



Technische Universität München

Department of Mathematics



Master's Thesis

# Phase Retrieval: a View from Optimal Transport

Arseniy Tsipenyuk

Supervisor: Prof. Gero Friesecke, Ph. D.

Advisor: Dr. Dominik Jüstel

Submission Date: 22.12.2014

I assure the single handed composition of this master's thesis only supported by declared resources.

Garching,

## Abstract

X-ray crystallography is the main tool for the structural analysis of molecules today. One of its main subtasks is the phase retrieval problem. In this thesis, the phase retrieval problem is discussed in the context of the optimal transport. Phase retrieval is formulated as an energy minimization problem on the quadratic Wasserstein space, and a corresponding gradient flow equation is presented. Under certain conditions, the considered energy functional exhibits signs of local  $\lambda$ -geodesic convexity.

## Zusammenfassung

Die Röntgenanalyse ist heutzutage die zentrale Methode zur Bestimmung der Struktur von Molekülen. Eine ihrer Zentralaufgaben ist das Phasenrekonstruktionsproblem. In dieser Arbeit, Phasenrekonstruktion ist im Kontext des optimalen Transportes vorgestellt. Phasenrekonstruktion ist als ein Energieminimierungsproblem auf dem quadratischen Wassersteinraum formuliert, und die zugehörige Gradientenflussgleichung ist hergeleitet. Unter gewissen Bedingungen, die betrachtete Energie zeigt die Eigenschaft der lokalen  $\lambda$ -geodätischen Konvexität.

# Contents

- Notation** **1**
  
- Introduction** **2**
  
- 1 Phase retrieval** **4**
  - 1.1 X-ray crystallography . . . . . 4
  - 1.2 Phase retrieval problem . . . . . 10
  - 1.3 Some approaches to phase retrieval . . . . . 13
  
- 2 Optimal transport and the quadratic Wasserstein space** **21**
  - 2.1 Classical Monge-Kantorovich problem . . . . . 21
  - 2.2 Time-dependent optimal transport . . . . . 28
  
- 3 Gradient flows on the quadratic Wasserstein space** **31**
  - 3.1 Riemannian structure of the quadratic Wasserstein space . . . . . 31
  - 3.2 Phase retrieval on the quadratic Wasserstein space . . . . . 37
  
- 4 Outlook** **43**
  
- References** **45**

# Acknowledgement

I would like to thank my supervisor Gero Friesecke for the possibility to work on this project, for his advice and for the mathematical insights he shared with me during my studies.

I especially thank Dominik Jüstel for his help, patience, and fruitful discussions during the work on this thesis. Without his guidance and insightful questions this thesis would lose a great deal in its quality.

I would like to thank Horst Osberger for helpful mathematical discussions on the numerical implementation of optimal transportation algorithms.

Finally, I thank my parents for their love, patience and support.

## Notation

$d$  dimension of an image

$x, y, \dots$  variables in object space  $\mathbb{R}^d$

$k, l, \dots$  variables in Fourier space  $\mathbb{R}^d$

$t$  real-valued variable (“time parameter”)

$f$  density of the investigated object

$A$  square root of the measured signal ( $A = |\hat{f}|$ )

$g, g_t$  densities of objects (approximations to  $f$ , functions from  $\mathbb{R}^d$  to  $\mathbb{R}$ )

$\mathcal{P}(\mathbb{R}^d)$  space of measures on  $\mathbb{R}^d$

$\mathcal{P}_{ac}(\mathbb{R}^d)$  space of absolutely continuous measures on  $\mathbb{R}^d$

$\mathcal{P}_w(\mathbb{R}^d)$  the quadratic Wasserstein space on  $\mathbb{R}^d$

$d_w(\cdot, \cdot)$  the quadratic Wasserstein metric on  $\mathbb{R}^d$

$\nu, \mu, \mu_t$  measures, elements of  $\mathcal{P}(\mathbb{R}^d)$  that generalize densities  $f, g, g_t$

$\mathcal{X}$  space that contains approximations  $g, g_t$

Fourier transforms are denoted by a circumflex

$\mathcal{F}()$ ,  $\mathcal{F}^{-1}()$  alternative notations for forward (resp. backward) Fourier transform

$P$  the set of non-negative functions

$M$  the set of the functions with Fourier transform modulus equal to  $A$

$P_P, P_M$  projections onto  $P$  and  $M$ , respectively

$\arg()$  phase of a complex-valued function (i.e.  $\hat{g} = |\hat{g}| \cdot e^{i \cdot \arg(\hat{g})}$ ); by convention,  $\arg(0) = 0$

Variables in object and Fourier spaces are sometimes written as indices, e.g.  $\hat{g}_k = \hat{g}(k)$

## Introduction

The X-ray crystallography is one of the major methods for studying molecular structures. Its methods have been developed for over 100 years, and the relatively recent discovery of quasicrystals has instigated further research in the area. However, the task of phase retrieval connected to the crystallographic image recovery is still far from being well understood.

One of the difficulties of phase retrieval is connected to the fact that the constraints that define its solutions are formulated the coordinate space (non-negativity constraint, but also support constraint, atomicity etc.) and in Fourier space (scattering amplitude constraint), making it difficult to satisfy these conditions simultaneously. The driving idea behind this thesis was to investigate phase retrieval as an energy minimization problem on a space that would naturally embed the non-negativity constraint. Thus, it would be enough if the energy would incorporate only the scattering amplitude condition.

A natural candidate for such a space is the quadratic Wasserstein space that has received a high interest in the last decades. The quadratic Wasserstein space is one of the key spaces in optimal transport theory; it is suited to describe transport of a measure from an initial to final location. This setting appears natural in the context of phase retrieval: it means that some initial random material distribution should be moved until it satisfies the scattering amplitude condition - with other words, until all atoms are in the right places. In optimal transport, the process of moving material around to minimize an energy functional can be written as a gradient flow. The relation between the energy minimization on the quadratic Wasserstein space and the equations corresponding to gradient flows was first discovered in [1]. The task of this thesis is to investigate how the energy functionals arising in phase retrieval are compatible to optimal transport theory. This thesis tries to present some basic concepts of phase retrieval and optimal transport, opting to explain the motivation and the interpretation behind the notions rather than presenting results in their most general form. For example, most results in the optimal transport description are discussed for absolutely continuous measures, although we try to mention the general results whenever it can be done without lengthy discussions. Therefore, this thesis is meant to be rather an introductory work. This choice was made during the initial writing periods and its significant drawback became apparent much later: to profit from the full power of optimal transport and gradient flow theory, it is insufficient to discuss the simplest cases. This is caused by the fact that the optimal transport theory “works especially well” for  $\lambda$ -geodesically convex energy functionals, which is not the case for phase retrieval. Still, the energy functions of phase retrieval appear to have enough structure to discuss their gradient flows on the quadratic Wasserstein space.

This thesis is structured in the following way. In the first section, we motivate the structure of phase retrieval constraints using X-ray crystallography (and geometrical plane wave scattering) as a starting point to formulate the phase retrieval problem. We present phase retrieval from two points of view: as the task of set intersection and as the task of energy minimization, illustrating these with some examples of existing phase retrieval algorithms. We also try to describe the difficulties that instigate considerations of phase retrieval on spaces different from the usual choices (such as  $L^2(\mathbb{R})$ ,  $\mathcal{S}(\mathbb{R})$  etc.). The essential aspects of the first section are based on [2–4].

In the second section, we describe the very basics of the optimal transport, such as optimal

maps, Monge and Kantorovich problems, the Brenier theorem and the time-dependent formulation of the Monge-Kantorovich problem. We try to motivate the reason why the quadratic Wasserstein space is so suited to describe not only the optimal transport, but also phase retrieval, and why the notions are defined one way or another. This section is mainly based on [5–7].

The last section introduces some machinery necessary to describe gradient flows based on [7, 8]. Its second subsection unifies the topics of phase retrieval and optimal transport and comprises the original results of this thesis. The example 3.14 presents an evolution equation that describes displacement of a random density towards a solution of the phase problem. The Proposition 3.17 is an indication of the connection between the phase problem energy and  $\lambda$ -geodesic convexity.



# 1 Phase retrieval

## 1.1 X-ray crystallography

The standard method for studying molecular structures is their reconstruction from X-ray diffraction patterns. The phase retrieval problem appears as one of the tasks during the reconstruction process. Namely, phase retrieval is a problem of finding an algorithm that uses the X-ray diffraction patterns to propose a possible approximated solution of the structure. Although in this thesis we discuss phase retrieval only in the context of X-ray crystallography, this mathematical problem also appears in various areas, such as astronomy, tomography, and others [9].

This subsection is structured in the following way. First, we mention some questions that are closely related to phase retrieval and briefly review some historical milestones of the X-ray crystallography. Second, we present a simplified case of geometric scattering to motivate how the scattered X-rays are connected to the structure of the sample. We use this connection to formulate the phase retrieval problem for structures described by density functions with compact support:  $f \in C_c^1(\mathbb{R}^d, \mathbb{R}_{\geq 0})$ ,  $d \in \mathbb{N}$ . Lastly, we discuss some ways in which phase retrieval can be generalized to other mathematical spaces (such as the space of tempered distributions).

### General remarks on image recovery and X-ray crystallography

The task of phase retrieval arises as a particular question in X-ray crystallography, but it also appears in other areas, all of which can be unified under the notion “image recovery”. The common ground for all these areas is the following setup: some measurement of an object  $f$  is performed, and data  $A$  is recovered. For example, in context of X-ray crystallography, the data  $A$  corresponds to the intensity of the scattered X-rays. In the case of coherent diffraction imaging (CDI), the data  $A$  corresponds to the intensity of the scattered X-rays measured on the two-dimensional screen by a detector far beyond the studied object.

Using the data  $A$  (and some additional knowledge of the measurement), one tries to make certain statements about the source  $f$  that produced  $A$ . There are various particular questions connected to the reconstruction process: how is the data  $A$  connected to the source  $f$ ? Is it possible to recover  $f$  from  $A$ , or, at least, recover an estimation  $g$  that is “close” to  $f$ ? To which degree is it possible to recover  $f$  if the data is contaminated by noise? Is there an additional (a priori) information that can be used during the reconstruction of  $f$ ? The phase problem concerns itself with the following question: assuming that at least one solution exists and the noise can be ignored, how to construct an algorithm that would recover an approximation  $g$  describing the studied structure?

To answer this question, it is necessary to know the details of the measurement process. We shall discuss how these details look in the case of X-ray crystallography. (Many of the notions described below are also relevant for methods related to the X-ray crystallography, such as Coherent Diffraction Imaging [10], Ptychography [11], or others [2]).

X-ray crystallography was pioneered in 1912, when the diffraction images of crystals were obtained in Wilhelm Röntgen’s laboratory in Munich by Max von Laue, Walter Friedrich and Paul Knippig [12] (M. v. Laue received the Nobel Prize in Physics 1914 for this discovery). In the following year, William Henry Bragg and William Lawrence Bragg

used X-ray scattering to resolve the first crystal structures [13]. These discoveries started what turned out to be one of the major methods for studying molecular structures (see, e.g., [14] for a detailed description of the development of crystallography). Among more recent notable contributions in that area one could name the mathematical work on phase retrieval by Herbert Hauptman and Jerome Karle [15] (Nobel Prize for Chemistry 1985) and the discovery of quasicrystals by Dan Shechtman [16] motivated by X-ray diffraction data (Nobel Prize for chemistry 2011).

Let us briefly motivate the reason behind the success of X-ray crystallography. With other words, we shall try to explain why the X-rays are a very suitable choice for the scattering beam, and why the crystalline structures can be studied particularly well.

Assume that we want to determine the molecular structure of an object, that is, to find out the position of its constituents relative to each other. This structure can be thought of as a three-dimensional image with the resolution of about 1 Ångstrom (which is the same as the diameter of the hydrogen atom or the distance between atoms in a water molecule). To observe effects at this resolution, it is necessary to use beams of comparable wavelength (cf. example 1.1). The first known beams of this wavelength were photon beams, namely the X-rays, discovered in 1895 by Wilhelm Röntgen<sup>1</sup>.

Crystals are particularly suited to be studied due to the sharp peaks exhibited by the scattered X-rays. For a long time, crystals were thought of as structures that exhibit certain internal symmetry; this symmetry leads to the appearance of the bright spots (known as “Bragg peaks”) in the diffraction pattern. However, with the discovery of quasicrystals it was shown that there exist structures whose diffraction pattern shows symmetry that is incompatible with symmetric three-dimensional sample lattices. According to the current definition of the International Union of Crystallography (IUCr) [18], a material is a crystal if it has essentially a sharp diffraction pattern (meaning that most of the intensity of the diffraction is concentrated in relatively sharp Bragg peaks, besides the always present diffuse scattering). With other words, the contemporary definition of crystals implies that they are especially suited to be resolved by means of studying their diffraction patterns. The usage of computational algorithms to solve crystalline structures is motivated by the fact that many object of interest possess a relatively complicated inner structure. For example, diffraction data is used as complementary information in resolution of some protein molecules; these molecules often consist of several thousands of atoms and present a challenge to phase retrieval algorithms (examples can be found in [19]).

Let us now motivate the connection between the sample density  $f$  and the scattering intensity registered by a detector far beyond the probe. Strictly speaking, the following example motivates coherent diffraction imaging rather than the crystallographic approach; however, they show some similarities that are important for phase retrieval. It is also important to keep in mind that the following example of geometric scattering is a rather heuristic chain of arguments. The generalized result can be derived for plane wave scat-

---

<sup>1</sup>Of course, photons are not the only option for the incident particles; in transmission electron microscopy (TEM), electrons with the appropriate energy are scattered on the source. Having de Broglie wavelength of the order of several picometers (one hundredth of the X-ray wavelength), electrons provide relatively good resolution. Another choice for the beam particles are neutrons. Due to the different scattering properties, all techniques have their specific details and can be used to complement each other. A beam is characterized by such properties as specimen damage, contrast and resolution for given specimen thickness, beam energy and beam precision. For more details on the applied imaging techniques, see, e.g., [17].

