathcal{D}(f_a)$, $q > 0$ and $x \in (0,1]$, the following holds: The function $f_a$ is*

a) strictly increasing *in the ingoing pressure, i.e.,*

$$\partial_p f_a(x;p,q) > 0,$$

b) strictly decreasing *in the flow, i.e.,*

$$\partial_q f_a(x;p,q) < 0,$$

c) strictly increasing *in flow direction, i.e.,*

$$\partial_x f_a(x;p,q) > 0.$$

*Proof.* We denote

$$G(Q) := (RTz_m)^{-1}\frac{\kappa-1}{\kappa}F(Q)+1 \qquad (1.78)$$

and obtain using the product rule

$$\partial_p f_a(x;p,q) = -x + \partial_p\left[G(Q)^{\frac{\kappa}{\kappa-1}}\right]px + G(Q)^{\frac{\kappa}{\kappa-1}}x$$

(using the chain rule for the second term)

$$= [G(Q)^{\frac{\kappa}{\kappa-1}} - 1]\, x + \frac{\kappa}{\kappa - 1} G(Q)^{\frac{\kappa}{\kappa-1}-1}\, \partial_Q G(Q)\, \partial_p Q\, p\, x$$

(by the definition (1.78) of $G$ and the simplification $\frac{\kappa}{\kappa-1} - 1 = \frac{1}{\kappa-1}$)

$$= [G(Q)^{\frac{\kappa}{\kappa-1}} - 1]\, x + (RT z_m)^{-1}\, G(Q)^{\frac{1}{\kappa-1}} \partial_Q F(Q)\, \partial_p Q\, p\, x$$

(as $\partial_p Q = -\frac{Ac^2 q}{p^2}$)

$$= [G(Q)^{\frac{\kappa}{\kappa-1}} - 1]\, x - (RT z_m)^{-1}\, G(Q)^{\frac{1}{\kappa-1}} \partial_Q F(Q)\, \frac{Ac^2 q}{p}\, x > 0,$$

because $G(Q) > 1$ and $\frac{\kappa}{\kappa-1} > 1$ imply that the first term is positive, while $\partial_Q F(Q) < 0$ ensures that the second term is positive.

For the derivative with respect to the flow, we obtain

$$\partial_q f_a(x; p, q) = \frac{\kappa}{\kappa - 1} G(Q)^{\frac{1}{\kappa-1}} \partial_Q G(Q)\, \partial_q Q$$

(by the definition (1.78) of $G$ and $\partial_q Q = \frac{Ac^2}{p}$)

$$= (RT z_m)^{-1} G(Q)^{\frac{1}{\kappa-1}} \partial_Q F(Q)\, \frac{Ac^2}{p} < 0,$$

because $\partial_Q F(Q)$ is negative, while the remaining terms are positive. For the space derivative, direct calculation shows

$$\partial_x f_a(x; p, q) = (G(Q)^{\frac{\kappa}{\kappa-1}} - 1)\, p > 0,$$

since $\kappa > 1$, which implies $\frac{\kappa}{\kappa-1} > 1$, and $G(Q) > 1$ hold. $\qquad\square$

With these results at hand, we can include turbo compressors and resistors in the gas network and obtain a similar existence-and-uniqueness result as in the case of passive networks. However, due to the construction of the solution, it is not possible to include turbo compressors in fundamental cycles, because they increase the pressure in flow direction.

**Theorem 5** (Existence and Uniqueness Including Active Elements)**.**
*Consider the set of arcs $\mathcal{A} = \mathcal{A}_{Pi} \,\dot\cup\, \mathcal{A}_{Re} \,\dot\cup\, \mathcal{A}_{Co}$, where the functions $f_a$ are given by* (Ideal) *for $a \in \mathcal{A}_{Pi}$, by* (R) *for $a \in \mathcal{A}_{Re}$ and by* (C) *for $a \in \mathcal{A}_{Co}$. We assume that all active elements have a bypass mode that treats them as*

*normal pipes for negative flow with the corresponding function like in* (Ideal).
*Let the states at the nodes $u \in \mathcal{V}$ be coupled by the pressure continuity* (Pr).
*Assume that the arcs $\mathcal{A}_{Co}$ are not contained in any fundamental cycle. Let
a pressure value $p_r > 0$ be prescribed at a simple node $r \in \mathcal{V}$.*

*Then there exists a constant $C(p_r, G) > 0$ such that for all $m^{out} \in \mathbb{R}^{|\mathcal{V}|}$
with*

$$\mathbf{1}^T m^{out} = 0, \quad \|m^{out}\| < C(p_r, G) \quad and \quad (0, p_r, m_r^{out}) \in \mathcal{D}(f_a), \ \forall a \in \mathcal{A}$$

*a unique solution $(p(0), p(L), q) \in \mathbb{R}^{|\mathcal{A}|} \times \mathbb{R}^{|\mathcal{A}|} \times \mathbb{R}^{|\mathcal{A}|}$ exists on the graph
$G = (\mathcal{V}, \mathcal{A})$. This solution fulfills*

(KirchhoffM), (Arc Coupling) *and* (Node Coupling).

*Proof.* On arcs that are no compressors, i.e., $a \in \mathcal{A}_{Pi} \cup \mathcal{A}_{Re}$, we see by
Lemma 18, Lemma 19 and Lemma 21 that the assumptions of Theorem 2
are fulfilled. For compressor arcs, we see by Lemma 22 that all assumptions
except a)(iii) and b)(vii) are fulfilled. Knowing that these assumptions
hold, we can follow the proof of Theorem 2 up to the point, where we need
to find a flow value $\lambda_-$ with $H(\lambda_-) < 0$ and a $\lambda_+$ with $H(\lambda_+) > 0$. Here
the monotonicity in space induced by Assumptions a)(iii) and b)(vii) is
required for all arcs lying on the constructed path $\mathcal{P}$ between $v_{cl}$ and $v_{cr}$.
We need to show that the construction of the path of unidirectional flow is
possible without using compressor arcs $a \in \mathcal{A}_{Co}$. By assumption the arcs
$a \in \mathcal{A}_{Co}$ are not part of a fundamental cycle. However, since the arc $c$ is
part of a fundamental cycle, the nodes $v_{cl}$ and $v_{cr}$ remain connected even
after the removal of the arcs in $\mathcal{A}_{Co}$. Also, the flows on the compressor
arcs do not depend on $\lambda$, since by Lemma 12, the redefined outflow vectors
sum to zero. To be precise: For an arc $a = (u, v) \in \mathcal{A}_{Co}$, let $v$ be the node
that lies in the same subgraph $\mathcal{G}_c = (\mathcal{V}_c, \mathcal{A}_c)$ with redefined outflows $f^c$ as
$v_{cl}$ and $v_{cr}$ after the removal of $a$. Then, by Lemma 12, the flow $f_v$ can be
calculated as

$$f_v = -\sum_{u \in \mathcal{V}_c} f_u^c = -\sum_{u \in \mathcal{V}_c \setminus \{v_{cl}, v_{cr}\}} f_u^c + \lambda - \lambda = -\sum_{u \in \mathcal{V}_c \setminus \{v_{cl}, v_{cr}\}} f_u^c.$$

The case of removing multiple compressor arcs is analogous with the tools
of Lemma 12. After completing the reduction, we end up with a subgraph
$\mathcal{G}_c = (\mathcal{V}_c, \mathcal{A}_c)$ with the new outflows $f^c$ as defined in Lemma 12. We apply
Lemma 14 to deduce that the choice

$$\lambda^+ = \sum_{u \in \mathcal{V}_c} |f_u^c|$$

leads to a path $\mathcal{P}$ from $v_{\text{cl}}$ to $v_{\text{cr}}$ that does not contain arcs $a \in \mathcal{A}_{\text{Co}}$ such that

$$q_a \geq 0 \quad \forall a \in \mathcal{P}$$

with $q_a > 0$ for at least one $a \in \mathcal{P}$. The choice

$$\lambda^- = - \sum_{u \in \mathcal{V}_c} |f_u^c|$$

leads to the same result with reversed signs. The rest of the proof poses no further problems. Choosing the prescribed flows $m^{\text{out}}$ sufficiently small ensures that the pressure and flow values stay inside the open domains of the functions $f_a$. $\qquad\square$

## 1.5. Examples

Throughout this section, we will demonstrate how to explicitly solve some example networks and we will give numerical results for realistic data.

### 1.5.1. Real Gas with Pressure Continuity

Consider $f_a$ as given by (Real) and the pressure coupling conditions (Pr). Throughout this section, we work under the assumptions of Theorem 3. We will use the examples given in Gugat, Schultz, and Wintergerst 2018, Section 4.1–4.3. The first network is a graph with no cycles, a tree. The second network is a graph with one cycle and the third example is a diamond graph, which is a network with two interconnected cycles.

**Tree Network**

Consider the topology shown in Figure 1.16. Let the pressure $p_r > 0$ be prescribed in the node $r$ and let outflows be prescribed at the boundary nodes. Let the outflow vector

$$q^{\text{out}} = (q_r^{\text{out}}, \, 0, \, q_v^{\text{out}}, \, q_w^{\text{out}}, \, q_s^{\text{out}})$$

be given. We assume that $r$ is a source and $v, w, s$ are sinks, hence $q_r^{\text{out}} < 0$ holds. The flow values on the arcs are given by

$$q_a = -q_r^{\text{out}} = \frac{D_b^2}{D_a^2} q_b^{\text{out}} + \frac{D_c^2}{D_a^2} q_c^{\text{out}} + \frac{D_d^2}{D_a^2} q_d^{\text{out}},$$
$$q_b = q_b^{\text{out}}, \quad q_c = q_c^{\text{out}}, \quad q_d = q_d^{\text{out}},$$

Figure 1.16.: Tree network



Figure 1.17.: The solution on a tree with the pressure values at the nodes (in $10^5$ Pa) and the flow values at the arcs (in $\mathrm{kg\,s^{-1}\,m^{-2}}$)

where $D_a > 0$ denotes the diameter of pipe $a$. Due to the pressure continuity through the node we will denote $p_u = p_a(x_a(u))$, $\forall a \in \delta(u)$. Since the flow along every arc is known, the pressure values can be calculated starting at the node $r$:

$$p_u = f_a(L_a; q_a, p_r), \quad p_v = f_b(L_b; q_b, p_u),$$
$$p_w = f_c(L_c; q_c, p_u), \quad p_s = f_d(L_d; q_d, p_u).$$

The solution for specific data can be seen in Figure 1.17.

**One Circle**

Let us now consider a network with two parallel pipes as shown in Figure 1.18. The node $r$ is a source and the node $w$ is a sink, i.e.,

$$q^{\mathrm{out}} = (q_r^{\mathrm{out}}, 0, 0, q_w^{\mathrm{out}}), \quad \text{with} \quad q_r^{\mathrm{out}} < 0,\ q_w^{\mathrm{out}} > 0,$$

Figure 1.18.: Two parallel pipes

where $D_a^2 q_r^{\text{out}} = -D_d^2 q_w^{\text{out}}$. We prescribe the pressure at the node $r$. Unlike in the case of a tree network it is not obvious how to calculate the flow value on the edges $b$ and $c$. Due to the Kirchhoff condition

$$D_a^2 q_a = D_b^2 q_b + D_c^2 q_c$$

in the node $u$, we make the ansatz

$$q_b(\lambda) = \lambda D_a^2 D_b^{-2} q_a \quad \text{and} \quad q_c(\lambda) = (1-\lambda) D_a^2 D_c^{-2} q_a, \quad \text{for } \lambda \in [0,1],$$

with the flow $q_a = -q_r^{\text{out}} > 0$. This is equivalent to taking the linear combination of the specific mass flow solution $m_s = (m_w^{\text{out}}, 0, m_w^{\text{out}}, m_w^{\text{out}})$ and a vector of the kernel $m_k = (0, m_w^{\text{out}}, -m_w^{\text{out}}, 0)$ and set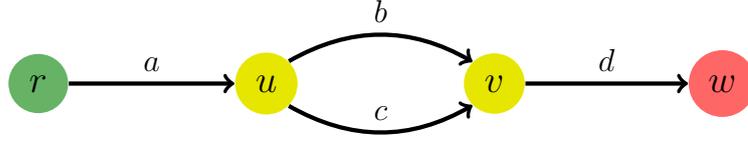ting $m(\lambda) = m_s + \lambda m_k$. By the pressure coupling condition in the node $v$, we have to choose $\lambda$ such that

$$f_b(L_b; q_b(\lambda), p_u) = f_c(L_c; q_c(\lambda), p_u),$$

where $p_u = f_a(L_a; q_w^{\text{out}}, p_r)$. It is equivalent to find a root of the function

$$H(\lambda) := f_b(L_b; q_b(\lambda), p_u) - f_c(L_c; q_c(\lambda), p_u).$$

By definition, $q_b$ is strictly increasing in $\lambda$ and $q_c$ is strictly decreasing. The functions $f_b(L_b; \cdot, p_u)$ and $f_c(L_b; \cdot, p_u)$ are strictly decreasing. Therefore, the function $H$ is strictly decreasing and continuous.

For $\lambda = 0$, we obtain

$$H(0) = f_b(L_b; 0, p_u) - f_c(L_c; D_a^2 D_c^{-2} q_a, p_u) > p_u - p_u = 0,$$

because $f_b(L_b; q, p_u) = p_u$ for $q = 0$ and $f_b(L_c; D_a^2 D_c^{-2} q_a, p_u) < p_u$ as $f_b$ strictly decreases in flow direction.

For $\lambda = 1$, we obtain

$$H(1) = f_b(L_b; D_a^2 D_b^{-2} q_a, p_u) - f_c(L_c; 0, p_u) < p_u - p_u = 0,$$

because $f_b(L_b; D_a^2 D_b^{-2} q_a, p_u) < p_u$ and $f_c(L_c; 0, p_u) = p_u$. Consequently, Bolzano's intermediate value theorem states that there is a unique solution

Figure 1.19.: Diamond graph

$\lambda^* \in (0,1)$ such that $H(\lambda^*) = 0$. Consequently, the flow values on the edges can be calculated via

$$q_a = -q_r^{\text{out}}, \quad q_b = q_b(\lambda^*), \quad q_c = q_c(\lambda^*), \quad q_d = q_w^{\text{out}}$$

and the pressure values can be calculated iteratively starting from the root node

$$\begin{aligned}
p_u &= f_a(L_a; q_a, p_r), \\
p_v &= f_b(L_b; q_b, p_u) = f_c(L_c; q_c, p_u), \\
p_w &= f_d(L_d; q_d, p_v).
\end{aligned}$$

### Diamond Graph

The diamond graph (Figure 1.19) is an example of a network with two interconnected cycles. We prescribe the pressure in the node $r$ and have an inflow through node $r$ and an outflow through node $t$, i.e.,

$$q^{\text{out}} = (q_r^{\text{out}},\ 0,\ 0,\ 0,\ 0,\ q_t^{\text{out}}) \quad \text{with} \quad q_r^{\text{out}} < 0,\ q_t^{\text{out}} > 0,$$

and $D_a^2 q_r^{\text{out}} = -D_g^2 q_t^{\text{out}}$. To obtain a more direct solution approach, we want to express the flow value at the arc $d$ as a function of the pressure in node $v$ and the pressure in node $w$. The following Lemma was formulated in Gugat, Schultz, and Wintergerst 2018, Lemma 6.

**Lemma 23** (Compensatory flow).
*On the arc $a = (u, v)$, let the pressure values in $u$ and $v$ be given by*

$p_u, p_v \in (0, |\alpha|^{-1})$. *Furthermore let the nonzero distance between $p_u$ and $p_v$ be small enough in the sense that*

$$|p_u - p_v| < p_u + p_v - 2 \frac{\sqrt{\alpha^{-1}(p_v - p_u) - \alpha^{-2}\ln\left(\frac{z(p_v)}{z(p_u)}\right)}}{\sqrt{\ln\left(\frac{p_v}{p_u}\frac{z(p_u)}{z(p_v)}\right) - \frac{1}{2}\operatorname{sign}(p_0 - p_1)\int_0^{L_a}\theta_a(s)\,\mathrm{d}s}}. \quad (1.79)$$

*Then, the sign of the flow is determined by*

$$\operatorname{sign}(q_a) = \operatorname{sign}(p_u - p_v) \quad (1.80)$$

*s and the flow value on the arc a is given by*

$$q_a = \operatorname{sign}(p_u - p_v)\sqrt{\frac{\alpha^{-1}(p_v - p_u) - \alpha^{-2}\ln\left(\frac{z(p_v)}{z(p_u)}\right)}{RT\left[\ln\left(\frac{p_v}{p_u}\frac{z(p_u)}{z(p_v)}\right) - \frac{1}{2}\operatorname{sign}(p_u - p_v)\int_0^{L_a}\theta_a(s)\,\mathrm{d}s\right]}} \quad (1.81)$$

*In the case of $p_u = p_v$, the flow $q_a$ is zero.*

*Proof.* We will assume subsonic flow $|q_a|\sqrt{RT} < \min\{p_u, p_v\}$ for the moment and show later that it indeed holds. By Equation (1.7), we have

$$F(p_v, q_a) - F(p_u, q_a) = -\frac{1}{2}q_a|q_a|RT\int_0^{L_a}\theta_a(s)\,\mathrm{d}s.$$

The monotonicity of $F$ with respect to the first argument stated in Lemma 1 shows

$$\operatorname{sign}(q_a) = \operatorname{sign}(F(p_u, q_a) - F(p_v, q_a)) = \operatorname{sign}(p_u - p_v),$$

which is Equation (1.80). Inserting the definition of $F$, see Equation (1.6), in (1.7) yields

$$-\frac{1}{2}q_a|q_a|RT\int_0^{L_a}\theta_a(s)\,\mathrm{d}s = \alpha^{-1}p_v - \alpha^{-2}\ln(z(p_v)) + RTq_a^2\ln\left(\frac{z(p_v)}{p_v}\right)$$
$$- \alpha^{-1}p_u + \alpha^{-2}\ln(z(p_u)) - RTq_a^2\ln\left(\frac{z(p_u)}{p_u}\right)$$

which is equivalent to

$$RTq_a^2\left[\ln\left(\frac{z(p_u)}{p_u}\right) - \ln\left(\frac{z(p_v)}{p_v}\right) - \frac{1}{2}\operatorname{sign}(q_a)\int_0^{L_a}\theta_a(s)\,\mathrm{d}s\right] =$$
$$\alpha^{-1}(p_v - p_u) - \alpha^{-2}\ln\left(\frac{z(p_v)}{z(p_u)}\right)$$

Hence,

$$q_a^2 = \frac{\alpha^{-1}(p_v - p_u) - \alpha^{-2}\ln\left(\frac{z(p_v)}{z(p_u)}\right)}{RT\left[\ln\left(\frac{z(p_u)}{p_u}\right) - \ln\left(\frac{z(p_v)}{p_v}\right) - \frac{1}{2}\operatorname{sign}(q_a)\int_0^{L_a}\theta_a(s)\,\mathrm{d}s\right]}$$

holds, which implies (1.81). It remains to show that the flow stays subsonic under condition (1.79), i.e., $\sqrt{RT}q_a < \min\{p_u, p_v\}$. We have

$$\min\{p_u, p_v\} = \frac{1}{2}(p_u + p_v - |p_u - p_v|)$$

$$> \frac{\sqrt{\alpha^{-1}(p_v - p_u) - \alpha^{-2}\ln\left(\frac{z(p_v)}{z(p_u)}\right)}}{\sqrt{\ln\left(\frac{p_v}{p_u}\frac{z(p_u)}{z(p_v)}\right) - \frac{1}{2}\operatorname{sign}(p_0 - p_1)\int_0^{L_a}\theta_a(s)\,\mathrm{d}s}}$$

$$= |q_a|\sqrt{RT}.$$

This concludes the proof. □

Additionally, for the construction of a solution on the diamond graph, the monotonicity of $q_a$ with respect to $p_u$ and $p_v$ is required; see Gugat, Schultz, and Wintergerst 2018, Lemma 7.

**Lemma 24** (Monotonicity of the Compensatory Flow).
*Under the assumptions of Lemma 23, the derivatives of $q_a$ satisfy*

$$\partial_{p_u}q_a > 0 \quad and \quad \partial_{p_v}q_a < 0.$$

*This means $q_a$ is strictly increasing as a function of $p_u$ and strictly decreasing as a function of $p_v$.*

*Proof.* Implicit differentiation of Equation (1.7) with respect to $p_u$ leads to

$$-\partial_p F(p_u, q_a) = -RT|q_a|\partial_{p_u}q_a\int_0^{L_a}\theta_a(s)\,\mathrm{d}s$$

$$\partial_{p_u}q_a = \frac{\partial_p F(p_u, q)}{RT\int_0^{L_a}\theta_a(s)\,\mathrm{d}s|q_a|}$$

and differentiation with respect to $p_v$ leads to

$$\partial_p F(p_v, q_a) = -RT|q_a|\partial_{p_v}q_a\int_0^{L_a}\theta_a(s)\,\mathrm{d}s$$

$$\partial_{p_v}q_a = -\frac{\partial_p F(p_v, q)}{RT|q_a|\int_0^{L_a}\theta_a(s)\,\mathrm{d}s}.$$

By Lemma 1, $\partial_p F(p_u, q_a)$ and $\partial_p F(p_v, q_a)$ are positive. Hence,

$$\partial_{p_u} q_a > 0 \quad \text{and} \quad \partial_{p_v} q_a < 0$$

follows. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

With the two Lemmas, we are equipped to derive a solution on the diamond graph. We make the same ansatz as in the case of one circle and set

$$q_b(\lambda) = \lambda D_a^2 D_b^{-2} q_a \quad \text{and} \quad q_c(\lambda) = (1 - \lambda) D_a^2 D_c^{-2} q_a, \quad \text{for } \lambda \in [0, 1].$$

Then, the pressure values in $u, v$ and $w$ can be calculated as

$$\begin{aligned} p_u &= f_a(L_a; p_r, q_a), \\ p_v(\lambda) &= f_b(L_b; p_u, q_b(\lambda)), \\ p_w(\lambda) &= f_c(L_c; p_u, q_c(\lambda)). \end{aligned}$$

By Lemma 23, the flow through the arc $d = (w, v)$ has to be

$$q_d(\lambda) = Q_d(p_w(\lambda), p_v(\lambda)),$$

where $Q_d(p_w(\lambda), p_v(\lambda))$ is given by the right hand side of Equation (1.81). Consequently, by the Kirchhoff conditions in $u$ and $v$, we obtain

$$\begin{aligned} q_e(\lambda) &= D_b^2 D_e^{-2} q_b(\lambda) - D_d^2 D_e^{-2} q_d(\lambda), \\ q_f(\lambda) &= D_c^2 D_f^{-2} q_c(\lambda) + D_d^2 D_f^{-2} q_d(\lambda). \end{aligned}$$

The pressure continuity in the node $s$ requires $\lambda$ to be chosen such that the function value of

$$H(\lambda) := f_e(L_e; p_v(\lambda), q_e(\lambda)) - f_f(L_f; p_w(\lambda), q_f(\lambda))$$

becomes zero. For a better overview, we note the monotonicity properties with respect to $\lambda$ in Table 1.2. The arguments for the monotonicity of $q_b, q_c, p_v$ and $p_w$ are straightforward. By Lemma 24, $q_d$ is strictly increasing as a function of $p_v$ and strictly decreasing as a function of $p_w$. Hence $q_d$ is strictly decreasing in $\lambda$. This directly implies that $q_e$ is strictly increasing and $q_f$ is strictly decreasing. Both $f_e$ and $f_f$ are strictly increasing in the pressure argument and strictly decreasing in the flow argument. Therefore, the monotonicity of $p_v$ and $q_e$ implies that $f_e$ strictly decreases in $\lambda$ and the monotonicity of $p_w$ and $q_f$ implies that $f_f$ strictly increases in $\lambda$.

Table 1.2.: Monotonicity with respect to $\lambda$

| Flow | Monotonicity (Flow) | Pressure | Monotonicity (Pressure) |
|------|--------------------|----------|------------------------|
| $q_b$ | strictly increasing | $p_v$ | strictly decreasing |
| $q_c$ | strictly decreasing | $p_w$ | strictly increasing |
| $q_d$ | strictly decreasing | | |
| $q_e$ | strictly increasing | $f_e(L_e; p_v(\cdot), q_e(\cdot))$ | strictly decreasing |
| $q_f$ | strictly decreasing | $f_f(L_f; p_w(\cdot), q_f(\cdot))$ | strictly increasing |



Figure 1.20.: Flow values in $\mathrm{kg\,s^{-1}\,m^{-2}}$ at each arc.



Figure 1.21.: Pressure values in $10^5\,\mathrm{Pa}$ in each node.

Consequently, $H$ is a strictly decreasing function. As in the example with one circle, we consider the cases $\lambda = 0$ and $\lambda = 1$. For $\lambda = 0$, we obtain

$$q_b(0) = 0 \quad \text{and} \quad q_c(0) = D_a^2 D_c^{-2} q_a,$$

which implies

$$p_v(0) = p_u \quad \text{and} \quad p_w(0) < p_u.$$

Consequently, $q_d(0) > 0$, leading to

$$q_e(0) = -D_d^2 D_e^{-2} q_d(0) < 0 \quad \text{and} \quad q_f(0) = D_a^2 D_f^{-2} q_a + D_d^2 D_f^{-2} q_d(0) > 0.$$

Subsequently,

$$\begin{aligned} H(0) &= f_e(L_e; p_v(0), q_e(0)) - f_f(L_f; p_w(0), q_f(0)) \\ &> p_v(0) - p_w(0) > p_u - p_u = 0 \end{aligned}$$

holds. For $\lambda = 1$, we obtain

$$q_b(1) = D_a^2 D_b^{-2} q_a \quad \text{and} \quad q_c(1) = 0,$$

which implies

$$p_v(1) < p_u \quad \text{and} \quad p_w(1) = p_u.$$

Consequently, $q_d(1) < 0$, leading to

$$q_e(1) = D_a^2 D_e^{-2} q_a - D_d^2 D_e^{-2} q_d(1) > 0 \quad \text{and} \quad q_f(1) = D_d^2 D_f^{-2} q_d(1) < 0.$$

Subsequently,

$$\begin{aligned} H(1) &= f_e(L_e; p_v(1), q_e(1)) - f_f(L_f; p_w(1), q_f(1)) \\ &< p_v(1) - p_w(1) < p_u - p_u = 0. \end{aligned}$$

Therefore, as $H$ is strictly decreasing and continuous, Bolzano's intermediate value theorem implies the existence of a unique $\lambda^* \in (0, 1)$ such that $H(\lambda^*) = 0$. The flow solution on the diamond graph can then be calculated as

$$(-q_r^{\text{out}}, \ q_b(\lambda^*), \ q_c(\lambda^*), \ q_d(\lambda^*), \ q_e(\lambda^*), \ q_f(\lambda^*), \ q_t^{\text{out}})$$

and the pressure values can be obtained by

$$(p_r, \ f_a(L_a; p_r, q_a), \ p_v(\lambda^*), \ p_w(\lambda^*), \ p_s(\lambda^*), \ f_g(L_g; p_s(\lambda^*), q_g)).$$

The flow solution for specific data is shown in Figure 1.20 and the pressure solution is shown in Figure 1.21.

## 1.5.2. Ideal Gas with Bernoulli Invariant Coupling

Consider $f_a$ as given by (Ideal) and the pressure coupling conditions (BI). Throughout this section, we work under the assumptions of Theorem 4. We will consider the first and second example of the previous subsection. The first network is a graph with no cycles, a tree. The second network is a graph with one cycle.

**Tree Network**

Consider the topology shown in Figure 1.16. Let the outflow vector

$$q^{\text{out}} = (q_r^{\text{out}}, \, 0, \, q_v^{\text{out}}, \, q_w^{\text{out}}, \, q_s^{\text{out}}),$$

be given. We assume that $r$ is a source and $v, w, s$ are sinks, hence $q_r^{\text{out}} < 0$ holds. The flow values on the arcs are given by

$$q_a = -q_r^{\text{out}} = \frac{D_b^2}{D_a^2} q_b^{\text{out}} + \frac{D_c^2}{D_a^2} q_c^{\text{out}} + \frac{D_d^2}{D_a^2} q_d^{\text{out}}$$
$$q_b = q_b^{\text{out}}, \quad q_c = q_c^{\text{out}}, \quad q_d = q_d^{\text{out}}.$$

Unlike in the case, where we worked with pressure continuity coupling, now, the pressure values can change over a node. This means $p_a(L_a)$ is not equal to $p_b(0)$ in general. Denote the function that solves

$$h_u(p_a, q_a) = h_u\big(P_b(p_a, q_a, q_b), q_b\big)$$

for prescribed $p_a, q_a$ and $q_b$ by $P_b(p_a, q_a, q_b)$. As discussed in the proof of Lemma 20, it can be explicitly calculated as

$$P_b(p_a, q_a, q_b) = q_b c \sqrt{-\left[W_0\left(-\frac{q_b^2}{c^2 \rho_a^2} \exp(-\eta_a)\right)\right]^{-1}}. \qquad (1.82)$$

Then, the pressure values can be calculated iteratively by

$$p_a(0) = p_r, \qquad\qquad p_a(L_a) = f_a(L_a; p_r, q_a),$$
$$p_b(0) = P_b(p_a(L_a), q_a, q_b), \qquad p_b(L_b) = f_b(L_b; p_b(0), q_b),$$
$$p_c(0) = P_c(p_a(L_a), q_a, q_c), \qquad p_c(L_c) = f_c(L_c; p_c(0), q_c),$$
$$p_d(0) = P_d(p_a(L_a), q_a, q_d), \qquad p_d(L_d) = f_d(L_d; p_d(0), q_d).$$

The simulation for the same data as in Section 1.5.1 is depicted in Figure 1.22.

Figure 1.22.: Solution on a tree. The pressure values at the beginning and the end of each pipe (in $10^5$ Pa) and the flow values at the arcs (in $\mathrm{kg\,s^{-1}\,m^{-2}}$) are shown. The value left of the middle node is the pressure value $p_a(L_a)$. The value right of the middle node is the value $p_b(0) = p_c(0) = p_d(0)$. The values are identical, because the flow values are on the arcs are identical.

**One Circle**

Consider a network with two parallel pipes as shown in Figure 1.18. For this example, we consider the case of equal pipe diameters $D_a$, $a \in \mathcal{A}$. The network has one source $r$ and one sink $w$, i.e.,

$$q^{\mathrm{out}} = (q_r^{\mathrm{out}}, 0, 0, q_w^{\mathrm{out}}), \quad \text{with} \quad q_r^{\mathrm{out}} < 0, \; q_w^{\mathrm{out}} > 0,$$

where $q_r^{\mathrm{out}} = -q_w^{\mathrm{out}}$. We proceed similarly to Section 1.5.1, yet we have to keep the coupling conditions in the nodes $u$ and $v$ in mind. Again, we make the ansatz

$$q_b(\lambda) = \lambda q_a \quad \text{and} \quad q_c(\lambda) = (1 - \lambda)q_a, \quad \text{for } \lambda \in [0, 1].$$

Following the idea of the proof of Theorem 2, we cut the arc $b$ in half and obtain the artificial nodes $u_l$ and $v_r$ and the corresponding arcs $b_l$ and $b_r$ with the lengths $L_{b_l} = L_b/2$ and $L_{b_r} = L_b/2$; see Figure 1.23. The flow on the arcs $b_l$ and $b_r$ is equal to $q_b(\lambda)$. For the pressure values on the new arcs, we have

$$p_{b_l}(L_{b_l}; \lambda) = f_b(L_{b_l}; p_{b_l}(0; \lambda), q_{b_l}(\lambda)),$$
$$p_{b_r}(0; \lambda) = f_b(L_{b_r}; p_{c_r}(L_{b_r}; \lambda), -q_{b_r}(\lambda)).$$

Figure 1.23.: The circle of Figure 1.18 after performing the cut.

The preceding pressure values are determined by

$$
\begin{aligned}
p_a(0) &= p_r, \\
p_a(L_a) &= f_a(L_a; p_r, q_a), \\
p_{b_l}(0; \lambda) &= P_b(p_a(L_a), q_a, q_b(\lambda)), \\
p_c(0; \lambda) &= P_c(p_a(L_a), q_a, q_c(\lambda)), \\
p_c(L_c; \lambda) &= f_c(L_c; p_c(0; \lambda), q_c(\lambda)), \\
p_{b_r}(L_{b_r}; \lambda) &= P_{b_r}(p_c(L_c; \lambda), -q_b(\lambda), q_c(\lambda)), \\
p_{b_r}(0; \lambda) &= f_b(L_{b_r}; p_{b_r}(0; \lambda), -q_b(\lambda)).
\end{aligned}
$$

The functions $P_b, P_c$ and $P_{b_r}$ are defined analogously to Equation (1.82) as long the flow argument of the outgoing arc is not zero. We need to determine the scalar $\lambda$ as a root of the function

$$
H(\lambda) := p_{b_l}(L_{b_l}; \lambda) - p_{b_r}(0; \lambda).
$$

We have (ommiting the argument $\lambda$ for readability), by the chain rule,

$$
\partial_\lambda p_{b_l}(L_{b_l}) = \partial_p f_b(L_{b_l}; p_{b_l}(0), q_b)\, \partial_{q_b} p_{b_l}(0)\, \partial_\lambda q_b + \partial_q f_b(L_{b_l}; p_{b_l}(0), q_b)\, \partial_\lambda q_b
$$

(by Lemma 15, Equation (1.46a))

$$
\begin{aligned}
= &-\partial_p f_b(L_{b_l}; p_{b_l}(0), q_b)\, \frac{\partial_q h_u(p_{b_l}(0), q_b)}{\partial_y h_u(p_{b_l}(0), q_b)}\, \partial_\lambda q_b \\
&+ \partial_q f_b(L_{b_l}; p_{b_l}(0), q_b)\, \partial_\lambda q_b < 0,
\end{aligned} \tag{1.83}
$$

since Assumption b)(v) of Theorem 2 holds due to Lemma 20 and $\partial_\lambda q_{b_l}$ is positive.