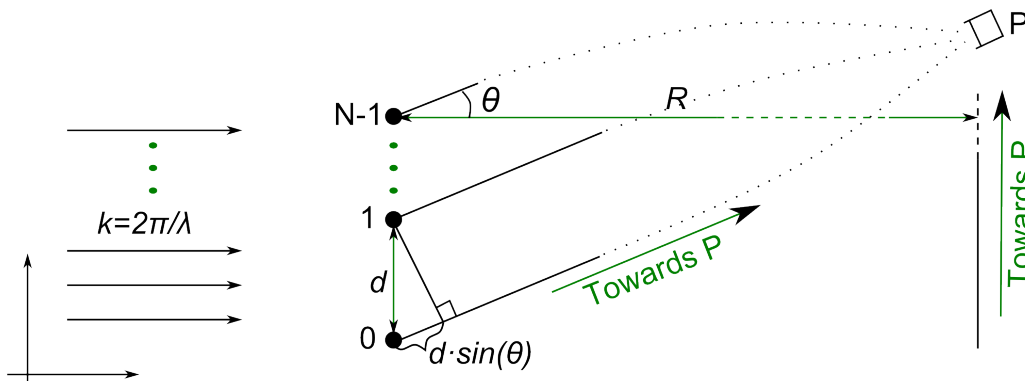
tering using the Maxwell equations, which was shown by Gero Friesecke, Richard James, and Dominik Jüstel in [2, 20, 21] (cf. theorems 1.6, 1.7). The following example is based on [22].

**Example 1.1** (Geometric scattering, field produced by identical oscillators). In this example, we derive how the amplitude  $A_P$  of the scattered wave depends on the angle  $\Delta\phi$  between the incoming wave and the line from source to the detector. Consider the following situation in the two-dimensional plane: a plane wave with the wavevector  $k = (k_0, 0)$  is scattered by  $N$  identical atoms placed on a line with equal distance  $d$  from each other (Fig. 1). The grid is placed perpendicularly to the wave propagation direction: the position of  $n$ -th atom is given by the coordinate  $x_n = (0, dn)$ ,  $n \in \mathbb{N}$ . Each of the atoms acts as an oscillator driven by the incoming wave. Since the plane wave is perpendicular to the grid, all oscillators have the same phase: for atoms  $n_1, n_2 \in \{0, \dots, N-1\}$  the amplitude  $A_{n_1}(x_{n_1}, t)$  of the wave created by the atom  $n_1$  at the point  $x_{n_1}$  is equal to (has the same phase as) the amplitude  $A_{n_2}(x_{n_2}, t)$  created by the atom  $n_2$  at  $x_{n_2}$ :

$$A_{n_1}(x_{n_1}, t) = A_{n_2}(x_{n_2}, t).$$

Assume that a screen is placed in the distance  $R \gg d$  (far-field limit) beyond the detector.

Figure 1: Object density and the resulting X-ray image.



The incident plane wave with wavevector  $k$  excites  $N$  identical atoms. The superposition of the emitted waves is measured at the point  $P$  far beyond the grid.

Consider a point on the screen  $P = (R, x_P)$  which has the angle  $\theta$  to the normal of the grid:  $\tan(\theta) = x_P/R$ . We would like to know the phase difference at the point  $P$  between the waves emitted by two neighbouring atoms - namely, the atoms placed at  $x_0 = (0, 0)$  and  $x_1 = (0, d)$ . For  $R \gg d$ , the distance from the atom placed at  $(0, 0)$  to  $P$  is by  $d \sin \theta$  longer than the distance from the atom with coordinates  $(0, d)$ . The number of waves in the segment with length  $d \sin \theta$  is equal to  $d \sin \theta / \lambda$ , where  $\lambda = 2\pi/k$  is the wavelength of the incoming wave. The phase difference of the two waves at the point  $P$  is given by

$$\Delta\phi = 2\pi \frac{d \sin \theta}{\lambda}. \quad (1.1)$$

Constructive interference is achieved when the phase difference is a multiple of  $2\pi$  (this is exactly when the maxima of the two waves overlap at  $P$ ):

$$2\pi m = \frac{2\pi d \sin \theta}{\lambda}, \quad m \in \mathbb{Z},$$

or

$$\lambda m = d \sin \theta, \quad m \in \mathbb{Z}. \quad (1.2)$$

Now, we can calculate the total amplitude at the point  $P$ . Namely, it is the superposition of amplitudes of all the atoms:

$$A = \sum_{n=0}^{N-1} A_{x_0,P} \cos(n\Delta\phi) = A_{x_0,P} \frac{\sin(N\Delta\phi/2)}{\sin\Delta\phi/2}, \quad (1.3)$$

where  $A_{x_0,P}$  is the amplitude at the point  $P$  of the wave emitted by the atom placed at  $x_0$ . (The amplitude dependencies on the time  $t$  or distance to the detector  $R$  are hidden in the amplitude  $A_{x_0,P}$ , since these quantities impact the absolute values of the measured signal and not the relative intensity or position of its constituents.) The intensity  $I$  of the total wave at the point  $P$  is given (up to a constant) by the following expression<sup>2</sup>:

$$I(\Delta\phi) \propto A^2(\Delta\phi) = A_{x_0,P}^2 \frac{\sin^2(N\Delta\phi/2)}{\sin^2\Delta\phi/2}. \quad (1.4)$$

This function is plotted in Fig. 2. In the interval  $(0, 2\pi)$ , it exhibits sharp peaks at

$$\Delta\phi \in \{0\} \cup \left\{ \frac{(2m+1)\pi}{N}, m \in \mathbb{N} \right\}.$$

Using the Taylor expansion of the intensity in  $\frac{\Delta\phi}{N}$ , one can obtain approximate intensity values; for the first peak,

$$I(0) = A_{x_0,P}^2 \cdot N^2;$$

for the next peak the intensity is

$$I\left(\frac{3\pi}{N}\right) = A_{x_0,P}^2 \left(\frac{2N}{3\pi}\right)^2 + \mathcal{O}\left(\left(\frac{\Delta\phi}{N}\right)^3\right).$$

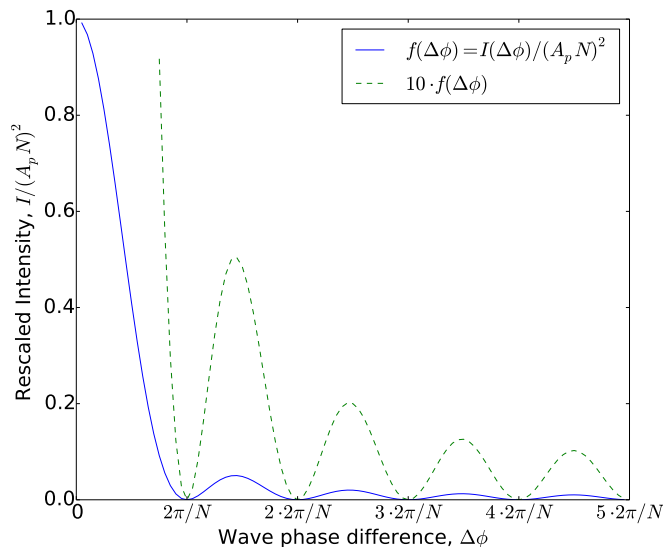
Therefore, the intensity of the second maximum constitutes about  $4/9\pi^2 \approx 0.045$  of the intensity at  $\Delta\phi = 0$  (for crystals described by the gitter with an infinite support, these secondary peaks are not present at all). The intensity  $I$  is  $2\pi$  periodic in  $\Delta\phi$ . Hence, the most notable peaks appear at  $\Delta\phi = 0, 2\pi, \dots$  etc., which (inserted in (1.1)) reaffirms equation (1.2). We see that for  $N$  oscillators — very similar to the case two oscillators — the most notable peaks occur when the angle  $\theta$  to the detector satisfies equation (1.2). These sharp intensity peaks can be experimentally seen as bright spots (also known as “Bragg peaks” [23]). These spots, their relative position and brightness provide the information which makes reconstruction of the sample possible. Depending on the modeling, these peaks can be described as Dirac delta distributions or as sharply peaked density functions.

Note that the trivial peak at  $\theta = 0$  is always present; it is the beam emitted by the source in the forward (and backward) direction of the incoming wave. This beam can be thought of as transmitted (reflected) part of the incoming wave; considered only by itself, it can not help us resolve the structure.

---

<sup>2</sup>The sign  $\propto$  is used with the meaning “proportional to”.

Figure 2: Wave intensity measured in the far field limit.



The intensity of the wave measured at the point  $P$  (scaled so that its maximal value equal to 1;  $N = 10$ ).

**Remark 1.2.** The presented example motivates the following observations. First, it is, indeed, necessary to use rays with a wavelength that is comparable to the scale of the sample. Indeed, for  $d < \lambda$ , equation (1.2) is satisfied only for  $m = 0$ ; there is only one peak in the forward direction (there is also a peak in the backward direction), and the information about the inner structure is not present in the form of multiple peaks.

Second, one can anticipate that diffraction patterns obtained by plane wave scattering are connected with the Fourier transform: in equation (1.4), the sum contains the term  $\cos(n\Delta\phi) = \cos(nkd \sin \theta)$ , and it is taken over the points in space where the atoms are placed. With other words, this sum has the form of a Fourier series that transforms the variable  $n$  (describing the gitter structure) to the variable  $\Delta\phi = kd \sin(\theta)$  (describing the scattered intensity). In the case when the density is described by a continuous function, a Fourier transform appears in place of the sum. The conditions for the points in space, where the scattered intensity shows maxima (points of the constructive interference of the scattered waves) are known as the von Laue condition. This condition can be derived from the Maxwell equations for plane waves scattered on crystallic structures, which was shown in [2, 20, 21]. To present this result, we first recall some definitions and results from the Fourier analysis and crystallography [2, 20, 24, 25].

**Definition 1.3** (Fourier transform and inversion theorem on  $L^1$ ). Let  $f \in L^1(\mathbb{R}^d)$ . Then its Fourier transform  $\hat{f}: \mathbb{R}^d \rightarrow \mathbb{C}$  is defined as<sup>3</sup>

$$\hat{f}(k) = \mathcal{F}(f)(k) := \int_{\mathbb{R}^d} f(x) e^{-ikx} dx, \quad k \in \mathbb{R}^3. \quad (1.5)$$

<sup>3</sup>Here and further, for two vectors  $k, x \in \mathbb{R}^d$ , expressions  $kx$  and  $k \cdot x$  are both used to denote scalar product of the vectors.

When  $\hat{f} \in L^1(\mathbb{R}^d)$ , then  $f$  can be reconstructed pointwise a.e. by

$$f(x) = \mathcal{F}^{-1}(f)(x) := \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} \hat{f}(k) e^{ikx} dk, \text{ for a.e. } x \in \mathbb{R}^d.$$

The transformation  $\hat{f} \mapsto \mathcal{F}^{-1}\hat{f}$  is called the inverse Fourier transform.

**Definition 1.4** (Bravais lattice, reciprocal lattice). A  $d$ -dimensional Bravais lattice  $\mathcal{B} \subset \mathbb{R}^d$  is a set of the form

$$\mathcal{B} = A\mathbb{Z}^d \quad \text{with } A \in GL_d(\mathbb{R}). \quad (1.6)$$

The reciprocal lattice  $\mathcal{B}^\perp \subset \mathbb{R}^d$  is defined as

$$\mathcal{B}^\perp = \{k \in \mathbb{R}^d \mid k \cdot b \in 2\pi\mathbb{Z} \text{ for all } b \in \mathcal{B}\}. \quad (1.7)$$

By  $\delta_{\mathcal{B}}$  we denote the tempered distribution associated to the Bravais lattice:

$$\delta_{\mathcal{B}} = \sum_{b \in \mathcal{B}} \delta_b, \quad (1.8)$$

where  $\delta_b$  is the Dirac delta distribution.

**Definition 1.5** (Fourier transform and convolution on tempered distributions). Recall the space of tempered distributions  $\mathcal{S}'(\mathbb{R}^d)$  that is defined as the dual of the Schwartz space  $\mathcal{S}(\mathbb{R}^d)$ . For a tempered distribution  $T \in \mathcal{S}'(\mathbb{R}^d)$  and  $\xi \in \mathcal{S}(\mathbb{R}^d)$ , the Fourier transform  $\widehat{T} \in \mathcal{S}'(\mathbb{R}^d)$  is defined by

$$\widehat{T}(\xi) := T(\widehat{\xi}).$$

For  $\phi \in \mathcal{S}(\mathbb{R}^d)$ , the convolution  $T * \phi \in \mathcal{S}'(\mathbb{R}^d)$  is defined by

$$(T * \phi)(\xi) := T(\phi_- * \xi),$$

with the notation  $\phi_-(x) = \phi(-x)$ .

**Theorem 1.6** (von Laue condition [2, 20, 21]). Consider a structure with the density  $f = \delta_{\mathcal{B}} * \varphi \in \mathcal{S}'(\mathbb{R}^d)$ , with a Bravais lattice  $\mathcal{B}$  and a Schwartz function  $\varphi \in \mathcal{S}(\mathbb{R}^3)$  modeling the electron density of the atoms. Then, the scattering intensity is proportional to

$$I(x) \propto \left| (P(x^\perp)n) \widehat{f} \left( \frac{\omega}{c} \left( \frac{x}{|x|} - \frac{k_{inc}}{|k_{inc}|} \right) \right) \right|^2 \text{ for a.e. } x \in \mathbb{R}^3 \quad (1.9)$$

where  $\omega$  is the angular wave frequency of the incoming wave,  $c$  is the speed of light,  $x \in \mathbb{R}^3$  is the coordinate of the intensity and  $k_{inc}$  is the wavevector. The complex amplitude of the scattered plane wave is denoted by  $n \in \mathbb{R}^3$  (with  $k \cdot n = 0$ ), and  $P(x^\perp) : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  denotes the projection onto the plane perpendicular to the vector  $x$ .

With other words, this theorem states how the scattering intensity is related to the Fourier transform of the density. Note that the Fourier transform of a Bravais lattice is proportional to the reciprocal Bravais lattice [2, 20, 21], and the constructive interference in equation (1.9) occurs if and only if  $\frac{\omega}{c} \left( \frac{x}{|x|} - \frac{k_{inc}}{|k_{inc}|} \right)$  is an element of  $\mathcal{B}^\perp$ . In particular (which is even more surprising), destructive interference occurs off the reciprocal lattice.

**Theorem 1.7** (Coherent Diffraction Imaging [2,20,21]). Consider the situation described in the previous theorem, denote the object density by  $f_{3D}$ . Assume that the incoming plane wave propagates in the direction  $x_3$ , and a detector is placed in the plane that is perpendicular to the incoming plane wave and has the distance  $R$  to the sample. The two-dimensional measured intensity  $I(k_1, k_2)$  is related to the planar projection  $f(x_1, x_2) = \int_{\mathbb{R}^d} f_{3D}(x_1, x_2, x_3) dx_3$  of the density function. Namely, it holds

$$I(k_1, k_2) = \left| \widehat{f} \left( \frac{\omega k_1}{cR}, \frac{\omega k_2}{cR} \right) \right|^2 + \mathcal{O} \left( \frac{|k|^2}{R^2} \right)$$

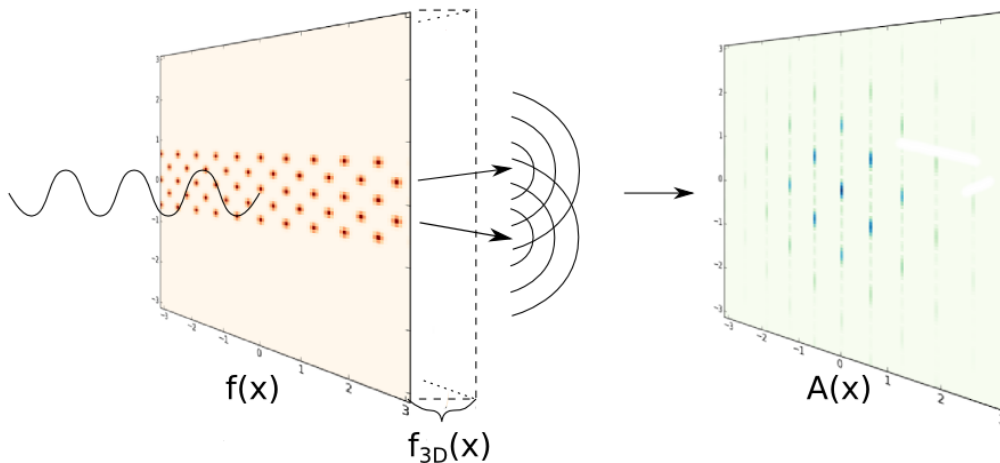
in the far field limit  $R \rightarrow \infty$ .

Let us recapitulate the most important points mentioned so far. The structure of studied objects can be described by a function or by a tempered distribution. The intensity of the scattered waves is related to the squared modulus of the Fourier transform of the sample. (From the algorithmical point of view, the physical constants can be rescaled to one). Note that in the context of this thesis, it shall be more convenient to work with the scattering amplitude, defined as the Fourier transform modulus of the source density:

$$A(k) \propto |\widehat{f}(k)|.$$

With these conventions for the scattering amplitude holds  $A(k) \propto \sqrt{I(k)}$ .

Figure 3: Object density and the resulting X-ray image.



The incident plane wave is scattered on an object with density  $f_{3D}(x)$ . Projecting this function onto the two-dimensional density  $f(x)$ , it is possible to show that in the far field limit for the scattering intensity holds  $I(k) \propto |\widehat{f}(k)|^2$ . Introduce the scattering amplitude  $A(k) = |\widehat{f}(k)|$ , which is proportional to the square root of the scattered intensity.

## 1.2 Phase retrieval problem

In this subsection we fix the mathematical context in which we shall discuss phase retrieval. First, we define phase retrieval for densities described by Schwarz functions. Then,

we describe how the phase problem is usually generalized to the space of measures and to the space of tempered distributions.

**Definition 1.8** (Phase Problem, Phase Retrieval). Given the scattering amplitude  $A \in \mathcal{S}(\mathbb{R}^d, \mathbb{R}_{\geq 0})$ , find a Lebesgue measurable function  $f \in \mathcal{S}(\mathbb{R}^d, \mathbb{R}_{\geq 0})$  such that

$$|\widehat{f}(k)| = A(k) \quad \text{for a.e. } k \in \mathbb{R}^d \quad (1.10)$$

w.r.t. the Lebesgue measure on  $\mathbb{R}^d$ . We call  $f$  the *solution density* to  $A$ . Consider the following questions:

- *Phase Problem*: under which conditions does  $f$  exist? Under which conditions is it unique?
- *Phase Retrieval*: assuming that at least one solution  $f$  exists, how to find it?

In this thesis, we discuss only phase retrieval; some information about the phase problem can be found, e.g., in [2, 9].

Effectively, it is necessary to satisfy two conditions when searching for the function  $f$ :

[M] Modulus constraint, defined by eq. (1.10);

[P] Non-negativity constraint, i.e., the target set of the density  $f$  is  $\mathbb{R}_{\geq 0}$ .

The modulus condition itself is too ambiguous, and even in conjunction with the non-negativity the problem remains rather underconstrained [2, 26]. In practice, one often makes use of some additional information about  $f$  such as known support or the expected number of atoms in the desired image [4].

The ultimate goal of this thesis is to investigate the connection between the phase retrieval and energy minimization on the quadratic Wasserstein space, which shall be defined in the next section. Since quadratic Wasserstein space is a space of measures, it is important to discuss how phase retrieval is generalized to measures, which is done in the remark below. Afterwards, we also mention how to generalize phase retrieval to crystalline and quasicrystalline structures with infinite support.

First, we fix the following

**Notation 1.9** (Measure spaces). Let  $\mathcal{P}(\mathbb{R}^d)$  ( $\mathcal{P}_{ac}(\mathbb{R}^d)$ ) denote the space of Lebesgue measures on  $\mathbb{R}^d$  (absolutely continuous Lebesgue measures on  $\mathbb{R}^d$ , respectively).

**Remark 1.10** (Generalization to measures). The goal of this remark is to generalize phase retrieval to the case when the density of the structure is described by the measure  $\nu \in \mathcal{S}'(\mathbb{R}^d; \mathbb{R}_{\geq 0})$ . To do so, one requires an appropriate generalization of the statement

“(modulus of the Fourier transform of the structure)<sup>2</sup> = measured scattering amplitude”.

For finite measures, the Fourier transform can be defined as the function

$$\hat{\nu}(k) := \int_{\mathbb{R}^d} e^{-ikx} d\nu(x), \quad k \in \mathbb{R}^d,$$



for any  $\nu(k) \in \mathcal{P}(\mathbb{R}^d)$ , since  $\nu(\mathbb{R}^d) = 1$  and  $|e^{-ikx}| = 1$ , and the generalization of the equation (1.10) is straight-forward. For non-finite measures (whose importance for description of crystalline structures can be seen in the definition 1.4), one can proceed as it is done usually in the case of generalizing functions to distributions. Assume that the structure to be resolved is described by the measure  $\nu \in \mathcal{S}'(\mathbb{R}^d; \mathbb{R}_{\geq 0})$  that generalizes the density function  $f$ . Then, it is necessary to search for a distribution  $\hat{\gamma}_\nu: \mathcal{S}(\mathbb{R}^d) \rightarrow \mathbb{C}$  such that

$$\gamma_\nu(\varphi) = \int_{\mathbb{R}^d} \varphi(k) |\hat{f}(k)|^2 dk, \quad \varphi \in \mathcal{S}(\mathbb{R}^d; \mathbb{C}), \quad (1.11)$$

in the regular case (i.e., in the case when  $\nu \in \mathcal{P}_{ac}(\mathbb{R}^d)$  with the density  $f \in L^1(\mathbb{R}^d) \cap L^2(\mathbb{R}^d)$ ). If such an operator is found, the condition (1.10) can be rewritten as

$$\hat{\gamma}_\nu(\varphi) = \int_{\mathbb{R}^d} \varphi(k) A(k)^2 dk, \quad \text{for all } \varphi \in \mathcal{S}(\mathbb{R}^d; \mathbb{C}), \quad (1.12)$$

and the goal of phase retrieval is to find the corresponding  $\gamma_\nu$ . To establish the relationship between  $\hat{\gamma}_\nu$  and  $f$  in the regular case, consider the measure  $\nu \in \mathcal{P}_{ac}(\mathbb{R}^d)$  with the density  $f$ . Use the following notation for the “flipped” measure:  $\tilde{\nu}(\phi) = \nu(\tilde{\varphi})$ , where  $\tilde{\varphi}(x) = \varphi^*(x) := \varphi^*(-x)$  for any (complex-valued) functions  $\varphi \in \mathcal{S}(\mathbb{R}^d)$ . Then, the equality

$$\begin{aligned} \hat{\gamma}_\nu(\varphi) &:= \widehat{\nu * \tilde{\nu}}(\varphi) = \nu * \tilde{\nu}(\hat{\varphi}) = \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \hat{\varphi}(x+y) d\nu(x) d\tilde{\nu}(y) = \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \hat{\varphi}(x+y) f(x) f_-^*(y) dx dy \\ &= \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} \hat{\varphi}(z) f(x) f_-^*(z-x) dx dz = \int_{\mathbb{R}^d} \hat{\varphi}(z) (f * f_-^*)(z) dz \\ &= \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} e^{-ikz} \varphi(k) (f * f_-^*)(z) dk dz = \int_{\mathbb{R}^d} \varphi(k) \widehat{f * f_-^*}(k) dk \\ &= \int_{\mathbb{R}^d} \varphi(k) \hat{f}(k) \widehat{f_-^*}(k) dk = \int_{\mathbb{R}^d} \varphi(k) \hat{f}(k) (\hat{f})^*(k) dk = \int_{\mathbb{R}^d} \varphi(k) |\hat{f}(k)|^2 dk \end{aligned} \quad (1.13)$$

holds for any test functions  $\varphi \in \mathcal{S}(\mathbb{R}^d)$ . Since  $f$  is interpreted as a density (non-negative and real-valued), the use of the complex conjugation  $f^*$  may seem somewhat redundant. However, when we construct various functions  $g$  approximating  $f$ , it can be profitable to allow the most general form of approximations  $g$ . The measure  $\gamma_\nu = \nu * \tilde{\nu}$  is known as *the autocorrelation measure* of  $\nu$ .

**Remark 1.11** (Generalization to crystals and quasicrystals). Canonical crystal structures are usually described as tempered distributions associated to a Bravais lattice with a well-defined Fourier transform given (up to scaling factors) by the reciprocal Bravais lattice (cf. Thm 1.6 or [2, App. C]).

The discovery of quasicrystals [16] encouraged research in diffraction of aperiodic structures (see, e.g., [27] for an informal review). Quasicrystals can be characterized by the following property: their diffraction spectrum contains “Bragg peaks”, but the symmetry of the spectrum cannot occur in three-dimensional periodic structures (the symmetry is “crystallographically forbidden”). The mathematical treatment of quasicrystals was pioneered in [23] and has been substantially developed since that time [26]. A common way to describe a quasicrystal is to consider a so-called Delone set  $\Lambda \subset \mathbb{R}^d$  with the following properties: i) there is a minimum distance  $d$  between any pair of points of  $\Lambda$ ; ii) there

exists  $\delta > 0$  such that any sphere with radius  $\delta$  contains at least one point of  $\Lambda$  (with other words, the points of  $\Lambda$  are neither too close nor too sparse). Then the structure of the crystal is described by the weighted Dirac comb

$$\nu = \sum_{x \in \Lambda} f(x) \delta_x, \quad (1.14)$$

where  $\delta_x$  is the Dirac delta measure and  $f(x) \in \mathbb{C}$  represents the scattering weight of an atom at position  $x$  (these conditions may be relaxed [23]). The autocorrelation measure of  $\nu$  is formally defined as the limit

$$\gamma_{\nu, \text{comb}} = \nu \circledast \tilde{\nu} := \lim_{D \rightarrow \infty} \frac{\nu|_D * \tilde{\nu}|_D}{\text{vol}(B_D)}, \quad (1.15)$$

where  $\nu|_D$  denotes the restriction of the measure  $\nu$  to the open ball  $B_R(0)$ , and  $\tilde{\nu}$  is the “flipped” measure defined in the previous remark. With other words, to generalize measures to lattices, it is necessary to check that the new autocorrelation measure  $\gamma_{\nu, \text{comb}}$  is well-defined for the Dirac comb  $\nu$ . The operator  $\circledast$  (also known as the *Eberlein convolution*) is the volume-averaged analogue of the ordinary convolution. Assuming that the autocorrelation measure  $\gamma_{\nu, \text{comb}}$  exists, its Fourier transform  $\hat{\gamma}_{\nu, \text{comb}}$  is a well-defined (in the distributional sense) positive measure, called the *diffraction measure*. It describes the scattering intensity observed in the experiment, and it has a unique decomposition in the form

$$\hat{\gamma}_{\nu, \text{comb}} = \hat{\gamma}_p + \hat{\gamma}_{ac} + \hat{\gamma}_{sc},$$

where  $\gamma_p$  denotes the point part (consisting of countably many Bragg peaks),  $\gamma_{ac}$  denotes the absolutely continuous part, and  $\gamma_{sc}$  denotes the singular continuous part (comprising everything that is left). This decomposition may show some surprising properties; for example, there are ordered (created by a deterministic rule) Dirac combs that have purely absolutely continuous spectrum (with the exception of the trivial Dirac peak at zero) [26].

### 1.3 Some approaches to phase retrieval

It is possible to describe phase retrieval from several perspectives [3, 4, 28], each having its advances and drawbacks. There are two ways especially important for our further discussion: set intersection formulation and gradient flow formulation of phase retrieval.

This subsection is structured in the following way. First, we formulate phase retrieval as a set intersection problem. We use this notation to write down some of the classical phase retrieval algorithms. Then, we define the notion of an energy functional that shall describe whether an approximation is close to the phase problem solution density or not. Using this definition, it is possible to formulate phase retrieval as an energy minimization problem. In particular, one of the phase retrieval algorithms (the error reduction algorithm) can be rewritten as a method of steepest descent for an appropriate energy functional.