Differentiating the pressure $p_{b_r}(0)$ at the other end of the cut leads to

$$
\begin{aligned}
\partial_\lambda p_{b_r}(0) = &\, \partial_p f_b\big(L_{b_r}; p_{b_r}(L_{b_r}), -q_b\big) \partial_\lambda p_{b_r}(L_{b_r}) \\
&- \partial_q f_b\big(L_{b_r}; p_{b_r}(L_{b_r}), -q_b\big) \partial_\lambda q_b,
\end{aligned} \tag{1.84}
$$

where derivative of the pressure $p_{b_r}(L_{b_r})$ is given by

$$\partial_\lambda p_{b_r}(L_{b_r}) = \partial_{p_c} P_{b_r}\big(p_c(L_c), -q_b, q_c\big)\partial_\lambda p_c(L_c) - \partial_{q_b} P_{b_r}\big(p_c(L_c), -q_b, q_c\big)\partial_\lambda q_b$$
$$+ \partial_{q_c} P_{b_r}\big(p_c(L_c), -q_b, q_c\big)\partial_\lambda q_c. \tag{1.85}$$

Furthermore, we obtain

$$\partial_\lambda p_c(L_c) = \partial_p f_c(L_c; p_c(0), q_c)\,\partial_\lambda p_c(0) + \partial_q f_c(L_c; p_c(0), q_c)\,\partial_\lambda q_c. \tag{1.86}$$

Similarly to the proof of Lemma 16, we denote

$$T_{p_c} := \partial_p f_b\big(L_{b_r}; p_{b_r}(L_{b_r}), -q_b\big)\partial_{p_c} P_{b_r}\big(p_c(L_c), -q_b, q_c\big)\partial_p f_c(L_c; p_c(0), q_c)$$
$$T_{q_b} := -\partial_p f_b\big(L_{b_r}; p_{b_r}(L_{b_r}), -q_b\big)\partial_{q_b} P_{b_r}\big(p_c(L_c), q_b, q_c\big)$$
$$- \partial_q f_b\big(L_{b_r}; p_{b_r}(L_{b_r}), -q_b\big)$$
$$T_{q_c} := \partial_{p_c} P_{b_r}\big(p_c(L_c), -q_b, q_c\big)\partial_q f_c(L_c; p_c(0), q_c) + \partial_{q_c} P_{b_r}\big(p_c(L_c), -q_b, q_c\big)$$

Inserting Equation (1.86) into Equation (1.85) leads to

$$\partial_\lambda p_{b_r}(L_{b_r}) = \partial_{p_c} P_{b_r}\big(p_c(L_c), -q_b, q_c\big)\Big[\partial_p f_c(L_c; p_c(0), q_c)\,\partial_\lambda p_c(0)$$
$$+ \partial_q f_c(L_c; p_c(0), q_c)\,\partial_\lambda q_c\Big] - \partial_{q_b} P_{b_r}\big(p_c(L_c), -q_b, q_c\big)\partial_\lambda q_b$$
$$+ \partial_{q_c} P_{b_r}\big(p_c(L_c), -q_b, q_c\big)\partial_\lambda q_c$$
$$= \partial_{p_c} P_{b_r}\big(p_c(L_c), -q_b, q_c\big)\partial_p f_c(L_c; p_c(0), q_c)\,\partial_\lambda p_c(0) + T_{q_c}\,\partial_\lambda q_c$$
$$- \partial_{q_b} P_{b_r}\big(p_c(L_c), -q_b, q_c\big)\partial_\lambda q_b. \tag{1.87}$$

Inserting Equation (1.87) into Equation (1.84) yields

$$\partial_\lambda p_{b_r}(0) = \partial_p f_b\big(L_{b_r}; p_{b_r}(L_{b_r}), -q_b\big)\Big[\partial_{p_c} P_{b_r}\big(p_c(L_c), -q_b, q_c\big)\partial_p f_c(L_c; p_c(0), q_c)$$
$$\partial_\lambda p_c(0) + T_{q_c}\,\partial_\lambda q_c - \partial_{q_b} P_{b_r}\big(p_c(L_c), -q_b, q_c\big)\partial_\lambda q_b\Big]$$
$$- \partial_q f_b\big(L_{b_r}; p_{b_r}(L_{b_r}), -q_b\big)\partial_\lambda q_b$$
$$= T_{p_c}\partial_\lambda p_c(0) + \partial_p f_b\big(L_{b_r}; p_{b_r}(L_{b_r}), -q_b\big)T_{q_c}\,\partial_\lambda q_c + T_{q_b}\partial_\lambda q_b \tag{1.88}$$

We know that $\partial_\lambda q_c < 0$ and $\partial_\lambda q_b > 0$. For the derivative of $p_c(0)$, we obtain

$$\partial_\lambda p_c(0) = \partial_{q_c} P_c(p_a(L_a), q_a, q_c)\partial_\lambda q_c$$

(by Lemma 15, Equation (1.47b))

$$= -\frac{\partial_q h_u(P_c(p_a(L_a), q_a, q_c), q_c)}{\partial_p h_u(P_c(p_a(L_a), q_a, q_c), q_c)}\partial_\lambda q_c > 0, \tag{1.89}$$

because the flows $q_a$ and $q_c$ are positive, Assumption b)(iii) of Theorem 2 holds and $\partial_\lambda q_c < 0$ is true. Hence, to show the assertion $\partial_\lambda p_{b_r}(0) > 0$, we have to verify $T_{p_c} > 0, T_{q_c} < 0$ and $T_{q_b} > 0$. The first term expands to

$$T_{p_c} = \partial_p f_b\big(L_{b_r}; p_{b_r}(L_{b_r}), -q_b\big)\partial_{p_c} P_{b_r}\big(p_c(L_c), -q_b, q_c\big)\partial_p f_c(L_c; p_c(0), q_c)$$

(by Lemma 15, Equation (1.47a))

$$= \partial_p f_b\big(L_{b_r}; p_{b_r}(L_{b_r}), -q_b\big)\partial_p f_c(L_c; p_c(0), q_c)\frac{\partial_p h_v(p_c(L_c), q_c)}{\partial_p h_v(P_{b_r}(p_c(L), q_c, -q_b), -q_b)}$$
$$> 0,$$

since each derivative with respect to $p$ is positive. The second term can be written as

$$T_{q_b} = -\partial_p f_b\big(L_{b_r}; p_{b_r}(L_{b_r}), -q_b\big)\partial_{q_b} P_{b_r}\big(p_c(L_c), q_b, q_c\big)$$
$$- \partial_q f_b\big(L_{b_r}; p_{b_r}(L_{b_r}), -q_b\big)$$

(by Lemma 15, Equation (1.47b))

$$= \partial_p f_b\big(L_{b_r}; p_{b_r}(L_{b_r}), -q_b\big)\frac{\partial_q h_v(P_{b_r}(p_c(L_c), q_c, q_b), q_b)}{\partial_p h_v(P_{b_r}(p_c(L_c), q_c, q_b), q_b)}$$
$$- \partial_q f_b\big(L_{b_r}; p_{b_r}(L_{b_r}), -q_b\big) > 0$$

as a consequence of Assumption b)(iii) of Theorem 2, which holds due to Lemma 20. The last term fulfills

$$T_{q_c} = \partial_{p_c} P_{b_r}\big(p_c(L_c), -q_b, q_c\big)\partial_q f_c(L_c; p_c(0), q_c) + \partial_{q_c} P_{b_r}\big(p_c(L_c), -q_b, q_c\big)$$

(by Lemma 15, Equations (1.47a) and (1.47c))

$$= \frac{\partial_p h_v(P_{b_r}(p_c(L_c), -q_b, q_c), -q_b)}{\partial_p h_v(p_c(L_c), q_c)}\partial_q f_c(L_c; p_c(0), q_c)$$
$$+ \frac{\partial_q h_v(p_c(L_c), q_c)}{\partial_p h_v(P_{b_r}(p_c(L_c), -q_b, q_c), -q_b)} < 0$$

by Assumption b)(vi) of Theorem 2, which holds due to Lemma 20. Consequently, we obtain

$$\partial_\lambda p_{b_r}(0) = T_{p_c}\partial_\lambda p_c(0) + \partial_p f_b\big(L_{b_r}; p_{b_r}(L_{b_r}), -q_b\big)T_{q_c}\partial_\lambda q_c + T_{q_b}\partial_\lambda q_b > 0.$$

This shows, together with Inequality (1.83), that the function $H$ is strictly decreasing.

Now, we find value for $\lambda$ such that

$$p_{b_l}(L_{b_l}; \lambda) - p_{b_r}(0; \lambda) = H(\lambda) < 0$$

and a value such that $H(\lambda) > 0$ holds. The choice $\lambda = 1$, i.e.,

$$q_b(\lambda) = q_a, \quad q_c(\lambda) = 0$$

leads to (omitting the argument $\lambda$ from now on)

$$p_{b_l}(L_{b_l}) < p_b(0), \tag{1.90}$$

since $f_b$ is strictly decreasing in flow direction.

The initial pressure on the arc $c$, results from the relation

$$h_u(p_c(0), 0) = h_u(p_b(0), q_b).$$

Inserting the definition of the Bernoulli invariant yields

$$\log\left(\frac{p_c(0)}{c^2}\right) = \log\left(\frac{p_b(0)}{c^2}\right) + \frac{1}{2}\eta_{b_l}.$$

Resolving for $p_c(0)$ leads to

$$p_c(0) = c^2 \exp\left(\log\left(\frac{p_b(0)}{c^2}\right) + \frac{1}{2}\eta_{b_l}\right) \geq p_b(0). \tag{1.91}$$

As the flow on the arc $c$ is zero, we obtain $p_c(0) = p_c(L_c)$. We denote by $P_{vc}(p_{b_r}(L_{b_r}), q_b, q_c)$ the function that solves the coupling condition

$$h_v(P_{vc}(p_{b_r}(L_{b_r}), q_b, q_c), q_c) = h_v(p_{b_r}(L_{b_r}), q_b)$$

in $v$. Therefore, we obtain

$$p_c(L_c) = P_{vc}(p_{b_r}(L_{b_r}), q_b, q_c)$$

(by the Taylor series with $\xi \in [q_c, q_b]$)

$$= P_{vc}(p_{b_r}(L_{b_r}), q_b, q_b) + (q_c - q_b)\, \partial_{q_c} P_{vc}(p_{b_r}(L_{b_r}), q_b, \xi)$$

(by Lemma 15, Equation (1.46a))

$$= p_{b_r}(L_{b_r}) - (q_c - q_b)\frac{\partial_q h_v(P_{vc}(p_{b_r}(L_{b_r}), q_b, \xi), \xi)}{\partial_p h_v(P_{vc}(p_{b_r}(L_{b_r}), q_b, \xi), \xi)}$$

$$= f_{b_r}(0; p_{b_r}(0), q_b) - (q_c - q_b)\frac{\partial_q h_v(P_{vc}(p_{b_r}(L_{b_r}), q_b, \xi), \xi)}{\partial_p h_v(P_{vc}(p_{b_r}(L_{b_r}), q_b, \xi), \xi)} < p_{b_r}(0),$$

$$\tag{1.92}$$

because Lemma 20 states that Assumption b)(vii) of Theorem 2 holds. Note that in the example $q_b$ has the role of $q_a$ in the assumption and $q_c$ has the role of $q_b$.

Combining (1.90), (1.91) and (1.92) leads to $p_{b_l}(L_{b_l}) < p_{b_r}(0)$, which shows $H(1) < 0$. With an analogous argumentation, we obtain $H(0) > 0$, which ensures that there is a unique solution $\lambda^* \in (0,1)$ such that $H(\lambda^*) = 0$. Hence all flow values on the network are known and the pressure values can be calculated, starting from the root, as in the tree network example.

## Conclusion

We constructed explicit stationary solutions of the quasilinear isothermal Euler equations on a single pipe and showed essential monotonicity properties. By introducing node coupling conditions, the states can be extended to networks. The existence theory is general enough to allow different node coupling conditions, the inclusion of active elements and even the combination of different models. The proof is constructive and therefore allows the computation of the state values, which is demonstrated by the examples. The requirements of the existence theorem were explicitly shown for real and ideal gas with the continuity of pressure and the continuity of the Bernoulli invariant as node coupling conditions. The values on the arcs can also be coupled by active elements including turbo compressors and resistors, which was thoroughly discussed.

Although only gas was considered in this thesis, the result can be applied to other balance laws with friction in an equilibrium that fulfill the required monotonicity properties. This includes the flow through networks of canals, modeled by the Saint-Venants Equations (see Gugat and Leugering 2009) and the flow of blood through an arterial network (see Bressan, Čanić, et al. 2014, Section 3.6).

# 2. Chance Constrained Optimization

> *As far as the laws of mathematics refer to reality, they are not certain, and as far as they are certain, they do not refer to reality.*
>
> (Albert Einstein)

In this chapter, we discuss optimization problems under uncertainty. The uncertainty is handled by using *chance constraints* (also known as *probabilistic constraints*). For a prescribed parameter $z \in \mathbb{R}^m$, we look at the optimization problem

$$\min_{x \in \mathbb{R}^n} \quad f(x)$$
$$\text{s.t.} \quad g(x, z) \leq 0.$$

We will refer to this problem as the *deterministic optimization problem.* If, however, the parameter $z$ is not precisely known but replaced by a random vector $\xi$ with known probability distribution, we can consider the constraint that $g(x, \xi) \leq 0$ holds with at least probability $p^{\mathrm{bd}} \in (0, 1)$. The resulting optimization problem reads

$$\min_{x \in \mathbb{R}^n} \quad f(x)$$
$$\text{s.t.} \quad \mathbb{P}(g(x, \xi) \leq 0) \geq p^{\mathrm{bd}}.$$

We will refer to this problem as the *stochastic optimization problem* despite the fact that it is actually completely deterministic. This can be easily seen by denoting $G(x) := \mathbb{P}(g(x, \xi) \leq 0)$ as the evaluation of the probability function is merely an integration over the set of $z \in \mathbb{R}^m$ fulfilling $g(x, z) \leq 0$ weighted with the probability density. The evaluation of high dimensional integrals over a complicated integration domain calls for the application of special integration techniques, namely *quasi-Monte Carlo* methods. Quasi-Monte Carlo methods are based on the same principle as *Monte Carlo* methods, but use low discrepancy sequences instead of random sequences. We will give an overview of the methods and known error estimates. If

the random vector $\xi$ is gaussian, we can apply a parameterization of the integral known as *spherical radial decomposition*. For this representation, we derive subgradients and gradients of the probability function, which are used in a gradient based optimization. We propose a multilevel approach that significantly reduces the computational time for solving the stochastic optimization problem and apply it to gas network optimization under uncertain demands.

**Literature Survey**

The book Niederreiter 1992 offers a well written introduction to quasi-Monte Carlo methods. The more recent book Lemieux 2009 on Monte Carlo and quasi-Monte Carlo is more focused on computational aspects and has a more extensive part on Monte Carlo methods. The focus in Dick and Pillichshammer 2010 lies on digital nets and sequences for quasi-Monte Carlo methods. A comparison of Monte Carlo methods, lattice rules and other low discrepancy point sets was provided in Lemieux and L'Ecuyer 1999.

Complexity theory of quasi-Monte Carlo algorithms and the weighted Koksma-Hlawka inequality are elaborated in Sloan and Woźniakowski 1998. Quasi-Monte Carlo methods in the weighted Hilbert space setting are discussed in Kuo, Schwab, and Sloan 2011. Weighted discrepancies for high dimensional integration are examined in Larcher, Pillichshammer, and Scheicher 2003.

Quasi-Monte Carlo methods for finance applications were investigated in Giles et al. 2008. Practical stopping criteria for Monte Carlo and randomized quasi-Monte Carlo methods that complement the theory are presented therein.

In Brauchart and Dick 2012 sphere integrals were solved by quasi-Monte Carlo integration. In Marques et al. 2013 spherical illumination integrals were treated by quasi-Monte Carlo method and the Fibonacci point set on the sphere was discussed as low discrepancy set.

Spherical radial decomposition of integrals of normal distribution density functions is treated in Genz and Bretz 2009 and Monahan and Genz 1997. Gradient representations for spherical radial decomposed probability functions were derived in Ackooij and Henrion 2014 and 2017.

Stochastic programming and, more specific, chance constrained programming is discussed in Shapiro, Dentcheva, and Ruszczyński 2009. Regularization and optimality conditions for chance constrained programs are examined in Adam and Branda 2016. The doctoral thesis Ackooij 2013 deals with chance constrained optimization and its application in energy

management. The compressor control under probabilistic constraints on stationary gas networks was recently discussed in Gugat and Schuster 2018. Explicit representations of derivatives of probability functions were first considered in Uryas'ev 1994 and Kibzun and Uryasev 1998. Safe tractable approximations of chance constraints were elaborated in Nemirovski 2012. In Curtis, Wächter, and Zavala 2018 chance constrained problems were solved by a specialized sequential quadratic programming solver with exact penalization.

## 2.1. Monte Carlo Methods

Throughout this section we follow the classical book *Random Number Generation and Quasi-Monte Carlo Methods*, Niederreiter 1992, Section 1.2. The motivation for using Monte Carlo methods instead of traditional methods like the Newton-Cotes formulas lies in the scaling with the dimension. Consider the trapezoidal rule that approximates the one dimensional integral

$$\int_0^1 f(x)\,\mathrm{d}x \quad \text{by} \quad I_h = \sum_{n=0}^{m} w_k f\left(\tfrac{k}{m}\right)$$

with weights $w_0 = w_m = 1/(2m)$ and $w_n = 1/m$, for $n = 1, \ldots, m-1$. If $f$ is two times continuously differentiable, the integration error is of order $\mathcal{O}(m^{-2})$. The multidimensional case is analogous: The integral

$$\int_{[0,1]^d} f(x)\,\mathrm{d}x$$

can be approximated by

$$I_h = \sum_{n_1=0}^{m} \cdots \sum_{n_d=0}^{m} w_{n_1} \cdots w_{n_d} f\left(\frac{n_1}{m}, \ldots, \frac{n_s}{m}\right)$$

with the weights defined as in the one dimensional case. Once again, the error is of order $\mathcal{O}(m^{-2})$ if the function $f$ is two times continuously differentiable. However, the number of evaluation points needed increases dramatically with the size of the dimension $s$, because the number of points is $N = (m+1)^s$ and thus the error bound in terms of evaluation points is given by $\mathcal{O}(N^{-2/s})$. To obtain an accuracy of order 0.01 in one dimension, 10 evaluation points are needed. For the same accuracy in dimension

100, the number of required evaluation points is $10^{100}$. This phenomenon has received the name the *curse of dimensionality*. The first approach to overcome this issue has been the Monte Carlo method.

**Definition 2** ($L^p$-Spaces)**.** Let $\mu$ be a measure on the measurable space $(\Omega, \Sigma)$ and let $(\Omega', \Sigma')$, $\Omega' \subset \mathbb{R}$ be a measurable space.

For $1 \le p < \infty$ and $A \subset \Omega$ define

$$L^p(A; \mu) := \left\{ f : \Omega \to \Omega' \mid f \text{ is measurable and } \int_A |f|^p \, \mathrm{d}\mu < \infty \right\}.$$

In the case, where $\mu = \lambda_s$ is the $s$-dimensional Lebesgue measure, we write $L^p(A)$ instead of $L^p(A; \lambda_s)$.

The idea of the Monte Carlo integration is simple: To calculate the integral $\int_\Omega f(x) \, \mathrm{d}x$, consider the integration domain $\Omega \subset \mathbb{R}^s$ with $0 < \lambda_s(\Omega) < \infty$, where $\lambda_s$ denotes the $s$-dimensional Lebesgue measure. The transformation $\mathrm{d}\mu = \mathrm{d}x/\lambda_s(\Omega)$ turns $\Omega$ into a probability space with probability measure $\mu$. For $f \in L^1(\Omega; \mu)$, we obtain

$$\int_\Omega f(x) \, \mathrm{d}x = \lambda_s(\Omega) \int_\Omega f \, \mathrm{d}\mu = \lambda_s(\Omega) \, \mathbb{E}(f), \tag{2.1}$$

where $\mathbb{E}(f)$ is the expected value of $f$ with respect to the probability measure $\mu$. The expected value can be estimated by the mean value of a sufficiently large number of samples.

**Definition 3** (Monte Carlo Estimate)**.** Let $f$ be a random variable on a probability space $(\Omega, \Sigma, \mu)$. Take $N$ independent $\mu$-distributed samples $x_1, \ldots, x_n \in \Omega$. The *Monte Carlo estimate* for the expected value $\mathbb{E}(f)$ is given by

$$\mathbb{E}_N(f) := \frac{1}{N} \sum_{n=1}^N f(x_n).$$

For $\mu = \mathrm{d}x/\lambda_s(\Omega)$, let $f$ be in $L^1(\Omega; \mu)$ and $0 < \lambda_s(A) < \infty$. The *Monte Carlo estimate* for the integral $\int_\Omega f(x) \, \mathrm{d}x$ is given by

$$I_N(f) := \frac{\lambda_s(\Omega)}{N} \sum_{n=1}^N f(x_n).$$

This definition is a consequence of replacing the expected value in Equation (2.1) by its Monte Carlo estimate.

The convergence of $\mathbb{E}_N(f)$ to $\mathbb{E}(f)$ in the sense of probability is a consequence of the strong law of large numbers.

**Theorem 6** (Strong Law of Large Numbers, Klenke 2013, Satz 5.17).
*Let $X_1, X_2, \ldots \in L^1(\Omega; \mu)$ be pairwise independent and identically distributed random variables. Then, the sequence $(X_n)_{n \in \mathbb{N}}$ obeys the strong law of large numbers, i.e.,*

$$\lim_{N \to \infty} \frac{1}{N} \sum_{n=1}^{N} X_n = \mathbb{E}(X_1) \qquad \mu^\infty\text{-almost everywhere,}$$

*where $\mu^\infty$ is product measure of countable many copies of $\mu$.*

This directly leads to the convergence of the Monte-Carlo estimates.

**Corollary 2** (Convergence of the Monte-Carlo Estimates).
*Consider the assumptions of Definition 3. Then*

$$\lim_{N \to \infty} \mathbb{E}_N(f) = \mathbb{E}(f) \qquad \mu^\infty\text{-almost everywhere}$$

*and*

$$\lim_{N \to \infty} I_N(f) = \int_\Omega f(x) \, \mathrm{d}x \qquad \mu^\infty\text{-almost everywhere.}$$

While this statement assures the convergences in probability, it says nothing about the speed of convergence. The variance gives us information about the average integration error.

**Definition 4** (Variance). Let $f$ be a random variable on a probability space $(\Omega, \Sigma, \mu)$. The *variance* of $f$ is defined by

$$\sigma^2(f) := \int_\Omega (f - \mathbb{E}(f))^2 \, \mathrm{d}\mu.$$

The variance is finite for $f \in L^2(\Omega; \mu)$.

**Theorem 7** (Average Integration Error).
*Let $f \in L^2(\Omega; \mu)$ and $N \geq 1$. Then, the integral average of the integration error is*

$$\int_\Omega \cdots \int_\Omega \left( \frac{1}{N} \sum_{n=1}^{N} f(x_n) - \mathbb{E}(f) \right)^2 \mathrm{d}\mu(x_1) \ldots \mathrm{d}\mu(x_N) = \frac{\sigma^2(f)}{N}.$$

*Proof.* See Niederreiter 1992, proof of Theorem 1.1. $\qquad\square$

For a probabilistic error bound of the Monte-Carlo method, we need the central limit theorem, for which the following definitions are required.

**Definition 5** (Product measure). Let $(\Omega, \Sigma, \mu)$ be a measure space and $A_1, \ldots, A_N \in \Sigma$. The extension of the measure

$$\mu^{\otimes N}(A_1, \ldots, A_N) := \prod_{n=1}^{N} \mu(A_n)$$

to the *product sigma-algebra* $\sigma(\Sigma^N)$ is called the *product measure.* We also denote it by $\mu^{\otimes N}$.

**Definition 6** (Weak Convergence of Measures, Klenke 2013, Def. 13.12). Let $\mu_1, \mu_2, \ldots$ be a sequence of finite measures on the measurable space $(\Omega, \Sigma)$. We say, $(\mu_n)_{n \in \mathbb{N}}$ *converges weakly* to $\mu$ and write $\mu_n \rightharpoonup \mu$, for $N \to \infty$ if

$$\lim_{n \to \infty} \int_\Omega f \, \mathrm{d}\mu_n = \int_\Omega f \, \mathrm{d}\mu \text{ for every } f \in \mathcal{C}_b(\Omega),$$

where $\mathcal{C}_b(\Omega)$ is the space of continuous bounded functions, i.e.,

$$\mathcal{C}_b(\Omega) := \{f : \Omega \to \mathbb{R} \mid f \text{ is continuous and bounded}\}.$$

**Theorem 8** (Central Limit Theorem, Klenke 2013, Satz 15.37). *Let $X_1, X_2, \ldots$ be independent and identically distributed random variables with finite variance $\sigma^2$. Define*

$$S_N(\omega_1, \ldots, \omega_N) := \sigma^{-1} N^{-1/2} \sum_{n=1}^{N} (X_n(\omega_n) - \mathbb{E}(X_1)).$$

*Then, the distribution converges weakly against the standard normal distribution, i.e.,*

$$\mu^{\otimes N} \circ S_N^{-1} \rightharpoonup \mathcal{N}_{0,1}, \quad \text{for } N \to \infty.$$

*For the cumulative distribution functions, we have, for $-\infty \le a < b \le +\infty$,*

$$\lim_{N \to \infty} \mu^{\otimes N} \big\{ (\omega_1, \ldots, \omega_N) \in \Omega^N \mid S_N(\omega_1, \ldots, \omega_N) \in [a, b] \big\}$$

$$= (2\pi)^{-1/2} \int_a^b e^{-t^2/2} \, \mathrm{d}t.$$

*This equation is sometimes abbreviated as*

$$\lim_{N \to \infty} \mathbb{P}\left( a \le \sigma^{-1} N^{-1/2} \sum_{n=1}^{N} (X_n - \mathbb{E}(X_1)) \le b \right) = (2\pi)^{-1/2} \int_a^b e^{-t^2/2} \, \mathrm{d}t.$$

As a direct consequence of the central limit theorem we obtain a probabilistic error bound for the Monte-Carlo estimate.

**Corollary 3** (Probabilistic Error Bound for the Monte Carlo Estimate)**.** *Assume $f$ has finite variance $\sigma^2(f)$. Then, the probabilistic error bound*

$$\lim_{N \to \infty} \mathbb{P}\left( a \, \frac{\sigma(f)}{\sqrt{N}} \leq \frac{1}{N} \sum_{n=1}^{N} f(x_n) - \mathbb{E}(f) \leq b \, \frac{\sigma(f)}{\sqrt{N}} \right) = (2\pi)^{-1/2} \int_a^b e^{-t^2/2} \, \mathrm{d}t$$

*holds, for any constants $-\infty \leq a < b \leq \infty$.*

**Example 2.** For a normal distributed variable, around 99.73% of all realizations lie inside a $3\sigma$ interval around the expected value. Hence, since $\sigma = 1$ for a standard normal distribution, we obtain by choosing $a = -3$ and $b = 3$ that

$$(2\pi)^{-1/2} \int_{-3}^{3} e^{-t^2/2} \, \mathrm{d}t \approx 99.73\%.$$

Take a random variable with variance $\sigma(f) = 2$. Then

$$\left| \frac{1}{N} \sum_{n=1}^{N} f(x_n) - \mathbb{E}(f) \right| \leq \frac{6}{\sqrt{N}}$$

with 99.73% probability. If the error should be smaller than 0.01, we have to choose

$$\frac{6}{\sqrt{N}} \leq 0.01 \iff \sqrt{N} \geq 600 \iff N \geq 360000.$$

At first this seems like a extraordinary large number of evaluation points, especially in the one dimensional case. However, unlike traditional integration rules, this number does not scale with the dimension of the integration domain.

Therefore, we obtain a probabilistic error bound of order $\mathcal{O}(N^{-1/2})$ compared to $\mathcal{O}(N^{-2/s})$ for the trapezoidal rule. The error bound shows that it can be advantageous to transform $f$ into a function $\tilde{f}$ with lower variance. Since this thesis focuses on the quasi-Monte Carlo approach rather than the classical Monte Carlo method, we do not provide further details on the subject of variance reduction techniques.

For implementation the sample standard variation

$$\sigma_N(f) := \left( \sum_{n=1}^{N} \frac{(f(x_n) - \mathbb{E}_N(f))^2}{n-1} \right)^{1/2}$$

can be used as an approximation of $\sigma(f)$; see Lemieux 2009, p. 10.

*Remark* 12 (Integrating over complicated domains). In the case, where our domain $\Omega$ is too complicated to compute $\lambda_s(\Omega)$, we use the following slight adjustment: Assume $\Omega \subset [0,1]^s$ (since $\Omega$ is bounded it is certainly contained in a cube $Q \subset \Omega$, which can be transformed to the unit cube). Then, we can rewrite the integral as

$$\int_\Omega f(x)\,\mathrm{d}x = \int_{[0,1]^s} f(x)\chi_\Omega(x)\,\mathrm{d}x$$

with the characteristic function of $\chi_\Omega$ of $\Omega$. The Monte-Carlo estimate is given by

$$I_h(f) := \frac{1}{N}\sum_{n=1}^{N} f(x_n)\chi(x_n) = \frac{1}{N}\sum_{\substack{n=1 \\ x_n \in \Omega}}^{N} f(x_n).$$

The results above still apply for the function $\chi_\Omega \cdot f$, because only very weak regularity is required.

We have seen that the Monte Carlo method provides decisive advantages over classical integration rules if the dimension is large. In Lemieux 2009, p. 10–11 an example is provided, where this is the case for dimensions greater than four. However, the drawbacks are not to be underestimated: The error bound is only probabilistic, higher regularity of the function to integrate does not improve the error bound and truly random samples are not easy to generate (which is why pseudorandom numbers are used in practice). The quasi-Monte Carlo methods seeks to overcome all of this problems by using deterministic point sets instead of random numbers. This achieves a deterministic error bound that is superior for medium dimensions.

## 2.2. Quasi-Monte Carlo Methods

The quasi-Monte Carlo is based on the same concept as the Monte Carlo method with the difference that it uses deterministic point sets instead of random numbers as integration points. This leads to a deterministic error bound of order $\mathcal{O}(N^{-1}[\log(N)]^{s-1})$ compared to the probabilistic error bound of order $\mathcal{O}(N^{-1/2})$. We will see that the Koksma-Hlawka inequality estimates the integration error by the Hardy-Krause variation of the function and the discrepancy of the point set. The discrepancy measures, how evenly the points are spread in the integration domain. This coincides with intuition: A numerical integration is easy if the function is no too "jumpy" and the evaluation points are evenly spread.

In view of Remark 12, we restrict ourselves to integration over the unit cube $[0,1]^s$. Like in the previous section we follow Niederreiter 1992, Chapter 2.

**Definition 7** (Quasi-Monte Carlo Estimate)**.** Let $f \in L^1([0,1]^s; \lambda_s)$ and the points $x_1, \ldots, x_N \in [0,1]^s$ be given. The *quasi-Monte Carlo estimate* of the integral

$$\int_{[0,1]^s} f(x) \, \mathrm{d}x$$

is given by

$$I_N(f) := \frac{1}{N} \sum_{n=1}^{N} f(x_n).$$

## 2.2.1. Classical Koksma-Hlawka Inequality

The definition of the integral estimate looks similar to the Monte Carlo method, yet we obtained more freedom in choosing the points $x_1, \ldots, x_N$. To obtain convergence of $I_N$ against the integral for a reasonable class of functions, the set of evaluation points should be evenly distributed across the unit cube. To measure the "evenness", we need a concept called discrepancy of the point set.

**Definition 8** (Discrepancy)**.** Consider the point set

$$P := \{x_1, \ldots, x_N\} \subset [0,1]^s$$

and an arbitrary set $A \subset [0,1]^s$. We denote the number of points of $P$ in $A$ by

$$N(A; P) := \sum_{n=1}^{N} \chi_A(x_n),$$

where $\chi_A$ is the characteristic function of $A$. Let $\mathcal{A}$ be a family of nonempty Lebesgue-measurable subsets of $[0,1]^s$. The *discrepancy* of the point set $P$ is given by

$$D_N(\mathcal{A}; P) := \sup_{A \in \mathcal{A}} \left| \frac{N(A; P)}{N} - \lambda_s(A) \right|.$$

The *star discrepancy* of the point set $P$ is defined as

$$D_N^*(P) := D_N(\mathcal{I}^*; P),$$

(a) Uniform grid

(b) Halton point set



(c) Pseudorandom points

Figure 2.1.: Discrepancy

where $\mathcal{I}^*$ is the family of all cuboids of the form $\prod_{i=1}^{s}[0, a_i) \subset [0,1]^s$. In the case of only one set $A$, one also speaks of the *local discrepancy*

$$D_N(A; P) := \left| \frac{N(A; P)}{N} - \lambda_s(A) \right|.$$

In the case of a set $[0, x)$ and fixed points $P$, we also abbreviate the local discrepancy as

$$D_N(x) := D_N([0, x); P).$$

*Remark* 13. Dividing by $\lambda_s([0,1])^s = 1$ shows that the discrepancy

$$D_N(\mathcal{A}; P) := \sup_{A \in \mathcal{A}} \left| \frac{N(A; P)}{N} - \frac{\lambda_s(A)}{\lambda_s([0,1]^s)} \right|$$

measures the deviation of the number of points that are expected by the relative volume from the actual number of points in the test set.

**Example 3.** Consider the set $A = [0, 1/5] \times [0, 1] \subset [0,1]^2$. It is shown in blue in Figure 2.1a. The point set $P$ is given by a uniform $6 \times 6$ grid.

The area of $A$ is given by $\lambda_2(A) = 1/5$. The number of points in $A$ is $N(A; P) = 12$. Consequently the discrepancy of $P$ on the set $A$ is given by

$$D_N(A; P) = \left| \frac{N(A; P)}{N} - \lambda_s(A) \right| = \left| \frac{12}{36} - \frac{1}{5} \right| = \frac{2}{15}.$$

It can be seen that no points are contained in the open cuboid $B = (0, 1/5) \times (0, 1)$. This leads to

$$D_N(B; P) = \frac{1}{5}$$

and it is the cuboid with the largest area and no points in it. Another extreme set is $C = \{0\} \times [0, 1]$, (red) e.g. a set which is only a line and has area $\lambda_2(C) = 0$, but six points in it. Hence

$$D_N(C; P) = \frac{1}{6}.$$

Taking the discrepancy over the family $\mathcal{A} := \{A, B, C\}$ leads to the largest value of the three namely

$$D_N(\mathcal{A}; P) = \frac{1}{5}.$$

We will later discuss low discrepancy sequences. Figure 2.1b shows 36 points of a Halton set. We denote the point set by $Q$. It can be seen that eight points are contained in the set $A$ from above. This leads to a discrepancy of

$$D_N(A; Q) = \left| \frac{N(A; P)}{N} - \lambda_s(A) \right| = \left| \frac{8}{36} - \frac{1}{5} \right| = \frac{1}{45},$$

which is significantly lower than the discrepancy $D_N(A; P)$ of the uniform grid. Figure 2.1c shows the typical clustering of pseudorandom points. The discrepancy of the 36 points $R$ on the sets $E = [1/10, 4/10] \times [0, 3/10]$ (orange) and $F = [0, 1/2] \times [1/3, 1/2]$ (green) is given by

$$D_N(E; R) = \frac{9}{100} \quad \text{and} \quad D_N(F; R) = \left| \frac{8}{36} - \frac{1}{12} \right| = \frac{5}{36}.$$

The error estimate depends not only on the point set, but also the behavior of the function to integrate. A higher variation make the function harder to integrate.

**Definition 9** (Variation)**.** Consider the function $f : [a, b] \to \mathbb{R}$ in one dimension. The *variation* of $f$ is defined by

$$V(f) := \sup\left\{ \sum_{n=0}^{N-1} \left| f(y_{n+1}^{(N)}) - f(y_n^{(N)}) \right| \mid N \in \mathbb{N}, a \le y_0^{(N)} < \cdots < y_N^{(N)} \le b \right\}.$$

We say that $f$ has bounded variation if $V(f) < \infty$.

In the one dimensional setting the *Koksma inequality* provides an error bound for the quasi-Monte Carlo integration; see Koksma 1942.

**Theorem 9** (Koksma Inequality)**.**
*Consider a function $f : [0, 1] \to \mathbb{R}$ of bounded variation $V(f)$ and evaluation points $x_1, \ldots, x_N \in [0, 1]$. Then the error bound*

$$\left| \frac{1}{N} \sum_{n=1}^{N} f(x_n) - \int_0^1 f(x) \, \mathrm{d}x \right| \le V(f) \, D_N^*(x_1, \ldots, x_N)$$

*holds.*

The generalization of the Koksma inequality to multiple dimensions requires a generalized concept of the variation of a function. A good overview that we follow here is presented in Owen 2005. We notate $x^I := (x_i)_{i \in I}$ for an index set $I \subset \{1, \ldots, s\}$ and for disjoint index sets $I$ and $J \subset \{1, \ldots, s\}$, we use the colon operator $x^I{:}y^J$ to denote the vector $z = x^I{:}y^J$ such that $z_i = x_i$, for $i \in I$ and $z_i = y_i$ for $i \in J$. Consider, for example, $I = \{1, 3, 5\}$ and $J = \{2, 4\}$. Then

$$x^I : y^J = (x_1, \ y_2, \ x_3, \ y_4, \ x_5)^T.$$

For the complement with respect to $\{1, \ldots, s\}$ we notate

$$-I := \{1, \ldots, s\} \setminus I.$$

**Definition 10** (Ladder)**.** A finite set $\mathcal{Y}$ of ordered points on $[a, b]$ containing $a$ and not containing $b$ except when $a = b$ is called a *ladder*. The points $y_i$ in $\mathcal{Y}$ are ordered by

$$a = y_0 < y_1 < \ldots < y_m.$$

We denote the *successor* $y_+$ of $y$. It is given by $(y_k)_+ = y_{k+1}$ for $k < m$ and $(y_m)_+ = b$. The set of all ladders on $[a, b]$ is denoted by $\mathbb{Y}$. The total

variation of a one dimensional function $f : [a, b] \to \mathbb{R}$ as in Definition 9 can be written as

$$V(f) = \sup_{\mathcal{Y} \in \mathbb{Y}} \sum_{y \in \mathcal{Y}} |f(y_+) - f(y)|.$$

In multiple dimensions, we call the product of ladders $\mathcal{Y}_i$ on $[a_i, b_i]$ a ladder on the cuboid $[a, b] := \{x \in \mathbb{R}^s \mid a \leq x \leq b\}$ and write $\mathcal{Y} := \prod_{i=1}^{s} \mathcal{Y}_i$. The successor $y_+$ is defined by taking the successor in each component of $y$.

**Definition 11** (Alternating Sum). The *s-fold alternating sum* of $f$ over $[a, b] \subset \mathbb{R}^s$ is given by

$$\Delta(f; [a, b]) := \sum_{I \subset \{1, \dots, s\}} (-1)^{|I|} f(a^I {:} b^{-I}).$$

**Definition 12** (Vitali Variation). Let $[a, b] \subset \mathbb{R}^s$ and $f : [a, b] \to \mathbb{R}$ and denote set of all ladders on $[a, b]$ by $\mathbb{Y}$. The *variation in the sense of Vitali* of $f$ is defined as

$$V_{\mathrm{V}}(f; [a, b]) := \sup_{\mathcal{Y} \in \mathbb{Y}} \sum_{y \in \mathcal{Y}} |\Delta(f; [y, y_+])|.$$

**Definition 13** (Hardy-Krause Variation). Let $[a, b] \subset \mathbb{R}^s$ and $f : [a, b] \to \mathbb{R}$ and denote set of all ladders on $[a, b]$ by $\mathbb{Y}$. The *variation in the sense of Hardy-Krause* of $f$ is defined as

$$V(f) := \sum_{I \subsetneq \{1, \dots, s\}} V_{\mathrm{V}}\big(f([\cdot]^{-I}; b^I); [a^{-I}, b^{-I}]\big),$$

where the function $f(x^{-I}; b^I) := f(x^{-I} {:} b^I)$ has the components in $-I$ as argument, while the components $I$ are set constant to the values in $b^I$.

Having these definitions at hand, we can generalize the Koksma inequality to the *s*-dimensional case. The resulting inequality is due to Hlawka 1961a and is known as *Koksma-Hlawka inequality*.

**Theorem 10** (Koksma-Hlawka Inequality).
*Consider $f : [0, 1]^s \to \mathbb{R}$ with bounded variation in the sense of Hardy-Krause and an arbitrary point set $P := \{x_1, \dots, x_N\} \subset (0, 1)^s$. Then,*

$$\left| \frac{1}{N} \sum_{n=1}^{N} f(x_n) - \int_{[0,1]^s} f(x) \, \mathrm{d}x \right| \leq V(f) D_N^*(P)$$

*holds.*

For convex integration domains in $[0,1]^s$ a similar result is available with the isotropic discrepancy.

**Definition 14** (Isotropic Discrepancy)**.** The *isotropic discrepancy* of a point set $P$ is given by

$$J_N(P) := D_N(\mathcal{C}; P),$$

where $\mathcal{C}$ is the family of all convex subsets over $[0,1]^s$.

**Theorem 11** (Quasi-Monte Carlo Error for Convex Domains)**.**
*Let $A \subset [0,1]^s$ be convex and let $f : [0,1]^s \to \mathbb{R}$ be a function of bounded variation in the sense of Hardy-Krause. Then, for any point set $P = \{x_1, \ldots, x_n\}$, we have*

$$\left| \frac{1}{N} \sum_{\substack{n=1 \\ x_n \in A}}^{N} f(x_n) - \int_A f(x)\,\mathrm{d}x \right| \le (f(\mathbf{1}) + V(f)) J_N(P).$$

## 2.2.2. Low Discrepancy Sequences

The construction of sequences with low discrepancy is thoroughly discussed in Niederreiter 1992. For a sequence $S = (x_n)_{n\in\mathbb{N}}$ in one dimension, we have by a result in Béjan 1982 and Niederreiter 1992, Proposition 2.4

$$\limsup_{N\to\infty} D_N^*(S)\,\frac{N}{\log(N)} \ge 0.06.$$

This shows that we cannot hope for a better bound of the star discrepancy than order $\mathcal{O}(N^{-1}\log(N))$. For the numerical computations in this thesis, Halton and Sobol point sets were used. We focus on discussing the basics for the Halton sequence.

**Definition 15** (Radical-Inverse Function)**.** For an given integer $b \ge 2$, set $Z_b := \{0, 1, \ldots, b-1\}$. Then any number $n \in \mathbb{N}$ has unique digit expansion in the base $b$ given by

$$n = \sum_{j=1}^{\infty} a_j(n)\,b^j,$$

where $a_j(n) \in Z_b$ and $a_j(n) = 0$ for sufficiently large $j$, this means, the expansion is finite.

The *radical-inverse function* $\phi_b : \mathbb{N} \to [0, 1)$ is defined by

$$\phi_b = \sum_{j=0}^{\infty} a_j(n)\,b^{-j-1}.$$

**Example 4.** Choose the base $b = 4$ and the integer $n = 27$. Then, $Z_b = \{0, 1, 2, 3\}$ and the digit expansion of $n$ is given by

$$n = 27 = 3 \cdot 1 + 2 \cdot 4 + 1 \cdot 16 = 3 \cdot 4^0 + 2 \cdot 4^1 + 1 \cdot 4^2$$
$$= a_0(27) \, b^0 + a_1(27) \, b^1 + a_2(27) \, b^2.$$

One also writes $123_4$. The radical-inverse function is the symmetric reflection at the "decimal point" resulting in

$$\phi_4(n) = 3 \cdot 4^{-1} + 2 \cdot 4^{-2} + 1 \cdot 4^{-3} = 0.890625$$

or in digit notation $\phi_4(123_4) = 0.321_4$.

The radical-inverse function can be used to construct a low discrepancy sequence (or quasirandom sequence) in one dimension.

**Definition 16** (Van der Corput Sequence). Choose an integer $b \geq 2$. The *van der Corput sequence* in base $b$ is defined by setting

$$x_n := \phi_b(n) \quad \text{for all } n \in \mathbb{N}.$$

The multidimensional generalization of the van der Corput sequence is the Halton sequence.

**Definition 17** (Halton Sequence). In dimension $s \geq 1$ choose $s$ bases $b_1, \ldots, b_s$ that are integers and greater than two. The *Halton sequence* in the bases $b_1, \ldots, b_s$ is defined as

$$x_n := \big(\phi_{b_1}(n), \ldots, \phi_{b_s}(n)\big) \in [0, 1)^s \quad \text{for all } n \in \mathbb{N}.$$

A example of the two dimensional Halton sequence in the bases $2, 3$ can be seen in Figure 2.1b.

**Theorem 12** (Discrepancy Bound for the Halton Sequence).
*Let $S$ be the Halton sequence in the pairwise relatively prime bases $b_1, \ldots, b_s$. Then,*

$$D_N^*(S) < \frac{s}{N} + \frac{1}{N} \prod_{i=1}^{s} \left( \frac{b_i - 1}{2 \log(b_i)} \log(N) + \frac{b_i + 1}{2} \right) \quad \text{for all } N \geq 1.$$

The coefficient

$$C(b_1, \ldots, b_s) := \prod_{i=1}^{s} \frac{b_i - 1}{2 \log(b_i)}$$

is minimized by choosing $b_1, \ldots, b_s$ as the first $s$ primes $p_1, \ldots, p_s$, which yields the bound

$$D_N^*(S) \leq C(p_1, \ldots, p_s) N^{-1} \log(N)^s + \mathcal{O}(N^{-1} \log(N)^{s-1}) \quad \text{for all } N \geq 2.$$

### 2.2.3. Weighted Koksma-Hlawka Inequality

There are other useful inequalities of similarly to the classical Koksma-Hlawka inequality that allow to shift the exponents between the norm of the function and the weighted discrepancy. We follow the article Sloan and Woźniakowski 1998. The error estimates hold in a weighted Sobolev space setting.

**Definition 18** (Sobolev Spaces)**.** The *Sobolev space* $W^{k,p}(\Omega)$ is defined by

$$W^{k,p} := \{f \in L^p(\Omega) \mid \text{the weak derivatives of } f \text{ up to order } k \text{ lie in } L^p(\Omega)\}.$$

The Hilbert space case $p = 2$ is also denoted by $H^k(\Omega) := W^{k,2}(\Omega)$. Furthermore, for a ordered vector of weights

$$\gamma_1 \geq \ldots \geq \gamma_s > 0$$

define the weighted Hilbert space norm on $[0,1]^s$ by

$$\|f\|_\gamma^2 := \sum_{I \subset J} \gamma_I^{-1} \int_{[0,1]^{|I|}} \left| \frac{\partial^{|I|}}{\partial x^I} f(x^I \colon \mathbf{1}^{-I}) \right|^2 \mathrm{d}x^I,$$

where $\gamma_I$ denotes the product weight

$$\gamma_I := \prod_{i \in I} \gamma_i.$$

and the *weighted Sobolev space*

$$H_\gamma([0,1]^s) := \{f \in [H^1([0,1])]^s \mid \text{the mixed weak partial}$$
$$\text{derivatives } \tfrac{\partial^{|I|}}{\partial x^I} f, \ I \subset J \text{ exist and } \|f\|_\gamma < \infty\},$$

where $[H^1([0,1])]^s$ denotes the $s$-times tensor product of the space of absolutely continuous functions with derivative in $L^2([0,1])$.

**Definition 19** (Weighted $L^2$-Discrepancy)**.** The *weighted $L^2$-discrepancy* of a point set $P$ is defined as

$$(D_\gamma(P))^2 := \sum_{\emptyset \neq I \subset J} \gamma_I \int_{[0,1]^{|I|}} \left(D_N(x^I \colon \mathbf{1}^{-I})\right)^2 \mathrm{d}x^I.$$

The basis to develop weighted Koksma-Hlawka inequalities is Zaremba's identity; see Zaremba 1968 or Hlawka 1961b. Remember the notion of the local discrepancy in Definition 8.

**Theorem 13** (Zaremba's Identity)**.**
*Let $f \in H_\gamma([0,1]^s)$ and $S = (x_n)_{n \in \mathbb{N}}$ be a sequence in $[0,1]^s$. Denote*

$$J := \{1, \ldots, s\} \quad \text{and} \quad I_N(f) := \frac{1}{N} \sum_{n=1}^{N} f(x_n).$$

*Then*

$$\int\limits_{[0,1]^s} f(x) \, \mathrm{d}x - I_N(f) = \sum_{\emptyset \neq I \subset J} (-1)^{|I|} \int\limits_{[0,1]^{|I|}} D_N\big(x^I \colon \mathbf{1}^{-I}\big) \frac{\partial^{|I|}}{\partial x^I} f(x^I \colon \mathbf{1}^{-I}) \, \mathrm{d}x^I.$$

With Zaremba's identity, the weighted Koksma-Hlawka inequality can be derived. It provides an error estimate for quasi-Monte Carlo integration in the weighted Hilbert space setting.

**Theorem 14** (Weighted Koksma-Hlawka Inequality)**.**
*Consider $f \in H_\gamma([0,1]^s)$ and a point set $P = \{x_1, \ldots, x_N\}$. Then the error bound*

$$\left| \int_{[0,1]^s} f(x) \, \mathrm{d}x - \frac{1}{N} \sum_{n=1}^{N} f(x_n) \right| \leq \|f\|_\gamma D_\gamma(P)$$

*holds.*

*Proof.* We follow the proof of Sloan and Woźniakowski 1998, p. 9–10. By expanding Zaremba's identity (Theorem 13) with the factor $\gamma_I^{1/2}$, we obtain

$$\int_{[0,1]^s} f(x) \, \mathrm{d}x - I_N(f) =$$

$$\sum_{\emptyset \neq I \subset J} (-1)^{|I|} \int_{[0,1]^{|I|}} \gamma_I^{1/2} D_N\big(x^I \colon \mathbf{1}^{-I}\big) \gamma_I^{-1/2} \frac{\partial^{|I|}}{\partial x^I} f(x^I \colon \mathbf{1}^{-I}) \, \mathrm{d}x^I.$$

Using the Cauchy-Schwartz inequality for the $L^2$ inner product leads to

$$\left| \int_{[0,1]^s} f(x) \, \mathrm{d}x - I_N(f) \right| \leq$$

$$\sum_{\emptyset \neq I \subset J} \sqrt{\gamma_I \int_{[0,1]^{|I|}} [D_N(x^I : \mathbf{1}^{-I})]^2 \, \mathrm{d}x^I} \sqrt{\gamma_I^{-1} \int_{[0,1]^{|I|}} \left[ \frac{\partial^{|I|}}{\partial x^I} f(x^I \colon \mathbf{1}^{-I}) \right]^2 \, \mathrm{d}x^I}$$

$$\leq D_\gamma(P) \|f\|_\gamma,$$

where the last step is due to the application of the Cauchy-Schwartz inequality for the euclidean scalar product. Note that the empty set is not excluded in the definition of the weighted Sobolev norm $\|\cdot\|_\gamma$. Adding it to the sum corresponds to adding the term $f(\mathbf{1})$, which makes the sum certainly larger and makes $\|\cdot\|_\gamma$ a norm instead of a semi-norm. $\qquad\square$

Weighted Koksma-Hlawka inequalities can be formulated in for functions in non-Hilbert Sobolev spaces using the $L^p$ norm instead of the $L^2$ norm.

**Definition 20** ($L^p$-Discrepancy)**.** The $L^p$-*discrepancy* of the point set $P$, for $p \in [1, \infty)$ is defined by

$$[D_{\gamma, p}(P)]^p := \sum_{\emptyset \neq I \subset J} \gamma_I^{p/2} \int_{[0,1]^{|I|}} \left[ D_N(x^I \colon \mathbf{1}^{-I}) \right]^p \mathrm{d}x^I$$

and for $p = \infty$ by

$$D_{\gamma, \infty}(P) := \sup_{x \in [0,1]^s} \max_{\emptyset \neq I \subset J} \gamma_I^{1/2} |D_N(x^I \colon \mathbf{1}^{-I})|.$$

**Definition 21** (Weighted Sobolev spaces)**.** Define for $q \in [1, \infty)$, the norm

$$\|f\|_{\gamma, q}^q := \sum_{I \subset J} \gamma_I^{-q/2} \int_{[0,1]^{|I|}} \left| \frac{\partial^{|I|}}{\partial x^I} f(x^I \colon \mathbf{1}^{-I}) \right|^q \mathrm{d}x^I$$

and for $q = \infty$

$$\|f\|_{\gamma, \infty} := \sup_{x \in [0,1]^s} \max_{I \subset J} \gamma_I^{-1/2} \left| \frac{\partial^{|I|}}{\partial x^I} f(x^I \colon \mathbf{1}^{-I}) \right|.$$

The *weighted Sobolev space* on $[0, 1]^s$ is defined by

$$H_{\gamma, q}([0, 1]^s) := \{ f \in [W^{1,q}([0, 1])]^s \mid \text{the mixed weak partial derivatives}$$
$$\tfrac{\partial^{|I|}}{\partial x^I} f, \ I \subset J \text{ exist and } \|f\|_\gamma < \infty \}.$$

**Theorem 15** (Weighted Koksma-Hlawka Inequality; $L^p$-Version)**.**
*Let $p \in [0, \infty]$ and $q$ the Hölder conjugate of $p$, i.e., $1/p + 1/q = 1$. Let $f$ be in $H_{\gamma, q}([0, 1]^s)$ and let $P = \{x_1, \ldots, x_N\}$ be a point set on $[0, 1]^s$. Then,*

$$\left| \int_{[0,1]^s} f(x) \, \mathrm{d}x - \frac{1}{N} \sum_{n=1}^N f(x_n) \right| \leq \|f\|_{\gamma, q} D_{\gamma, p}(P)$$

*holds.*

*Proof.* Use the Hölder inequality instead of the Cauchy-Schwartz inequality in the proof of Theorem 14. $\qquad\square$

## 2.3. Spherical Radial Decomposition

The quasi-Monte Carlo methods explained in the previous sections can be combined with other integration techniques. The integration over the whole space can be split into an integral over the unit sphere and an integral over the radius. The one dimensional radius integral can be treated by classical integration techniques or analytically and the high dimensional sphere integral by quasi-Monte Carlo integration. This proves to be advantageous, because the sphere integral has a lower dimension compared to the whole space, which reduces the variance; see Ackooij and Henrion 2014, Equation (1.5). Furthermore, we will discuss a gradient representation for the spherical radial decomposition of a probabilistic constraint function, which can be used for optimization. For the derivation of the transformation, we follow Genz and Bretz 2009, Section 4.1.1.

We start by sketching the basic idea. For $x \in \mathbb{R}^m$, consider a integrable function $f : \mathbb{R}^m \to \mathbb{R}$. By Stein and Shakarchi 2003, p. 180, the integral over $f$ can be transformed by writing $x = rv$, where $v$ is a vector on the unit sphere

$$\mathbb{S}^{m-1} := \{v \in \mathbb{R}^m \mid \|v\| = 1\}$$

and $r = \|x\|$. The transformation is given by

$$\int_{\mathbb{R}^m} f(x)\,\mathrm{d}x = \int_{\mathbb{S}^{m-1}} \int_0^\infty r^{m-1} f(rv)\,\mathrm{d}r\,\mathrm{d}\sigma(v),$$

where $\sigma$ is the surface measure on the sphere. The integral over a cuboid $[a,b] \subset \mathbb{R}^m$ instead of the whole space can be calculated by transforming the boundaries for the radius integral via

$$r_l(v) := \min\{r \geq 0 \mid a \leq rv \leq b\}, \tag{2.2a}$$

$$r_u(v) := \max\{r \geq 0 \mid a \leq rv \leq b\}. \tag{2.2b}$$

The surface measure on the sphere can be normalized with the area of the surface of the unit sphere

$$\mathcal{H}^{m-1}(\mathbb{S}^{m-1}) = \frac{2\pi^{m/2}}{\Gamma(\frac{m}{2})},$$

to obtain the uniform distribution on the sphere by $u = \sigma/\mathcal{H}^{m-1}(\mathbb{S}^{m-1})$, which verifies

$$\int_a^b f(x)\,\mathrm{d}x = \frac{2\pi^{m/2}}{\Gamma(\frac{m}{2})} \int_{\mathbb{S}^{m-1}} \int_{r_l(v)}^{r_u(v)} r^{m-1} f(rv)\,\mathrm{d}r\,\mathrm{d}u(v). \tag{2.3}$$

We apply the method to the density function of a multivariate normal distribution. Consider the positive definite covariance matrix $\Sigma \in \mathbb{R}^{m \times m}$, which has the Cholesky decomposition $\Sigma = LL^T$. The matrix $L$ is a lower triangular matrix. We discuss the case for the expected value $\mu = 0$. A nonzero expectation $\mu \in \mathbb{R}^m$ for $x$ is treated by the simple shift $y = x - \mu$, where $y$ has an expected value of zero. The multivariate normal density function is defined by

$$\vartheta(x) := \frac{\exp\left(-\frac{1}{2}x^T\Sigma^{-1}x\right)}{\sqrt{(2\pi)^m \det(\Sigma)}}.$$

The probability for $x \sim \mathcal{N}(0, \Sigma)$ to lie in the box $[a, b]$ is given by $\int_a^b \vartheta(x)\,\mathrm{d}x$. To make $\vartheta$ only dependent on the norm of $x$, we use the transformation $Ly = x$. This makes $y \sim \mathcal{N}(0, \mathbb{1})$ a standard normal distributed random variable. The Jacobian is given by $\det(L) = \sqrt{\det(\Sigma)}$. The integral becomes

$$(2\pi)^{-m/2} \det(\Sigma)^{-1/2} \int_a^b \exp\left(-\tfrac{1}{2}x^T\Sigma^{-1}x\right)\,\mathrm{d}x = (2\pi)^{-m/2} \int_{a \leq Ly \leq b} \exp\left(-\tfrac{1}{2}y^Ty\right)\,\mathrm{d}y.$$

Using the transformation (2.3), the probability function can be rewritten as

$$\int_a^b \vartheta(x)\,\mathrm{d}x = \frac{2^{1-m/2}}{\Gamma(\frac{m}{2})} \int_{\|v\|=1} \int_{r_l(Lv)}^{r_u(Lv)} r^{m-1} \exp\left(-r^2/2\right)\,\mathrm{d}r\,\mathrm{d}u(v) \qquad (2.4)$$

with $r_l$ and $r_u$ as in Equations (2.2).

The integrand in Equation (2.4) is the probability density function of the chi distribution with $m$ degrees of freedom given by
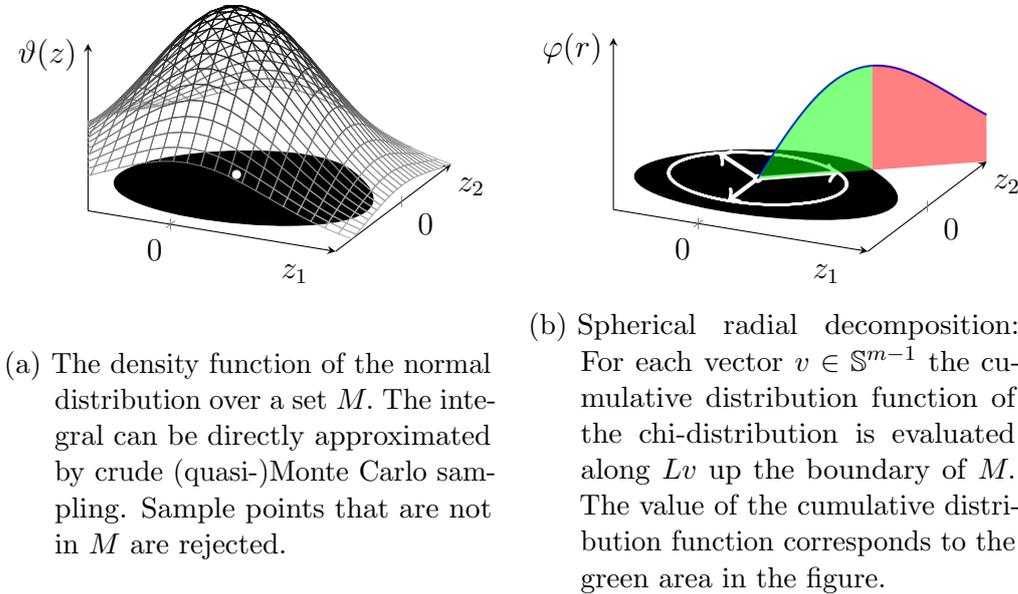
$$f_\chi(x) = \frac{2^{1-m/2}}{\Gamma(\frac{m}{2})} x^{m-1} \exp\left(-x^2/2\right), \quad \text{for } x \geq 0$$

and $f_\chi(x) = 0$, for $x < 0$. This allows us to write (2.4) as

$$\int_a^b \vartheta(x)\,\mathrm{d}x = \int_{\|v\|=1} \int_{r_l(Lv)}^{r_u(Lv)} f_\chi(r)\,\mathrm{d}r\,\mathrm{d}u(v). \qquad (2.5)$$

See Figure 2.2 for the idea.

Figure 2.2.: Evaluation of the probabilistic constraint by a direct approach and by spherical radial decomposition

(a) The density function of the normal distribution over a set $M$. The integral can be directly approximated by crude (quasi-)Monte Carlo sampling. Sample points that are not in $M$ are rejected.

(b) Spherical radial decomposition: For each vector $v \in \mathbb{S}^{m-1}$ the cumulative distribution function of the chi-distribution is evaluated along $Lv$ up the boundary of $M$. The value of the cumulative distribution function corresponds to the green area in the figure.

## 2.3.1. Spherical Radial Decomposition in Chance Constrained Programming

The spherical radial decomposition is a useful tool to treat the chance constrained optimization problem of the introduction of this chapter. A preliminary version of the following sections has been made publicly available online; see Wintergerst 2017.

Let $f : \mathbb{R}^n \to \mathbb{R}$ and $g : \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}$ be continuously differentiable functions, and let $\xi \sim \mathcal{N}(0, \Sigma)$ be a normal distributed random vector with realizations in $\mathbb{R}^m$. The positive definite covariance matrix $\Sigma \in \mathbb{R}^{m \times m}$ has the Cholesky decomposition $\Sigma = LL^T$. Consider the optimization problem

$$
\begin{aligned}
\min_{x \in \mathbb{R}^n} \quad & f(x) \\
\text{s.t.} \quad & \mathbb{P}(g(x, \xi) \leq 0) \geq p^{\mathrm{bd}}.
\end{aligned}
$$

Denote the chi distribution by $\chi$ and define

$$
e(x, v) := \chi(\{r \geq 0 : g(x, rLv) \leq 0\}). \tag{2.6}
$$

Then, a representation of $G(x) := \mathbb{P}(g(x,\xi) \leq 0)$ is given by

$$G(x) = \int_{v \in \mathbb{S}^{m-1}} e(x,v) \, \mathrm{d}u(v), \tag{2.7}$$

where $u$ is the uniform distribution over the sphere.

**Definition 22** (Set of Feasible Realizations)**.** For a fixed optimization variable $x \in \mathbb{R}^n$, denote the *set of feasible realizations* by

$$M(x) := \{z \in \mathbb{R}^m \mid g(x,z) \leq 0\}$$

and its boundary by

$$N(x) := \{z \in \mathbb{R}^m \mid g(x,z) = 0\}.$$

## 2.3.2. Gradient Representation for a Single Constraint

To apply efficient optimization algorithms, an analytic representation of the gradient of the probability function is desirable. In the articles Ackooij and Henrion 2014 and Ackooij and Henrion 2017 such a representation is given, under the strong assumption that the function $g$ is convex in the second argument. We generalize this to the case, where it is not assumed that the function $g$ is convex. Instead, we only require the sets of feasible realizations $M(x)$ to be convex and bounded. We follow the ideas of Ackooij and Henrion 2014, while adapting the argumentation for the new assumptions. Throughout this section, we pose the following assumptions to the point $x \in \mathbb{R}^n$, where we want to differentiate the probability function:

(A1) The sets $M(y)$ are convex in a neighborhood $U$ of $x$.

(A2) The set $M(x)$ is bounded.

(A3) The function $g : \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}$ is continuously differentiable.

(A4) The gradient $\nabla_z g(x,z)$ is nonzero on the set $N(x)$.

(A5) The point $(x,0)$ is a Slater-point, i.e., $g(x,0) < 0$.

*Remark* 14. The assumption (A1) is significantly weaker than imposing the convexity of the function $g(x,\cdot)$. The function $g(x,\cdot)$ does not even have to be quasiconvex for (A1) to be fulfilled. For example, consider $g : \mathbb{R} \times \mathbb{R} \to \mathbb{R}$ with

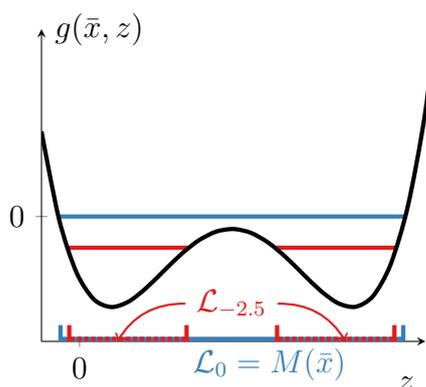$$g(x,z) := x\,z(z-1)(z-3)(z-4) - 5x$$

Figure 2.4.: The function fulfills (A1) at the point $\bar{x}$ as the level set $\mathcal{L}_0 := M(\bar{x})$ (in blue) is convex. However, the function is not quasiconvex in the second argument, because the level set $\mathcal{L}_{-2.5} := \{z \in \mathbb{R} \mid g(\bar{x}, z) \leq -2.5\}$ (in red) is not convex.

at the point $\bar{x} = 1$. Its graph is shown in Figure 2.4. The function $g(\bar{x}, \cdot)$ is not quasiconvex, yet (A1) is fulfilled in the point $\bar{x}$.

*Remark* 15. The nonzero gradient assumption

$$\nabla_z g(x, z) \neq 0, \quad \text{for } z \in N(x)$$

is fulfilled if $g(x, \cdot)$ is convex. This can be shown as follows: Suppose $\nabla_z g(x, z) = 0$, then using, in this order, the convexity of $g(x, \cdot)$, the equation $g(x, z) = 0$ and the Slater condition (A5) yields

$$0 = \nabla_z g(x, z)^T (0 - z) \leq g(x, 0) - g(x, z) = g(x, 0) < 0,$$

which is a contradiction.

*Remark* 16. Assumption (A5) is no severe restriction for a symmetrical probability distribution with expected value 0. If 0 is a point on the boundary of $M(x)$, the convexity of $M(x)$ implies that $M(x)$ is contained in a halfspace defined by a separating plane through 0. This, however implies that $G(x) \leq 0.5$.

For applications, we are typically interested in points $G(x) \geq p^{\mathrm{bd}}$ for a $p^{\mathrm{bd}}$ close to 1 or, to phrase it differently: Points that do not fulfill (A5) are not feasible with respect to the constraint $G(x) \geq p^{\mathrm{bd}}$ for $p^{\mathrm{bd}} > 0.5$.

Our goal is to show a gradient representation of the form

$$\nabla G(x) = -\int_{v \in \mathbb{S}^{m-1}} \frac{f_\chi(r(x, v))}{\nabla_z g(x, r(x, v)Lv)^T Lv} \, \nabla_x g(x, r(x, v)Lv) \, \mathrm{d}u(v), \quad (2.8)$$

where $r(x, v) > 0$ is chosen such that $g(x, r(x, v)Lv) = 0$ and $f_\chi$ is the probability distribution function of the chi distribution. The key points

to show are the following: The denominator in Equation (2.8) does not approach zero, and a locally unique continuously differentiable function $r$ with $g(x, r(x, v)Lv) = 0$ exists in a neighborhood of a given point $(x, v)$.