### Phase retrieval as set intersection problem

**Definition 1.12** (Fourier modulus set, non-negative set). Consider the phase problem with the density  $f$  and the scattering amplitude  $A$ . Define the following sets:

$$\begin{aligned} M &= \{g \in \mathcal{S}(\mathbb{R}^d) \mid |\widehat{g}| = A \text{ a.e.}\}, \\ P &= \{g \in \mathcal{S}(\mathbb{R}^d) \mid g \geq 0\}. \end{aligned} \quad (1.16)$$

(The Schwartz space  $\mathcal{S}(\mathbb{R}^d)$  is one of the possible choices; the definitions of  $M$  and  $P$  apply mutatis mutandis for any other space where the Fourier transform is well defined.) With this notation, the phase retrieval problem can be formulated as follows: find an element of  $M \cap P$ . This formulation is very compact, but does not provide any guidelines on how to find such elements. To make it more specific, one could try to find some properties of the defined sets; avoiding a thorough discussion, we would like to point out a specific relation of phase retrieval with convexity: namely,  $M \cap P$  is “almost” an intersection of two convex sets.

It is easy to see that the set  $P$  of non-negative functions is convex. The set  $M$  is itself not convex, but it is a boundary of a convex set: let  $M_{\leq A} = \{g \in L^2(\mathbb{R}^d) \mid |\widehat{g}| \leq A \text{ a.e.}\}$ . The set  $M_{\leq A}$  is, indeed, convex: for two elements  $f, g \in M_{\leq A}$  and a parameter  $t \in (0, 1)$ ,

$$\begin{aligned} |\mathcal{F}((1-t)f(k) + tg(k))| &= |(1-t)\widehat{f}(k) + t\widehat{g}(k)| \leq (1-t)|\widehat{f}(k)| + t|\widehat{g}(k)| \\ &\leq (1-t)A(k) + tA(k) = A(k) \end{aligned}$$

holds for almost all  $k \in \mathbb{R}^d$ .

Using this observation, one could try to approach phase retrieval using methods of convex optimization (the task of finding the intersection of two convex sets). Indeed, there is a correspondence between many algorithms in phase retrieval and convex optimization [28]. The set intersection formulation allows us to write down some phase retrieval algorithms in a brief and elegant way; it also provides some helpful starting points to solve phase retrieval. Intuitively, we can expect that the best way to “land” in the intersection set  $M \cap P$  is to guess some initial density and consecutively apply projections onto sets  $P$  and  $M$ . These projection operators can be written out explicitly: for  $g \in L^2(\mathbb{R}^d)$ ,

$$P_P g(x) = \frac{1}{2}(g(x) + |g(x)|) = \begin{cases} g(x), & \text{if } g(x) \geq 0; \\ 0 & \text{otherwise;} \end{cases} \quad (1.17)$$

$$P_M g(x) = \mathcal{F}^{-1} \left( A \cdot e^{i \cdot \arg(\widehat{g})} \right) = \begin{cases} \mathcal{F}^{-1} \left( A \cdot \frac{\widehat{g}}{|\widehat{g}|} \right) (x), & \text{if } \widehat{g}(x) \neq 0; \\ \mathcal{F}^{-1}(A)(x) & \text{otherwise.} \end{cases} \quad (1.18)$$

The argument function  $\arg()$  returns the phase of a complex function; by convention,  $\arg(0) = 0$ . One can easily check that  $P_P$  and  $P_M$  are, indeed, projections [28].

**Example 1.13.** Using the defined projection operators, we now write down some of the phase retrieval algorithms:

$$\text{(ER)} \quad g_{n+1} = P_P \circ P_M g_n, \quad (1.19)$$

$$\text{(BIO)} \quad g_{n+1} = (1 - \beta P_M + \beta P_P \circ P_M) g_n, \quad (1.20)$$

$$\text{(HIO)} \quad g_{n+1} = (1 - P_P - \beta P_M + (1 + \beta) P_P \circ P_M) g_n, \quad (1.21)$$

with parameter  $\beta \in \mathbb{R}$ ;  $g_n$  is a sequence of approximations with random  $g_0 \in L^2(\mathbb{R}^d)$ . These iteration methods are known as error reduction (also Gerchberg-Saxon), basic input-output and hybrid input-output algorithms, respectively. The error reduction algorithm is one of the simplest phase retrieval algorithms. Each approximation produced by (ER) is non-negative. The drawback of the error reduction algorithm is its poor convergence: after several iterations, the iteratives often stay very close to each other (stagnation behaviour) despite being far from the convergence point (see below for precise definition of “proximity”). Such stagnation may persist for a very long time. The algorithms (BIO) and (HIO) were designed to “disturb” the approximations of (ER) in an attempt to force the iterations out of the stagnation region. The initial density  $g_0$  can be a random non-negative function. See [4, 28] and the references therein for the detailed discussion and other examples of phase retrieval algorithms.

**Example 1.14.** A very promising phase retrieval algorithm is the difference map algorithm [4]. Consider the update step

$$(DM) \quad g_{n+1} = (1 + \beta D)g_n, \quad \beta \in \mathbb{R}, \quad (1.22)$$

where  $\beta \geq 0$  is the step length, and  $D$  is a difference operator of the form

$$D = P_1 \circ \psi_2 - P_2 \circ \psi_1. \quad (1.23)$$

with projections  $P_1 = P_P, P_2 = P_M: \mathcal{S}(\mathbb{R}^d) \rightarrow \mathcal{S}(\mathbb{R}^d)$ . The operators  $\psi_1, \psi_2: \mathcal{S}(\mathbb{R}^d) \rightarrow \mathcal{S}(\mathbb{R}^d)$  are some operators that must be tuned in a certain way. Namely, for appropriately chosen  $\psi_1, \psi_2$ , the fixed points of equation (DM) are attractive. One of the possible choices for the functions  $\psi_i, i \in \{1, 2\}$ , is the following:

$$\psi_i(g) = (1 + \gamma_i)P_i(g) - \gamma_i g, \quad \gamma_i \in \mathbb{R}. \quad (1.24)$$

Then,  $\psi_i(g)$  lies on the line connecting current approximation  $g$  and its projection  $P_i g$ , and

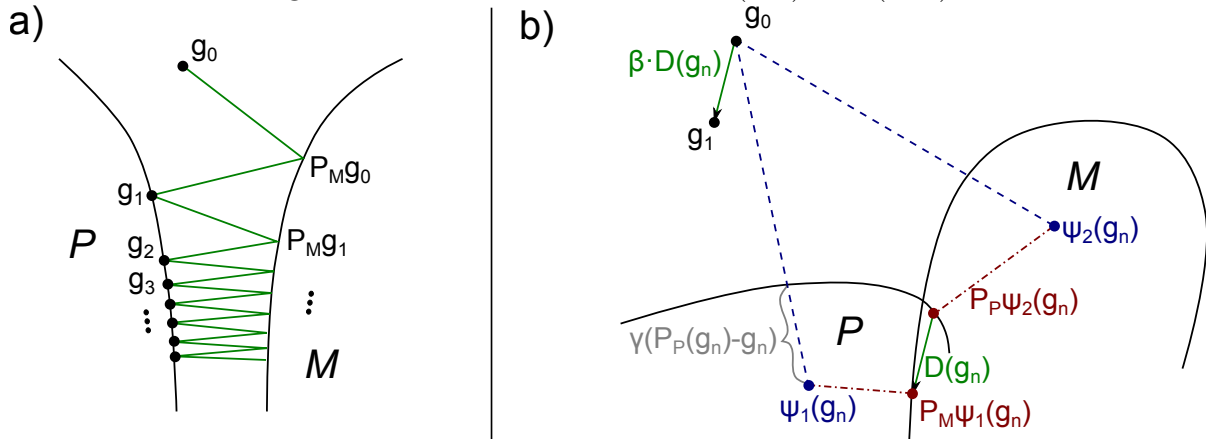
$$(DM_{sp}) \quad g_{n+1} = (1 - \beta\gamma_2 P_P + \beta\gamma_1 P_M + \beta(1 + \gamma_2)P_P \circ P_M - \beta(1 + \gamma_1)P_M \circ P_P)g_n, \quad (1.25)$$

Referring for the detailed discussion of the difference map to [4], we remark that the (HIO) algorithm is a particular case of (DM) for  $P_1, P_2, \psi_1, \psi_2$  as defined above, and  $\gamma_1 = -1, \gamma_2 = \beta^{-1}$ .

## Energy functional

To determine whether an algorithm is any good, we need an estimate  $E$  to tell us how “close” the current approximation  $g_n$  is to the desired density  $f$ . The quantity  $E$  should be an error of some sort, i.e. it should be minimal if the solution  $g_n = f$  is found. With other words, it should incorporate the conditions [M] and [P] in some way. The domain of such an error functional should contain all possible approximations generated by a given algorithm.

Figure 4: Heuristic illustrations of (ER) and (DM).



a) Sketch of the (ER) stagnation behaviour; b) Sketch of the (DM) step.

**Notation 1.15** (Approximation set). Consider a phase retrieval algorithm defined by some operator  $H: \mathcal{S}(\mathbb{R}^d) \rightarrow \mathcal{S}(\mathbb{R}^d)$  with the following prescription:

$$g_{n+1} = H(g_n). \quad (1.26)$$

We call the set  $\mathcal{X} \subseteq \mathcal{S}(\mathbb{R}^d)$  an *approximation set* to this phase retrieval algorithm, if for all  $g \in \mathcal{X}$  holds  $H(g) \in \mathcal{X}$ . With other words, the set  $\mathcal{X}$  is invariant under the action of the operator  $H$ .

This definition allows us to characterize properties of some algorithms. For example, for the error reduction algorithm (1.19), the set  $\mathcal{X} = P$  is an approximation set: at the end of each step, a projection  $P_P$  is applied, so the algorithm produces only non-negative approximations. Since a sensible algorithm should be able to produce the solution density  $f$ , it is important to check the inclusion  $f \in \mathcal{X}$ . (But this is, of course, not strictly necessary: the algorithm may be viable if it produces an approximation  $\tilde{f} \in \mathcal{X}$  that is sufficiently close to  $f \notin \mathcal{X}$ .)

Note that sometimes we shall first define an approximation set  $\mathcal{X}$  and only later opt for an algorithm that shall produce approximations in  $\mathcal{X}$  (this does not cause any contradictions as long as we check that the chosen algorithm does indeed leave  $\mathcal{X}$  invariant).

**Definition 1.16** (Energy functional). Consider phase retrieval with density  $f$ , scattering amplitude  $A$ , and approximation set  $\mathcal{X}$ . We call a functional  $E: \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$  an *energy functional* (or an *error functional*), if

$$P \cap M = \arg \min_{g \in \mathcal{X}} E[g]. \quad (1.27)$$

If only

$$P \cap M \subseteq \arg \min_{g \in \mathcal{X}} E[g] \quad (1.28)$$

is true, we call  $E$  a *relaxed energy functional*.

**Example 1.17.** As already mentioned, the error reduction algorithm keeps the approximates  $g_n$  always non-negative; therefore, an energy functional for this algorithm can

ignore the non-negativity constraint. A possible choice could be

$$E_M[g] = \frac{1}{2} \|(1 - P_M)g\|_2^2$$

for  $\mathcal{X} = \mathcal{S}(\mathbb{R}^d; \mathbb{R}_{\geq 0})$ ;  $\|\cdot\|_2$  denotes the  $L^2$ -norm. The functional  $E_M$  is equal to zero exactly when the modulus constraint is satisfied. We can easily verify that  $E_M$  is an energy functional:

$$\begin{aligned} \subseteq: \quad & f \in P \cap M \implies P_M f = f \implies E_M[f] = 0; \\ \supseteq: \quad & \begin{cases} f \in X \\ E_M[f] = 0 \end{cases} \implies \begin{cases} f \in P \\ P_M f = f \end{cases} \implies f \in P \cap M. \end{aligned}$$

For an algorithm that does not produce strictly non-negative functions, an energy functional can incorporate the constraint **[P]** as well. The energy functional  $E_M$  defined above is one of the simplest energy functionals describing the modulus constraint; it shall appear multiple times during further discussion and deserves a proper notation.

**Definition 1.18** (Modulus energy functional). Consider phase retrieval with density  $f$ , signal  $A$  and approximation set  $\mathcal{X} \subseteq L^2(\mathbb{R}^d)$ . The error functional

$$E_M[g] = \frac{1}{2} \|(1 - P_M)g\|_2^2 \quad (1.29)$$

is called *the modulus energy functional*. Using Plancherel's theorem, one can rewrite  $E_M$  in a more convenient form: for  $g \in \mathcal{X}$ ,

$$\begin{aligned} E_M[g] &= \frac{1}{2} \|g - P_M g\|_2^2 = \frac{1}{2} \int_{\mathbb{R}^d} |\mathcal{F}^{-1}(\widehat{g}(k) - A(k)e^{i \arg(\widehat{g}(k))})|^2 dk \\ &= \frac{1}{2(2\pi)^d} \int_{\mathbb{R}^d} |\widehat{g}(k)e^{i \arg(\widehat{g}(k))} - A(k)e^{i \arg(\widehat{g}(k))}|^2 dk = \frac{1}{2(2\pi)^d} \|\widehat{g} - A\|_2^2. \end{aligned} \quad (1.30)$$

The notion of the energy functional is an important tool to analyze the behaviour of various algorithms. For example, one can show that the modulus energy functional is non-increasing for the error reduction algorithm. This is stated in the following lemma.

**Lemma 1.19.** Consider phase retrieval with  $f$ ,  $A \in \mathbb{L}^2(\mathbb{R}^d)$  and the non-negative approximation set  $\mathcal{X} = P$ . For  $g_0 \in P$ , let the sequence  $(g_n)_{n \in \mathbb{N}}$  be generated by the error reduction algorithm (1.19). Then, for all  $n \in \mathbb{N}$  holds

$$E_M[g_{n+1}] \leq E_M[g_n]. \quad (1.31)$$

*Proof.* For a proof, see [3].

### Energy minimization formulation

Instead of starting with an algorithm and then choosing an energy to check its behaviour, one can first choose the error functional and use it as a starting point to generate a phase retrieval algorithm [3]. Consider an approximation set  $\mathcal{X}$  and an error functional  $E[g]$  that reaches its minimum if some combination of conditions **[P]**, **[M]** is satisfied. Then the

solution of phase problem would correspond to a minimizer of the functional  $E$ . One way to construct an algorithm knowing the error  $E$  is to use the method of steepest descent:

$$g_{n+1}(x) = g_n(x) - h_n(x)\partial_{g_n}E(x), \quad (1.32)$$

where  $h_n$  is the step size and  $\partial_{g_n}E(x)$  is the functional derivative of  $E$  with respect to the approximation  $g_n$ . In the example below, we show that the error reduction algorithm is connected to the method of steepest descent for the modulus energy functional. Unfortunately, the modulus energy functional alone is not enough to generate the error reduction algorithm; in example 1.21, we discuss some other energy functionals that aim to generate the “whole” error reduction algorithm, but fail to do so due to inner properties of phase retrieval. The discussion of the underlying difficulties concludes this subsection and motivates the introduction of the Wasserstein space.