**Lemma 25** (Nonzero directional derivative).
*For a direction $v$, let $r_v > 0$ be a scalar fulfilling*

$$g(x, r_v Lv) = 0, \qquad (2.9)$$

*which exists, because $M(x)$ is bounded by Assumption (A2).*
*   Then, there is a constant $C > 0$ such that*

$$\nabla_z g(x, r_v Lv)^T Lv > C \quad \text{for all } v \in \mathbb{S}^{m-1}. \qquad (2.10)$$

*Proof.* First, note that

$$\nabla_z g(x, r_v Lv)^T Lv \geq 0, \qquad (2.11)$$

as by (2.9) and the convexity of $M(x)$ stated in (A1),

$$g(x, r_v Lv + tLv) - g(x, r_v Lv) = g(x, (r_v + t)Lv) \geq 0,$$

holds. Suppose that

$$\nabla_z g(x, r_v Lv)^T Lv = 0.$$

Define the affine linear (half-)spaces

$$H^=(x) := \{r_v Lv + z \mid z \in \mathbb{R}^m : \nabla_z g(x, r_v Lv)^T z = 0\},$$
$$H^\leq(x) := \{r_v Lv + z \mid z \in \mathbb{R}^m : \nabla_z g(x, r_v Lv)^T z \leq 0\},$$
$$H^>(x) := \{r_v Lv + z \mid z \in \mathbb{R}^m : \nabla_z g(x, r_v Lv)^T z > 0\}.$$

Because of $\nabla_z g(x, r_v Lv)^T(-r_v Lv) = 0$, the point 0 is in $H^=(x)$. Since the set $M(x)$ is convex, the inclusion

$$M(x) \subset H^\leq(x) \qquad (2.12)$$

holds. In a ball

$$B_\epsilon(0) := \{y \in R^m \mid \|y\| < \epsilon\}$$

around 0, we have $g(x, y) < 0$ for $y \in B_\epsilon(0)$, because of the generalized Slater condition (A5). Choose a point $y \in H^>(x) \cap B_\epsilon(0)$. Because the origin lies in $H^=(x)$ and the gradient $\nabla_z g(x, rLv)$ is nonzero, the set $H^>(x) \cap B_\epsilon(0)$ is nonempty. The strict inequality

$$g(x, y) < 0,$$

Figure 2.5.: If the directional derivative is zero, while the z-gradient at $rLv$ is not, then either the Slater point 0 is not an interior point or the set $M(x)$ is not convex.

holds, since $y \in B_\epsilon(0)$. However, as $y$ lies in $H^>(x) = \mathcal{C}\big(H^\le(x)\big)$, the Inclusion (2.12) implies that $y \notin M(x)$. This is a contradiction to

$$y \in B_\epsilon(x) \subset M(x).$$

Therefore, $\nabla_z g(x, r_v Lv)^T Lv \ne 0$ holds and, by Equation (2.11), the inequality

$$\nabla_z g(x, r_v Lv)^T Lv > 0$$

follows. The idea of this step is illustrated in Figure 2.5. The boundedness of $M(x)$ implies that the maximum possible radius $r_v$ is bounded from above by a constant $\bar{r}$, while the Slater condition (A5) ensures it is bounded from below with a positive lower bound $\underline{r}$. Furthermore, by the continuity of $g$ the set $N(x)$ is closed. Hence, the minimum of the continuous function

$$(v, r_v) \mapsto \nabla_z g(x, r_v Lv)^T Lv$$

is attained and we can choose

$$C = \min\{\nabla_z g(x, r_v Lv)^T Lv \mid v \in \mathbb{S}^{m-1}, r_v \in [\underline{r}, \bar{r}] : r_v Lv \in N(x)\} > 0. \quad \square$$

Before we begin the construction of the radius function $r(\cdot, \cdot)$, we show that Assumption (A4) extends to a neighborhood of $x$. This prevents the situation of Figure 2.6 for small changes of $x$.

**Lemma 26** (Nonzero Gradient in a Neighborhood).
*Let*

$$\nabla_z g(x, z) \ne 0 \quad \forall z \in N(x)$$

Figure 2.6.: In the Figure, the function $g$ is locally constant, which means that the boundary $N(x)$ of $M(x)$ is given by an area instead of a line. This situation is avoided by Assumption (A4).

*Then, there exists a neighborhood $U$ of $x$ such that for all $y \in U$*

$$\nabla_z g(y, z) \neq 0 \quad \forall z \in N(y).$$

*Proof.* Observe that $\nabla_z g(x, z) \neq 0$ near the set $N(x)$, i.e.,

$$\nabla_z g(x, z) \neq 0 \text{ on } N_\epsilon(x) := \{z \in \mathbb{R}^m \mid \text{dist}(z, N(x)) < \epsilon\}, \qquad (2.13)$$

for $\epsilon > 0$ small enough. This is true, because $N(x)$ is compact, which implies the existence of a $c > 0$ such that

$$\|\nabla_z g(x, z)\| > c \quad \text{on } N(x).$$

Therefore, Equation (2.13) holds by the continuity of the gradient. By choosing $U$ small enough to guarantee $N(y) \subset N_\epsilon(x)$ for $y \in U$ we obtain

$$\nabla_z g(y, z) \neq 0 \quad \text{for all } z \in N(y) \subset N_\epsilon(x). \qquad \square$$

The next step is to show the existence of a continuously differentiable radius function, that fulfills $g(x, r(x, v)Lv) = 0$. This is a consequence of the implicit function theorem (Theorem A). The results of the following Lemma are analogous to Ackooij and Henrion 2014, Lemma 3.2 and parts of its proof are similar.

**Lemma 27** (Implicit Function for the Radius)**.**
*There exist neighborhoods $U$ of $x$, $V$ of $v$ and a continuously differentiable function $r : U \times V \to \mathbb{R}_{\geq 0}$ such that*

 *a) For all $(y, w, s) \in U \times V \times \mathbb{R}_{\geq 0}$, we have*

$$g(y, sLw) = 0 \iff s = r(y, w).$$

*b) For all $(y, w) \in U \times V$ the gradient representation*

$$\nabla_x r(y, w) = -\frac{1}{\nabla_z g(y, r(y, w)Lw)^T Lw} \, \nabla_x g(y, r(y, w)Lw)$$

*holds.*

*Proof.* The inequality

$$\nabla_z g(x, rLv)^T Lv \geq C > 0$$

from Lemma 25 allows us to apply the implicit function theorem (Theorem A) to the Equation

$$g(x, r_v Lv) = 0,$$

to derive the existence of neighborhoods $U$ of $x$, $V$ of $v$ and $W$ of $r_v$ such that a continuously differentiable function $r : U \times V \to W$ exists and fulfills the equivalence

$$g(y, sLw) = 0, \ (y, w, s) \in U \times V \times W \iff s = r(y, w), \ (y, w) \in U \times V.$$

By choosing $U$ sufficiently small, it is guaranteed that $r$ maps to $\mathbb{R}_{\geq 0}$ and, due to the continuity of the function $g$, that

$$g(y, 0) < 0 \quad \text{for all } y \in U.$$

Because of Lemma 26, by further shrinking $U$, we obtain

$$\nabla_z g(y, z) \neq 0 \quad \text{for all } z \in N(y) \subset N_\epsilon(x)$$

if $y$ lies in $U$. Consequently, choosing

$$0 < s = r(y, w), \quad (y, w) \in U \times V$$

implies

$$g(y, sLw) = 0, \text{ for all } (y, w, s) \in U \times V \times W, \text{ where } W \subset \mathbb{R}_+.$$

For the converse implication in a), we have to show the uniqueness of the root of the function $g(y, \cdot Lw)$ in $\mathbb{R}_{>0}$. We assume the contrary. Suppose there is a $r^* \in \mathbb{R}_{\geq 0}$ not equal to $r(y, w)$ fulfilling $g(y, r^*Lw) = 0$, for $(y, w) \in U \times V$. Consider the case $r^* > r(y, w)$. By the convexity of $M(y)$, the point

$$[r(y, w) + \lambda(r^* - r(y, w))]Lw$$

is feasible for all $\lambda \in [0,1]$. This yields

$$0 \geq g(y, [r(y,w) + \lambda(r^* - r(y,w))]Lw)$$

(by the Taylor formula)

$$= g(y, r(y,w)Lw) + \lambda(r^* - r(y,w))\nabla_z g(y, \theta Lw)^T Lw, \qquad (2.14)$$

where $\theta \in [r(y,w), r(y,w) + \lambda(r^* - r(y,w))]$. By Lemma 25, the bound

$$\nabla_z g(x, r(x,v)Lv)^T Lv \geq C > 0$$

holds and we have, due to the continuity of the gradient,

$$\nabla_z g(y, r(y,w)Lw)^T Lw \geq D \text{ for a } D \in (0, C).$$

Hence, it is possible to choose a neighborhood $\widetilde{W}$ of $r(y,w)$, where

$$\nabla_z g(y, rLw)^T Lw \geq \tfrac{D}{2} \text{ holds for } r \in \widetilde{W}.$$

For $\lambda > 0$ small enough, $\theta$ lies in $\widetilde{W}$. Continuing (2.14), we obtain

$$0 \geq g(y, r(y,w)Lw) + \lambda(r^* - r(y,w))\nabla_z g(y, \theta Lw)^T Lw$$
$$\geq \lambda(r^* - r(y,w))\tfrac{D}{2} > 0,$$

which is a contradiction. The case $r(y,w) > r^*$ can be treated analogously by exchanging the roles of $r^*$ and $r(y,w)$. For an illustration of this idea see Figure 2.7.

Equation b) holds by the implicit function theorem. $\qquad \square$

With the radius function $r$, we can represent the function $e$ defined in (2.6) as well as the gradient of $e$. This allows us to deduct (2.8) from (2.7).

**Lemma 28** (Characterization of $e(x,v)$)**.**
*Remember the Definition (2.6) given by*

$$e(x,v) := \chi(\{r \geq 0 : g(x, rLv) \leq 0\}),$$

*where $\chi$ is the measure of the chi-distribution. There exist neighborhoods $U$ of $x$ and $V$ of $v$ such that for $y \in U$ and $w \in V$,*

$$e(y,w) = F_\chi(r(y,w)) \qquad (2.15)$$

*holds, where $F_\chi$ is the cumulative distribution function of the chi-distribution.*

Figure 2.7.: A non unique root of $g(x, \cdot Lv)$ (marked by the blue dots) implies either that the directional derivative becomes zero or that the feasible set (depicted by the green line) is nonconvex.

*Proof.* Choose the neighborhoods $U$ and $V$ according to Lemma 27. Then, by Lemma 27 a), we obtain

$$
\begin{aligned}
e(y, w) &:= \chi\big(\{r \geq 0 : g(y, rLw) \leq 0\}\big) \\
&= \chi\big([0, r(y, w)]\big) \\
&= F_\chi\big(r(y, w)\big) \text{ for all } (y, w) \in U \times V,
\end{aligned}
$$

because the chi-distribution is zero for negative values. $\qquad\square$

This leads to the gradient representation of $e$.

**Corollary 4** (Gradient of $e(x, v)$)**.**
*The function $e : \mathbb{R}^n \times \mathbb{S}^{m-1} \to \mathbb{R}$ is continuously differentiable with respect to $x$ and its gradient is given by*

$$
\nabla_x e(x, v) = -\frac{f_\chi(r(x, v))}{\nabla_z g(x, r(x, v)Lv)^T Lv} \, \nabla_x g(x, r(x, v)Lv),
$$

*where $f_\chi$ is the density of the chi-distribution with $m$ degrees of freedom and $r$ is the function of Lemma 27.*

*Proof.* The application of the chain rule to Equation (2.15) together with Lemma 27 b) and

$$
F_\chi'(r(x, v)) = f_\chi(r(x, v))
$$

yields the formula for the partial derivatives. Because

$$\nabla_z g(x, r(x,v)Lv)^T Lv \geq C > 0$$

holds and all functions in the formula are continuous, the function $\nabla_x e(\,\cdot\,, v)$ is continuous. $\square$

Now we have the necessary tools to prove the gradient representation (2.8).

**Theorem 16** (Gradient Formula for the Probability Function).
*Let Assumptions (A1)-(A5) hold in $x \in \mathbb{R}^n$.*

*Then, $G$ is continuously differentiable in a neighborhood $U$ of $x$ and its gradient is given by*

$$\nabla G(y) = - \int\limits_{v \in \mathbb{S}^{m-1}} \frac{f_\chi(r(y,v))}{\nabla_z g(y, r(y,v)Lv)^T Lv} \, \nabla_x g(y, r(y,v)Lv) \, \mathrm{d}u(v),$$

*for all $y \in U$.*

*Proof.* By Lemma 26 and assumption (A4), we have

$$\nabla_z g(y, z) \neq 0 \quad \forall z \in N(y),$$

provided we choose $y$ in a sufficiently small neighborhood $U$ of $x$. The continuity of $g$ together with (A2) implies $M(y)$ is bounded in each neighborhood of $x$. Using the representation (2.7) we have

$$\nabla G(y) = \nabla \int_{v \in \mathbb{S}^{m-1}} e(y,v) \, \mathrm{d}u(v). \tag{2.16}$$

To differentiate under the integral, we want to apply Lebesgue's dominated convergence theorem. Therefore, we have to show that $\|\nabla_x e(y,v)\|$ is bounded for all directions $v \in \mathbb{S}^{m-1}$. By Corollary 4 we have

$$\|\nabla_x e(y,v)\| = \frac{|f_\chi(r(y,v))|}{|\nabla_z g(y, r(y,v)Lv)^T Lv|} \|\nabla_x g(y, r(y,v)Lv)\|$$

$$\leq \frac{1}{C} \max_{v \in \mathbb{S}^{m-1}} \|\nabla_x g(y, r(y,v)Lv)\|,$$

where we used Lemma 25 and the fact that $f_\chi$ is bounded by 1. The occurring maximum exists as a consequence of the Weierstrass theorem, because the sphere $\mathbb{S}^{m-1}$ is compact and $\nabla_x g(y, r(y, \cdot)L\,\cdot)$ is continuous

(the continuity of $r$ is a result of Lemma 27). Hence, we can differentiate under the integral in (2.16) and obtain

$$\nabla G(y) = \int_{v \in \mathbb{S}^{m-1}} \nabla_x e(y, v) \, \mathrm{d}u(v)$$

(by Corollary 4)

$$= - \int_{v \in \mathbb{S}^{m-1}} \frac{f_\chi(r(y, v))}{\nabla_z g(y, r(y, v)Lv)^T Lv} \, \nabla_x g(y, r(y, v)Lv) \, \mathrm{d}u(v). \quad \square$$

### 2.3.3. Gradient Representation for Multiple Constraints

In the previous section, we have treated the case of a real-valued constraint function $g : \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}$. In this section, we consider a vector-valued constraint function $g : \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}^p$ and the *joint chance constraint* or *joint probabilistic constraint*

$$G(x) := \mathbb{P}(g(x, \xi) \leq 0) \leq p^{\mathrm{bd}}.$$

This case was considered in Ackooij and Henrion 2017 under the assumption that $g(x, \cdot)$ is convex. We follow the basic line of argumentation therein, but adapt it to our assumptions. Define the maximum function $g^{\mathrm{m}} : \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}$

$$g^{\mathrm{m}}(x, z) := \max_{i=1,\dots,p} g_i(x, z). \tag{2.17}$$

Note that $g^m$ is not differentiable in general.

**Definition 23** (Set of Feasible Realizations)**.** Define the *set of feasible realizations for all components*, the *set of feasible realizations for the i-th component* and its boundary for the *i*-th component as

$$M(x) := \{z \in \mathbb{R}^m \mid g(x, z) \leq 0\},$$
$$M_i(x) := \{z \in \mathbb{R}^m \mid g_i(x, z) \leq 0\},$$
$$N_i(x) := \{z \in \mathbb{R}^m \mid g_i(x, z) = 0\}.$$

**Definition 24** (Sets of Finite and Infinite Directions)**.** We denote the sets of *finite* and *infinite* directions by

$$\mathbb{F}_i(x) := \{v \in \mathbb{S}^{m-1} \mid \exists r > 0 : g_i(x, rLv) = 0\},$$
$$\mathbb{I}_i(x) := \{v \in \mathbb{S}^{m-1} \mid \forall r > 0 : g_i(x, rLv) < 0\}.$$

Throughout this section, we assume that in the point $x \in \mathbb{R}^n$, where we want to differentiate, the following holds:

(B1) The sets $M_i(y)$ are convex for $y$ in a neighborhood $U$ of $x$.

(B2) The set $M(x)$ is bounded.

(B3) The functions $g_i : \mathbb{R}^n \times \mathbb{R}^m \to \mathbb{R}$ are continuously differentiable.

(B4) The gradients $\nabla_z g_i(x, z)$ are nonzero on the sets $N_i(x)$.

(B5) The point $(x, 0)$ is a Slater-point, i.e., $g^m(x, 0) < 0$.

The indices $i$ are in $\{1, \ldots, p\}$. Under this assumptions, directions in a finite direction $v \in \mathbb{F}_i(x)$ stay finite. We formulate this in a Lemma.

**Lemma 29** (The Set of Finite Directions is Open)**.**
*For a fixed index $i \in \{1, \ldots, p\}$, the set $\mathbb{F}_i(x)$ is open.*

*Proof.* Because $\mathbb{F}_i(x)$ and $\mathbb{I}_i(x)$ are complements of each other, the equality

$$\mathbb{F}_i(x) = \mathbb{S}^{m-1} \setminus \mathbb{I}_i(x)$$

holds and it is equivalent to show that $\mathbb{I}_i(x)$ is closed. Let us assume it is not. Then, there is a sequence

$$(v^k)_{k \in \mathbb{N}} \subset \mathbb{I}(x) \quad \text{with} \quad \lim_{k \to \infty} v^k = v$$

such that $v \in \mathbb{F}_i(x)$. By the definition of $\mathbb{F}_i(x)$, there exists a $\bar{r} > 0$ with $g_i(x, \bar{r}Lv) = 0$. Since $M_i(x)$ is convex and $\nabla_z g_i(x, z) \neq 0$ on $N_i(x)$, we have

$$g_i(x, (\bar{r} + t)Lv) > 0 \quad \text{for } t > 0$$

and moreover, by the continuity of $g_i$, we have

$$g_i(x, z) > 0 \quad \text{for all } z \in B_\epsilon((\bar{r} + t)Lv)$$

if the ball $B_\epsilon((\bar{r} + t)Lv)$ around $(\bar{r} + t)Lv$ is sufficiently small. However, since $v^k$ converges to $v$, the vector $(\bar{r} + t)Lv^k$ lies in $B_\epsilon((\bar{r} + t)Lv)$ for $k$ sufficiently large. This shows

$$g_i(x, (\bar{r} + t)Lv^k) > 0,$$

whereas the Slater condition states

$$g_i(x, 0) < 0.$$

Now Bolzano's intermediate value theorem implies the existence of a scalar $r \in (0, \bar{r} + t)$ such that

$$g_i(x, rLv^k) = 0.$$

This contradicts $v^k \in \mathbb{I}_i(x)$ and proves the Lemma. $\qquad\square$

Next, we show that the directional derivative in direction $Lv$ is locally bounded from below.

**Corollary 5** (Locally Bounded Directional Derivative)**.**
*Consider the case, where $M_i(x)$ is not bounded. Then, for a given $v \in \mathbb{F}_i(x)$, there is a neighborhood $V$ of $v$ and a constant $C > 0$ such that*

$$\nabla_z g_i(x, r_v Lv)^T Lv > C \text{ for all } v \in V.$$

*Proof.* By Lemma 29 the set $\mathbb{F}(x)$ is open, hence it is possible to choose a neighborhood of $v$ fulfilling $\widetilde{V} \subset \mathbb{F}(x)$. Choose a sufficiently small subset

$$V \subset \widetilde{V} \quad \text{such that} \quad \mathrm{cl}(V) \subset \widetilde{V}.$$

One can completely follow the proof of Lemma 25 up to its last paragraph, where it is necessary to replace $\mathbb{S}^{m-1}$ by $\mathrm{cl}(V)$ and choose

$$C = \min\{\nabla_z g_i(x, r_v Lv)^T Lv \mid v \in \mathrm{cl}(V), r_v \in [\underline{r}, \bar{r}] : r_v Lv \in N_i(x)\} > 0,$$

which shows the claim. $\qquad\square$

For each active constraint, we retain results similar to Section 2.3.2. This will allow us to characterize the Clarke subdifferential of the probability function.

**Definition 25** (Set of Active Constraints)**.** We denote the *set of active constraints* by

$$\mathcal{A}(x, v) := \{i \in \{1, \dots, p\} \mid \exists r > 0 : g_i(x, rLv) = 0\}.$$

**Corollary 6** (Implicit Function for the Radius for the Active Constraints)**.**
*For any active index $i \in \mathcal{A}(x, v)$, there exist neighborhoods $U$ of $x$, $V$ of $v$ and a continuously differentiable function $r_i : U \times V \to \mathbb{R}_{\geq 0}$ such that*

   *a) For all $(y, w, s) \in U \times V \times \mathbb{R}_{\geq 0}$, we have*

$$g_i(y, sLw) = 0 \quad \Longleftrightarrow \quad s = r_i(y, w).$$

*b) For all $(y, w) \in U \times V$ the gradient formula*

$$\nabla_x r_i(y, w) = -\frac{1}{\nabla_z g_i(y, r_i(y, w)Lw)^T Lw} \; \nabla_x g_i(y, r_i(y, w)Lw)$$

*holds.*

*Proof.* Lemma 29 states that for a finite direction $v \in \mathbb{F}_j(x)$, there is a neighborhood $V$, such that $w \in \mathbb{F}_j(x)$ for all $w \in V$. By the continuity of $g_i$, the direction $w \in \mathbb{F}_j(x)$ also fulfills $w \in \mathbb{F}_j(y)$ for $y$ in a sufficiently small neighborhood $U$ of $x$. By Corollary 5, the directional derivative is bounded by $C > 0$ in the following way:

$$\nabla_z g_i(x, r_w Lw)^T Lw > C \text{ for all } w \in V.$$

This brings us in the situation of Lemma 27 and its proof applies. $\qquad\square$

With the result for the active components of $g$, the Clarke-subdifferential of the minimal radius function can be characterized.

**Definition 26** (Clarke-Subdifferential, Clarke 1990, p. 10)**.** Let the function $f : \mathbb{R}^n \to \mathbb{R}$ be locally Lipschitz continuous. The *Clarke-directional derivative* in direction $d \in \mathbb{R}^n$ is defined as

$$f^\circ(x; d) := \limsup_{\substack{y \to x \\ h \searrow 0}} \frac{f(y + hd) - f(y)}{h}.$$

The *Clarke-subdifferential* is defined as the set of linear functionals that are dominated by the Clarke-directional derivative for all directions, i.e.,

$$\partial f(x) := \{l \in \mathbb{R}^n \mid l^T d \leq f^\circ(x; d) \text{ for all } d \in \mathbb{R}^n\}.$$

We choose the Clarke-subdifferential, because it is convenient. However, our results can be adapted to other subdifferentials, for example by Demyanov, Rubinov and Mordukhovich (see Knossalla 2015 for an overview).

**Lemma 30** (Implicit Function for the Radius for Multiple Constraints)**.** *Let the point $(x, v)$ be in $\mathbb{R}^n \times \mathbb{S}^{m-1}$. Apply Corollary 6 for each component $g_i$ of $g$ with $i \in \mathcal{A}(x, v)$ to get the functions $r_i$ defined on the neighborhoods $U_i$ of $x$ and $V_i$ of $v$. Furthermore, define the intersections*

$$\widetilde{U} := \bigcap_{i=1,\dots,p} U_i \quad and \quad \widetilde{V} := \bigcap_{i=1,\dots,p} V_i.$$

*The following results hold: There exist neighborhoods $U \subset \widetilde{U}$ of $x$ and $V \subset \widetilde{V}$ of $v$ such that*

*a) the mapping $r : \widetilde{U} \times \widetilde{V} \to \mathbb{R}_{\geq 0}$ defined as*

$$r(y, w) := \min\{r_i(y, w) \mid i \in \mathcal{A}(x, v)\}$$

*fulfills the equivalence*

$$g^m(y, sLw) = 0 \iff s = r(y, w),$$

*for all $(y, w, s) \in U \times V \times \mathbb{R}_{\geq 0}$.*

*b) the partial Clarke-subdifferential of $r$ with respect to $x$ is given by*

$$\partial_x r(y, w) = \text{conv}\{\nabla_x r_i(y, w) : i \in \mathcal{I}^{x,v}(y, w)\}$$

*for $(y, w) \in U \times V$, with the index set*

$$\mathcal{I}^{x,v}(y, w) := \text{argmin}\{r_i(y, w) \mid i \in \mathcal{A}(x, v)\}.$$

*Proof.* The existence of the functions $r_i(y, w)$ for $i \in \mathcal{A}(x, v)$ is due to Corollary 6.

The forward direction of the implication

$$s = r(y, w) \implies g^m(y, sLw) = 0$$

can be shown as follows: If $j \in \mathcal{A}(x, v)$ is such that $r(y, w) = r_j(y, w)$, we have

$$0 = g_j(y, r_j(y, w)Lw) = g_j(y, r(y, w)Lw).$$

If $g^m(y, r(y, w)Lw) = 0$ holds, the implication is true. Suppose

$$g^m(y, r(y, w)Lw) > 0.$$

The continuity of the functions $g_i$ implies the continuity of the function $g^m$. Due to Assumption (B5), we have $g^m(y, 0) < 0$, for $y$ in a sufficiently small neighborhood $U$. By Bolzano's intermediate value theorem, there exists a radius

$$r^* \in (0, r(y, w)) \quad \text{such that} \quad g^m(y, r^*Lw) = 0$$

holds. Let $k \in \mathcal{A}(x, v)$, $k \neq j$ be an index that fulfills the equation

$$0 = g^m(y, r^*Lw) = g_k(y, r^*Lw).$$

By Corollary 6 this is equivalent to $r^* = r_k(y, w)$. However, as $r^*$ lies in $(0, r(y, w))$, we obtain

$$\min\{r_i(y, w) \mid i \in \mathcal{A}(x, v)\} = r(y, w) > r^* = r_k(y, w),$$

which contradicts the minimality in the choice of $r(y, w)$.

For the converse implication

$$g^{\mathrm{m}}(y, sLw) = 0 \implies s = r(y, w),$$

we start with $(y, w, s) \in U \times V \times \mathbb{R}_{\geq 0}$ such that $g^{\mathrm{m}}(y, sLw) = 0$. For an index $j$ that fulfills

$$g^{\mathrm{m}}(y, sLw) = g_j(y, sLw) = 0,$$

we have by Corollary 6 that $s = r_j(y, sLw)$. Suppose that $r_j$ is not minimal, i.e.,

$$r(y, w) < r_j(y, w)$$

and that the minimal radius is attained for the index $k \neq j$, i.e.,

$$g^{\mathrm{m}}(y, r(y, w)Lw) = g_k(y, r(y, w)Lw) = 0.$$

The convexity of $M(x)$ from Assumption (B1) implies that each convex combination $r(y, w) + \lambda(s - r(y, w))$ for $\lambda \in [0, 1]$ is feasible. Hence,

$$0 \geq g^{\mathrm{m}}(y, r(y, w) + \lambda(s - r(y, w))Lw)$$

(by the definition $g^{\mathrm{m}}$)

$$\geq g_k(y, r(y, w) + \lambda(s - r(y, w))Lw)$$

(by the Taylor theorem)

$$= g_k(y, r(y, w)) + \lambda(s - r(y, w))\nabla_z g_k(y, \theta Lw)^T Lw,$$

(since $g_k(y, r(y, w)) = 0$ holds)

$$= \lambda(s - r(y, w))\nabla_z g_k(y, \theta Lw)^T Lw,$$

where $\theta \in [r(y, w), \lambda(s - r(y, w))]$. Now we are in the situation of the proof of Lemma 27 and can follow it from Equation (2.14) onwards under the use of Corollary 5 to obtain a contradiction. Therefore, we have $s = r(y, w)$, which shows a).

For the statement in b), we use Corollary 6b), which states that for an index $i \in \mathcal{A}(x, v)$ and $(y, w) \in U \times V$, the gradient formula

$$\nabla_x r_i(y, w) = -\frac{1}{\nabla_z g_i(y, r(y, w)Lw)^T Lw} \nabla_x g_i(y, r(y, w)Lw)$$

holds. Since the functions $r_i(y, w)$, $i \in \mathcal{A}(x, v)$ are continuously differentiable, they are regular[1] in the sense of Clarke. Therefore by Clarke 1990, Proposition 2.3.12, we have

$$\partial_x r(y, w) = \partial_x(\min\{r_i(y, w) \mid i \in \mathcal{A}(x, v)\})$$
$$= \text{conv}\{\nabla_x r_i(y, w) \mid i \in \mathcal{I}^{x,v}(y, w)\},$$

where $I^{x,v}(y, w)$ is the set of indices that attain the minimal radius, i.e.,

$$\mathcal{I}^{x,v}(y, w) := \text{argmin}\{r_i(y, w) \mid i \in \mathcal{A}(x, v)\}.$$

This shows b) and concludes the proof. $\qquad\square$

The next steps are straightforward: We characterize the function $e$ in the case of multiple constraints and obtain its subdifferential as a consequence of the subdifferential representation of the radius in Lemma 30. This leads to a subdifferential characterization of the probability function $G$. As in the case of one constraint, an argument for exchanging the subdifferential with the integral is needed.

**Definition 27.** Define the function $e : \mathbb{R}^n \times \mathbb{S}^{m-1} \to \mathbb{R}$ as

$$e(x, v) := \chi\big(\{r \geq 0 : g^{\mathrm{m}}(x, rLv) \leq 0\}\big).$$

**Corollary 7** (Characterization of $e$ for Multiple Constraints)**.**
*There exist neighborhoods $U$ of $x$ and $V$ of $v$ such that*

$$e(y, w) = F_\chi(r(y, w)), \quad y \in U, \; w \in V \tag{2.18}$$

*where $F_\chi$ is the cumulative distribution function of the chi-distribution and $r$ is defined as in Lemma 30b).*

*Proof.* Follow the proof of Lemma 28. $\qquad\square$

**Corollary 8** (Partial subdifferential of $e$)**.**
*The partial Clarke subdifferential of the function $e : \mathbb{R}^n \times \mathbb{S}^{m-1} \to \mathbb{R}$ is given by*

$$\partial_x e(x, v) = f_\chi(r(x, v)) \, \text{conv}\{\nabla_x r_i(x, v) \mid i \in \mathcal{I}^{x,v}(x, v)\}.$$

---

[1]the directional derivative exists and is identical with the Clarke-directional derivative

*Proof.* We use the representation $e(x, v) = F_\chi\big(r(x, v)\big)$ from Corollary 7 and the chain rule for the Clarke-subdifferential Clarke 1990, Theorem 2.3.10 as well as $F'_\chi = f_\chi$. This yields

$$\partial_x e(x, v) = f_\chi(r(x, v)) \, \partial_x r(x, v)$$

and Lemma 30b) implies the result. $\qquad\qquad\square$

Similar to the gradient representation of Theorem 16 for one constraint, we obtain a subdifferential representation for multiple constraints.

**Theorem 17** (Subdifferential of $G$).
*Let Assumptions (A1)–(A5) hold in $x \in \mathbb{R}^n$. Then, the Clarke-subdifferential of $G$ is nonempty for all $y$ in a neighborhood $U$ of $x$ and characterized by*

$$\partial G(y) \subset \int\limits_{v \in \mathbb{S}^{m-1}} f_\chi\big(r(y, v)\big) \operatorname{conv}\{\nabla_x r_i(y, v) \mid i \in \mathcal{I}^{y,v}(y, v)\} \, du(v),$$

*where*

$$\nabla_x r_i(y, v) = -\frac{1}{\nabla_z g_i(y, r(y, v)Lv)^T Lv} \, \nabla_x g_i(y, r(y, v)Lv),$$

*for $i \in \mathcal{I}^{y,v}(y, v)$.*

*Remark* 17. The integral is to be understood as the set of integrals over all measurable selections (see Kuratowski and Ryll-Nardzewski 1965) of

$$\partial_x e(y, v) = f_\chi\big(r(y, v)\big) \operatorname{conv}\{\nabla_x r_i(y, v) \mid i \in \mathcal{I}^{y,v}(y, v)\}.$$

This means to every subgradient $l \in \partial G(y)$ corresponds a measurable mapping

$$h_l : \mathbb{S}^{m-1} \to \mathbb{R}^n, \quad v \mapsto l_v$$

such that $l_v \in \partial_x e(y, v)$ holds $u$-almost everywhere, $h_l$ is in $L^1(\mathbb{S}^{m-1}, \mathbb{R}^n; u)$ and

$$\int_{v \in \mathbb{S}^{m-1}} \partial_x e(y, v) \, du(v) = \left\{ \int_{v \in \mathbb{S}^{m-1}} h_l(v) \, du(v) \mid l \in \partial G(y) \right\}.$$

*Proof of Theorem 17.* We have

$$\partial G(y) = \partial \int_{v \in \mathbb{S}^{m-1}} e(y, v) \, du(v). \tag{2.19}$$

To interchange the subdifferential with the integral, we use Clarke 1990, Theorem 2.7.2. Clearly, the space $\mathbb{R}^n$ is separable and by Corollary 7 there is a neighborhood $V$ of $v$ for each $v \in \mathbb{S}^{m-1}$ such that

$$e(y, w) = F_\chi(r(y, w)),$$

which, noting that $r(y, \cdot)$ is the minimum over continuous functions, is a continuous function and thus measurable. It remains to assert the Lipschitz continuity of $e(\cdot, v)$ on $U$. As its subdifferential is known for each $y \in U$, the Lipschitz constant can be estimated by the norm of the largest subgradient due to a generalized version of the mean value theorem, see Lebourg 1979, Theorem 1.7. By Corollary 8 each subgradient $l \in \partial_x e(y, v)$ can be written as

$$l = f_\chi(r(y,v)) \sum_{i \in \mathcal{I}^{x,v}(y,v)} \lambda_i \nabla_x r_i(y,v), \quad \lambda_i \in [0,1], \quad \sum_{i \in \mathcal{I}^{x,v}(y,v)} \lambda_i = 1. \quad (2.20)$$

Because the set $M(y)$ is bounded due to Assumption (B2), the function $r(y,v)$ is bounded as a consequence of Lemma 30a). The continuity of $f_\chi$ implies the boundedness of $f_\chi(r(y,v))$. Therefore, the boundedness of each gradient

$$\nabla_x r_i(y,v), \quad i \in \mathcal{I}^{y,v}(y,v)$$

is sufficient to prove the boundedness of $l$. Corollary 6 provides the representation

$$\nabla_x r_i(y,v) = -\frac{1}{\nabla_z g_i(y, r_i(y,v)Lv)^T Lv} \nabla_x g_i(y, r_i(y,v)Lv).$$

Clearly, $\nabla_x g_i(y, r_i(y,v)Lv)$ is bounded on $U \times \mathbb{S}^{m-1}$ since $g$ is continuously differentiable due to (B3) and $r_i(y,v)$ is bounded. We need to show that $\nabla_z g(x, r_i(y,v)Lv)^T Lv$ is bounded away from zero. By an analogous argumentation as in the proof of Lemma 25, it is sufficient to show that the sets of directions

$$\mathcal{V}_i(y) := \{v \in \mathbb{S}^{m-1} \mid i \in I^{y,v}(y,v)\},$$

such that the $i$-th radius function $r_i(y,w)$ stays minimal, are compact for $y \in U$. By Lemma 30a) and Corollary 6, we have

$$i \in I^{y,v}(y,v) \iff g^m(y, r_i(y,v)Lv) = 0.$$

The continuity of $g_i$ and $r_i(y, \cdot)$ imply that for a convergent sequence $\{w_k\} \subset \mathcal{V}_i(y)$ with limit $w$ the equality

$$g^m(y, r_i(y,w)Lw) = 0$$

holds true. This shows $r(y, w) = r_i(y, w)$ and hence $i \in I^{y,w}(y, w)$. Consequently $w \in \mathcal{V}_i(y)$ holds, which asserts that $\mathcal{V}_i(y)$ is closed. Together with its boundedness implied by $\|v\| = 1$, the compactness of $\mathcal{V}_i(y)$ follows. Hence, we have to consider the continuous functions

$$C_{ij}(y) := \min_{v \in \mathcal{V}_i(y)} \nabla_z g_j(y, r(y, v)Lv)^T Lv$$

for $i \in \{1, \ldots, p\}$ and $j \in I^{y,v}(y, v)$. The construction with the sets $\mathcal{V}_i(y)$ is needed, because the sets $M_i(y)$ do not need to be bounded, but the sets

$$\{rLv \mid v \in \mathcal{V}_i, r > 0\} \cap M_i(y)$$

are bounded. This allows to adapt the argument of Lemma 25. Thus, we know that

$$\min_{v \in \mathcal{V}_i(y)} \nabla_z g_j(y, r(y, v)Lv)^T Lv$$

is bounded away from zero for each admissible combination of $i$ and $j$ and hence

$$\inf_{y \in U} C_{ij}(y) > C > 0.$$

Consequently, the inequality $|l| \leq K$ for each $l \in \partial_x e(y, v)$ for a $K > 0$ follows from the representation (2.20) and subsequently $e(\,\cdot\,, v)$ is Lipschitz continuous in $U$. Finally, the exchange of subdifferential and integral in (2.19) is justified and we obtain

$$\partial G(y) \subset \int_{v \in \mathbb{S}^{m-1}} \partial_x e(y, v) \, \mathrm{d}\mu_\zeta(v).$$

Application of Corollary 8 yields the claim. $\qquad\square$

## 2.3.4. Differentiability of the Probability Function

To obtain differentiability of $G$, it is sufficient to show that the subdifferential of $G$ contains only one subgradient. This is the case if the right hand side of the formula of Theorem 17 contains only one element. Therefore a sufficient condition is given by the assumption that the set of vectors $v \in \mathbb{S}^{m-1}$ with more than one active constraint is of measure zero, because then the convex hull is taken over only one gradient. Let us formalize this in a Theorem.