**Example 1.20** (ER as gradient flow). In this example, we calculate one step of the steepest descent for the modulus energy functional. Consider phase retrieval with the modulus energy  $E_M$  and the non-negative approximation set  $\mathcal{X} = P$ . We want to calculate  $\partial_g E_M$  and the resulting form of steepest descent method (1.32). To calculate the gradient of  $E_M$ , recall the functional derivative of the Fourier transform. The Fourier transform functional is defined as

$$g \mapsto \hat{g}_k := \hat{g}(k),$$

for any  $k \in \mathbb{R}^d$ . For any  $g, \delta\varphi \in \mathcal{S}(\mathbb{R}^d)$

$$\widehat{(g + \delta\varphi)}(k) = \int_{\mathbb{R}^d} g(k)e^{-ikx} dk + \int_{\mathbb{R}^d} \delta\varphi e^{-ikx} dx.$$

Therefore, the functional derivative of  $\hat{g}_k$  by  $\delta\varphi$  evaluated at the point  $x \in \mathbb{R}^d$  is given by

$$\partial_{\delta\varphi}\hat{g}_k(x) = e^{-ik \cdot x}, \quad (1.33)$$

Another property of the Fourier transform we want to point out is

$$\hat{g}_k = \widehat{g}_{-k}^*, \quad (1.34)$$

which holds for any real-valued functions  $g \in \mathcal{S}(\mathbb{R}^d; \mathbb{R}_{\geq 0})$ . Here and in the future we shall write arguments of functions as their indices (e.g.,  $\widehat{g}_k = \widehat{g}(k)$ ,  $A_k = A(k)$ ), where it can not cause confusion with the iteration index  $n$ . Calculate the functional derivative of  $E_M$  with respect to  $g \in \mathcal{X}$  at the point  $y \in \mathbb{R}^d$ :

$$\begin{aligned} \partial_g E_M[g](y) &= \partial_g \left( \frac{1}{2(2\pi)^d} \int (|\widehat{g}_k| - A_k)^2 dk \right) (y) \\ &= \frac{1}{2(2\pi)^d} \int 2(|\widehat{g}_k| - A_k) \partial_g \sqrt{\widehat{g}_k^* \widehat{g}_k}(y) dk \\ &= \frac{1}{(2\pi)^d} \int (|\widehat{g}_k| - A_k) \frac{1}{2|\widehat{g}_k|} (\widehat{g}_k^* e^{-ik \cdot y} + \widehat{g}_k e^{ik \cdot y}) dk \\ &= \frac{1}{(2\pi)^d} \int (|\widehat{g}_k| - A_k) \frac{2\widehat{g}_k e^{ik \cdot y}}{2|\widehat{g}_k|} dk \\ &= \mathcal{F}^{-1}(\widehat{g} - A e^{i \arg(\widehat{g})})(y) = g(y) - P_M g(y). \end{aligned} \quad (1.35)$$

This means that method of steepest descent on  $E_M[g] = \frac{1}{2} \|(1 - P_M)g\|_2^2$  with step length  $h_n(y) = 1$  yields the preliminary result

$$\tilde{g}_{n+1}(y) = g_n(y) - h_n(y) \partial_{g_n} E(y) = P_M g_n(y), \quad (1.36)$$

which is “a half” of the error reduction algorithm. In fact, since we have assumed the non-negativity of the approximations  $g_n \in P$ , we must impose the non-negativity condition on the preliminary approximation  $\tilde{g}_{n+1}$ . The simplest way to do so is to set:

$$g_{n+1} = P_P \tilde{g}_{n+1},$$

thus regaining the error reduction algorithm.

Using the gradient flow approach we can investigate the properties of phase retrieval algorithms: all we need is to find the corresponding energy functional. Vice versa, one can also think of various energy functionals and check the corresponding algorithms.

**Example 1.21.** As demonstrated in the example above, the modulus energy functional can not reproduce the whole error reduction algorithm. In this example, we write down the gradients of two functionals with the similar structure for illustrative purposes. As above, one can easily verify that the derivative of the functional  $E_P[g] = \frac{1}{2} \|(1 - P_P)g\|_2^2$  is given by

$$\partial_g E_P[g](y) = g(y) - P_P g(y), \quad (1.37)$$

and the method of steepest descent with step length  $h_n(y) = 1$  yields

$$g_{n+1}(y) = P_P g_n(y). \quad (1.38)$$

As before, we do not gain the complete error reduction algorithm; the modulus projection is missing. It is tempting to try out the natural generalization of this and the energy modulus functional, namely the functional

$$E_{P \circ M}[g] = \frac{1}{2} \|(1 - P_P \circ P_M)g\|_2^2.$$

This functional has the following derivative:

$$\begin{aligned} \partial_g E_{P \circ M}[g](y) &= g(y) - P_P \circ P_M g(y) \\ &+ \beta_k \cdot \mathbb{1}_{\{P_M g(x) \geq 0\}} (g(x) - P_P \circ P_M g(x)) \Big|_{x=-y}, \end{aligned} \quad (1.39)$$

with  $\beta_k = \int dk |A(k)| \frac{\hat{g}(k)^2}{|\hat{g}(k)|^3}$ . The correction terms in the last row arise because  $P_P$  and  $P_M$  are coupled via the Fourier transform. When considering only the modulus energy  $E_M$ , one can evade the Fourier transform with Plancherel’s theorem (as in (1.30)), but this trick does not work when both constraints are coupled together in the energy functional.

The underlying difficulty (that makes it complicated to find an energy functional that would produce the error reduction algorithm) is that our constraints are defined in different spaces - namely, in the object space and in the Fourier space. However, it is possible to separate these constraints. In the Example 1.20, we have considered the energy functional



that depended only on the modulus and the set of approximations that only accepted non-negative functions. Unfortunately, after each gradient step (1.36) we have lost the non-negativity property, and we had to impose it in an explicit way (1.20). As long as we consider the gradient on the space  $L^2(\mathbb{R}^d)$ :

$$E[g + t\delta g] = E[g_0] + t\langle \partial_g E, \delta g \rangle_{L^2} + \mathcal{O}(t^2), \quad (1.40)$$

we can not expect  $\partial_g E$  to keep the non-negativity intact. This issue can be resolved if the gradient is taken over the space of non-negative functions with an appropriate metric and the gradient definition suited to this space. Then, the non-negativity constraint would be embedded in the underlying space, and the modulus constraint would be satisfied by minimizing an appropriate energy functional — for example, the modulus energy functional.

A perfect candidate for such space is the quadratic Wasserstein space: it preserves positivity and has enough structure to define gradient flows. The investigation of the connection between energy minimization problems and the corresponding gradient flows was encouraged by the seminal paper [1] and has been substantially developed in the recent decade [8, 29]. In the next section, we review how the quadratic Wasserstein space is motivated by the optimal transport problems and present some of its basic properties.

## 2 Optimal transport and the quadratic Wasserstein space

This section describes some basics of the optimal transport theory based on [5–8]. The emphasis is often kept on the case of absolutely continuous measures for illustrative purposes. We still try to present general results in cases where it can be done without lengthy digressions, since the restriction to absolutely continuous measures is relatively severe to the results of optimal transport.

This section is organized as follows. In the first subsection, the basic notions of optimal transport are introduced and illustrated with some examples (we discuss push-forward, measure coupling, Monge and Kantorovich problem). Thereafter, we recall the definition of the quadratic Wasserstein space and Brenier’s theorem, sketching some important details of its proof. We conclude the first subsection with an example that compares optimal transport with respect to the  $L^2$  metric to the optimal transport with respect to the quadratic Wasserstein metric. In the second subsection, we formulate the Monge-Kantorovich problem in the time-dependent context and recall the Benamou-Brenier theorem, describing some steps of its proof using [5].

### 2.1 Classical Monge-Kantorovich problem

The topic of optimal transport motivates the introduction of the quadratic Wasserstein space and provides some beneficial ideas for an interpretation of phase retrieval. Speaking informally, optimal transportation investigates how to relocate some material (“a pile of sand”) to the desired location in a cost-efficient way. To formulate this problem in a precise way, one needs mathematical notions that describe the following elements: a material distribution, a method to transport a given distribution to another one, and a cost corresponding to the method of transportation.

It is convenient to describe the initial and the final distribution of a material with probability densities. This way, such properties as positivity and mass conservation follow immediately from the definition of the material distribution. Sometimes it is more convenient to use probability measures instead of probability densities (for example, when the material distribution is discrete).

There is a special notion for a function  $F$  that allows us to transform the given measure  $\mu_0$  into another measure  $\mu_1$ .

**Definition 2.1** (Push-forward). Let  $\mu_0 \in \mathcal{P}(\mathbb{R}^d)$ . Given a measurable function  $F: \mathbb{R}^d \rightarrow \mathbb{R}^d$ , define for all Borel measurable sets  $A \in \mathcal{B}(\mathbb{R}^d)$ :

$$\mu_1(A) = F_{\#}\mu_0(A) := \mu_0(F^{-1}(A)). \quad (2.1)$$

Then,  $\mu_1$  is called a *push-forward measure* of  $\mu_0$  by  $F$ ; we call  $F$  a *transport map* from  $\mu_0$  to  $\mu_1$ .

**Remark 2.2.** It is easy to check that  $\mu_1 \in \mathcal{P}(\mathbb{R}^d)$ . Also, if  $\mu_0, \mu_1 \in \mathcal{P}_{ac}(\mathbb{R}^d)$  (absolutely continuous measures) with densities  $g_0, g_1 \in L^1(\mathbb{R}^d)$  respectively, one can rewrite (2.1) as

$$\int_A g_1(x) dx = \int_{F^{-1}(A)} g_0(y) dy. \quad (2.2)$$

For bijective functions  $F$  with the Jacobian  $DF$ , this equation is equivalent to

$$|\det(DF(y))| \cdot g_1(F(y)) = g_0(y) \quad (2.3)$$

or, in other form,

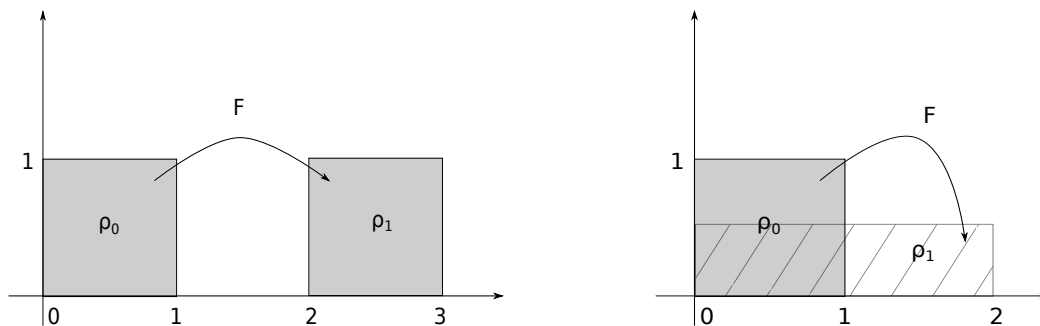
$$g_1(y) = \frac{1}{|\det(DF(y))|} \cdot g_0(F^{-1}(y)) \quad (2.4)$$

$$(2.5)$$

for almost all  $y \in \mathbb{R}^d$ .

**Example 2.3** (Push-forward). It is easy to check that a push-forward does exactly what one could anticipate by its name. For example,  $F(x) = x - c$  translates a measure by  $c \in \mathbb{R}$  to the right;  $F(x) = ax$  dilates it by  $a \in \mathbb{R}_+$  (see Fig. 5). A transport map between

Figure 5: Examples of push-forward measures.



a) Push-forward of  $\mu_0$  with  $\rho_0(x) = \mathbb{1}_{[0;1]}(x)$  by  $F(x) = x - 2$ .

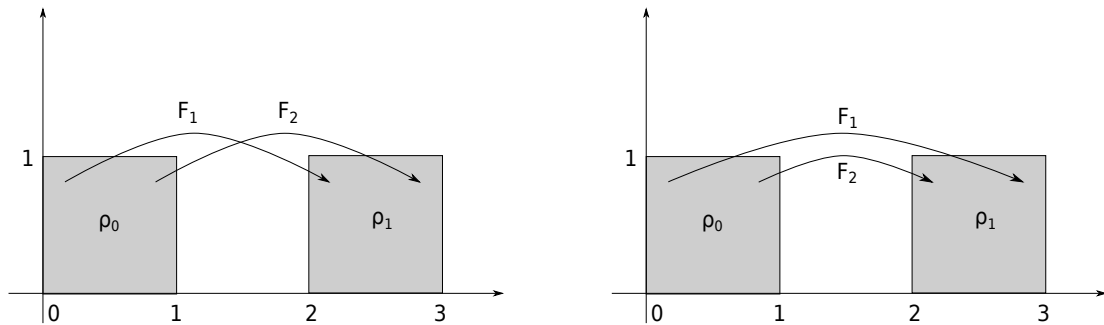
b) Push-forward of  $\mu_0$  with  $\rho_0(x) = \mathbb{1}_{[0;1]}(x)$  by  $F(x) = 2x$ .

two measures is, of course, not unique (see Fig. 6).

**Remark 2.4.** It is important to remember that a push forward is a very specific way to describe material relocation: even the mere fact that  $F$  is a function implies that splitting of the material is not allowed. For example, the situation when all material is concentrated at the point  $x \in \mathbb{R}^d$  is described by the Dirac delta distribution  $\delta_x$ . Any push forward of this measure will still remain a Dirac delta distribution. This is a significant restriction, especially important for discrete distributions. Jumping ahead, this restriction constitutes the difference between the Monge problem and its relaxed generalization, the Kantorovich problem.