**Theorem 18** (Sufficient Condition for Differentiability, Ackooij and Henrion 2017, Theorem 4.1)**.**
*Let Assumptions (B1)-(B5) hold in $x \in \mathbb{R}^n$. Furthermore, assume that*

$$u\left(\left\{v \in \mathbb{S}^{m-1} \mid |\mathcal{I}^{y,v}(y,v)| > 1\right\}\right) = 0,$$

*for all $y \in U$.*

*Then, the function $G$ is differentiable in $U$ and its gradient is given by*

$$\nabla G(y) = \int\limits_{\substack{v \in \mathbb{S}^{m-1}: \\ |\mathcal{I}^{y,v}(y,v)|=1}} f_\chi(r_{i(v)}(y,v)) \, \nabla_x r_{i(v)}(y,v) \, \mathrm{d}u(v),$$

*where*

$$\nabla_x r_i(y,v) = -\frac{1}{\nabla_z g_i(y, r_i(y,v)Lv)^T Lv} \, \nabla_x g_i(y, r_i(y,v)Lv)$$

*and $i(v)$ is the single index in $I^{y,v}(y,v)$.*

*Proof.* The Lipschitz continuity of $e(\,\cdot\,, v)$ was shown in the proof of Theorem 17. It ensures that the subdifferential of $G$ is nonempty. Because the set

$$\left\{v \in \mathbb{S}^{m-1} \mid |\mathcal{I}^{y,v}(y,v)| > 1\right\}$$

has measure zero, it can be excluded from integration. Therefore,

$$\partial G(y) \subset \int\limits_{v \in \mathbb{S}^{m-1}} f_\chi(r(y,v)) \operatorname{conv}\{\nabla_x r_i(y,v) \mid i \in \mathcal{I}^{y,v}(y,v)\} \, \mathrm{d}u(v)$$

$$= \int\limits_{\substack{v \in \mathbb{S}^{m-1}: \\ |\mathcal{I}^{y,v}(y,v)|=1}} f_\chi(r_{i(v)}(y,v))\{\nabla_x r_{i(v)}(y,v)\} \, \mathrm{d}u(v).$$

The facts that the left hand side is nonempty and the right hand side contains only one element lead us to the conclusion that equality holds and $G$ is indeed differentiable. $\square$

It is not easy to verify the zero measure condition directly. The *rank 2 constraint qualification* (R2CQ) for the gradients of $g$ is easier to verify and implies the zero measure condition.

**Definition 28** (Rank 2 Constraint Qualification (R2CQ), Ackooij and Henrion 2017, Section 4). The functions $g_i$ satisfy the *rank 2 constraint qualification* (R2CQ) if

$$\text{rank}\{\nabla_z g_i(x,z), \nabla_z g_j(x,z)\} = 2,$$

for all pairs $i, j \in \mathcal{J}(x,z)$ with $i \neq j$, and all $z \in M(x)$, with the index set $\mathcal{J}(x,z)$ given by

$$\mathcal{J}(x,z) := \{i \in \{1, \ldots, p\} \mid g_i(x,z) = 0\}.$$

**Lemma 31** (Zero Measure Condition under (R2CQ), Ackooij and Henrion 2017, Lemma 4.3).
*If $g$ satisfies the (R2CQ) at the point $x$, then*

$$u\big\{v \in \mathbb{S}^{m-1} \mid |\mathcal{I}^{x,v}(x,v)| > 1\big\} = 0.$$

*Proof.* See the proof of Ackooij and Henrion 2017, Lemma 4.3. □

Therefore Theorem 18 can be formulated under (R2CQ).

**Corollary 9** (Theorem 18 under (R2CQ), Ackooij and Henrion 2017, Corollary 4.4).
*Theorem 18 holds if the zero measure condition*

$$u\big\{v \in \mathbb{S}^{m-1} \mid |\mathcal{I}^{y,v}(y,v)| > 1\big\} = 0$$

*is replaced by (R2CQ). Additionally, the gradient $\nabla_x G(y)$ is continuous in the neighborhood $U$ of $x$.*

## 2.4. Chance Constrained Optimization on Gas Networks

In this section, we turn our attention to gas networks. The goal is to control the pressure in the entry node of a tree network such that the pressure values in the nodes stay inside prescribed bounds with a certain probability. The demands at the exit nodes are uncertain and given by a multivariate normal distributed random variable.

To apply the results of Section 2.3.3 and Section 2.3.4, we formulate the feasible set for gas network optimization on a tree, simplify the system such that it can be written with only inequality constraints of the form $g(x,z) \leq 0$ and show that the necessary assumptions (B1)–(B5) are fulfilled.

We consider a tree network described by the graph $\mathcal{T} = (\mathcal{V}, \mathcal{A})$, where $\mathcal{V}$ is the set of nodes and $\mathcal{A}$ is the set of arcs. The network has one entry, which is the root node $r$ of the tree. The arcs point away from the root node. For a given node $v \in \mathcal{V}$ the set of leaf nodes that are reachable by the directed arcs from $v$ is denoted by

$$\mathcal{L}(v) := \{u \in \mathcal{V} \mid |\delta(u)| = 1 \text{ and there is a directed path from } v \text{ to } u\}.$$

We denote $|\mathcal{L}(r)| = m$. The pressure function on each edge $a \in \mathcal{A}$ as defined in (Real) on page 58 is denoted by $f_a$. The node-arc incidence matrix is denoted by $A$.

Coupling the flow values by the Kirchhoff conditions and the pressure values by the continuity condition yields the system

$$\begin{aligned}
Aq &= q^{\text{out}}, \\
p_v &= f_a(L_a; p_u, q_a), \quad \forall\, a = (u, v) \in \mathcal{A}, \\
p^{\text{lb}} &\leq p,
\end{aligned} \qquad \text{(Stat-Net-Full)}$$

for $p \in \mathbb{R}^{|\mathcal{V}|}$ and $q \in \mathbb{R}^{|\mathcal{A}|}$, where $p^{\text{lb}} \in \mathbb{R}_{>0}^{|\mathcal{V}|}$ is the vector of pressure bounds for $p$. In the terms of Section 1.3 this is the system

$$\text{(Kirchhoff)}, \ \text{(Arc Coupling)}, \ \text{(Node Coupling)}$$

with pressure continuity as node coupling condition. Therefore, it is sufficient to assign one variable

$$p_v = p_a(x_a(v)) \quad \text{for all arcs } a \in \delta(v)$$

to the pressure in each node. However, by Theorem 3, we know that the solution to the first two Equations of (Stat-Net-Full) is in fact completely determined if the pressure in the root node is prescribed. Therefore, we can reduce the system: First, we eliminate the equality constraints. Since the graph is a tree, the system $Aq = q^{\text{out}}$ has a unique solution

$$q = A^+ q^{\text{out}},$$

for balanced outflows, that is $\mathbf{1}^T q^{\text{out}} = 0$. The matrix $A^+$ is the pseudoinverse of $A$. Since $|\mathcal{A}| = |\mathcal{V}| - 1$ and $A \in \mathbb{R}^{m \times |\mathcal{A}|}$, it is given by

$$A^+ = (A^T A)^{-1} A^T.$$

The demand vector $z$ describes the outflows in the leave nodes $\mathcal{L}(r)$. For inner nodes the outflow is zero and the inflow in the rootnode balances the outflows, i.e.,

$$\mathbf{1}^T q^{\text{out}} = q_r^{\text{out}} + \mathbf{1}^T z = 0,$$

where

$$q^{\text{out}} = \begin{pmatrix} q_r^{\text{out}} \\ z \end{pmatrix} \quad q_r^{\text{out}} \in \mathbb{R}_-.$$

Once the flow values on every edge are known, the system can be parameterized by one pressure value $p_r$ in the root node. Furthermore, for a node $v$, there exists a unique path from $r$ to $v$. Denote the set of arcs on this path by $[r, v]$. The pressure in $v$ can be calculated by the function $p_v(p_r, z)$ that is the composition of the functions $f_a(L_a; \cdot, q_a)$ along the path $[r, v]$. The flow $q_a$ is only depending on $q^{\text{out}}$. For example, in the case $[r, v] = \{(r, u), (u, v)\} = \{a, b\}$, we obtain

$$p_v(p_r, z) = f_b(L_b; f_a(L_a; p_r, q_a), q_b),$$

where $q = A^+ q^{\text{out}}$.

**Definition 29** (Pressure Function in the Node). Formally, on the path $[r, v]$, we have the recursion

$$p_v(p_r, z) := f_a(L_a; p_u(p_r, z), q_a),$$

for the arc $a \in \delta_{\text{in}}(v) \cap [r, v]$. The pressure $p_u$ is calculated along the path $[r, u] \subsetneq [r, v]$. The recursion is well defined, because the pressure at the node $r$ is given as a argument of the function and the flow $q = A^+ q^{\text{out}}$ is uniquely defined by the argument $z$.

The system (Stat-Net-Full) reduces to

$$p_v^{\text{lb}} \leq p_v(p_r, z) \quad \forall v \in \mathcal{V}. \tag{Stat-Net}$$

We are interested in the set of admissible outflows for a given $p_r > 0$.

**Definition 30** (Set of Feasible Outflows). Define the *set of feasibly outflows* by

$$Q(p_r) := \left\{ z \in \mathbb{R}^m \mid p_v^{\text{lb}} \leq p_v(p_r, z) \text{ for all } v \in \mathcal{V} \setminus \{r\}, \ z \geq q^{\text{lb}} \right\} \tag{2.21}$$

and the *set of feasible outflows for the $v$-th constraints* by

$$Q_v(p_r) := \left\{ z \in \mathbb{R}_{\geq 0}^m \mid p_v^{\text{lb}} \leq p_v(p_r, z) \right\}.$$

The flow bound $q^{\text{lb}} \in \mathbb{R}^m_{>0}$ pays respect to the assumption that every exit node has a minimal positive demand. We assume $z \sim \mathcal{N}(\mu, \Sigma)$ to be normally distributed with expected value $\mu \geq q^{\text{lb}}$ and a positive definite covariance matrix $\Sigma$. The joint probabilistic constraint with respect to system (Stat-Net) is given by

$$\mathbb{P}(z \in Q(p_r)) \geq p^{\text{bd}} \tag{2.22}$$

for a probability bound $p^{\text{bd}} \in (0,1)$, which is usually close to one. With the probability density function $f_\nu(z)$, this can be rewritten as

$$G(p_r) := \mathbb{P}(z \in Q(p_r)) = \int\limits_{Q(p_r)} f_\nu(z)\, \mathrm{d}z. \tag{2.23}$$

We will now show assumptions (B1)–(B5) for the functions describing (2.21). The describing functions are continuously differentiable. Therefore (B3) is fulfilled. For $p^{\text{bd}} > 0.5$, there is no need to discuss (B5) as a consequence of Remark 16. Hence we have to show (B1), (B2) and (B4).

We begin by showing the monotonicity of the functions $p_v$ with respect to a single outflow.

**Lemma 32** (Monotonicity of the functions $p_v$ in $z_w$)**.**
*For $z \in \mathbb{R}^m_{\geq 0}$, the functions $p_v(p_r, z)$ are strictly decreasing as a function of $z_w$ for each leaf node $w \in \mathcal{L}(r)$.*

*Proof.* On the arc $a = (u, v)$, the pressure function $f_a$ strictly decreases as a function of $q_a$ and strictly increases as a function of $p_u$ for $q_a > 0$. For $u = r$, which means we are on the first arc of the tree starting from the root, the conjecture is clear, because the flow $q_a$ can be calculated as

$$q_a = \sum\nolimits_{w \in \mathcal{L}(r)} z_w.$$

Hence also $p_v$ strictly decreases as a function of $q_a$. For the following vertices pursue inductively. $\qquad\square$

We assume the convexity of the sets of feasible flow. The assumption seems to hold in all practical cases; see Figure 2.8. A tree network with three arcs and two exits was used for the computation. The outflows were evaluated on a $500 \times 500$ grid between 0 and $400\,\text{kg}\,\text{s}^{-1}$. For the data see Table C.1 and Table C.2 in the appendix.

*Assumption* 2 (Convexity of $Q_v(p_r)$). Assume that the sets $Q_v(p_r)$ are convex.

(a) $p_r = 61 \cdot 10^5 \, \text{Pa}$

(b) $p_r = 62 \cdot 10^5 \, \text{Pa}$

(c) $p_r = 63 \cdot 10^5 \, \text{Pa}$

(d) $p_r = 64 \cdot 10^5 \, \text{Pa}$

Figure 2.8.: The sets of feasible outflows (in green) for a tree network with one entry and two exits for different pressure values at the root node.

**Lemma 33** (Convexity of $B_w(p_r)$).
*The sets*

$$B_w(p_r) := \{z \in \mathbb{R}^m \mid z_w \geq q_w^{lb}\}.$$

*are convex.*

*Proof.* The restriction in $B_w(p_r)$ is linear. This implies that $B_w(p_r)$ is a convex set. □

**Lemma 34** (Boundedness of $Q(p_r)$).
*The set $Q(p_r)$ is bounded.*

*Proof.* The pressure functions $f_a$ on each arc $a \in \mathcal{A}$ are strictly decreasing in the flow value $q_a$. Consider the first node $u$ connected to the root node $r$.

By the strict monotonicity of the function $f_a(p_r, \cdot)$, see Lemma 5, there is a unique flow[2] $\bar{q}_a$ such that

$$f_a(p_r, \bar{q}_a) = p_u^{\text{lb}}.$$

The flow on each arc $a = (u, v)$ is bounded from below by

$$q_a^{\text{lb}} = \sum_{v \in \mathcal{L}(u)} q_v^{\text{lb}}.$$

Therefore, the flow on each arc is bounded by

$$q_a^{\text{lb}} \leq q_a \leq \bar{q}_a.$$

For the next layer of arcs, note that the flow on each arc may be larger if the ingoing pressure $p_u$ is larger, because a higher pressure drop is allowed before attaining the lower pressure bound. To be precise: If $p_l < p_h$ and $f_a(p_l, q_l) = f_a(p_h, q_h)$ holds, then $q_l$ has to be greater than $q_h$. Consequently, choosing a lower ingoing pressure only relaxes the obtained upper bound for the flow. Set

$$p_u := f_a(p_r, q_a^{\text{lb}}) \geq f_a(p_r, q_a)$$

on the arc $a = (r, u)$ and consider the arc $b = (u, v)$. The strict monotonicity of $f_b$ with respect to the flow allows the calculation of a unique flow $\bar{q}_b$ such that

$$f_b(p_u, \bar{q}_b) = p_v^{\text{lb}}.$$

Proceeding inductively gives flow bounds $q_a^{\text{lb}}$ and $\bar{q}_a$ for every arc $a \in \mathcal{A}$, which are written in the vector $\bar{q} \in \mathbb{R}^{|\mathcal{A}|}$. Hence, we have

$$q^{\text{lb}} \leq q \leq \bar{q}.$$

Multiplying by the node arc incidence matrix $A$ yields that the vector of outflows $Aq = q^{\text{out}}$ is bounded. Consequently, the set $Q(p_r)$ is bounded. $\square$

It remains to show assumption (B4). We define the sets that correspond to the sets $N_i(x)$ of Definition 23.

**Definition 31.** Define the sets

$$\text{NQ}_v(p_r) := \{z \in \mathbb{R}^m \mid p_v^{\text{lb}} = p_v(p_r, z)\}$$
$$\text{NB}_w(p_r) := \{z \in \mathbb{R}^m \mid z_w = q_w^{\text{lb}}\}.$$

---

[2]under the reasonable assumption that $p_u^{\text{lb}}$ lies in the image of $f_a(p_r, \cdot)$, which is the case for subsonic flow, i.e., if $p_u^{\text{lb}}$ is not too small

**Lemma 35** (Nonzero Flow Gradient at the Boundaries of $Q_v(p_r), B_w(p_r)$)**.**
*The flow gradient with respect to $z$ of the single restriction in $Q_v(p_r)$ or*
*$B_w(p_r)$ respectively is nonzero on the sets $\mathrm{NQ}_v(p_r)$ and $\mathrm{NB}_w(p_r)$.*

*Proof.* It is easy to see that the gradient of the function

$$g_w : z \mapsto z_w - q_w^{\mathrm{lb}}$$

is the unit vector $e_w$, which is one at the component $w$. Therefore the
gradient is nonzero on $\mathrm{NB}_w(p_r)$. Let us now discuss the mappings

$$h_v : z \mapsto p_v^{\mathrm{lb}} - p_v(p_r, z).$$

By Lemma 32 the functions $p_v$ are strictly decreasing in $z_w$ for $w \in \mathcal{L}(v)$.
This implies $\partial_{z_w} h_v(z) > 0$ and therefore $\nabla_z h_v(z) \neq 0$. $\qquad\square$

**Corollary 10** (Subdifferential of $G$ for the Gas Network Problem)**.**
*Let Assumption 2 hold. Then, the Clarke subdifferential of $G$ fulfills*

$$\partial G(p_r) \subset \int_{d \in \mathbb{S}^{m-1}} f_\chi(r(p_r, d)) \operatorname{conv} \left\{ \mathbf{g}_v(p_r, d), \mathbf{g}_w(p_r, d) \mid v, w \in \mathcal{I}(p_r, d) \right\} \mathrm{d}u(d),$$

*where the derivatives $\mathbf{g}_v$ and $\mathbf{g}_w$ are given by*

$$\mathbf{g}_v(p_r, d) = \frac{\partial_{p_r} p_v(p_r, r_v(p_r, d)Ld)}{\nabla_z p_v(p_r, r_v(p_r, d)Ld)^T Ld}, \qquad \text{for } v \in \mathcal{V} \setminus \{r\}$$

$$\text{and} \quad \mathbf{g}_w(p_r, d) = 0, \qquad\qquad\qquad \text{for } w \in \mathcal{L}(r).$$

*The index set $\mathcal{I}(p_r, d)$ is defined as*

$$\mathcal{I}(p_r, d) := \operatorname{argmin}\{r_i(p_r, d) \mid i \in \mathcal{A}(p_r, d)\},$$

*compare Lemma 30.*

*Proof.* Assumption 2 and Lemmas 33–35 show that the assumptions re-
quired for Theorem 17 are fulfilled. $\qquad\square$

Under (R2CQ) (see Definition 28) the function $G$ is differentiable. This
allows the application of gradient based optimization algorithms for the
numerical treatment.

**Corollary 11** (Derivative of $G$ for the gas network problem)**.**
*Let assumption 2 be fulfilled and assume that (R2CQ) holds. Then, $G$ is continuously differentiable and its derivative is given by*

$$\partial_{p_r} G(p_r) = \int\limits_{\substack{d \in \mathbb{S}^{m-1} \\ |\mathcal{I}(p_r,d)|=1}} f_\chi(r(p_r,d)) \, \mathbf{g}(p_r,d) \, \mathrm{d}u(d),$$

*where*

$$\mathbf{g}(p_r,d) = \begin{cases} \mathbf{g}_v(p_r,d) & \text{if } \mathcal{I}(p_r,d) = \{v\}, \\ \mathbf{g}_w(p_r,d) & \text{if } \mathcal{I}(p_r,d) = \{w\}. \end{cases}$$

*Proof.* This is a direct consequence of Corollary 10 and Corollary 9. $\qquad\square$

## 2.5. Numerical Results

### 2.5.1. Algorithmic Treatment

The evaluation of $G$ is extremely time consuming for higher sample numbers. In the case of gas network optimization each function evaluation of $G$ might mean simulating the whole network 10000 times. The high computational cost of the evaluation of the chance constraint is a challenge of the chance constrained optimization problems of the form

$$\begin{aligned} \min_{x \in \mathbb{R}^n} \quad & f(x) \\ \text{s.t.} \quad & \mathbb{P}(g(x,\xi) \le 0) \ge p^{\mathrm{bd}}. \end{aligned}$$

This makes an approach that uses fewer number of samples to generate a warmstart for the highest sample number attractive. The pseudo code is shown in Algorithm 1. Like in previous sections, we denote

$$G(x) := \mathbb{P}(g(x,\xi) \le 0),$$

the function $f_\chi$ is the density of the chi distribution and $e$ is defined as in Definition 27. Moreover, we define the (approximated) Lagrangian function using the definition of $G_l$ in Algorithm 1.

**Definition 32** (Lagrangian Function)**.** The *Lagrangian function*

$$L : \mathbb{R}^n \times \mathbb{R} \to \mathbb{R}$$

is defined as

$$L(x, \lambda) := f(x) + \lambda(p^{\mathrm{bd}} - G(x))$$

and its *approximation on level l* is given by

$$L_l(x, \lambda) := f(x) + \lambda(p^{\mathrm{bd}} - G_l(x)).$$

The samples $v_i$ on the sphere that are required for Algorithm 1 are obtained sampling random or quasirandom uniformly distributed points in the unit cube and using the inversion method to obtain a standard normal distributed sequence. Normalizing the vector of standard normal distributed values yields a vector on the unit sphere. The samples are uniform distributed. See Weisstein 2018 for the method. The procedure is shown in Algorithm 2.

*Remark* 18. It is essential that the approximations $G_l$ and $\nabla G_l$ in Algorithm 1 are continuous in $x$ even for low sample numbers. Compare them to the crude quasi Monte Carlo approximation

$$H_l(x) := \frac{1}{s_l} \sum_{i=1}^{s_l} \mathbb{1}_{\{g(x, z_i) \leq 0\}}(x),$$

which is discontinuous in $x$ as a sum over characteristic functions and needs high sample numbers to be "approximately" continuous.

For Algorithm 2, we need the definition of the error function.

**Definition 33** (Error Function)**.** The *error function* $\mathrm{erf} : \mathbb{R} \to (-1, 1)$ is defined as

$$\mathrm{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x \exp(-t^2)\,\mathrm{d}t$$

and its inverse is notated by $\mathrm{erf}^{-1}$. We also use the vector notation

$$\mathrm{erf}(x) = \begin{pmatrix} \mathrm{erf}(x_1) \\ \vdots \\ \mathrm{erf}(x_m), \end{pmatrix} \quad \text{and} \quad \mathrm{erf}^{-1}(x) = \begin{pmatrix} \mathrm{erf}^{-1}(x_1) \\ \vdots \\ \mathrm{erf}^{-1}(x_m), \end{pmatrix}$$

for $x \in \mathbb{R}^m$ (or $x \in (-1, 1)^m$ for the inverse).

*Remark* 19 (Wrong Approach to Sphere Sampling)*.* The approach to sample two angles $\theta \in [0, 2\pi]$, and $\phi \in [0, \pi]$ uniformly and set

$$x = \begin{pmatrix} \cos(\theta)\sin(\phi) \\ \sin(\theta)\sin(\phi) \\ \cos(\phi) \end{pmatrix}$$

---

**Algorithm 1** Successive Subsampling Solver (SSS)

---

**Require:** A vector $s$ of $\mathcal{L}$ refinement levels, all $s_{\mathcal{L}}$ uniform distributed samples $\{v_i\}_{1,\dots,s_{\mathcal{L}}}$ on the sphere for the finest level, the constraint function $g$ used in the joint probabilistic constraint, an objective function $f$, a probability bound $p^{\mathrm{bd}}$, a starting point $x_0$
  For $l = 1, \dots, \mathcal{L}$, define

$$G_l(x) := \frac{1}{s_l} \sum_{i=1}^{s_l} e(x, v_i) \quad \text{and}$$

$$\nabla G_l(x) := \frac{1}{s_l} \sum_{i=1}^{s_l} \frac{f_\chi(r(x, v_i))}{\nabla_z g_j(x, r(x, v_i) L v_i)^T L v_i} \nabla_x g_j(x, r(x, v_i) L v_i),$$

  for $\{j\} = \mathcal{I}(x, v_i)$.
  Set $\lambda = 1$.
  **for** $l = 1, \dots, \mathcal{L}$ **do**
    **if** $l < \mathcal{L}$ **then**
      Set the constraint tolerance $\tau_{\mathrm{con}}$ to $|G_{l+1}(x_0) - G_l(x_0)|$.
      Set the optimality tolerance $\tau_{\mathrm{opt}}$ to $|\lambda| \, \|\nabla G_{l+1}(x_0) - \nabla G_l(x_0)\|$.
    **else**
      Set tolerances reasonable fine.
    **end if**
    Solve the optimization problem

$$\min_{x \in \mathbb{R}^n} f(x), \quad \text{s.t. } G_l(x) \geq p^{\mathrm{bd}}$$

    with starting point $x_0$ to obtain the optimal solution $x^*$ and the corresponding Lagrange multiplier $\lambda^*$. The stopping criteria for the optimization is given by

$$p^{\mathrm{bd}} - G_l(x) \leq \tau_{\mathrm{con}} \quad \text{and} \quad \|\nabla_x L_l(x, \lambda)\| \leq \tau_{\mathrm{opt}}.$$

    Set $x_0 \leftarrow x^*$ and $\lambda \leftarrow \lambda^*$.
  **end for**
  **return** $x^*$

---

---

**Algorithm 2** Sphere Point Picking

---

**Require:** Number of desired samples $N$, dimension $m$

Generate uniform distributed samples $x^i$ for $i = 1, \ldots, N$ on $[0, 1]^m$.

Use the inversion method to obtain standard normal distributed points $y^i$ via

$$y^i = \sqrt{2} \mathrm{erf}^{-1}(2x^i - \mathbf{1}), \quad \text{for } i = 1, \ldots, N.$$

The normalized vectors

$$v^i = \frac{1}{\|y^i\|} y^i, \quad \text{for } i = 1, \ldots, N,$$

are uniformly distributed on the sphere $\mathbb{S}^{m-1}$.

**return** $\{v^i\}_{i=1,\ldots,N}$

---

to generate a uniform distribution on the sphere $\mathbb{S}^2$ is faulty as it leads to a accumulation of points around the poles. This is due to the nonlinear surface element $\mathrm{d}S = \sin(\phi) \, \mathrm{d}\theta \, \mathrm{d}\phi$. See Figure 2.9 for the comparison between this method and Algorithm 2.

It is not desirable to solve the optimization problem to small tolerances as long as the approximation of $G$ is very coarse. The optimization stops, when

$$p^{\mathrm{bd}} - G_l(x) \leq \tau_{\mathrm{con}} \quad \text{and} \quad \|\nabla_x L_l(x, \lambda)\| \leq \tau_{\mathrm{opt}}.$$

The difference

$$\tau_{\mathrm{con}} := |G_l(x_0) - G_{l+1}(x_0)|$$

provides an estimate of the error $|G_l(x) - G(x)|$. Similarly, the difference of the gradient of the Lagrangian on the coarse and the next finer level is given by

$$\|\nabla_x L_l(x_0, \lambda) - \nabla_x L_{l+1}(x_0, \lambda)\|$$
$$= \|\nabla f(x_0) - \lambda \nabla G_l(x_0) - \nabla f(x_0) + \lambda \nabla G_{l+1}(x_0)\|$$
$$= |\lambda| \|\nabla G_l(x_0) - \nabla G_{l+1}(x_0)\| =: \tau_{\mathrm{opt}},$$

where we can estimate the Lagrange multiplier by the Lagrange multiplier obtained in the last optimization with the constraint approximation $G_{l-1}$. It is useful to project the tolerances back to a box $[\tau_l, \tau_r]$ to avoid making not enough "cheap" steps while not demanding unachievable tolerances if $G_l$ and $G_{l+1}$ happen to (almost) coincide at the evaluation point $x_0$. In our
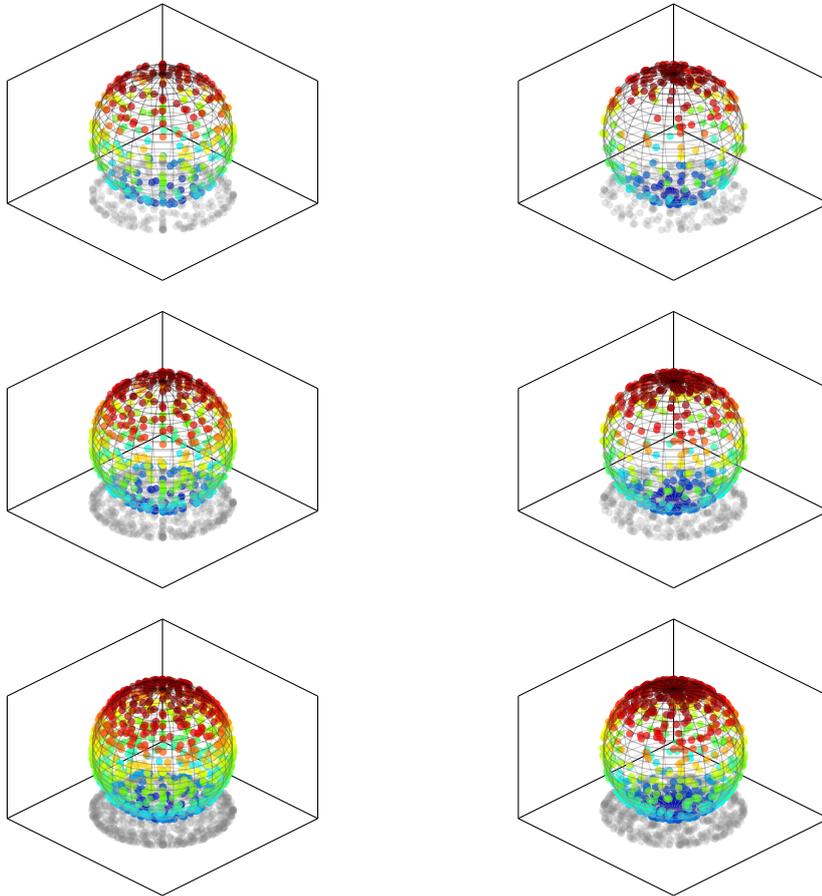
Figure 2.9.: Sampling of points on the sphere for $N = 300, 500$ and $800$ (top to bottom). The samples in the left column are generated Algorithm 2 based on Sobol points on $[0, 1]^3$. The right column shows the faulty approach of sampling two angles uniformly and using spherical coordinates. This leads to an accumulation of points around the poles of the sphere.

computations we did choose $\tau_l = 0.1$ and $\tau_r = 10^{-3}$. If the approximation on the fine level is close to $G$ and the pair $(x_0, \lambda)$ is in a neighborhood of the optimal solution $(x^*, \lambda^*)$ of the exact problem, then the error estimate provides a useful bound.

**Lemma 36** (Error Estimates)**.**
*Let $(x^*, \lambda^*)$ be a solution of the problem*

$$\min_{x \in \mathbb{R}^n} f(x), \quad s.t. \ G(x) \geq p^{bd}.$$

*Assume that*

$$\|x_0 - x^*\| \leq \epsilon_1, \quad \|\lambda - \lambda^*\| \leq \epsilon_2$$

*and*

$$|G_{l+1}(x^*) - G(x^*)| \leq \epsilon_3, \quad \|\nabla G_{l+1}(x^*) - \nabla G(x^*)\| \leq \epsilon_4,$$

*for $\epsilon_i > 0$, $i = 1, \ldots, 4$. Furthermore, let $G$ be continuously differentiable on $B_{\epsilon_1}(x^*)$ and let the Lipschitz conditions (the conditions on the left hand side are no additional assumptions)*

$$
\begin{aligned}
|G(x) - G(y)| &\leq K\|x - y\| & \|\nabla G(x) - \nabla G(y)\| &\leq K\|x - y\| \\
|G_l(x) - G_l(y)| &\leq K\|x - y\| & \|\nabla G_l(x) - \nabla G_l(y)\| &\leq K\|x - y\| \\
|G_{l+1}(x) - G_{l+1}(y)| &\leq K\|x - y\| & \|\nabla G_{l+1}(x) - \nabla G_{l+1}(y)\| &\leq K\|x - y\|
\end{aligned}
$$

*hold for all $x, y \in B_{\epsilon_1}(x^*)$.*
  *Then, the following holds:*

*a)*
$$|G_l(x^*) - G(x^*)| \leq 2K\epsilon_1 + \epsilon_3 + \tau_{con},$$

*b)*
$$\|\nabla_x L_l(x^*, \lambda^*) - \nabla_x L(x^*, \lambda^*)\| \leq$$
$$(\epsilon_2 + |\lambda|)(2K\epsilon_1 + \epsilon_4 + \|\nabla G_l(x_0) - \nabla G_{l+1}(x_0)\|).$$

*For $\epsilon_i \searrow 0$, $i = 1, \ldots, 4$ the right hand side of a) converges to $\tau_{con}$ and the right hand side of b) converges to $\tau_{opt} = |\lambda| \|\nabla G_l(x_0) - \nabla G_{l+1}(x_0)\|$.*

*Proof.* The function values fulfill

$$|G_l(x^*) - G(x^*)| = |G_l(x^*) - G_l(x_0) + G_l(x_0) - G_{l+1}(x_0) + G_{l+1}(x_0)$$
$$- G_{l+1}(x^*) + G_{l+1}(x^*) - G(x^*)|$$

(by the triangle inequality)

$$\leq |G_l(x^*) - G_l(x_0)| + |G_l(x_0) - G_{l+1}(x_0)|$$
$$+ |G_{l+1}(x_0) - G_{l+1}(x^*)| + |G_{l+1}(x^*) - G(x^*)|$$

(by the Lipschitz continuity and $|G_{l+1}(x^*) - G(x^*)| \leq \epsilon_3$)

$$\leq K\|x^* - x_0\| + \tau_{\text{con}} + K\|x_0 - x^*\| + \epsilon_3$$
$$\leq 2K\epsilon_1 + \epsilon_3 + \tau_{\text{con}}.$$

For the gradients of the Lagrangians, we obtain

$$\|\nabla_x L_l(x^*,\lambda^*) - \nabla_x L(x^*, \lambda^*)\| = |\lambda^*|\|\nabla G_l(x^*) - \nabla G(x^*)\|$$
$$= |\lambda^*|\|\nabla G_l(x^*) - \nabla G_l(x_0) + \nabla G_l(x_0) - \nabla G_{l+1}(x_0)$$
$$+ \nabla G_{l+1}(x_0) - \nabla G_{l+1}(x^*) + \nabla G_{l+1}(x^*) - \nabla G(x^*)\|$$

(using the triangle inequality)

$$\leq |\lambda^*|(\|\nabla G_l(x^*) - \nabla G_l(x_0)\| + \|\nabla G_l(x_0) - \nabla G_{l+1}(x_0)\|$$
$$+ \|\nabla G_{l+1}(x_0) - \nabla G_{l+1}(x^*)\| + \|\nabla G_{l+1}(x^*) - \nabla G(x^*)\|)$$

(by the Lipschitz continuity and $\|\nabla G_{l+1}(x^*) - \nabla G(x^*)\| \leq \epsilon_4$)

$$\leq |\lambda^*|(K\|x^* - x_0\| + \|\nabla G_l(x_0) - \nabla G_{l+1}(x_0)\|$$
$$+ K\|x_0 - x^*\| + \epsilon_4)$$

(using the triangle inequality on $|\lambda^* - \lambda + \lambda|$)

$$\leq (|\lambda^* - \lambda| + |\lambda|)(2K\epsilon_1 + \epsilon_4 + \|\nabla G_l(x_0) - \nabla G_{l+1}(x_0)\|)$$
$$\leq (\epsilon_2 + |\lambda|)(2K\epsilon_1 + \epsilon_4 + \|\nabla G_l(x_0) - \nabla G_{l+1}(x_0)\|).$$

$\square$

## 2.5.2. Example 1: Control of the Pressure in the Root Node of a Tree

**Problem Description**

We consider a tree network with one entry and six exits as depicted in Figure 2.10. The outflows are normally distributed with expected value $\mu \in \mathbb{R}^6$ and positive definite covariance matrix $\Sigma \in \mathbb{R}^{6\times 6}$. Our objective is to find the minimal pressure such that the lower pressure bounds in all nodes are fulfilled with high probability. As a consequence of the discussion in Section 2.4, the pressure values in the network depend only on the

Figure 2.10.: A tree network with six exits (red) with uncertain outflows. The pressure at the green entry node can be controlled.

pressure in the root node and the outflows and the problem can be written as

$$\min_{p_r \in \mathbb{R}} \quad \tfrac{1}{2} p_r^2$$
$$\text{s.t.} \quad G(p_r) \geq p^{\text{bd}},$$
$$p_r \in [p_r^{\text{lb}}, p_r^{\text{ub}}]. \tag{Ex1}$$

The pressure bounds are set to $p_v^{\text{lb}} = 40 \cdot 10^5 \, \text{Pa}$ for all nodes $v \in \mathcal{V}$ and the upper pressure bound in the root node is set to $p_r^{\text{ub}} = 60 \cdot 10^5 \, \text{Pa}$. For the normal distribution the expected value

$$\mu = \begin{pmatrix} 100 & 110 & 120 & 130 & 140 & 150 \end{pmatrix}^T \rho_0 \cdot 1000 \, \text{m}^3 \, \text{h}^{-1},$$

where $\rho_0$ is the norm density, was chosen. The covariance matrix was set to

$$\Sigma = 0.9\mathbb{1} + 0.1\mathbb{E},$$

where $\mathbb{1}$ is the identity in $\mathbb{R}^{6 \times 6}$ and $\mathbb{E}$ the matrix of only ones. The probability bound was set to $p^{\text{bd}} = 0.75$. A similar problem is discussed in Gotzes et al. 2016. However, in the mentioned article the set of feasible realizations is defined as the outflows, such that there exists a pressure in the root node that keeps the system within the prescribed thresholds. In our case, we choose the pressure before the realization of the demands.

We provide results for both the common Weymouth model (see e.g. Koch et al. 2015, Lemma 2.2) and the more accurate but computationally more expensive stationary isothermal Euler equation with nonconstant compressibility factor as described by (Real). On the arc $a = (u, v)$, the pressure loss described by the Weymouth equation is given by

**Weymouth formula**

$$p_v = f_a(p_u, q_a) := \sqrt{p_u^2 - \Lambda q_a |q_a|} \tag{Wey}$$

with a constant $\Lambda > 0$.

**Data and Methods Used for the Numerical Computations**

We solve the optimization problem with a maximal sample number of $20\,000$ and the starting point $p_r^0 = 50 \cdot 10^5\,\mathrm{Pa}$. For the optimization problem on each level, Matlab's SQP solver (part of the fmincon routine) was used with the constraint tolerance $\tau_{\mathrm{con}}$, the optimality tolerance $\tau_{\mathrm{opt}}$ and the step size tolerance $10^{-8}$ (practically avoiding abortion due to low step sizes). An analytic gradient for the objective was provided, while for the chance constraint the gradient approximation described in Algorithm 1 was used with finite difference gradients for the functions $g_j$. In our example the gradient of $g_i$ is the gradient of the negative pressure functions $p_v$ in each node. The sampling of points on the sphere was implemented by Algorithm 2 with Sobol points (see Sobol' 1967) as initialization[3]. All computations were carried out on a desktop PC with a Intel(R) Core(TM) i7-2600 CPU @ 3.40GHz processor and 16 GB RAM.

**Comparison of the Two Models for Different Numbers of Refinement Levels**

The results in Table 2.1 show that by solving the optimization on the coarse intermediate levels one quickly approaches the optimal solution. In comparison the Weymouth model is much cheaper to evaluate; see Table 2.2. The optimal solution $p_r^* = 46.8513 \cdot 10^5\,\mathrm{Pa}$ is lower than for the model with the isothermal Euler equation. A lower optimal pressure in the root node means that it is easier to fulfill the lower pressure bounds and therefore shows a tendency of the Weymouth formula to overestimate the pressure. The seemingly low difference in time between the first and the second level can be explained by the increasing benefits of the parallelization for the evaluation of the chance constraint. Table 2.3 shows the results for a direct computation on the finest level of 20000 samples with starting point $p_r^* = 50 \cdot 10^5\,\mathrm{Pa}$. The calculation time needed is significantly higher compared to Algorithm 1. The results in Table 2.4 and 2.5, where only two levels were used, show that the run times are even lower than with five levels. This indicates that the intermediate levels may be to time consuming.

---

[3]The Matlab implementation found in Burkardt and Fox 2010 was used

Table 2.1.: Results for the (Real)-Model with 5-Level-(SSS)

| $l$ | $s_l$ | SQP iterations | $p_r^*$ in $10^5$ Pa | Time in s |
|---|---|---|---|---|
| 1 | 50 | 20 | 46.8870 | 38 |
| 2 | 250 | 6 | 46.8854 | 73 |
| 3 | 1000 | 5 | 46.8900 | 282 |
| 4 | 5000 | 4 | 46.8900 | 1026 |
| 5 | 20000 | 6 | 46.8894 | 3767 |
| Total Time | | | | 5186 |

Table 2.2.: Results for the (Wey)-Model with 5-Level-(SSS)

| $l$ | $s_l$ | SQP iterations | $p_r^*$ in $10^5$ Pa | Time in $s$ |
|---|---|---|---|---|
| 1 | 50 | 18 | 46.8489 | 11 |
| 2 | 250 | 6 | 46.8473 | 15 |
| 3 | 1000 | 5 | 46.8519 | 49 |
| 4 | 5000 | 4 | 46.8519 | 201 |
| 5 | 20000 | 6 | 46.8513 | 734 |
| Total Time | | | | 1010 |

Table 2.3.: Results without (SSS)

| | $s_{\mathcal{L}}$ | SQP steps | $p_r^*$ in $10^5$ Pa | Time in s |
|---|---|---|---|---|
| Weymouth | 20000 | 16 | 46.8521 | 1660 |
| isothermal Euler | 20000 | 20 | 46.8900 | 7399 |

Table 2.4.: Results for the (Real)-Model with 2-Level-(SSS)

| $l$ | $s_l$ | SQP iterations | $p_r^*$ in $10^5$ Pa | Time in s |
|---|---|---|---|---|
| 1 | 50 | 20 | 46.8870 | 140 |
| 2 | 20000 | 7 | 46.8893 | 4013 |
| Total Time | | | | 4153 |

Table 2.5.: Results for the (Wey)-Model with 2-Level-(SSS)

| $l$ | $s_l$ | SQP iterations | $p_r^*$ in $10^5$ Pa | Time in s |
|---|---|---|---|---|
| 1 | 50 | 18 | 46.8489 | 33 |
| 2 | 20000 | 7 | 46.8512 | 835 |
| Total Time | | | | 868 |

Table 2.6.: Results for the (Real)-model with 2-Level-(SSS) and Weymouth warmstart

| Model | $l$ | $s_l$ | SQP iterations | $p_r^*$ in $10^5$ Pa | Time in $s$ |
|---|---|---|---|---|---|
| Weymouth | 1 | 50 | 18 | 46.8489 | 12 |
| isothermal Euler | 1 | 50 | 8 | 46.8863 | 293 |
| | 2 | 20000 | 7 | 46.8893 | 4034 |
| Total Time | | | | | 4339 |

## Using Weymouth as a Warmstart

Different model hierarchies can be coupled in the (SSS) approach. The cheap Weymouth model with a low number of samples can be used to generate a warmstart for the accurate but expensive isothermal Euler equations. The results displayed in Table 2.6 are not very encouraging. The difference in time needed on the coarse level does not make up for the worse warmstart, when compared to a coarse approximation with the more accurate (Real) model in Table 2.4.

## Comparison to Optimization without Uncertainty

Consider the optimization problem with the (Real)-model, where the outflow is not stochastic, but simply replaced by the expected demand. For the remaining data, the values of Section 2.5.2 were used. The resulting problem reads (again with the constant arguments of $p_v$ omitted).

$$\min_{p_r \in \mathbb{R}} \quad \frac{1}{2} p_r^2$$
$$\text{s.t.} \quad p_v(p_r, \mu) \geq p_v^{\mathrm{lb}}, \quad \forall v \in \mathcal{V} \setminus \{r\}, \qquad \text{(Ex1-Det)}$$
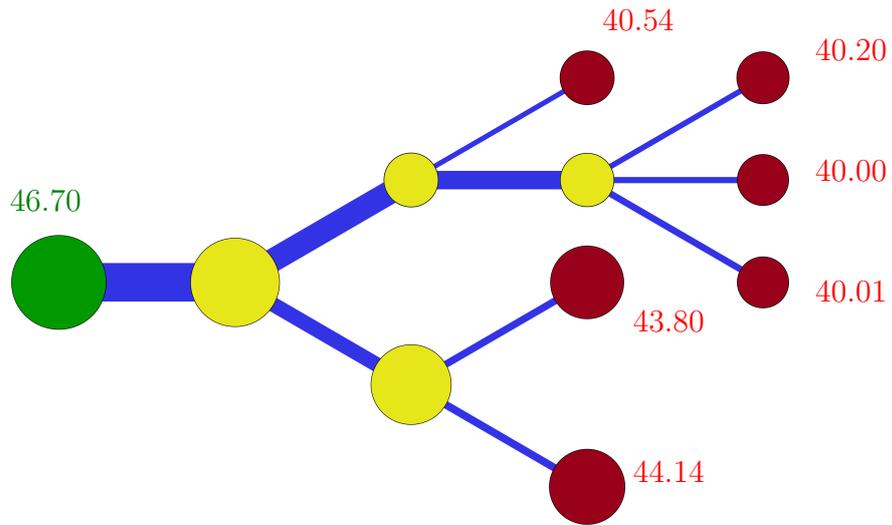$$p_r \in [p_r^{\mathrm{lb}}, p_r^{\mathrm{ub}}].$$

149

Figure 2.11.: Network simulation for the optimal pressure in the deterministic problem for the expected demand. The size of each circle corresponds to the pressure in the node and the width of each edge to the flow on the edge. The pressure values in $10^5$ Pa are shown at the source and sink nodes only.

Optimizing yields an optimal pressure of $p^*_{\text{det}} = 46.6964 \cdot 10^5$ Pa, which is clearly lower than the optimal pressures obtained from the computations including uncertainty. The network simulation for the optimal pressure is shown in Figure 2.11. The lower pressure bound of $40 \cdot 10^5$ Pa is exactly met at one node. The optimal pressure of $p_{\text{stoch}} = 46.8893 \cdot 10^5$ Pa for the stochastic problem calculated with two levels (see Table 2.4) was used for the network simulation in Figure 2.12. Here the lowest sink node still has some slack to the lower pressure bound. The approximate probability on the finest level is

$$G_{\mathcal{L}}(p^*_{\text{det}}) = 0.4837 \quad \text{and} \quad G_{\mathcal{L}}(p^*_{\text{stoch}}) = 0.7500,$$

i.e. the chance constraint is active in the solution of the stochastic optimization problem. However, the solution of (Ex1-Det) is feasible with less than 50% probability in the stochastic formulation.

## 2.5.3. Example 2: Controlling Multiple Compressors

We consider a problem similar to the one above, but with two additional controls. Two arcs in the tree graph correspond two compressor stations

Figure 2.12.: Network simulation with the expected demand for the optimal pressure in the stochastic problem. The size of each circle corresponds to the pressure in the node and the width of each edge to the flow on the edge. The pressure values in $10^5$ Pa are shown at the source and sink nodes only.

with controllable compressionrates $\gamma \in [1, 2]^2$. To obtain reasonable scaling of the problem, we use the variable $\beta := p_r \cdot (40 \cdot 10^5 \, \text{Pa})^{-1}$ for the pressure in the root node. The node pressures in the nodes $u, v$ connected by the arc $a = (u, v)$ corresponding to the compressor $i$ satisfy

$$p_v(\beta, \gamma, z) = \gamma_i \, p_u(\beta, \gamma, z).$$

The goal is to find the minimal compression rates $\gamma_i$ and a pressure

$$40 \cdot 10^5 \, \text{Pa} \cdot \beta = p_r \in [50 \cdot 10^5 \, \text{Pa}, 60 \cdot 10^5 \, \text{Pa}]$$

in the root node such that the pressures in the nodes stay inside the prescribed bounds $[40 \cdot 10^5 \, \text{Pa}, 60 \cdot 10^5 \, \text{Pa}]$ with high probability. Similarly to above, we define

$$G(\beta, \gamma) := \mathbb{P}\big(z \in \mathbb{R}^m_{\geq 0} \mid p_v(\beta, \gamma, z) \geq p_v^{\text{lb}} \quad \forall v \in \mathcal{V} \setminus \{r\}\big),$$

with the obvious adjustments to the functions $p_v$. It is advantageous to rather use the two compressor stations than demanding a higher pressure in the root node, which corresponds to a possibly complex control of a

Table 2.7.: Results for the (Real)-Model with 2-Level-(SSS)

| $l$ | $s_l$ | SQP iterations | $\beta^*$ | $\gamma_1^*$ | $\gamma_2^*$ | Time in s |
|---|---|---|---|---|---|---|
| 1 | 50 | 6 | 1.2500 | 1.0818 | 1.0757 | 844 |
| 2 | 5000 | 6 | 1.2500 | 1.0815 | 1.0725 | 11603 |
| Total Time | | | | | | 12447 |

connected network. Therefore the pressure in the root node is assigned a larger constant in the objective function. The resulting optimization problem reads

$$\min_{(\beta,\gamma)\in\mathbb{R}^3} \quad \tfrac{1}{2}\big(10\,\beta^2 + \gamma_1^2 + \gamma_2^2\big)$$
$$\text{s.t.} \quad G(\beta,\gamma) \geq p^{\mathrm{bd}},$$
$$\beta \in [1.25, 1.5],$$
$$\gamma \in [1, 2]^2. \tag{Ex2}$$

The value $p^{\mathrm{bd}} = 0.90$ was used. For the multivariate normally distributed outflow, the expected values

$$\big(100 \quad 125 \quad 500 \quad 150 \quad 175 \quad 200 \quad 225\big)^T \rho_0 \cdot 1000\,\mathrm{m}^3\,\mathrm{h}^{-1},$$

and the covariance matrix

$$\Sigma = 0.9\mathbb{1} + 0.1\mathbb{E}$$

were used. The results for two levels can be seen in Table 2.7. The relatively high computation times can be explained by two effects: The optimization method needs more function evaluations for each iteration. In regions of the optimization variable, where the expected demand is not feasible, i.e., Assumption (B5) is not fulfiled, a smoothened direct quasi Monte-Carlo method with higher sample numbers was used as a backup routine to evaluate the probability function. The state of the gas network for the optimal solution

$$\big(\beta^* \quad \gamma_1^* \quad \gamma_2^*\big) = \big(1.2500 \quad 1.0815 \quad 1.0725\big)$$

under the expected demand is depicted in Figure 2.13.

## Conclusion

We have shown that the gradient formula for the spherical radial decomposition of the chance constraint is applicable for gas network optimization
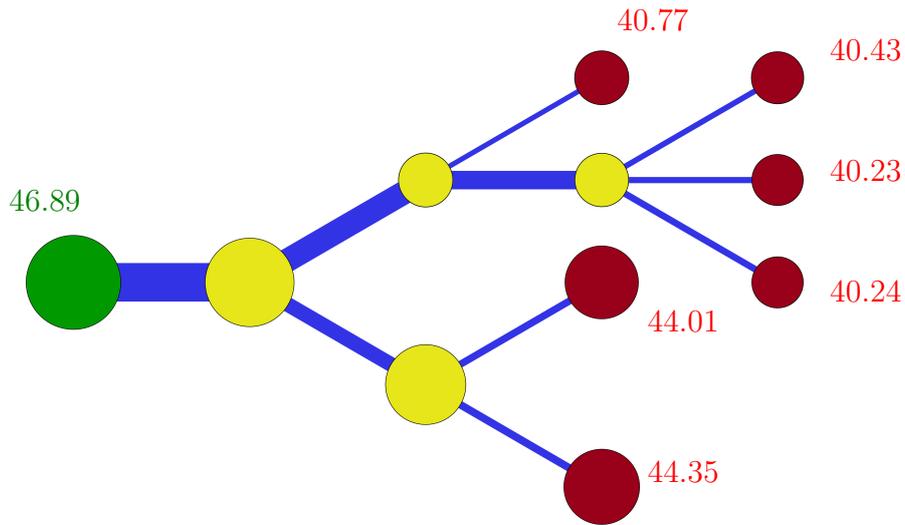
Figure 2.13.: Network simulation for the expected demand with optimal controls of the stochastic problem. The size of each circle corresponds to the pressure in the node and the width of each edge to the flow on the edge. The pressure values in $10^5\,\mathrm{Pa}$ are shown at the source and sink nodes only. The green edges correspond to compressors that raise the pressure by 7% and 8% respectively.

on trees. With a multilevel approach medium sized networks are solved in reasonable time with the quasilinear isothermal Euler equations with non-constant compressibility factor as a model for the gas flow. The numerical results show the advantage of optimization with chance constraints over the naive optimization using the expected values for the uncertain parameters.

How to extend the framework to the case, where the sets of feasible realizations are nonconvex, remains an open question. Furthermore, it is not clear if the theoretical results remain applicable for gas networks containing cycles. Approaches that do not depend on spherical radial decomposition include direct quasi-Monte Carlo methods with smoothening, kernel density approximation (see Caillau et al. 2017) or exact penalization (see Curtis, Wächter, and Zavala 2018).

# 3. Feedback Stabilization of the Linear Wave Equation under Uncertain Initial-Boundary Data

> *Have you considered that if you don't make waves, nobody including yourself will know that you are alive?*
>
> (Theodore Isaac Rubin)

Until now, we have discussed the case of the gas being in an equilibrium. The instationary case is more complex, which makes simplifications of the model desirable. For ideal gas, the pressure $p$ and the density $\rho$ of the gas obey the linear dependency $p = c^2\rho$, where $c > 0$ is the speed of sound in the gas. See also Table 1.1 for the notation. The system (Iso) from Chapter 1 can be written as

$$\begin{cases} \partial_t\rho + \partial_x q = 0, \\ \partial_t q + \partial_x\left(c^2\rho + \dfrac{q^2}{\rho}\right) = -\dfrac{\theta}{2}\dfrac{q|q|}{\rho}. \end{cases} \tag{Iso}$$

With the velocity $v = q/\rho$, we can reformulate the term in the space derivative of the momentum equation

$$c^2\rho + \frac{q^2}{\rho} = c^2\rho\left(1 + \frac{v^2}{c^2}\right),$$

see Domschke et al. 2017, Section 4.1. For small gas velocities $|v| \ll c$ this term is close to 1 and we arrive at the semilinear model on the acoustic timescale

$$\begin{cases} \partial_t\rho + \partial_x q = 0, \\ \partial_t q + c^2\partial_x\rho = -\dfrac{\theta}{2}\dfrac{q|q|}{\rho}. \end{cases} \tag{Iso-semi}$$

A rigorous derivation of this model by asymptotic analysis is given in Brouwer, Gasser, and Herty 2011, Section 3.2.2. We can show that the density $\rho$ and the mass flow $q$ obey a semilinear wave equation. Taking the time derivative of the first equation and the space derivative of the second equation yields

$$\partial_{tt}^2\rho + \partial_x\partial_t q = 0,$$

$$\partial_x\partial_t q + c^2\partial_{xx}^2\rho = -\frac{\theta}{2}\partial_x\left(\frac{q|q|}{\rho}\right).$$

Inserting the second equation into the first equation leads to

$$\partial_{tt}^2\rho = c^2\partial_{xx}^2\rho + \frac{\theta}{2}\partial_x\left(\frac{q|q|}{\rho}\right). \tag{3.1}$$

Taking the space derivative of the first equation and the time derivative of the second equation yields

$$\partial_t\partial_x\rho + \partial_{xx}^2 q = 0,$$

$$\partial_{tt}^2 q + c^2\partial_t\partial_x\rho = -\frac{\theta}{2}\partial_t\left(\frac{q|q|}{\rho}\right).$$

By inserting the first equation into the second equation, we arrive at the semilinear wave equation

$$\partial_{tt}^2 q = c^2\partial_{xx}^2 q - \frac{\theta}{2}\partial_t\left(\frac{q|q|}{\rho}\right). \tag{3.2}$$

Thus neglegting the friction in Equation (3.1) and Equation (3.2) leads to the linear wave equations

$$\partial_{tt}^2\rho = c^2\partial_{xx}^2\rho \tag{WaveEq-$\rho$}$$

and

$$\partial_{tt}^2 q = c^2\partial_{xx}^2 q. \tag{WaveEq-$q$}$$

**Literature Survey**

In Gugat and Leugering 2008, $L^\infty$-minimal control of the wave equation was considered. The article Gugat, Leugering, and Wang 2017 examined Neumann boundary feedback stabilization of a nonlinear wave equation using $H^2$-Lyapunov functions. Time optimal control of the wave equations has been considered in Kunisch and Wachsmuth 2013.

Stabilization of the wave equation on one dimensional networks was analyzed in Valein and Zuazua 2009. Dynamic control games on string networks governed by the wave equation were analyzed in Gugat and Steffensen 2017. Traveling wave solutions of the quasilinear isothermal Euler equations on networks and a corresponding optimal control problem were discussed in Gugat and Wintergerst 2018.

Stochastic optimal control for the semilinear isothermal Euler equations is discussed in Zavala 2014. The numerical results of different objective functions are compared. Chance constrained optimal control in aerospace via kernel density estimation is investigated in Caillau et al. 2017. The properties of chance constrained problems in infinite dimension were analyzed in Farshbaf-Shaker, Henrion, and Hömberg 2017.

## 3.1. Deterministic Wave Equation

The attentive reader may have guessed by now that in view of Equations (WaveEq-$\rho$), (WaveEq-$q$) and the quote at the beginning of this chapter, it makes sense to study the linear wave equation. Like in previous parts of this thesis, we want to include stochastic boundary data (and analogously stochastic initial data). We start, however, by discussing the deterministic case with fixed initial and boundary data. Consider the wave equation on the space interval $[0, L]$ and the time interval $[0, T]$ with prescribed initial data $(v_0, v_1) \in L^\infty(0, L) \times L^1(0, L)$, boundary data $\xi \in L^\infty(0, T)$ on the boundary $x = L$ and Neumann feedback with feedback parameter $\eta \in \mathbb{R}$ on the boundary $x = 0$. This leads to the system

$$\begin{cases} v(0, x) = v_0(x), & v_t(0, x) = v_1(x), & \text{(Initial Condition)} \\ v_x(t, 0) = \eta v_t(t, 0), & v(t, L) = \xi(t), & \text{(Boundary Condition)} \quad \text{(S)} \\ v_{tt}(t, x) = c^2 v_{xx}(t, x). & & \text{(PDE)} \end{cases}$$

An explicit representation of the generated state in terms of travelling waves (d'Alembert's solution) is given in Gugat 2015, Gugat and Leugering 2008. This allows the computation of the system state $v \in L^\infty((0, T) \times (0, L))$ without discretization errors. The feedback stabilization of a wave equation is discussed in Adelhütte et al. 2018 (In Revision), Section 5 for the case of completely absorbing Neumann feedback.

**Theorem 19** (Solution of system (S))**.**
*Consider system* (S) *with initial-boundary data*

$$\xi \in L^\infty(0, T) \quad and \quad (v_0, v_1) \in L^\infty(0, L) \times L^1(0, L)$$

*for the feedback parameter $\eta \neq -\frac{1}{c}$. Define the antiderivative of $v_1$ by*

$$V_1(x) := \int_0^x v_1(s) \, \mathrm{d}s$$

*and define*

$$\alpha(s) := \begin{cases} v_0(cs) + \frac{1}{c} V_1(cs) & \text{for } s \in \left[0, \frac{L}{c}\right) \\ 2\xi\left(s - \frac{L}{c}\right) - \beta\left(s - \frac{L}{c}\right) & \text{for } s \in \left[\frac{L}{c}, T + \frac{L}{c}\right] \end{cases}$$

*and*

$$\beta(s) := \begin{cases} v_0(L - cs) - \frac{1}{c} V_1(L - cs) & \text{for } s \in \left[0, \frac{L}{c}\right) \\ \gamma(\eta)\alpha\left(s - \frac{L}{c}\right) + C(\eta) & \text{for } s \in \left[\frac{L}{c}, T + \frac{L}{c}\right] \end{cases}$$

*with the constants*

$$C(\eta) := [1 - \gamma(\eta)]v_0(0),$$
$$\gamma(\eta) := \frac{1 - c\eta}{1 + c\eta}.$$

*Then, the function*

$$v(t, x) := \tfrac{1}{2}\alpha\left(t + \tfrac{x}{c}\right) + \tfrac{1}{2}\beta\left(t + \tfrac{L-x}{c}\right) \tag{3.3}$$

*solves system* (S) *and the solution $v$ lies in $L^\infty((0,T) \times (0,L))$.*

Note that for the feedback parameter $\eta = 1/c$ all energy is absorbed at the boundary $x = 0$ as $\gamma(\eta)$ becomes zero and the solution simplifies.

*Proof.* Because the initial-boundary data is not necessarily weakly differentiable, all upcoming derivatives are to be understood in the sense of distributions.

*Wave Equation:*
First, we see that $v$ satisfies the wave equation, because we have

$$\begin{aligned} v_{tt} &= \tfrac{1}{2}\alpha''\left(t + \tfrac{x}{c}\right) + \tfrac{1}{2}\beta''\left(t + \tfrac{L-x}{c}\right), \\ v_x &= \tfrac{1}{2c}\alpha'\left(t + \tfrac{x}{c}\right) - \tfrac{1}{2c}\beta'\left(t + \tfrac{L-x}{c}\right), \\ v_{xx} &= \tfrac{1}{2c^2}\alpha''\left(t + \tfrac{x}{c}\right) + \tfrac{1}{2c^2}\beta''\left(t + \tfrac{L-x}{c}\right). \end{aligned}$$

*Initial Conditions:*

At $t = 0$, we have for $x \in (0, L)$

$$v(0, x) = \tfrac{1}{2}\alpha\left(\tfrac{x}{c}\right) + \tfrac{1}{2}\beta\left(\tfrac{L-x}{c}\right)$$
$$= \tfrac{1}{2}\left[v_0(x) + \tfrac{1}{c}V_1(x)\right] + \tfrac{1}{2}\left[v_0(x) - \tfrac{1}{c}V_1(x)\right] = v_0(x).$$

The time derivative at $t = 0$, $x \in (0, L)$ fulfills

$$v_t(0, x) = \tfrac{1}{2}\alpha'\left(\tfrac{x}{c}\right) + \tfrac{1}{2}\beta'\left(\tfrac{L-x}{c}\right)$$
$$= \tfrac{1}{2}[v_0'(x) + v_1(x)] - \tfrac{1}{2}[v_0'(x) - v_1(x)] = v_1(x).$$

*Boundary Conditions:*

Now we prove that the Dirichlet boundary condition at $x = L$ is fulfilled for $t > 0$. We obtain

$$v(t, L) = \tfrac{1}{2}\alpha(t + \tfrac{L}{c}) + \tfrac{1}{2}\beta(t)$$
$$= \tfrac{1}{2}[2\xi(t) - \beta(t)] + \tfrac{1}{2}\beta(t) = \xi(t).$$

The feedback law at $x = 0$, implies

$$0 = 2[v_x(t, 0) - \eta v_t(t, 0)] = \tfrac{1}{c}\alpha'(t) - \tfrac{1}{c}\beta'(t + \tfrac{L}{c}) - \left[\eta\alpha'(t) - \eta\beta'(t + \tfrac{L}{c})\right]$$
$$= [\tfrac{1}{c} - \eta]\alpha'(t) - [\tfrac{1}{c} + \eta]\beta'(t + \tfrac{L}{c}).$$

This is equivalent to

$$\beta'(t + \tfrac{L}{c}) = \frac{1 - c\eta}{1 + c\eta}\alpha'(t).$$

Integrating yields

$$\beta(t + \tfrac{L}{c}) = \gamma(\eta)\alpha(t) + C,$$

with a constant of integration $C(\eta)$ that is given by

$$C(\eta) := \beta(\tfrac{L}{c}) - \gamma(\eta)\alpha(0) = v_0(0) - \tfrac{1}{c}V_1(0) - \gamma(\eta)\left[v_0(0) - \tfrac{1}{c}V_1(0)\right]$$
$$= [1 - \gamma(\eta)]v_0(0). \qquad \square$$

Now we show that $v$ lies in $L^\infty((0, T) \times (0, L))$. By the assumptions, we have $v_0 \in L^\infty(0, L)$ and $\xi \in L^\infty(0, L)$. The claim is true if $V_1$ is in $L^\infty(0, L)$. We know that $v_1$ is in $L^1(0, L)$. This implies

$$\|V_1\|_{L^\infty} = \operatorname*{ess\,sup}_{x \in (0, L)} \left| \int_0^x v_1(s)\, ds \right|$$
$$\leq \operatorname*{ess\,sup}_{x \in (0, L)} \int_0^x |v_1(s)|\, ds \leq \int_0^L |v_1(s)|\, ds = \|v_1\|_{L^1}.$$

This finishes the proof.

## 3.2. Karhunen-Loève Approximation

For uncertain initial-boundary data given by stochastic processes, we are interested in the stability of the system in the following sense: The $L^\infty((0,T)\times(0,L))$-norm of the solution should stay below a prescribed bound with high probability. This means that the probability

$$G(\eta) := \mathbb{P}(\|v(\eta)\|_\infty \leq v_{\text{ub}}) \qquad (3.4)$$

should be large. The solution $v(\eta)$ depends on the feedback parameter $\eta$ and the realization of the inital-boundary data. For each realization it is given by Theorem 19.