The cost of moving one distribution to another can be interpreted as a metric between  $\mu_0$  and  $\mu_1$ . For example, for regular measures  $\mu_0$  and  $\mu_1$  with respective densities  $g_0, g_1 \in L^2(\mathbb{R}^d)$  the  $L_2$ -metric  $\|g_0 - g_1\|_2$  would be a possible, but not a very good choice: this metric is distance-independent. For example, it yields the same values for initial distribution

Figure 6: Non-uniqueness of push-forwards.



a) Push-forward of  $\mu_0$  with  $\rho_0(x) = \mathbb{1}_{[0;1]}(x)$  by  $F_1(x) = x - 2$ .

b) Push-forward of  $\mu_0$  with  $\rho_0(x) = \mathbb{1}_{[0;1]}(x)$  by  $F_2(x) = 3 - x$ .

density  $g_0(x) := \mathbb{1}_{[0;1]}(x)$  and target densities  $g_1(x) := \mathbb{1}_{[1;2]}(x)$  and  $\tilde{g}_1(x) := \mathbb{1}_{[9;10]}(x)$ , respectively:

$$\|g_0 - g_1\|_2^2 = \int_0^1 dx + \int_1^2 dx = 2 = \int_0^1 dx + \int_9^{10} dx = \|g_0 - \tilde{g}_1\|_2^2. \quad (2.6)$$

Intuitively, the second target  $\tilde{g}_1$  is “farther away” from the initial distribution  $g_0$ . A better choice for the transportation distance between two measures should take into account the way the initial measure  $\mu_0$  was transformed to the target measure. In praxis, one knows some cost function  $c(x, y)$ , which describes how much needs to be spent in order to move a unit amount of material from  $x$  to  $y$ . Given the transport map  $F$ , the transportation cost between the points  $x$  and  $y = F(x)$  is  $c(x, F(x))$ .

**Definition 2.5** (The Monge problem). Consider two Borel measures  $\mu_0, \mu_1$  on  $\mathbb{R}^d$  and a cost function  $c: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$ . Find

$$(MP) \quad \inf \int_{\mathbb{R}^d} c(x, F(x)) d\mu_0(x), \quad (2.7)$$

where the infimum is taken over all transport maps  $F$  from  $\mu_0$  to  $\mu_1$ .

As mentioned before, since  $F$  is a function, splitting of the material is forbidden in this formulation: if  $\mu_0$  is a Dirac delta potential and  $\mu_1$  is not, there are no transport maps from  $\mu_0$  to  $\mu_1$ . To relax this restriction of (MP), one needs to allow multiple destination points for a single origin. Consider a Dirac delta distribution  $\mu_0 = \delta_x, x \in \mathbb{R}^d$  that is transported to the measure  $q_x$ . Assume that some material was transported to the set  $B \in \mathcal{B}(\mathbb{R}^d)$ ; the cost for this transportation is given by  $\int_B c(x, y) dq_x(y)$ . The measure  $q_x$  is sometimes known as the transport kernel; it describes how much of the material currently present in  $x$  shall be relocated to  $y$ . If the initial distribution  $\mu_0$  is not point-like, the cost of transport from the set  $A \in \mathcal{B}(\mathbb{R}^d)$  to the set  $B$  is

$$\int_A \int_B c(x, y) d\mu_0(x) dq_x(y).$$

Instead of speaking about two measures  $\mu_0$  and  $q_x$ , it is often simpler to define a general measure  $q: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$  called a *coupling*. Then, the generalized transport information is contained in the coupling  $q$  instead of the transport map  $F$ :  $q(A, B)$  tells us how much material from the set  $A$  is relocated to the destination  $B$ .

**Definition 2.6** (Measure coupling). Let  $\mu_0, \mu_1$  be Borel measures on  $\mathbb{R}^d$ . The measure  $q: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is called a *coupling* (or a transport plan) of  $\mu_0$  and  $\mu_1$ , if and only if  $q(A \times \mathbb{R}^d) = \mu_0(A)$ ,  $q(\mathbb{R}^d, B) = \mu_1(B)$  for all measurable sets  $A, B \subseteq \mathbb{R}^d$ .

**Definition 2.7** (The Kantorovich problem). Let  $\mu_0, \mu_1$  be Borel measures on  $\mathbb{R}^d$ . Given a cost function  $c: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$ , find

$$(KP) \quad \inf \int_{\mathbb{R}^d} c(x, y) dq(x, y), \quad (2.8)$$

where the infimum is taken over all couplings  $q$  of measures  $\mu_0$  and  $\mu_1$ .

**Remark 2.8** (Transport plans generalize transport maps). For each transport map  $F$  there exists a corresponding transport plan given by  $q(A, B) = \mu_0(A) \cdot F_{\#}\mu_0(B)$ .

**Remark 2.9** ((KP) always has a solution). In contrary to the Monge problem, the Kantorovich problem always has a solution: for all initial and final conditions  $\mu_0$  and  $\mu_1$  there exists at least one corresponding transport plan  $q$ , namely  $q(A, B) = \mu_0(A) \cdot \mu_1(B)$  for  $A, B \in \mathcal{B}(\mathbb{R}^d)$ . For regular measures, this plan tells us that “a handful of material” from the initial location must be distributed over the whole support of  $\mu_1$  proportional to the value of its density.

An important example of a transportation cost is the function

$$c(x, y) = |x - y|^p, \quad (2.9)$$

also known as the Wasserstein transportation cost. Two exponents  $p$  received special interest: the original Monge problem was formulated for  $p = 1$  (cf. [5] and references therein); in this report, we shall consider the other case, namely  $p = 2$  (quadratic Wasserstein cost). It is possible to show that equation (2.8) with cost function (2.9) defines a metric on appropriate measure space [30].

**Definition 2.10** (Quadratic Wasserstein space). Let

$$\mathcal{P}_w(\mathbb{R}^d) = \left\{ \mu \in \mathcal{P}(\mathbb{R}^d) \mid \int |x|^2 d\mu(x) < \infty \right\}. \quad (2.10)$$

The quadratic Wasserstein distance between  $\mu_0, \mu_1 \in \mathcal{P}_w(\mathbb{R}^d)$  is defined as

$$d_w(\mu_0, \mu_1) = \inf_q \int_{\mathbb{R}^d \times \mathbb{R}^d} |x - y|^2 dq(x, y); \quad (2.11)$$

the infimum is taken over all couplings  $q$  of  $\mu_0$  and  $\mu_1$ . The pair  $(\mathcal{P}_w(\mathbb{R}^d), d_w)$  is called the *quadratic Wasserstein space*.

**Corollary 2.11.** Quadratic Wasserstein space  $(\mathcal{P}_w(\mathbb{R}^d), d_w)$  is a complete separable metric space. For a proof, see [7, 30].

**Theorem 2.12** (Brenier [31]). Let  $\mu_0 \in \mathcal{P}_w(\mathbb{R}^d) \cap \mathcal{P}_{ac}(\mathbb{R}^d)$ ,  $\mu_1 \in \mathcal{P}_w(\mathbb{R}^d)$ . Then, for the cost function  $c(x, y) = \frac{1}{2}|x - y|^2$ , there exists the unique optimal coupling  $q = (Id, F)_\# \mu_0$  solving the Kantorovich problem (2.8); the function  $F$  solves the corresponding Monge problem (2.7). The optimal map  $F$  is unique and can be written in the form  $F = \nabla \varphi$  for some convex function  $\varphi: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ . For a proof, see [31].

**Remark 2.13.** Brenier's theorem states that the Monge problem is equivalent to the Kantorovich problem on quadratic Wasserstein space for regular initial metric; the solution exists, it is unique and possesses a certain form (gradient of a convex function). One can show that a solution of the Monge problem exists on general Wasserstein distances (see [31] and references therein). The form of the optimal map  $F = \nabla \varphi$  may be anticipated by the following argument: for the one-dimensional problem with cost increasing with distance rapidly enough, one expects  $F$  to be non-decreasing (otherwise,  $F$  would be not optimal); considering  $F = \nabla \varphi$  turns out to be an appropriate generalization to higher dimensions.

Let us briefly indicate which role in the proof of Brenier's theorem is played by the fact that the chosen Wasserstein metric is quadratic. Due to the Kantorovich duality theorem (cf. [6]), one can rewrite the Kantorovich problem (2.8) in the following way:

$$\inf \int c(x, y) dq(x, y) = \sup \int \varphi(x) d\mu_0(x) + \int \psi(y) d\mu_1(y); \quad (2.12)$$

the supremum is taken over all functions  $\varphi \in L^1(\mu_0)$ ,  $\psi \in L^1(\mu_1)$  satisfying

$$\varphi(x) + \psi(y) \leq c(x, y) \quad (2.13)$$

for  $\mu_0$  (resp.  $\mu_1$ )-almost all  $x$  (resp.  $y$ )  $\in \mathbb{R}^d$ . For  $c(x, y) = |x - y|^2$ , the latter inequality can be rewritten as

$$\underbrace{\left(\frac{1}{2}x^2 - \varphi(x)\right)}_{:=\varphi_1(x)} + \underbrace{\left(\frac{1}{2}y^2 - \psi(y)\right)}_{:=\psi_1(y)} \leq -xy,$$

or

$$\varphi_1(x) + \psi_1(y) \leq -xy. \quad (2.14)$$

Since we are trying to maximize the right-hand side of the equation (2.12) under condition (2.13), it is sensible to choose the functions  $\varphi$  and  $\psi$  in such a way that the inequality (2.13) is as sharp as possible (i.e., it is an equality most of the times). In terms of inequality (2.14), this choice corresponds to

$$\psi_1(y) = \inf_{x \in \mathbb{R}^d} (-xy - \varphi_1(x)) = - \sup_{x \in \mathbb{R}^d} (xy + \varphi_1(x));$$

with other words, the functions  $-\psi_1$  and  $-\varphi_1$  should be Legendre transforms of each other! Moreover, it is known (for example, cf. [32]) that for concave, differentiable  $\varphi_1(x)$  with Legendre transform  $\psi_1(y)$  holds

$$\varphi_1(x) + \psi_1(y) = -xy \quad (2.15)$$

with the additional properties

$$\begin{aligned} -\nabla\psi_1(-\nabla\varphi_1(x)) &= x, \\ -\nabla\varphi_1(-\nabla\psi_1(y)) &= y. \end{aligned}$$

If one considers  $x$  as a free variable and  $y = F(x)$  as a function of  $x$  defined through (2.15), then

$$y = F(x) = -\nabla\varphi_1(x).$$

Thus, quadratic Wasserstein cost leads to the specific form of  $F$ . Summing up, the optimal map  $F$  is the gradient of a convex function because the functions minimizing the quadratic Wasserstein cost (in Kantorovich dual formulation) are related through the Legendre transform (establishing connection of the Monge problem to convex functions).

**Example 2.14** (Optimal transport on  $L^2$  and  $\mathcal{P}_w$ ). As mentioned above, one can consider optimal transport with respect to the different metrics. The goal of this example is to compare optimal transport between two simple densities with respect to the  $L_2$ -metric and Wasserstein metric  $d_w$ , and to confirm that optimal transport with respect to the quadratic Wasserstein metric corresponds to the intuitive ideas of the optimal transport. Let  $\rho$  be a symmetric probability density with compact support  $A \subseteq \mathbb{R}$ ;  $\rho$  should be strictly increasing on  $A \cap \mathbb{R}_-$ . Define the following initial and target densities:

$$\begin{aligned} g_0(x) &= \rho\left(x + \frac{c}{2}\right), \\ g_1(x) &= \rho\left(x - \frac{c}{2}\right), \end{aligned}$$

where  $c \in \mathbb{R}$  is a constant such that  $(A + \frac{c}{2})$ ,  $(A - \frac{c}{2})$  and  $0$  are pairwise disjoint sets (see Fig. 7). One could look for the density  $g_{1/2, \mathcal{X}}$  describing the situation when the transport between  $g_0$  and  $g_1$  is “half way complete” with respect to the metric space  $(\mathcal{X}, d) \in \{(L^2, \|\cdot\|_2), (\mathcal{P}_w, d_w(\cdot))\}$ . In this context, “half way complete” means:

- The density  $g_{1/2, \mathcal{X}}$  is equidistant between  $g_0$  and  $g_1$ :

$$g_{1/2, \mathcal{X}} \in \mathcal{X}_{eq} := \{g \in \mathcal{X} \mid d(g_0, g_{1/2, \mathcal{X}}) = d(g_1, g_{1/2, \mathcal{X}})\}. \quad (2.16)$$

- Also,  $g$  should be “on the shortest way” with respect to  $d$ , i.e.,

$$d(g_0, g_{1/2, \mathcal{X}}) = \min_{g \in \mathcal{X}_{eq}} d(g_0, g). \quad (2.17)$$

One can show that

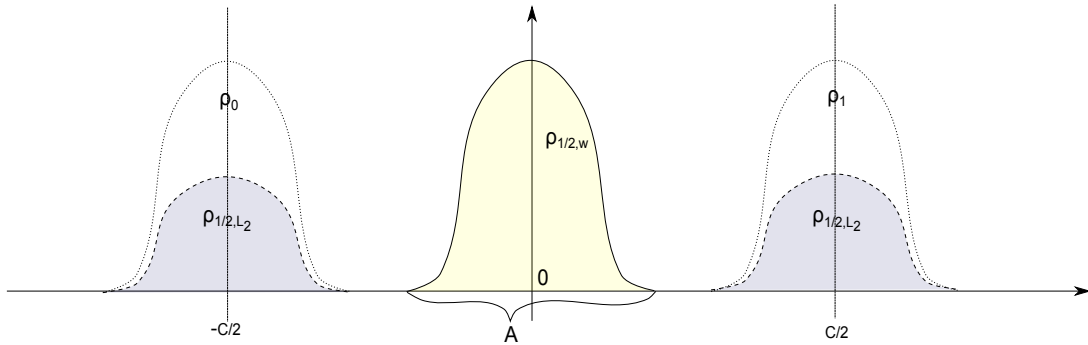
$$g_{1/2, L^2} = \frac{1}{2}(g_0 + g_1); \quad (2.18)$$

$$g_{1/2, \mathcal{P}_w} = \rho(x) = g_0\left(x + \frac{1}{2}\right); \quad (2.19)$$

the Wasserstein transportation map is  $F_{1/2}(x) = x - \frac{1}{2}$  (see Fig. 7). There are several possible interpretations of equations (2.18). First, one can see that transporting “a handful of material” with respect to the Wasserstein metric means translation by a certain

distance, while transport with respect to the  $L_2$ -metric means that material is distributed among the whole support of  $g_1$  proportional to the value of  $g_1$ . A second observation is that it makes sense to consider optimal transportation in a time-dependent context: then, the  $L_2$ -transport would mean that material little by little disappears at the initial location and reappears at destination; the  $d_w$ -transport would mean that the initial density is “gradually moved” to its destination.