Because the infinite dimensional case of arbitrary stochastic processes as boundary data cannot be handled numerically, we need a tool to break down the dimension to the finite dimensional case. This tool is the Karhunen-Loève series. It allows us to write an stochastic process as a series of random coefficients multiplied with the eigenfunctions of the covariance function[1] of the stochastic process. The series can be cut off after a finite number of summands, which leads to an approximation of the stochastic process that requires only a finite number of random variables and is optimal in the $L^2$-sense. The theorem is due to Karhunen 1947 and Loève 1978.

**Theorem 20** (Karhunen-Loève Theorem)**.**
*Let $\{X_t\}_{t\in[a,b]}$ be a zero-mean square integrable stochastic process on the probability space $(\Omega, \Sigma, \mu)$ over the finite time intervall $[a,b]$ with continuous covariance function*

$$C(s,t) := \text{cov}(X_s, X_t) = \mathbb{E}(X_s X_t).$$

*Define the Hilbert-Schmidt operator of the covariance function as*

$$T_C : L^2([a,b]) \to L^2([a,b]), \quad T_C f := \int_a^b C(s,\,\cdot\,)f(s)\,\mathrm{d}s.$$

*Denote the pairs of eigenvalues $\lambda_k$ and eigenfunctions $e_k$ of the operator $T_C$ that solve*

$$T_C e_k = \lambda_k e_k.$$

*Then, the eigenfunctions $e_k$ form an orthonormal basis of $L^2([a,b])$ and the stochastic process has the representation*

$$X_t = \sum_{k=1}^{\infty} a_k e_k(t)$$

---

[1]To be precise: the eigenfunctions of the Hilbert–Schmidt integral operator corresponding to the covariance function

*with the random variables*

$$a_k := \int_a^b X_t e_k(t)\, \mathrm{d}t.$$

*The convergence of the series is in $L^2(\Omega)$ and uniform in $t$. The random coefficients have zero-mean, are uncorrelated and have variance $\lambda_k$.*

The covariance function of a Wiener process is given by

$$C(t,s) = \mathrm{cov}(W_t, W_s) = \min(s,t),$$

which allows a direct calculation of the Karhunen-Loève series.

**Corollary 12** (Karhunen-Loève series of a Wiener process)**.**
*Let $\{W_t\}_{t\in[0,1]}$ be a Wiener process on $[0,1]$.*
*   *Then there is a sequence of independent standard normal distributed random variables $\{a_k\}_{k=1,\dots,\infty}$ such that*

$$W_t = \sqrt{2}\sum_{k=1}^{\infty} a_k \frac{\sin(\omega_k \pi t)}{\omega_k \pi}, \quad \omega_k = \left(k - \tfrac{1}{2}\right).$$

*The convergence is in $L^2(\Omega)$ and uniform in $t$.*

We use a finite approximation of the representation of a Wiener process as in Corollary 12 as initial boundary data. However, by the use of Theorem 20, the methodology can be applied to arbitrary zero-mean square integrable stochastic processes with the obvious adaptions. The finite dimensional approximation of a Wiener process on $[0,T]$ yields

$$\xi(t) = \sqrt{2T}\sum_{k=1}^{N} a_k \frac{\sin\left(\omega_k \pi \frac{t}{T}\right)}{\omega_k \pi}, \quad \omega_k = k - \tfrac{1}{2}, \text{ on } [0,T]. \qquad \text{(KL-bd)}$$

Analogously, we choose the compatible initial data on $[0,L]$

$$v_0(x) = \sqrt{2L}\sum_{k=1}^{N} b_k \frac{\sin\left(\omega_k \pi \frac{L-x}{L}\right)}{\omega_k \pi}, \quad \omega_k = k - \tfrac{1}{2}, \text{ on } [0,L]. \qquad \text{(KL-id)}$$

The random variables $a_k$ and $b_k$ are independently normal distributed. The compatibility condition $\xi(0) = v_0(L) = 0$ holds. Furthermore, set $v_1 = 0$. The value of $\|v\|_{L^\infty}$ is not easily expressed as an analytic function of the random variables. This means a sampling scheme based on spheric radial

Figure 3.1.: Realizations of the random initial-boundary data



(a) Different realizations (21) of the initial data for a Karhunen-Loève sum with 20 standard normal distributed coefficients

(b) Different realizations (21) of the boundary data for a Karhunen-Loève sum with 20 standard normal distributed coefficients

decomposition can not be directly be applied. We use a quasi-Monte Carlo method based on a Sobol sequence instead. Different realizations of the initial-boundary data are shown in Figure

For an approximation of the $L^\infty$-norm of the velocity by pointwise evaluation on a grid, the Lipschitz continuity of the velocity is required.

**Definition 34** (Space of Lipschitz Continuous Functions)**.** Denote the space of *Lipschitz continuous functions* on $X$ by

$$\mathcal{C}^{0,1}(X) := \{f : X \to \mathbb{R} \mid \text{there exists a constant } K > 0 \text{ such that}$$
$$|f(x) - f(y)| \le K\|x - y\|_X\}$$

**Theorem 21** (Lipschitz Continuity of the Solution)**.**
*Assume that the boundary data* $\xi \in \mathcal{C}^{0,1}([0,T])$ *and initial data* $v_0 \in \mathcal{C}^{0,1}([0,L])$ *are Lipschitz continuous and assume that Lipschitz compatibility over the edge holds, i.e.,*

$$|\xi(t) - v_0(L - x)| \le K|t - L + x|, \quad \text{for} \quad (t,x) \in [0,T] \times [0,L] \quad (3.5)$$

*with a Lipschitz constant $K > 0$. Furthermore, let the initial time derivative $v_1 \in L^\infty(0, L)$ be bounded. Then the solution $v$ of system (S) is Lipschitz continuous in time and space, i.e., $v \in \mathcal{C}^{0,1}([0, T] \times [0, L])$.*

*Proof.* The sum of Lipschitz continuous functions is Lipschitz continuous. It is therefore sufficient to show the Lipschitz continuity of $\alpha$ and $\beta$ defined in Theorem 19. Without loss of generality—by going to the maximum of the occurring Lipschitz constants—we assume that they are all the same and denote each of them by $K > 0$. First, we show the Lipschitz continuity of $V_1$. We have, for $x, y \in [0, L]$

$$
\begin{aligned}
|V_1(x) - V_1(y)| &= \left| \int_0^x v_1(s)\,\mathrm{d}s - \int_0^y v_1(s)\,\mathrm{d}s \right| \\
&= \left| \int_y^x v_1(s)\,\mathrm{d}s \right| \le |x - y|\, \|v_1\|_{L^\infty}.
\end{aligned}
$$

The functions $\alpha$ and $\beta$ are Lipschitz continuous on the interval $[0, \frac{L}{c})$ by the Lipschitz continuity of $v_0$ and $V_1$. On the interval $[\frac{L}{c}, T + \frac{L}{c}]$ we can use an induction argument: We show the Lipschitz continuity of $\alpha$ and $\beta$ for

$$
s \in I_k := \left[ \frac{kL}{c}, \frac{(k+1)L}{c} \right), \ k \ge 1
$$

under the assumption that $\alpha$ and $\beta$ are Lipschitz continuous on the interval

$$
I_{k-1} := \left[ \frac{(k-1)L}{c}, \frac{kL}{c} \right).
$$

By the definition of $\alpha$ in Theorem 19, we have for $s, r \in I_k$ using the triangle inequality

$$
\begin{aligned}
|\alpha(s) - \alpha(r)| &= \left| 2\xi\left(s - \tfrac{L}{c}\right) - \beta\left(s - \tfrac{L}{c}\right) - 2\xi\left(r - \tfrac{L}{c}\right) + \beta\left(r - \tfrac{L}{c}\right) \right| \\
&\le 2\left| \xi\left(s - \tfrac{L}{c}\right) - \xi\left(r - \tfrac{L}{c}\right) \right| + \left| \beta\left(s - \tfrac{L}{c}\right) - \beta\left(r - \tfrac{L}{c}\right) \right| \\
&\le 2K|s - r| + K|s - r| = 3K|s - r|.
\end{aligned}
$$

by the Lipschitz continuity of $\xi$ and the Lipschitz continuity of $\beta$ on $I_{k-1}$. For the function $\beta$, we obtain for $s, r \in \left[ \frac{kL}{c}, \frac{(k+1)L}{c} \right)$

$$
\begin{aligned}
|\beta(s) - \beta(r)| &= \left| \gamma(\eta)\alpha\left(s - \tfrac{L}{c}\right) - \gamma(\eta)\alpha\left(r - \tfrac{L}{c}\right) \right| \\
&\le K|\gamma(\eta)||s - r|,
\end{aligned}
$$

because $\alpha$ is Lipschitz continuous on $I_{k-1}$. Thus, we have shown the Lipschitz continuity of $\alpha$ and $\beta$ on the individual intervals $\left[0, \frac{L}{c}\right)$ and $\left[\frac{L}{c}, T + \frac{L}{c}\right]$

(the inclusion of the right boundary is due to the continuity of the functions on $\left[\frac{L}{c}, T + \frac{L}{c}\right]$). Clearly, the functions are Lipschitz continuous on $\left[0, T + \frac{L}{c}\right]$ if they are continuous in the point $\frac{L}{c}$. The Lipschitz constant can be chosen as the maximum of the Lipschitz constants on each interval. We calculate the left and right limits:

$$\lim_{s \nearrow \frac{L}{c}} \alpha(s) = \lim_{s \nearrow \frac{L}{c}} v_0(cs) + \tfrac{1}{c}V_1(cs) = v_0(L) + \tfrac{1}{c}V_1(L)$$

$$\lim_{s \searrow \frac{L}{c}} \alpha(s) = \lim_{s \searrow \frac{L}{c}} 2\xi\left(s - \tfrac{L}{c}\right) - \beta\left(s - \tfrac{L}{c}\right)$$

$$= 2\xi(0) - \beta(0) = 2\xi(0) - v_0(L) + \tfrac{1}{c}V_1(L)$$

Both limits are equal, because $\xi(0)$ equals $v_0(L)$ as a consequence of (3.5). For $\beta$, we have

$$\lim_{s \nearrow \frac{L}{c}} \beta(s) = \lim_{s \nearrow \frac{L}{c}} v_0(L - cs) - \tfrac{1}{c}V_1(L - cs) = v_0(0)$$

$$\lim_{s \searrow \frac{L}{c}} \beta(s) = \lim_{s \searrow \frac{L}{c}} \gamma(\eta)\alpha\left(s - \tfrac{L}{c}\right) + C(\eta) = \gamma(\eta)\alpha(0) + C(\eta)$$

$$= \gamma(\eta)[v_0(0) + \tfrac{1}{c}V_1(0)] + [1 - \gamma(\eta)]v_0(0) = v_0(0),$$

where $V_1(0) = 0$ was used. This shows the continuity of $\alpha$ and $\beta$ and concludes the proof. $\qquad\square$

The solution for different realizations of the initial-boundary data with the value of the $L^\infty$-norm is depicted in Figure 3.3 for the feedback parameter $\eta = 1$ and in Figure 3.4 for completely absorbing feedback $\eta = 1/c$. Note that in the case of feedback $\eta = 1/c$, the value of the $L^\infty$-norm can be calculated analytically in dependence of the initial and boundary data.

**Theorem 22** (Value of $\|v\|_{L^\infty}$ in terms of initial and boundary data).
*Let $v$ be a solution of system* (S) *for $\eta = 1/c$ under the assumptions of Theorem 19. For $(t, x) \in (0, T) \times (0, L)$, define*

$$m_1(t, x) := \tfrac{1}{2}\left[v_0(x + ct) + \tfrac{1}{c}V_1(x + ct)\right] + \tfrac{1}{2}\left[v_0(x - ct) - \tfrac{1}{c}V_1(x - ct)\right],$$

$$m_2(t, x) := \tfrac{1}{2}\left[v_0(ct + x) + \tfrac{1}{c}V_1(ct + x) + v_0(0)\right],$$

$$m_3(t, x) := \xi(t + \tfrac{x - L}{c}) - \tfrac{1}{2}\left[v_0(2L - x - ct) - \tfrac{1}{c}V_1(2L - x - ct)\right]$$
$$\qquad\qquad + \tfrac{1}{2}\left[v_0(x - ct) - \tfrac{1}{c}V_1(x - ct)\right],$$

$$m_4(t, x) := \xi(t + \tfrac{x - L}{c}) + \tfrac{1}{2}\left[\tfrac{1}{c}V_1(2L - x - ct) + v_0(0) - v_0(2L - x - ct)\right],$$

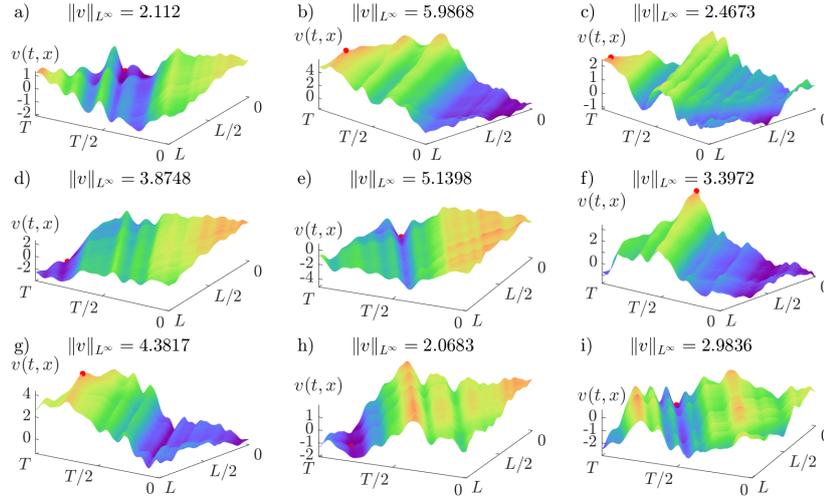$$m_5(t, x) := \xi(t + \tfrac{x - L}{c}),$$

Figure 3.3.: The solution $v$ of the wave equation with boundary and initial data given by the functions defined in (KL-bd) and (KL-id) for nine samples of the standard normal distributed random vector $(a, b)$ with realizations in $\mathbb{R}^{40}$, i.e. $N = 20$. The data $T = 8$, $L = 2$, $c = 0.5$, $\eta = 1$ was used and the value of the $L^\infty$-norm was evaluated on a $100 \times 100$ grid in time and space. The estimated probability after 10000 samples is 0.8101.

*denote $\Omega_T = (0, T) \times (0, L)$ and set*

$$\Omega_1 := \left\{ (t, x) \in \overline{\Omega}_T \mid t < \min\{\tfrac{L-x}{c}, \tfrac{x}{c}\} \right\},$$
$$\Omega_2 := \left\{ (t, x) \in \overline{\Omega}_T \mid \tfrac{x}{c} \leq t < \tfrac{L-x}{c} \right\},$$
$$\Omega_3 := \left\{ (t, x) \in \overline{\Omega}_T \mid \tfrac{L-x}{c} \leq t < \tfrac{x}{c} \right\},$$
$$\Omega_4 := \left\{ (t, x) \in \overline{\Omega}_T \mid \max\{\tfrac{L-x}{c}, \tfrac{x}{c}\} \leq t < \tfrac{L}{c} + \tfrac{L-x}{c} \right\},$$
$$\Omega_5 := \left\{ (t, x) \in \overline{\Omega}_T \mid t \geq \tfrac{L}{c} + \tfrac{L-x}{c} \right\};$$

*see Figure 3.5. For $i \in \{1, \ldots, 5\}$ denote*

$$M_i := \sup\{ |m_i(t, x)| : (t, x) \in \Omega_i \}.$$

*Then the $L^\infty$-norm of the solution $v$ is given by*

$$\|v\|_{L^\infty} = \max\{M_1, M_2, M_3, M_4, M_5\}.$$

*Proof.* By Theorem 19 the solution of system (S) is given by

$$v(t, x) := \tfrac{1}{2}\alpha\big(t + \tfrac{x}{c}\big) + \tfrac{1}{2}\beta\big(t + \tfrac{L-x}{c}\big).$$
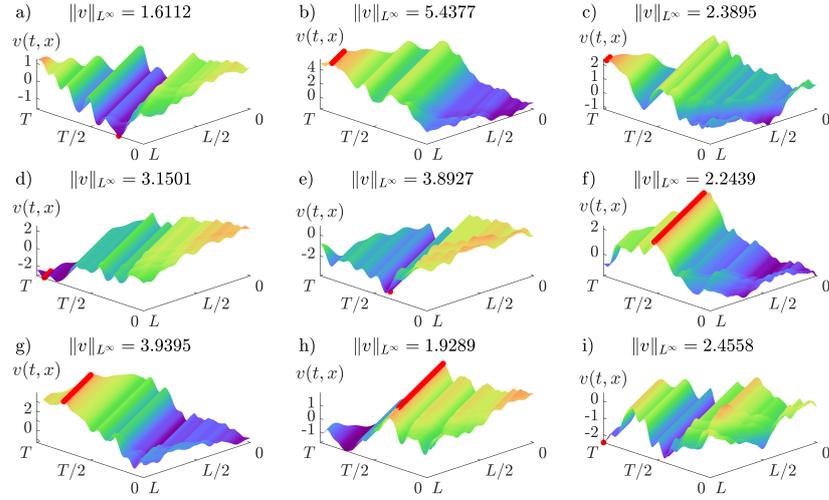
Figure 3.4.: The solution $v$ of the wave equation with boundary and initial data given by the functions in (KL-bd) and (KL-id) for nine samples of the standard normal distributed random vector $(a, b)$ with realizations in $\mathbb{R}^{40}$, i.e., $N = 20$. The constants $T = 6$, $L = 2$, $c = 0.5$ and the completely absorbing feedback $\eta = 2$ were used. The bound $v_{\max} = 5$ was chosen. The probability of $\|v\|_{L^\infty} \leq v_{\max}$ is 0.8808 with 10000 samples used. The value of the $L^\infty$-norm is approximated by evaluation on a $100 \times 100$ grid on $[0, T] \times [0, L]$. The points, where the value of the $L^\infty$-norm is attained are marked in red.

Remember that the definition of $\alpha$ was

$$\alpha(s) := \begin{cases} v_0(cs) + \frac{1}{c} V_1(cs) & \text{for } s \in \left[0, \frac{L}{c}\right) \\ 2\xi\left(s - \frac{L}{c}\right) - \beta\left(s - \frac{L}{c}\right) & \text{for } s \in \left[\frac{L}{c}, T + \frac{L}{c}\right]. \end{cases}$$

Because $\eta = 1/c$ implies $\gamma(\eta) = 0$ and $C(\eta) = v_0(0)$, the function $\beta$ simplifies to

$$\beta(s) := \begin{cases} v_0(L - cs) - \frac{1}{c} V_1(L - cs) & \text{for } s \in \left[0, \frac{L}{c}\right) \\ v_0(0) & \text{for } s \in \left[\frac{L}{c}, T + \frac{L}{c}\right]. \end{cases}$$

By the definition of $\alpha$ and $\beta$, there are four cases to consider. The last case is split into two subcases. The cases correspond to the different regions in the space time domain that are separated by the characteristic curves

$$t = \frac{x}{c}, \quad t = \frac{L-x}{c} \quad \text{and} \quad t = \frac{L}{c} + \frac{L-x}{c}.$$
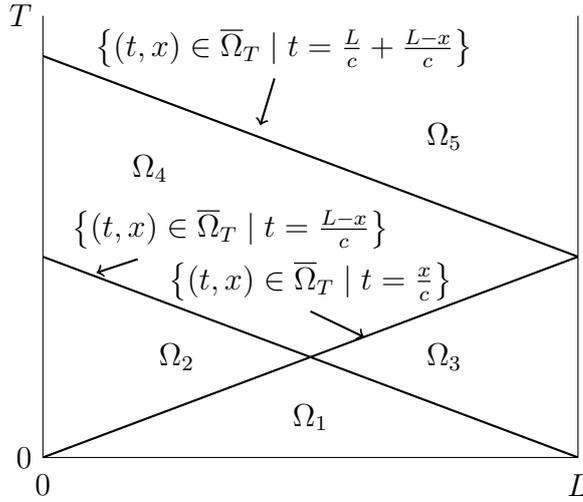
Figure 3.5.: The different cases to consider are determined by the sectioning of $\overline{\Omega}_T = [0, T] \times [0, L]$ by the characteristic curves.

The regions are depicted in Figure 3.5.

The first case $t < \min\{\frac{x}{c}, \frac{L-x}{c}\}$, which is $(t, x) \in \Omega_1$, corresponds to being in the first case for both $\alpha$ and $\beta$. Therefore, we obtain

$$v(t, x) = \tfrac{1}{2}\big[v_0(x + ct) + \tfrac{1}{c}V_1(x + ct)\big] + \tfrac{1}{2}\big[v_0(x - ct) - \tfrac{1}{c}V_1(x - ct)\big].$$

For $\frac{x}{c} \leq t < \frac{L-x}{c}$, which is $(t, x) \in \Omega_2$, we are in the first case for $\alpha$ and in the second case for $\beta$. Note that the interval for $t$ can only be nonempty for $x \in \big(0, \frac{L}{2}\big)$. Hence,

$$v(t, x) = \tfrac{1}{2}\big[v_0(ct + x) + \tfrac{1}{c}V_1(ct + x) + v_0(0)\big].$$

For $\frac{L-x}{c} \leq t < \frac{x}{c}$, which corresponds to $(t, x) \in \Omega_3$, we are in the second case for $\alpha$ and in the first case for $\beta$. Note that the interval for $t$ can only be nonempty for $x \in \big(\frac{L}{2}, L\big)$. Consequently,

$$\begin{aligned}
v(t, x) &= \tfrac{1}{2}\big[2\xi(t + \tfrac{x-L}{c}) - \beta(t + \tfrac{x-L}{c}) + \beta(t + \tfrac{L-x}{c})\big] \\
&= \xi(t + \tfrac{x-L}{c}) - \tfrac{1}{2}\beta(t + \tfrac{x-L}{c}) + \tfrac{1}{2}\big[v_0(x - ct) - \tfrac{1}{c}V_1(x - ct)\big] \\
&= \xi(t + \tfrac{x-L}{c}) - \tfrac{1}{2}\big[v_0(2L - x - ct) - \tfrac{1}{c}V_1(2L - x - ct)\big] \\
&\quad + \tfrac{1}{2}\big[v_0(x - ct) - \tfrac{1}{c}V_1(x - ct)\big],
\end{aligned}$$

since $t < \frac{x}{c} < \frac{L}{c}$ and $\frac{x-L}{c} < 0$, imply $t + \frac{x-L}{c} < \frac{L}{c}$. The last case to consider is $t \geq \max\{\frac{L-x}{c}, \frac{x}{c}\}$. It leads to

$$v(t, x) = \tfrac{1}{2}\big[2\xi(t + \tfrac{x-L}{c}) - \beta(t + \tfrac{x-L}{c}) + v_0(0)\big],$$

which is equal to

$$\xi(t + \tfrac{x-L}{c}) + \tfrac{1}{2}\left[\tfrac{1}{c}V_1(2L - x - ct) + v_0(0) - v_0(2L - x - ct)\right],$$

for $t < \tfrac{L}{c} + \tfrac{L-x}{c}$ and

$$\xi(t + \tfrac{x-L}{c}), \quad \text{for } t \geq \tfrac{L}{c} + \tfrac{L-x}{c}.$$

This corresponds to the domains $\Omega_4$ and $\Omega_5$. $\qquad\qquad\qquad\square$

## 3.3. Optimization of the Feedback Parameter

The feedback parameter $\eta$ can be chosen such that the probability $G(\eta)$ defined in Equation (3.4) is maximized. We consider the probability to



Figure 3.6.: The probability to stay under the bound $v_{\mathrm{ub}} = 5$ over the feedback parameter $\eta$ for the data $L = 2$, $c = 0.5$, $T = 2$ using 2000 samples. In this case, the maximum of the probability is reached for completely absorbing feedback $\eta = 1/c = 2$.



Figure 3.7.: The quasi-Monte Carlo approximation of the probability to stay under the bound $v_{\mathrm{ub}} = 5$ over the number of samples $N$ for the data $L = 2$, $c = 0.5$, $T = 2$ for the feedback parameter $\eta = 3$.
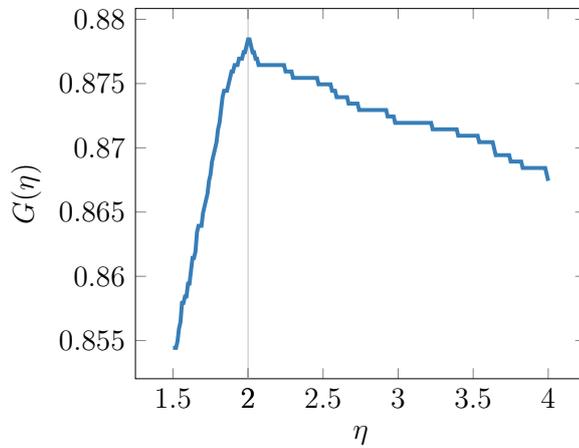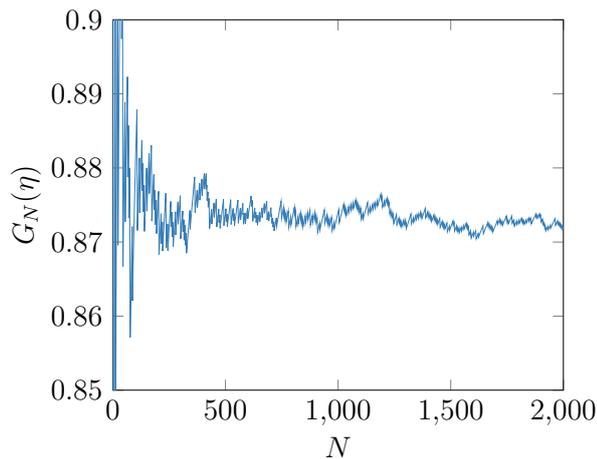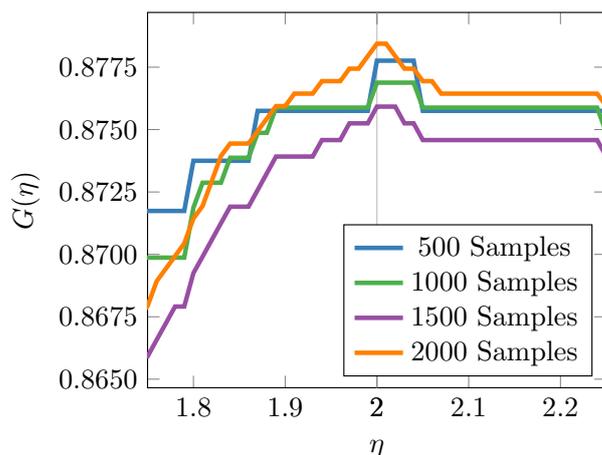
Figure 3.8.: The probability to stay under the bound $v_{\mathrm{ub}} = 5$ over the feedback parameter $\eta$ for the data $L = 2$, $c = 0.5$, $T = 2$ using $500, 1000, 1500$ and $2000$ samples.

stay under the bound $v_{\max} = 5$ for different feedback parameters $\eta > 0$ on a grid with stepsize 0.01 between 1.5 and 4. The data for the example has been chosen as $L = 2$, $T = 2$, $c = 0.5$. For the approximation of the probability, 2000 samples were used for each value of $\eta$. In Figure 3.6 the probability is plotted over the feedback parameter. The maximum of the probability is reached for completely absorbing feedback $\eta = 1/c = 2$. The peak in probability is very distinct. At the peak the probability function appears to be nonsmooth.

The Figure 3.7 shows the quasi-Monte Carlo approximation of the probability function over the number of samples. The variance decrease and thus the increasing quality of the approximation for higher sample numbers is apparent. The importance of sufficiently high sample numbers is highlighted in Figure 3.8. The lower sample numbers do not show the sharp peak at $\eta = 2$ as it is the case for $N = 2000$. The grid stepsize between $\eta = 1.75$ and $\eta = 2.25$ is 0.01 for all sample sizes.

## Conclusion

We introduced a way to model stochastic boundary conditions for partial differential equations (PDEs) using a Karhunen-Loève series. When including uncertainty into the system, the existing literature focuses on uncertainty of the dynamics, leading to stochastic PDEs. Although the uncertainty of the input data in the form of boundary conditions is important in applications such as gas pipeline control, it is a barely explored research topic. The robustness in the sense that the supremum of the absolute value of the state stays under a certain bound with high probability has been

considered. The system can be stabilized using a feedback stabilization to maximize this probability.

There are different directions to continue the research of this topic: Other norms than the $L^\infty$-norm may be used in the probability function. There is hope that the classical choice of a squared $L^2$-norm might lead to a probability function that is differentiable with respect to the feedback parameter. Analytical gradients or subgradients are desirable, yet traditional methods as in Uryas'ev 1994, Kibzun and Uryasev 1998 or Shapiro, Dentcheva, and Ruszczyński 2009, Theorem 4.77 fail due to the infinite dimensional nature of the problem (the solution $v$ lies in the space $L^\infty((0,T) \times (0,L))$). Furthermore, the feedback boundary may be replaced by a free control expanding the framework from a one dimensional optimization parameter to a infinite dimensional optimal control problem. Nontrivial optimality conditions for this problem are not known due to the lack of derivatives even for the feedback control problem. The extension to networks is open to further research.

# A. Mathematical Basics

**Implicit Function Theorem** (Rudin 1976, Theorem 9.28)**.**
*Let $U_1 \subset \mathbb{R}^k$ and $U_2 \subset \mathbb{R}^m$ be open subsets and let*

$$F : U_1 \times U_2 \to \mathbb{R}^m, \quad (x, y) \mapsto F(x, y),$$

*be continuously differentiable. Let a point $(a, b) \in U_1 \times U_2$ that fulfills*

$$F(a, b) = 0$$

*be given. Assume that the Jacobian $D_y F$ with respect to $y$ is invertible in the point $(a, b)$. Then, there exist open neighborhoods $V_1(a) \subset U_1$ of $a$ and $V_2(b) \subset U_2$ of $b$ and a continuously differentiable mapping $g : V_1 \to V_2 \subset \mathbb{R}^m$ with $g(a) = b$ such that*

$$F(x, g(x)) = 0 \quad \text{for all } x \in V_1.$$

*If $(x, y) \in V_1 \times V_2$ fulfills $F(x, y) = 0$, then $y = g(x)$ follows. The derivative of $g$ is given by*

$$g'(x) = -(D_y F(x, g(x)))^{-1} D_x F(x, g(x)).$$

# B. Data of Chapter 1

For the convenience of the reader, Tables B.1–B.7 below contain the data used in the numerical examples presented in Chapter 1. The constant $\alpha_a$ in the compressibility factor model is computed via (see e.g., Domschke et al. 2017)

$$\alpha = \frac{0.257}{p_c} - 0.533 \frac{T_c}{p_c\,T},$$

where $p_c > 0$ is the pseudocritical pressure and $T_c > 0$ is the pseudocritical temperature. We have used the formula of Chen (s. Cerbe 2008), which is an explicit estimate for the Colebrook-White law, to calculate realistic values for $\theta_a = \lambda_a/D_a$. The equations for the Reynolds number and the friction factor read as follows:

$$\mathrm{Re}_a := \frac{q_a D_a}{\eta_v},$$

$$\frac{1}{\sqrt{\lambda_a}} = -2\log_{10}\left(\frac{k_a/D_a}{3.7065} - \frac{5.0425}{\mathrm{Re}_a}\log_{10}\left[\frac{(k_a/D_a)^{1.1098}}{2.8257} + \frac{5.8506}{(\mathrm{Re}_a)^{0.8981}}\right]\right).$$

In practice the friction factor $\lambda_a$ is of the order $10^{-1} - -1$. Since the flow $q_a$ is constant along each pipe, for the stationary states, this leads to constant friction factors $\lambda_a$ along each pipe. In all examples, the dynamic viscosity is $\eta_v := 11.9 \cdot 10^{-6}\,\mathrm{kg\,m^{-1}\,s^{-1}}$. Note that our analysis does not directly cover the dependency of the friction factor on $q_a$, but it poses no problems for numerical calculations. For trees, there is no difference between a flow-dependent and a flow-independent friction law.

Table B.1 contains the constants used for the single pipe computations in Section 1.1 and 1.2. Table B.2 contains the globally used constants and Table B.3 the constants for each pipe for the example of a tree networks in Section 1.5.1 and Section 1.5.2. The boundary values for the tree example are presented in Table B.4. The globally used constants for the diamond graph (Section 1.5.1) are contained in Table B.5. The pipe data for this example is presented in Table B.6 and its boundary flow values in Table B.7.

Table B.1.: Values for single pipe examples

| Symbol | Name | Value | Unit |
|--------|------|-------|------|
| $R$ | Specific gas constant | 448.66 | $\mathrm{J\,kg^{-1}\,K^{-1}}$ |
| $T$ | Temperature | 290 | K |
| $p_c$ | Pseudocritical pressure | 46.70206 | $10^5\,\mathrm{Pa}$ |
| $T_c$ | Pseudocritical temperature | 202.43951 | K |
| $\alpha$ | Constant in the state equation | -2.46391 | $10^{-8}\,\mathrm{Pa^{-1}}$ |
| $k_a$ | Roughness | 0.06 | mm |

Table B.2.: Global values for the tree network example

| Symbol | Name | Value | Unit |
|--------|------|-------|------|
| $R$ | Specific gas constant | 447.80152 | $\mathrm{J\,kg^{-1}\,K^{-1}}$ |
| $T$ | Temperature | 273.15 | K |
| $p_c$ | Pseudocritical pressure | 45.92935 | $10^5\mathrm{Pa}$ |
| $T_c$ | Pseudocritical temperature | 188.54976 | K |
| $\alpha$ | Constant in the state equation | -2.41499 | $10^{-8}\mathrm{Pa^{-1}}$ |

Table B.3.: Pipe data for the tree network example

| Pipe name | From | To | Length $L_a$ | Diameter $D_a$ | Roughness $k_a$ |
|-----------|------|----|--------------|----------------|------------------|
| $a$ | $r$ | $v$ | 15.99078 km | 1.0 m | 0.05 mm |
| $b$ | $u$ | $v$ | 44.71898 km | 1.0 m | 0.05 mm |
| $c$ | $u$ | $w$ | 7.05313 km | 1.0 m | 0.05 mm |
| $d$ | $u$ | $s$ | 38.05409 km | 1.0 m | 0.05 mm |

Table B.4.: Given boundary data for the tree network example

| Node name | Node outflow $q_v^{\mathrm{out}}$ | Pressure $p_r$ |
|-----------|-----------------------------------|-----------------|
| $r$ | -130.83333 $\mathrm{kg\,s^{-1}}$ | $13.00000 \cdot 10^5\,\mathrm{Pa}$ |
| $u$ | 0 $\mathrm{kg\,s^{-1}}$ | – |
| $v, w, s$ | 43.61111 $\mathrm{kg\,s^{-1}}$ | – |

Table B.5.: Global values for the diamond graph example

| Symbol | Name | Value | Unit |
|---|---|---|---|
| $R$ | Specific gas constant | 460.66628 | $\mathrm{J\,kg^{-1}\,K^{-1}}$ |
| $T$ | Temperature | 288.15 | K |
| $p_c$ | Pseudocritical pressure | 46.70206 | $10^5\mathrm{Pa}$ |
| $T_c$ | Pseudocritical temperature | 202.43951 | K |
| $\alpha$ | Constant in the state equation | -2.51506 | $10^{-8}\mathrm{Pa^{-1}}$ |

Table B.6.: Pipe data for the diamond graph example

| Pipe name | From | To | Length $L_a$ | Diameter $D_a$ | Roughness $k_a$ |
|---|---|---|---|---|---|
| $a$ | $r$ | $u$ | 39.74748 km | 1.3 m | 0.01 mm |
| $b$ | $u$ | $v$ | 37.57120 km | 1.0 m | 0.01 mm |
| $c$ | $u$ | $w$ | 28.40076 km | 1.3 m | 0.01 mm |
| $d$ | $v$ | $w$ | 26.59033 km | 1.3 m | 0.01 mm |
| $e$ | $v$ | $s$ | 17.97404 km | 1.0 m | 0.01 mm |
| $f$ | $w$ | $s$ | 25.14802 km | 1.0 m | 0.01 mm |
| $g$ | $s$ | $t$ | 14.58364 km | 1.0 m | 0.01 mm |

Table B.7.: Given boundary data for the diamond graph example

| Node name | Node outflow $q_v^{\mathrm{out}}$ | Pressure $p_r$ |
|---|---|---|
| $r$ | -453.1495 $\mathrm{kg\,s^{-1}}$ | $60 \cdot 10^5$ Pa |
| $u, v, w, s$ | 0 $\mathrm{kg\,s^{-1}}$ | – |
| $t$ | 453.1495 $\mathrm{kg\,s^{-1}}$ | – |

# C. Data of Chapter 2

The data for the calculations in Chapter 2 is presented in Tables C.3–C.6. The friction $\lambda$ and the compressibility factor $z$ are calculated as described in Chapter B. The simple tree network (for Figure 2.8) uses the data of Table C.1 and Table C.2. The example of a tree with six exit nodes (see Section 2.5.2) uses the data of Table C.3 and Table C.4. The example with multiple compressors (see Section 2.5.3) uses the data of Tables C.5 and Table C.6.

Table C.1.: Global values for the simple tree network

| Symbol | Name | Value | Unit |
|--------|------|-------|------|
| $R$ | Specific gas constant | 461.91443 | $\mathrm{J\,kg^{-1}\,K^{-1}}$ |
| $T$ | Temperature | 300 | K |
| $p_c$ | Pseudocritical pressure | 46.7 | $10^5\,\mathrm{Pa}$ |
| $T_c$ | Pseudocritical temperature | 202 | K |
| $\alpha$ | Constant in the state equation | -2.18172734 | $10^{-8}\,\mathrm{Pa^{-1}}$ |

Table C.2.: Pipe data for the simple tree example

| Pipe name | From | To | Length $L_a$ | Diameter $D_a$ | Roughness $k_a$ |
|-----------|------|-----|--------------|----------------|-----------------|
| Entry1_Innode | Entry1 | Innode | 100 km | 0.8 m | 0.06 mm |
| Innode_Exit1 | Innode | Exit1 | 80 km | 0.8 m | 0.06 mm |
| Innode_Exit2 | Innode | Exit2 | 60 km | 0.8 m | 0.06 mm |

Table C.3.: Global values for the larger tree example

| Symbol | Name | Value | Unit |
|--------|------|-------|------|
| $R$ | Specific gas constant | 447.80152 | $J\,kg^{-1}\,K^{-1}$ |
| $T$ | Temperature | 273.15 | K |
| $p_c$ | Pseudocritical pressure | 45.92935 | $10^5\,Pa$ |
| $T_c$ | Pseudocritical temperature | 188.54976 | K |
| $\alpha$ | Constant in the state equation | -2.41599 | $10^{-8}\,Pa^{-1}$ |

Table C.4.: Pipe data for the larger tree example

| Pipe name | From | To | Length $L_a$ | Diameter $D_a$ | Roughness $k_a$ |
|-----------|------|------|--------------|----------------|-----------------|
| pipe_1 | source_1 | innode_1 | 13.07109 km | 1.0 m | 0.05 mm |
| pipe_2 | innode_1 | innode_2 | 76.89355 km | 0.8 m | 0.05 mm |
| pipe_3 | innode_2 | sink_1 | 21.55757 km | 1.0 m | 0.05 mm |
| pipe_4 | innode_2 | innode_3 | 6.99805 km | 1.0 m | 0.05 mm |
| pipe_5 | innode_3 | sink_2 | 58.21897 km | 0.8 m | 0.05 mm |
| pipe_6 | innode_3 | sink_3 | 86.69027 km | 0.8 m | 0.05 mm |
| pipe_7 | innode_3 | sink_4 | 16.57933 km | 0.6 m | 0.05 mm |
| pipe_8 | innode_1 | innode_4 | 10.02278 km | 0.6 m | 0.05 mm |
| pipe_9 | innode_4 | sink_5 | 35.21884 km | 0.6 m | 0.05 mm |
| pipe_10 | innode_4 | sink_6 | 20.32221 km | 0.6 m | 0.05 mm |

Table C.5.: Global values for the compressor tree example

| Symbol | Name | Value | Unit |
|--------|------|-------|------|
| $R$ | Specific gas constant | 447.80152 | $J\,kg^{-1}\,K^{-1}$ |
| $T$ | Temperature | 273.15 | K |
| $p_c$ | Pseudocritical pressure | 45.92935 | $10^5\,Pa$ |
| $T_c$ | Pseudocritical temperature | 188.54976 | K |
| $\alpha$ | Constant in the state equation | -2.41499 | $10^{-8}\,Pa^{-1}$ |

Table C.6.: Pipe data for the compressor tree example

| Pipe name | From | To | Length $L_a$ | $D_a$ | $k_a$ |
|---|---|---|---|---|---|
| pipe_1 | source_1 | innode_1 | 13.07109 km | 1.0 m | 0.01 mm |
| pipe_2 | innode_1 | innode_2 | 16.89355 km | 0.8 m | 0.05 mm |
| pipe_3 | innode_2 | sink_1 | 21.55757 km | 1.0 m | 0.05 mm |
| pipe_4 (comp) | innode_2 | innode_3 | – | – | – |
| pipe_5 | innode_3 | sink_2 | 88.21897 km | 0.8 m | 0.05 mm |
| pipe_6 | innode_3 | sink_3 | 86.69027 km | 0.8 m | 0.05 mm |
| pipe_7 | innode_3 | sink_4 | 16.57933 km | 0.6 m | 0.05 mm |
| pipe_8 | innode_11 | innode_4 | 10.02278 km | 0.6 m | 0.05 mm |
| pipe_9 | innode_4 | sink_5 | 15.21884 km | 0.6 m | 0.01 mm |
| pipe_10 (comp) | innode_4 | innode_5 | – | – | – |
| pipe_11 | innode_5 | sink_6 | 91.32221 km | 0.6 m | 0.05 mm |
| pipe_12 | innode_5 | sink_7 | 21.32221 km | 0.6 m | 0.05 mm |

# D. Equation Glossary

(Iso)
$$\begin{cases} \partial_t \rho + \partial_x q = 0 \\ \partial_t q + \partial_x \left( p + \dfrac{q^2}{\rho} \right) = -\dfrac{\theta}{2} \dfrac{q|q|}{\rho} \end{cases}$$

(IsoStat)
$$\begin{cases} \partial_x q = 0 \\ \partial_x \left( p + \dfrac{q^2}{\rho} \right) = -\dfrac{\theta}{2} \dfrac{q|q|}{\rho} \end{cases}$$

(StateEq)
$$p = RT z(p) \rho$$

(AGA)
$$z(p) = 1 + \alpha p$$

(param-IVP)
$$\begin{cases} u_x(x, z) = s(x; u(x, z), z) & x \in (0, L] \\ u(0, z) = u_0 \end{cases}$$

(Initial State)
$$y_a(0) = y_r$$

(Kirchhoff)
$$Aq = q^{\text{out}}$$

(Arc Coupling)
$$y_a(L_a) = f_a(L_a; y_a(0), q_a) \quad \forall a = (u, v) \in \mathcal{A}$$

(Node Coupling)
$$h_u(y_a(x_a(u)), q_a) = h_u(y_b(x_b(u)), q_b)$$

(KirchhoffM)
$$Am = m^{\text{out}}$$

(Real)
$$p_a(L_a) = f_a(L_a; p_a(0), q_a),$$
$$f_a(L; p, q) := F^{-1}\left( F(p, q) - RTq|q| \int_0^L \tfrac{\theta(s)}{2}\, \mathrm{d}s \right),$$
$$F(p, q) := \frac{p}{\alpha} - \frac{1}{\alpha^2} \ln(z(p)) + q^2 RT \ln\left( \frac{z(p)}{p} \right)$$

(Ideal)
$$p_a(L_a) = f_a(L_a; p_a(0), q_a),$$
$$f_a(L; p, q) := c\,|q|\Big[-W_{-1}\Big(-\exp\big(-C + \mathrm{sign}(q)\int_0^L \theta_a(s)\,\mathrm{d}s\big)\Big)\Big]^{1/2},$$
$$C := \tfrac{1}{\eta} + \ln(\eta), \quad \eta := \tfrac{q^2 c^2}{p^2}$$

(Pr)
$$h_u(p, q) = p$$

(MF)
$$h_u(p, q) = \frac{p}{c^2}\left(1 + \frac{q^2 c^2}{p^2}\right) = \rho(1 + \eta)$$

(BI)
$$h_u(p, q) = \log\left(\frac{p}{c^2}\right) + \frac{1}{2}\frac{q^2 c^2}{p^2} = \log(\rho) + \frac{1}{2}\eta$$

(R)
$$f_a(x; p, q) := p - \frac{1}{2}\zeta\frac{q^2 c^2}{p}x, \quad x \in [0, 1]$$

(C)
$$f_a(x; p, q) = p(1 - x) + \big[(RT z_m)^{-1}\tfrac{\kappa - 1}{\kappa}F(Q) + 1\big]^{\frac{\kappa}{\kappa - 1}}px$$

(Stat-Net-Full)
$$Aq = q^{\mathrm{out}},$$
$$p_v = f_a(L_a; p_u, q_a), \quad \forall\, a = (u, v) \in \mathcal{A},$$
$$p^{\mathrm{lb}} \leq p$$

(Stat-Net)
$$p_v^{\mathrm{lb}} \leq p_v(p_r, z) \quad \forall v \in \mathcal{V}$$

(Ex1)
$$\min_{p_r \in \mathbb{R}} \quad \tfrac{1}{2}p_r^2$$
$$\text{s.t.} \quad G(p_r) \geq p^{\mathrm{bd}},$$
$$p_r \in [p_r^{\mathrm{lb}}, p_r^{\mathrm{ub}}]$$

(Wey)
$$p_v = f_a(p_u, q_a) := \sqrt{p_u^2 - \Lambda q_a|q_a|}$$

(Ex1-Det)
$$\min_{p_r \in \mathbb{R}} \quad \tfrac{1}{2}p_r^2$$
$$\text{s.t.} \quad p_v(p_r, \mu) \geq p_v^{\mathrm{lb}}, \quad \forall v \in \mathcal{V} \setminus \{r\},$$
$$p_r \in [p_r^{\mathrm{lb}}, p_r^{\mathrm{ub}}]$$

(Ex2) 
$$\min_{(\beta,\gamma)\in\mathbb{R}^3} \quad \tfrac{1}{2}\left(10\,\beta^2 + \gamma_1^2 + \gamma_2^2\right)$$
$$\text{s.t.} \quad G(\beta,\gamma) \geq p^{\mathrm{bd}},$$
$$\beta \in [1.25, 1.5],$$
$$\gamma \in [1,2]^2$$

(Iso-semi)
$$\begin{cases} \partial_t\rho + \partial_x q = 0, \\ \partial_t q + c^2\partial_x\rho = -\dfrac{\theta}{2}\dfrac{q|q|}{\rho} \end{cases}$$

(WaveEq-$\rho$)
$$\partial_{tt}^2\rho = c^2\partial_{xx}^2\rho$$

(WaveEq-$q$)
$$\partial_{tt}^2 q = c^2\partial_{xx}^2 q$$

(S)
$$\begin{cases} v(0,x) = v_0(x), & v_t(0,x) = v_1(x), \\ v_x(t,0) = \eta v_t(t,0), & v(t,L) = \xi(t), \\ v_{tt}(t,x) = c^2 v_{xx}(t,x) \end{cases}$$

(Initial Condition)
$$v(0,x) = v_0(x), \quad v_t(0,x) = v_1(x)$$

(Boundary Condition)
$$v_x(t,0) = \eta v_t(t,0), \quad v(t,L) = \xi(t)$$

(PDE)
$$v_{tt}(t,x) = c^2 v_{xx}(t,x)$$

(KL-bd)
$$\xi(t) = \sqrt{2T}\sum_{k=1}^{N} a_k \frac{\sin\left(\omega_k\pi\frac{t}{T}\right)}{\omega_k\pi} \quad \text{on } [0,T],$$
$$\omega_k = k - \tfrac{1}{2}$$

(KL-id)
$$v_0(x) = \sqrt{2L}\sum_{k=1}^{N} b_k \frac{\sin\left(\omega_k\pi\frac{L-x}{L}\right)}{\omega_k\pi} \quad \text{on } [0,L],$$
$$\omega_k = k - \tfrac{1}{2}$$

# Bibliography

[1]  Wim van Ackooij. "Chance constrained programming: with applications in energy management". PhD thesis. Ecole Centrale Paris, 2013.

[2]  Wim van Ackooij and René Henrion. "(Sub-)Gradient Formulae for Probability Functions of Random Inequality Systems under Gaussian Distribution". In: *SIAM/ASA Journal on Uncertainty Quantification* 5.1 (2017), pp. 63–87.

[3]  Wim van Ackooij and René Henrion. "Gradient Formulae for Nonlinear Probabilistic Constraints with Gaussian and Gaussian-Like Distributions". In: *SIAM Journal on Optimization* 24.4 (2014), pp. 1864–1889.

[4]  Lukáš Adam and Martin Branda. "Nonlinear Chance Constrained Problems: Optimality Conditions, Regularization and Solvers". In: *Journal of Optimization Theory and Applications* 170.2 (2016), pp. 419–436.

[5]  Dennis Adelhütte, Denis Aßmann, Tatjana Gonzàlez Grandòn, Martin Gugat, Holger Heitsch, Renè Henrion, Frauke Liers, Sabrina Nitsche, Rüdiger Schultz, Michael Stingl, and David Wintergerst. "Joint model of probabilistic-robust (probust) constraints with application to gas network optimization". University of Duisburg-Essen, Weierstrass Institute Berlin, Humboldt University Berlin, Friedrich-Alexander Universität Erlangen-Nürnberg. Preprint. https://opus4.kobv.de/opus4-trr154/frontdoor/index/index/docId/215. 2018 (In Revision).

[6]  Nikolai Sergeevich Bakhvalov. "The existence in the large of a regular solution of a quasilinear hyperbolic system". In: *Zhurnal Vychislitel'noi Matematiki i Matematicheskoi Fiziki* 10.4 (1970), pp. 969–980.

[7]  Robert Béjan. "Minoration de la discrépance d'une suite quelconque sur T". In: *Acta Arithmetica* 41 (1982), pp. 185–202.

[8]   Johann S. Brauchart and Josef Dick. "Quasi-Monte Carlo rules for numerical integration over the unit sphere $\mathbb{S}^2$". In: *Numerische Mathematik* 121.3 (2012), pp. 473–502.

[9]   Alberto Bressan, Sunčica Čanić, Mauro Garavello, Michael Herty, and Benedetto Piccoli. "Flows on networks: recent results and perspectives". In: *EMS Surveys in Mathematical Sciences* 1.1 (2014), pp. 47–111.

[10]  Alberto Bressan, Geng Chen, Qingtian Zhang, and Shengguo Zhu. "No BV bounds for approximate solutions to p-system with general pressure law". In: *Journal of Hyperbolic Differential Equations* 12.04 (2015), pp. 799–816.

[11]  Jens Brouwer, Ingenuin Gasser, and Michael Herty. "Gas Pipeline Models Revisited: Model Hierarchies, Nonisothermal Models, and Simulations of Networks". In: *Multiscale Modeling & Simulation* 9.2 (2011), pp. 601–623.

[12]  Edward Burkard. *Introduction to Ordinary Differential Equations and Some Applications*. 2014.

[13]  John Burkardt and Bennett Fox. *SOBOL The Sobol Quasirandom Sequence*. 2010. URL: https://people.sc.fsu.edu/~jburkardt/m_src/sobol/sobol.html (visited on 06/13/2018).

[14]  Jean-Baptiste Caillau, Max Cerf, Achille Sassi, Emmanuel Trélat, and Hasnaa Zidani. "Solving chance constrained optimal control problems in aerospace via Kernel Density Estimation". IMB - Institut de Mathématiques de Bourgogne [Dijon]; McTAO - Mathematics for Control, Transport and Applications CRISAM - Inria Sophia Antipolis - Méditerranée; Airbus Safran Launchers; OC - Optimisation et commande; CaGE - Control And GEometry; Inria de Paris, LJLL - Laboratoire Jacques-Louis Lions; LJLL - Laboratoire Jacques-Louis Lions. Preprint. https://hal.inria.fr/hal-01507063. 2017.

[15]  Günter Cerbe. *Grundlagen der Gastechnik: Gasbeschaffung - Gasverteilung - Gasverwendung, 7. Auflage*. Hanser, 2008.

[16]  Francis Clarke. *Optimization and Nonsmooth Analysis*. Classics in Applied Mathematics. Philadelphia: Society for Industrial and Applied Mathematics, 1990.

[17]   Rinaldo M. Colombo, Graziano Guerra, Michael Herty, and Veronika Schleper. "Optimal Control in Networks of Pipes and Canals". In: *SIAM Journal on Control and Optimization* 48.3 (2009), pp. 2032–2050.

[18]   Rinaldo M. Colombo and Francesca Marcellini. "Smooth and discontinuous junctions in the p-system". In: *Journal of Mathematical Analysis and Applications* 361.2 (2010), pp. 440–456.

[19]   Robert M. Corless, Gaston H. Gonnet, David E. G. Hare, David J. Jeffrey, and Donald E. Knuth. "On the LambertW function". In: *Advances in Computational Mathematics* 5.1 (1996), pp. 329–359.

[20]   Frank E. Curtis, Andreas Wächter, and Victor M. Zavala. "A Sequential Algorithm for Solving Nonlinear Optimization Problems with Chance Constraints". In: *SIAM Journal on Optimization* 28.1 (2018), pp. 930–958.

[21]   Julio Cézar De Almeida, José Antonio Velásquez, and Renato Barbieri. "A methodology for calculating the natural gas compressibility factor for a distribution network". In: *Petroleum Science and Technology* 32.21 (2014), pp. 2616–2624.

[22]   Josef Dick and Friedrich Pillichshammer. *Digital Nets and Sequences.* Cambridge: Cambridge University Press, 2010.

[23]   Pia Domschke, Benjamin Hiller, Jens Lang, and Caren Tischendorf. *Modellierung von Gasnetzwerken: Eine Übersicht.* Tech. rep. `https://opus4.kobv.de/opus4-trr154/frontdoor/index/index/docId/191`. Technische Universität Darmstadt; Konrad-Zuse-Institut Berlin; Humboldt Universität Berlin, 2017.

[24]   M. Hassan Farshbaf-Shaker, René Henrion, and Dietmar Hömberg. "Properties of Chance Constraints in Infinite Dimensions with an Application to PDE Constrained Optimization". In: *Set-Valued and Variational Analysis* (2017).

[25]   Mauro Garavello and Benedetto Piccoli. "Conservation laws on complex networks". In: *Annales de l'I.H.P. Analyse non linéaire* 26.5 (2009), pp. 1925–1951.

[26]   Alan Genz and Frank Bretz. *Computation of Multivariate Normal and t Probabilities.* Vol. 195. Lecture Notes in Statistics 9. Berlin, Heidelberg: Springer Berlin Heidelberg, 2009, pp. 1689–1699.

[27]    Mike Giles, Frances Y. Kuo, Ian H. Sloan, and Benjamin J. Water-house. "Quasi-Monte Carlo for finance applications". In: *ANZIAM Journal* 50 (2008), pp. 308–323.

[28]    Claudia Gotzes, Holger Heitsch, René Henrion, and Rüdiger Schultz. "On the quantification of nomination feasibility in stationary gas networks with random load". In: *Mathematical Methods of Operations Research* 84.2 (2016), pp. 427–457.

[29]    Martin Gugat. *Optimal Boundary Control and Boundary Stabilization of Hyperbolic Systems*. SpringerBriefs in Electrical and Computer Engineering. Springer International Publishing, 2015.

[30]    Martin Gugat, Falk Hante, Markus Hirsch-Dick, and Günter Leuger-ing. "Stationary states in gas networks". In: *Networks and Heterogeneous Media* 10.2 (2015), pp. 295–320.

[31]    Martin Gugat and Michael Herty. "Existence of classical solutions and feedback stabilization for the flow in gas networks". In: *ESAIM: Control, Optimisation and Calculus of Variations* 17.1 (2011), pp. 28–51.

[32]    Martin Gugat, Michael Herty, Axel Klar, Günter Leugering, and Veronika Schleper. "Well-posedness of Networked Hyperbolic Systems of Balance Laws". In: *Constrained Optimization and Optimal Control for Partial Differential Equations*. Vol. 160. Internat. Ser. Numer. Math. Birkhäuser/Springer Basel AG, Basel, 2012, pp. 123–146.

[33]    Martin Gugat and Günter Leugering. "$L^\infty$-norm minimal control of the wave equation: on the weakness of the bang-bang principle". In: *ESAIM: Control, Optimisation and Calculus of Variations* 14.2 (2008), pp. 254–283.

[34]    Martin Gugat and Günter Leugering. "Global boundary controllability of the Saint-Venant system for sloped canals with friction". In: *Annales de l'Institut Henri Poincare (C) Analyse Non Lineaire* 26.1 (2009), pp. 257–270.

[35]    Martin Gugat, Günter Leugering, Alexander Martin, Martin Schmidt, Mathias Sirvent, and David Wintergerst. "MIP-based instantaneous control of mixed-integer PDE-constrained gas transport problems". In: *Computational Optimization and Applications* 70.1 (2018), pp. 267–294.

[36]  Martin Gugat, Günter Leugering, Alexander Martin, Martin Schmidt, Mathias Sirvent, and David Wintergerst. "Towards simulation based mixed-integer optimization with differential equations". In: *Networks* 72.1 (2018), pp. 60–83.

[37]  Martin Gugat, Günter Leugering, and Ke Wang. "Neumann boundary feedback stabilization for a nonlinear wave equation: A strict $H^2$-Lyapunov function". In: *Mathematical Control and Related Fields* 7.3 (2017), pp. 419–448.

[38]  Martin Gugat, Rüdiger Schultz, and David Wintergerst. "Networks of pipelines for gas with nonconstant compressibility factor: stationary states". In: *Computational and Applied Mathematics* 37.2 (2018), pp. 1066–1097.

[39]  Martin Gugat and Michael Schuster. "Stationary Gas Networks with Compressor Control and Random Loads: Optimization with Probabilistic Constraints". In: *Mathematical Problems in Engineering* 2018 (2018), pp. 1–17.

[40]  Martin Gugat and Sonja Steffensen. "Dynamic boundary control games with networks of strings". In: *ESAIM: Control, Optimisation and Calculus of Variations* (2017), pp. 1–23.

[41]  Martin Gugat and David Wintergerst. "Transient Flow in Gas Networks: Traveling Waves". In: *International Journal of Applied Mathematics and Computer Science* 28.2 (2018), pp. 341–348.

[42]  Edmund Hlawka. "Funktionen von beschränkter Variation in der Theorie der Gleichverteilung". In: *Ann. Mat. Pura Appl.* 54.2 (1961), pp. 325–333.

[43]  Edmund Hlawka. "Über die Diskrepanz mehrdimensionaler Folgen mod. 1". In: *Mathematische Zeitschrift* 77 (1961), pp. 273–284.

[44]  Kari Karhunen. *Über lineare Methoden in der Wahrscheinlichkeitsrechnung.* Annales Academiae scientiarum Fennicae: Mathematica - Physica. Universitat Helsinki, 1947.

[45]  Andrey Kibzun and Stanislav Uryasev. "Differentiability of probability function". In: *Stochastic Analysis and Applications* 16.6 (1998), pp. 1101–1128.

[46]  Gustav Kirchhoff. "Ueber die Auflösung der Gleichungen, auf welche man bei der Untersuchung der linearen Vertheilung galvanischer Ströme geführt wird". In: *Annalen der Physik* 148.12 (1847), pp. 497–508.

[47]   Achim Klenke. *Wahrscheinlichkeitstheorie*. Berlin, Heidelberg: Springer, 2013.

[48]   Martin Knossalla. "Stetige äußere Subdifferentiale und deren Anwendung zur Optimierung lokal Lipschitz-stetiger Funktionen". PhD thesis. Friedrich-Alexander-Universität Erlangen-Nürnberg, 2015.

[49]   Thorsten Koch, Benjamin Hiller, Marc E. Pfetsch, and Lars Schewe, eds. *Evaluating Gas Network Capacities*. Vol. 145. 1956. Philadelphia, PA: Society for Industrial and Applied Mathematics, 2015.

[50]   Jurien F. Koksma. "Een algemeene stelling uit de theorie der gelijkmatige verdeeling modulo 1". In: *Mathematica B (Zutphen)* 11.7-11 (1942), p. 43.

[51]   Karl Kunisch and Daniel Wachsmuth. "On Time Optimal Control of the Wave Equation and Its Numerical Realization as Parametric Optimization Problem". In: *SIAM Journal on Control and Optimization* 51.2 (2013), pp. 1232–1262.

[52]   Frances Y. Kuo, Christoph Schwab, and Ian H. Sloan. "Quasi-Monte Carlo methods for high-dimensional integration: the standard (weighted Hilbert space) setting and beyond". In: *The ANZIAM Journal* 53.01 (2011), pp. 1–37.

[53]   Kazimierz Kuratowski and Czesław Ryll-Nardzewski. "A general theorem on selectors". In: *Bull. Acad. Polon. Sci. Ser. Sci. Math. Astronom. Phys* 13.1 (1965), pp. 397–403.

[54]   Jens Lang and Pascal Mindt. "Entropy-preserving coupling conditions for one-dimensional Euler systems at junctions". In: *Networks and Heterogeneous Media* 13.1 (2018), pp. 177–190.

[55]   Gerhard Larcher, Friedrich Pillichshammer, and Klaus Scheicher. "Weighted discrepancy and high-dimensional numerical integration". In: *BIT Numerical Mathematics* 43 (2003), pp. 123–137.

[56]   Gérard Lebourg. "Generic differentiability of Lipschitzian functions". In: *Transactions of the American Mathematical Society* 256.December (1979), pp. 125–144.

[57]   Christiane Lemieux. *Monte Carlo and Quasi-Monte Carlo Sampling*. Springer Series in Statistics. New York, NY: Springer New York, 2009.

[58]   Christiane Lemieux and Pierre L'Ecuyer. "A Comparison of Monte Carlo, Lattice Rules and Other Low-Discrepancy Point Sets". In: *Recherche* September (1999), pp. 1–12.

[59]   Michel Loève. *Probability theory, vol. ii.* New York, Heidelberg: Springer-Verlag New York-Heidelberg, 1978.

[60]   Ricardo Marques, Christian Bouville, Mickaël Ribardière, Luis Paulo Santos, and Kadi Bouatouch. "Spherical Fibonacci Point Sets for Illumination Integrals". In: *Computer Graphics Forum* 32.8 (2013), pp. 134–143.

[61]   Alexander Martin, Markus Möller, and Susanne Moritz. "Mixed integer models for the optimisation of gas networks in the stationary case". In: *Mathematical Programming, Series B* 105 (2006), pp. 563–582.

[62]   E. Shashi Menon. *Gas Pipeline Hydraulics.* Boca Raton: Crc Press, 2005.

[63]   John Monahan and Alan Genz. "Spherical-Radial Integration Rules for Bayesian Computation". In: *Journal of the American Statistical Association* 92.438 (1997), pp. 664–674.

[64]   Arkadi Nemirovski. "On safe tractable approximations of chance constraints". In: *European Journal of Operational Research* 219.3 (2012), pp. 707–718.

[65]   Harald Niederreiter. *Random Number Generation and Quasi-Monte Carlo Methods.* Philadelphia: Society for Industrial and Applied Mathematics, 1992.

[66]   James M. Ortega and Werner C. Rheinboldt. "Convergence under partial ordering". In: *Iterative Solution of Nonlinear Equations in Several Variables.* Elsevier, 1970, pp. 432–472.

[67]   Art B. Owen. "Multidimensional variation for quasi-Monte Carlo". In: *Contemporary Multivariate Analysis and Design of Experiments* January (2005), pp. 49–74.

[68]   Gunhild A. Reigstad. "Numerical network models and entropy principles for isothermal junction flow". In: *Networks and Heterogeneous Media* 9.1 (2014), pp. 65–95.

[69] Gunhild A. Reigstad, Tore Flåtten, Nils Erland Haugen, and Tor Ytrehus. "Coupling constants and the generalized Riemann problem for isothermal junction flow". In: *Journal of Hyperbolic Differential Equations* 12.01 (2015), pp. 37–59.

[70] Roger Z. Rios-Mercado and Conrado Borraz-Sanchez. "Optimization problems in natural gas transportation systems: A state-of-the-art review". In: *Applied Energy* 147 (2015), pp. 536–555.

[71] Roger Z. Ríos-Mercado, Suming Wu, L. Ridgway Scott, and E. Andrew Boyd. "A Reduction Technique for Natural Gas Transmission Network Optimization Problems". In: *Annals of Operations Research* 117.1-4 (2002), pp. 217–234.

[72] Walter Rudin. *Principles of Mathematical Analysis*. 3rd ed. New York: McGraw-Hill, 1976.

[73] Martin Schmidt, Marc C. Steinbach, and Bernhard M. Willert. *High Detail Stationary Optimization Models for Gas Networks - Part 2: Validation and Results*. Tech. rep. http://www.optimization-online.org/DB_HTML/2014/10/4602.html. Friedrich-Alexander-Universität Erlangen-Nürnberg, Department Mathematik; Leibniz Universität Hannover, Institut für Angewandte Mathematik, 2014.

[74] Alexander Shapiro, Darinka Dentcheva, and Andrzej Ruszczyński. *Lectures on Stochastic Programming*. Philadelphia: Society for Industrial and Applied Mathematics, 2009.

[75] Ian H. Sloan and Henryk Woźniakowski. "When Are Quasi-Monte Carlo Algorithms Efficient for High Dimensional Integrals?" In: *Journal of Complexity* 14.1 (1998), pp. 1–33.

[76] Il'ya Meerovich Sobol'. "On the distribution of points in a cube and the approximate evaluation of integrals". In: *Zhurnal Vychislitel'noi Matematiki i Matematicheskoi Fiziki* 7.4 (1967), pp. 784–802.

[77] Claudia Stangl. "Modelle, Strukturen und Algorithmen für stationäre Flüsse in Gasnetzen". PhD thesis. Universität Duisburg-Essen, 2014.

[78] Kenneth E. Starling, Jeffrey L. Savidge, American Gas Association, et al. "Compressibility factors of natural gas and other related hydrocarbon gases". In: American Gas Association, Operating Section. 1992.

[79] Elias M. Stein and Rami Shakarchi. *Fourier Analysis: An Introduction*. Princeton: Princeton University Press, 2003, p. 326.

[80]  Gerard Teschl. *Ordinary differential equations and dynamical systems.* Vol. 140. Providence, Rhode Island: American Mathematical Society, 2012.

[81]  Stanislav Uryas'ev. "Derivatives of probability functions and integrals over sets given by inequalities". In: *Journal of Computational and Applied Mathematics* 56.1-2 (1994), pp. 197–223.

[82]  Julie Valein and Enrique Zuazua. "Stabilization of the Wave Equation on 1-d Networks". In: *SIAM Journal on Control and Optimization* 48.4 (2009), pp. 2771–2797.

[83]  Darko Veberic. "Having fun with Lambert W(x) function". University of Nova Gorica, Slovenia; IK, Forschungszentrum Karlsruhe, Germany; J. Stefan Institute, Ljubljana, Slovenia. Preprint. http://arxiv.org/abs/1003.1628. 2010.

[84]  Eric W. Weisstein. *Sphere Point Picking.* From MathWorld–A Wolfram Web Resource. 2018. URL: http://mathworld.wolfram.com/SpherePointPicking.html (visited on 06/07/2018).

[85]  David Wintergerst. "Application of Chance Constrained Optimization to Gas Networks". Friedrich-Alexander-Universität Erlangen-Nürnberg. Preprint. https://opus4.kobv.de/opus4-trr154/frontdoor/index/index/docId/158. 2017.

[86]  Stanisław K. Zaremba. "Some applications of multidimensional integration by parts". In: *Annales Polonici Mathematici* 1 (1968), pp. 85–96.

[87]  Victor M. Zavala. "Stochastic optimal control model for natural gas network operations". In: *Computers & Chemical Engineering* 64 (2014), pp. 1–22.