Figure 7: Interpretation of Example 2.14.



Intermediary states  $g_{1/2,w}$ ,  $g_{1/2,L_2}$  for initial and final densities  $g_0$ ,  $g_1$ .

We separate the proof into two claims. *Claim 1:* for metric  $d(\rho_0, \rho_1) = \|\rho_1 - \rho_0\|_2$  the intermediary state is uniquely defined by conditions (2.16), (2.17) and is given by  $\rho_{1/2,L_2} = \frac{1}{2}(\rho_0 + \rho_1)$ .

*Proof of claim 1.* Since  $\text{supp } \rho_0 \cap \text{supp } \rho_1 = \emptyset$ , for  $\rho \in B$

$$\begin{aligned} 2\|\rho_0\|_2^2 &= \|\rho_0\|_2^2 + \|\rho_1\|_2^2 = \|\rho_1 - \rho_0\|_2^2 \leq \|\rho_1 - \rho\|_2^2 + \|\rho - \rho_0\|_2^2 = \|\rho - \rho_0\|_2^2 \\ &\Rightarrow \|\rho_0 - \rho\|_2 \geq 1. \end{aligned} \quad (2.20)$$

Since for a small variation  $\delta\rho$

$$\|\rho_0 - (\rho + \delta\rho)\|_2^2 = \|\rho_0 - \rho\|_2^2 - 2 \int (\rho_0(x) - \rho(x))\delta\rho(x) dx + \int (\delta\rho(x))^2 dx, \quad (2.21)$$

the second variation of  $\|\rho_0 - \rho\|_2^2$  w.r.t.  $\rho$  is equal to 1 and positive for all densities  $\rho \in B$ . Hence, if minimizer  $\rho$  exists, it is unique. It is easy to check that  $\|\rho_{1/2,L_2}\|_2 = 1$ ; hence, it is a minimizer and is unique.

*Claim 2:* for metric  $d(\rho_0, \rho_1) = d_w(\rho_0, \rho_1)$  the intermediary state is uniquely defined by conditions (2.16), (2.17) and is given by  $\rho_{1/2,w} = g(x) = \rho_0(x + \frac{c}{2})$  with transportation map  $F_{1/2}(x) = x - \frac{c}{2}$ .

*Proof of claim 2.* We use without proof the fact that the optimal map from  $\rho_0$  to  $\rho_1$  is given by  $F = x - c$  (this fact follows from Brenier’s theorem [5]). Then, the Wasserstein



distance between  $\rho_0$  and  $\rho_1$  is equal to  $d_w(\rho_0, \rho_1) = c$ . One can proceed as in the previous case:

$$d_w(\rho_0, \rho_1) \leq 2d_w(\rho_0, \rho_{1/2}). \quad (2.22)$$

The second variation of  $d_w(\rho_0, \rho_{1/2})$  w.r.t.  $F_{1/2}$  is 1, so an optimal transport map is unique. The claim then follows as in  $L_2$ -case: it is easy to check that  $d_w(\rho_0, \rho_{1/2}) = c/2$ ; hence,  $\rho_{1/2}$  is a minimizer and it is unique.  $\square$

## 2.2 Time-dependent optimal transport

The previous example motivates a discussion of the optimal transport in time-dependent context this approach was pioneered in [5]). In this subsection, we discuss the time-dependent Monge-Kantorovich problem and the Benamou-Brenier theorem that establishes the relationship between the time-dependent Monge-Kantorovich problem and the original Kantorovich problem.

An appropriate setting includes an absolutely continuous curve  $t \mapsto \mu_t \in \mathcal{P}_w(\mathbb{R}^d)$  describing the material distribution at the time  $t \in [0; 1]$ , and a (smooth enough) velocity vector field  $v(t, x) = v_t(x) \in \mathbb{R}^d$ . The curve  $\mu_t$  should satisfy boundary conditions

$$\begin{aligned} \mu_0 &= \tilde{\mu}_0; \\ \mu_1 &= \tilde{\mu}_1 \end{aligned} \quad (2.23)$$

for some initial measures  $\tilde{\mu}_0, \tilde{\mu}_1 \in \mathcal{P}_w(\mathbb{R}^d)$ . The pair  $(\mu_t, v_t)$  should also fulfill the continuity equation

$$\partial_t \mu_t + \nabla \mu_t \cdot v = 0. \quad (2.24)$$

With other words, the underlying idea is to treat the material as an incompressible pressureless fluid, that moves from the initial distribution  $\mu_0$  to the target distribution  $\mu_1$ . The continuity equation then states that the material is not lost or gained during the relocation. The continuity equation above should be satisfied in the sense of distributions, i.e., for all  $u \in C_c^\infty([0; 1] \times \mathbb{R}^d)$  should hold

$$\int_0^1 \int \left( \frac{d}{dt} u(t, x) + \nabla_x u(t, x) \cdot v(x) \right) d\mu_t(x) dt = 0.$$

One can easily verify that for regular measures  $\mu_t$ ,  $t \in [0; 1]$  with corresponding densities  $g_t$  the equation above is equivalent to

$$\partial_t g_t + \nabla g_t \cdot v = 0. \quad (2.25)$$

Further, let the velocity field  $v$  be divergence-free (with other words, the material is treated as perfect incompressible fluid):

$$\nabla \cdot v = 0. \quad (2.26)$$

**Definition 2.15** (Time-dependent Monge-Kantorovich problem on  $\mathcal{P}_w(\mathbb{R}^d)$ ). Let  $\tilde{\mu}_0, \tilde{\mu}_1 \in \mathcal{P}_w(\mathbb{R}^d)$ . Let

$$\begin{aligned} \mathcal{X} = \left\{ (\mu, v) \mid \mu: [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}, \right. \\ (t, x) \mapsto \mu(t, x) = \mu_t(x), \\ t \mapsto \mu_t \in \mathcal{P}_w(\mathbb{R}^d) \text{ absolutely continuous;} \\ v: [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d \\ (t, x) \mapsto v(t, x) = v_t(x), \\ v_t \in L^2(\mu_t); \\ \left. \mu, v \text{ satisfy equations (2.23), (2.24), (2.26)} \right\}. \end{aligned} \quad (2.27)$$

Find the pair  $(\mu, v)$  that minimizes

$$(tMKP) \quad \int_0^1 \int_{\mathbb{R}^d} |v(t, x)|^2 d\mu_t(x) dt. \quad (2.28)$$

The next theorem shows how this problem is connected to the original Kantorovich problem.

**Theorem 2.16** (Benamou-Brenier). Let  $\tilde{\mu}_0, \tilde{\mu}_1 \in \mathcal{P}_w(\mathbb{R}^d)$ ,  $\mathcal{X}$  as described in the definition above. Then,

$$d_w^2(\mu_0, \mu_1) = \inf_{\mathcal{X}} \int_0^1 \int_{\mathbb{R}^d} |v(t, x)|^2 d\mu_t(x) dt, \quad (2.29)$$

which is sometimes referred to as the Benamou-Brenier formula for the Wasserstein distance.

**Remark 2.17.** This theorem was first shown in [5]; its generalized proof can be found, e.g., in [7].

**Remark 2.18** (Relation between the pair  $(\mu, v)$  and an optimal map  $F$ ). To determine the connection between the Monge problem and its time-dependent formulation, let us recall some arguments that were used in [5] to prove the Benamou-Brenier formula. Consider the situation described above; for simplicity, assume that measures  $\tilde{\mu}_0, \tilde{\mu}_1, \mu_t \in \mathcal{P}_w(\mathbb{R}^d)$  are regular with respective densities  $\tilde{g}_0, \tilde{g}_1, g_t \in L^1(\mathbb{R}^d)$ ; denote the optimal map from  $\mu_0$  to  $\mu_1$  with  $F = \nabla\phi$ . Introduce Lagrangian coordinates  $X_t(x) = X(t, x)$  defined by

$$\begin{aligned} X(0, x) &= x; \\ \partial_t X(t, x) &= v(t, X(t, x)). \end{aligned} \quad (2.30)$$

One can think about  $X$  as the coordinate of the infinitesimal amount of the transported material, or, alternatively, as the characteristic curve of the equation (2.24). During the transportation, the material density along the coordinate  $X$  remains constant:

$$\frac{d}{dt} g(t, X_t(x)) = \partial_t g(t, X_t(x)) + \nabla g(t, X_t(x)) \cdot \partial_t X_t(x) = -g(t, X_t(x)) \cdot \nabla v(t, x) = 0; \quad (2.31)$$

we used the continuity equation (2.25) and the assumption that  $v$  is divergence-free. Hence,

$$g(t, X(t, x)) = g(0, X(0, x)) = g(0, x) = g_0(x). \quad (2.32)$$

According to [5], we can rewrite this condition in terms of test functions:

$$\begin{aligned} \int_0^1 \int_{y \in \mathbb{R}^d} u(t, y) g(t, y) dy dt &= \int_0^1 \int_{z \in \mathbb{R}^d} u(t, z) g_0(X_t^{-1}(z)) dz dt \Rightarrow \\ \int_0^1 \int_{y \in \mathbb{R}^d} u(t, y) g(t, y) dy dt &= \int_0^1 \int_{y \in \mathbb{R}^d} u(t, X_t(y)) g_0(y) dz dt \end{aligned}$$

(changing to  $y = X_t^{-1}(z)$  on the r.h.s.). Using this equation, it is possible to show after some computations [5] that

$$\int_0^1 \int_{\mathbb{R}^d} |v(t, x)|^2 g(t, x) dx dt \geq \int_{\mathbb{R}^d} |\nabla \varphi(x) - x|^2 g_0(x) dx.$$

This inequality becomes an equality, if we choose

$$v_t(x) = \nabla \varphi(x) - x. \quad (2.33)$$

From the definition of Lagrangian coordinates (2.30), recover

$$X(t, x) = x + t(\nabla \varphi(x) - x). \quad (2.34)$$

Inserting this result into equation (2.32), obtain

$$g(t, x + tv_0(x)) = g_0(x), \quad (2.35)$$

or, equivalently (cf. eqs. (2.3), (2.4)),

$$g(t, x) = g_0\left((Id + tv_0)^{-1}(x)\right) \quad (2.36)$$

Note that  $F(x) = X(1, x)$  equations (2.2) and (2.34) coincide. Therefore, we define the transport map  $F_t$  corresponding to  $\mu_t$  via

$$F_t = X(t, \cdot) = Id + tv_0. \quad (2.37)$$

By the Benamou-Brenier formula, the map  $F_t$  is the optimal transport map from  $\mu_0$  to  $\mu_t$ . Let us sum up some results from this remark. If the initial distribution  $\mu_0$  is regular, there is a particular solution of the continuity equation (defined by eqs. (2.33), (2.35)) with constant velocity  $v_t(x) = v_0(x)$ . This is one of the solutions of the time-dependent Monge-Kantorovich problem; this solution is the unique solution for which the velocity is constant (follows from (2.33) and the fact that the optimal map  $F$  obtained by Brenier's theorem is unique). We call the solution  $(\mu, v)$  the constant speed solution of the time-dependent Monge-Kantorovich problem (tMKP).

### 3 Gradient flows on the quadratic Wasserstein space

In the previous section we have established some basic notions to work in the quadratic Wasserstein space. Our goal is to apply the steepest descent method to phase retrieval on the quadratic Wasserstein space. To do so, it is necessary to introduce the corresponding notion for a gradient. This, in turn, requires an appropriate definition of the tangent space to the Wasserstein space.

This section is structured in the following way. The aim of the first subsection is to define a gradient on a Wasserstein space. We start by discussing the form of certain geodesics (namely, constant speed geodesics) on Wasserstein space. The investigation of their behaviour will be used to motivate the definition of the tangent space to the quadratic Wasserstein space. After introducing the tangent space, we shall write down the definition of the gradient, and discuss a practical way to calculate it. The results describing geodesics and tangent spaces are based on [7]. The definition of the gradient is a simplified version of the gradient flow notions discussed in [6, 8].

In the second subsection, we shall apply the definition of the gradient to the energy modulus functional of phase retrieval (example 3.14) and discuss the corresponding results. We also mention the connection between the modulus energy functional and the geodesic convexity. This section comprises the original results of this thesis.

#### 3.1 Riemannian structure of the quadratic Wasserstein space

The quadratic Wasserstein space has many similarities to an infinite-dimensional Riemannian manifold; in particular, it is possible to define a gradient of an appropriate functional. Let us briefly review of some of these definitions following [7]. Often we shall use the assumption  $\mu_0 \in P_{ac}(\mathbb{R}^d) \cap P_w(\mathbb{R}^d)$ ; however, most results are also valid for general measures  $\mu_0 \in \mathcal{P}_w(\mathbb{R}^d)$  and can be generalized for quadratic Wasserstein spaces over manifolds [7, 8].

**Definition 3.1.** Let  $\mu_0, \mu_1 \in \mathcal{P}_w(\mathbb{R}^d)$ , let  $\mu_0$  be regular; let  $\nabla\varphi$  be the optimal map from  $\mu_0$  to  $\mu_1$ . The interpolation curve  $\mu: [0; 1] \rightarrow \mathcal{P}_w(\mathbb{R}^d)$  is defined as

$$t \mapsto \mu_t = \left( (1-t)Id + t\nabla\varphi \right)_{\#} \mu_0. \quad (3.1)$$

**Proposition 3.2.** A curve  $\gamma: [0; 1] \rightarrow \mathcal{P}_w(\mathbb{R}^d)$  from  $\mu_0$  to  $\mu_1$  is a constant speed geodesic if and only if it is an interpolation curve from  $\mu_0$  to  $\mu_1$ .

*Proof.* We demonstrate only the direction “geodesic  $\Leftarrow$  interpolation curve”; for the other direction, see [7].

“ $\Leftarrow$ ”. By definition of the Wasserstein metric  $d_w$ , for any measure  $\mu \in \mathcal{P}_w(\mathbb{R}^d)$  and any functions  $T, S \in L^2(\mu, \mathbb{R}^d)$  holds:

$$d_w(T_{\#}\mu, S_{\#}\mu) \leq \|T - S\|_{L^2(\mu, \mathbb{R}^d)}. \quad (3.2)$$

Let  $\mu_t$  be an interpolation curve from  $\mu_0$  to  $\mu_1$  with the optimal map  $F$ ; by remark 2.18, the optimal map from  $\mu_0$  to  $\mu_t$  is given by  $F_t(x) = x + t(F(x) - x)$  For every  $0 \leq s < t \leq 1$ ,

$$\begin{aligned} d_w^2(\mu_s, \mu_t) &\leq \|F(s) - F(t)\|_{L^2_{\mu_0}}^2 \\ &= \int |(1-s)x + sF(x) - (1-t)x + tF(x)|^2 d\mu_0(x) = (t-s)^2 d_w^2(\mu_0, \mu_1). \end{aligned} \quad (3.3)$$

Assume that this inequality were strict. Then, by the triangle inequality follows:

$$\begin{aligned} d_w(\mu_0, \mu_1) &\leq d_w(\mu_0, \mu_s) + d_w(\mu_s, \mu_t) + d_w(\mu_t, \mu_1) \\ &< s d_w(\mu_0, \mu_1) + (t - s) d_w(\mu_0, \mu_1) + (1 - t) d_w(\mu_0, \mu_1) = d_w(\mu_0, \mu_1). \end{aligned}$$

The claim follows by contradiction.

**Corollary 3.3.** A constant speed geodesic  $\mu_t$  from regular measure  $\mu_0$  to  $\mu_1$  is given by (3.1) and is unique.

**Reminder 3.4** (Some notions on Riemannian manifolds). Before introducing the tangent space for quadratic Wasserstein space, let us recall how certain notions are defined on Riemannian manifolds. Let  $(M, g)$  be a finite-dimensional Riemannian manifold, let  $g$  be the inner product on the tangent bundle  $TM$ . Further, let  $\gamma: [0; 1] \rightarrow M$  be a differentiable curve in  $M$  with  $\gamma(0) = p \in M$ ,  $\gamma'(0) = v \in T_p M$ , let  $\varphi \in C^\infty(\mathbb{R}^d)$ . The *directional derivative* of the function  $\varphi$  at point  $p$  in direction  $v$  is defined as

$$\partial_v \varphi(p) = \left. \frac{d}{dt} \right|_{t=0} \varphi(\gamma(t)).$$

The *gradient* of the function  $\varphi$  at  $p$  is defined through the following equation:

$$g_p((\nabla \varphi)_p, v_p) = \partial_v \varphi(p)$$

for all vector fields  $v \in T_p M$ . The *induced metric* on  $M$  is defined by

$$d(p, q) = \inf \int_0^1 |\gamma'(t)| dt,$$

where the infimum is taken over all differentiable curves  $\gamma \in M$  with  $\gamma(0) = p$ ,  $\gamma(1) = q$ . The *exponential map*  $\exp_p: T_p M \rightarrow M$  is defined in the neighbourhood of the point  $p$  by the equation

$$\exp_p(tv) = \gamma_v(t),$$

where  $\gamma \in C^1([0, 1]; M)$  is the unique geodesic with initial conditions  $\gamma(0) = p$ ,  $\gamma'(0) = v$ . If the function  $\exp_p$  defines a local diffeomorphism between the tangent space  $T_p M$  and the manifold  $M$ , the radius of the largest ball in  $T_p M$  diffeomorphically mapped onto  $M$  is called *the injectivity radius* of the exponential map at the point  $p$ .

To consider the Riemannian structure of the quadratic Wasserstein spaces, it is necessary to find an appropriate candidate for the tangent space  $T\mathcal{P}_w(\mathbb{R}^d)$ . Let  $\tilde{\mu}_0, \tilde{\mu}_1 \in \mathcal{P}_w(\mathbb{R}^d)$ , assume that  $\tilde{\mu}_0$  is regular. Let  $F$  be the optimal map from  $\mu_0$  to  $\mu_1$ . As shown in remark 2.18, for a curve  $\mu: t \mapsto \mu_t \in \mathcal{P}_w(\mathbb{R}^d)$  with boundary conditions  $\mu_0 = \tilde{\mu}_0$ ,  $\mu_1 = \tilde{\mu}_1$  the optimal map from  $\mu_0$  to  $\mu_t$  is given by

$$F_t(x) = x + tv_0(x),$$

where  $v_0(x) = F(x) - x$ , and it holds

$$d_w(\mu_t, \mu_0) = \left( \int_{\mathbb{R}^d} |v(x)|^2 d\mu_0(x) \right)^{1/2}. \quad (3.4)$$

The former equation motivates that elements of a possible tangent space  $T\mathcal{P}_w(\mathbb{R}^d)$  could be defined using optimal maps; the latter equation indicates these elements should be contained in  $L^2(\mu_0)$ .

**Definition 3.5** (Tangent space induced by optimal maps). At a measure  $\mu \in \mathcal{P}_w(\mathbb{R}^d)$ , the tangent space induced by optimal maps is defined by

$$\text{Tan}_\mu \mathcal{P}_w(\mathbb{R}^d) = \overline{\{v \in L^2_\mu(\mathbb{R}^d) : Id + tv \text{ is optimal for some } t \in [0, 1]\}}^{L^2_\mu(\mathbb{R}^d)}. \quad (3.5)$$

**Remark 3.6** (Other choices for tangent spaces). The definition above approaches elements of the tangent space via optimal maps; note, that the regularity of  $\mu_0$  is not required for this definition. As shown in remark 2.18, for regular initial measures  $\mu_0$  there is a particular solution  $(v, \mu)$  that corresponds to a constant speed geodesic; the velocities corresponding to these solutions can be used to define an alternative tangent space (*the regular tangent space* in notation of [7]):

$$\text{Tan}_\mu^* \mathcal{P}_w(\mathbb{R}^d) = \overline{\{v = \nabla \varphi \in C_c^\infty(\mathbb{R}^d)\}}^{L^2_\mu(\mathbb{R}^d)}.$$

It is possible to show that the two definitions are equivalent [7, Ch. 4.4]; this claim is easy to check using Brenier's theorem for regular measures  $\mu$ , the general case  $\mu \in \mathbb{R}^d$  is more complicated. Referring to [7] for a deeper discussion of these problems, we remark that the defined tangent space  $\text{Tan}_\mu \mathcal{P}_w(\mathbb{R}^d)$  is only a subset of a more general so-called *geometric tangent space* that is induced not by optimal maps, but by optimal plans and is a more appropriate (and more technical) setting (introducing which would take this discussion beyond the bounds of an overview).

The tangent space  $\text{Tan}_\mu \mathcal{P}_w$  can be endowed with an inner product in a natural way:

$$\langle \nabla \varphi, \nabla \psi \rangle_{\text{Tan}_\mu \mathcal{P}_w} = \langle \nabla \varphi, \nabla \psi \rangle_\mu := \int_{\mathbb{R}^d} \nabla \varphi(x) \nabla \psi(x) d\mu(x).$$

For regular measures  $\mu_0$ , the corresponding induced metric  $d_{ind}$  on  $\mathcal{P}_w(\mathbb{R}^d)$  coincides with the quadratic Wasserstein metric: for  $\mu \in \mathbb{R}^d$ ,

$$\begin{aligned} d_{ind}(\mu_0, \mu) &= \inf \int_0^1 |\dot{\mu}_t| dt = \inf \int_0^1 \lim_{\varepsilon \rightarrow 0} \frac{d_w(\mu_{t+\varepsilon}, \mu_t)}{t} dt \\ &= \inf \int_0^1 \lim_{\varepsilon \rightarrow 0} \frac{\varepsilon \cdot d_w(\mu_0, \mu_1)}{t} dt = d_w(\mu_0, \mu_1), \end{aligned}$$

where we used eq. (3.3). Despite all the similarities, the quadratic Wasserstein space is not a Riemannian manifold, which can be seen by the following argument. Define the exponential map

$$\begin{aligned} \exp_\mu : T_\mu \mathcal{P}_w(\mathbb{R}^d) &\rightarrow \mathcal{P}_w(\mathbb{R}^d) \\ v &\mapsto \exp_\mu(tv) := (Id + tv)_\# \mu. \end{aligned} \quad (3.6)$$

The injectivity radius of this map is zero:

$$\begin{aligned} t \mapsto \exp_\mu(tv) \text{ is a geodesic for } t \in [0, T] &\Leftrightarrow \\ Id + Tv \text{ is an optimal map from } \mu \text{ to } (Id + Tv)_\# \mu &\Leftrightarrow \\ \text{its Jacobian } \mathcal{I} + TDv \text{ is bounded from below.} & \end{aligned}$$

Since  $Dv$  is not bounded for  $v \in T_\mu \mathcal{P}_w(\mathbb{R}^d)$ , the radius of injectivity is equal to zero.

**Definition 3.7** (Gradient). Consider a functional  $E : \mathcal{P}_w(\mathbb{R}^d) \rightarrow \mathbb{R} \cup \{\infty\}$ , let  $\mu$  be a regular measure in the domain of  $E$ . An element  $\nabla_w E \in T_\mu \mathcal{P}_w$  is called the gradient of  $E$  at  $\mu$  if for every measure  $\nu \in \mathcal{P}_w(\mathbb{R}^d)$  with the optimal map  $F = Id + \nabla\varphi$  from  $\mu$  to  $\nu$  holds

$$\begin{aligned} E(\nu) &= E(\mu) + \langle \nabla_w E, F - Id \rangle_\mu + o(d_w(\mu, \nu)) \\ &= E(\mu) + \langle \nabla_w E, \nabla\varphi \rangle_\mu + o(d_w(\mu, \nu)). \end{aligned} \quad (3.7)$$

**Remark 3.8** (How to calculate a gradient). For the situation described above, consider a constant speed geodesic  $\mu_t$  from  $\mu_0 = \mu$  to  $\mu_1 = \nu$ . Equation (3.7) is equivalent (cf. equation (3.3)) to

$$E(\mu_t) = E(\mu_0) + \langle \nabla E, t\nabla\varphi \rangle_{\mu_0} + o(td_w(\mu_0, \mu_1)). \quad (3.8)$$

Therefore,

$$\left. \frac{\partial}{\partial t} E(\mu_t) \right|_{t=0} = \langle \nabla_w E, \nabla\varphi \rangle_{\mu_0}, \quad (3.9)$$

or, more explicitly, the equation

$$\left. \frac{\partial}{\partial t} E((Id + t\nabla\varphi)_\# \mu_0) \right|_{t=0} = \langle \nabla_w E, \nabla\varphi \rangle_{\mu_0} \quad (3.10)$$

should hold for all  $\nabla\varphi \in T_\mu \mathcal{P}_w(\mathbb{R}^d)$ .

**Remark 3.9.** To remain consistent with the definition 3.7, it is necessary to check for each particular case that the gradient  $\nabla_w E$  defined this way is indeed an element of  $T\mathcal{P}_w(\mathbb{R}^d)$ .

The next proposition verifies that the condition “gradient must vanish at an extremum” remains valid for the functionals on the Wasserstein space.

**Proposition 3.10** (Necessary condition for a minimum). Assume that the functional  $E : \mathcal{P}_w(\mathbb{R}^d) \rightarrow \mathbb{R} \cup \{\infty\}$  reaches its minimum at the regular measure  $\mu \in \mathcal{P}_w(\mathbb{R}^d)$  and the gradient  $\nabla_w E \in T_\mu \mathcal{P}_w(\mathbb{R}^d)$  exists. Then,  $\nabla_w E = 0$ .

*Proof.* For all  $\nabla\varphi \in T_\mu \mathcal{P}_w(\mathbb{R}^d)$  holds

$$E((Id + \nabla\varphi)_\# \mu) = E(\mu) + \langle \nabla_w E, \nabla\varphi \rangle_\mu + o(d_w((Id + \nabla\varphi)_\# \mu, \mu)). \quad (3.11)$$

Assume that  $\nabla_w E \neq 0$ ; choose  $\nabla\varphi = -\nabla_w E$ . For  $t \in \mathbb{R}$ ,

$$E((Id + t\nabla\varphi)_\# \mu) = E(\mu) - t \int_{\mathbb{R}^d} \|\nabla_w E\|^2 d\mu + o(d_w((Id + \nabla\varphi)_\# \mu, \mu)). \quad (3.12)$$

Therefore, for all  $\epsilon > 0$  there exists  $t > 0$  s.t. for all  $\nabla\varphi \in T_\mu \mathcal{P}_w(\mathbb{R}^d)$  satisfying the condition  $d_w((Id + \nabla\varphi)_\# \mu, \mu) < \epsilon$  holds  $E((Id + t\nabla\varphi)_\# \mu) < E(\mu)$ ; hence,  $\mu$  is not a minimizer of  $E$ . The claim of the lemma follows by contradiction.

**Remark 3.11** (Steepest descent: functional approach). Using the definition of the gradient, one can establish the relationship between a functional and its possible minimizer. This remark shows that for each functional  $E$  one can define a functional  $H$ , so that the Wasserstein gradient of  $H$  imposes a dynamic that minimizes the energy  $E$ .

For a given functional  $E: \mathcal{P}_w(\mathbb{R}^d) \rightarrow \mathbb{R} \cup \{\infty\}$ , define the functional

$$\begin{aligned} H: \mathcal{P}_w(\mathbb{R}^d) \times \mathcal{P}_w(\mathbb{R}^d) &\rightarrow \mathbb{R} \cup \{\infty\}; \\ (\mu, \nu) &\mapsto E[\nu] + d_w(\mu, \nu). \end{aligned} \quad (3.13)$$

Fix  $\mu_0 \in \mathcal{P}_w(\mathbb{R}^d)$ ; assume that the minimum of  $H(\mu_0, \cdot)$  exists and is reached at some  $\mu_1 \in \mathcal{P}_w(\mathbb{R}^d)$ , i.e.

$$\mu_1 = \arg \min_{\nu \in \mathcal{P}_w(\mathbb{R}^d)} H(\mu_0, \nu). \quad (3.14)$$

Let  $F$  be the optimal map from  $\mu_0$  to  $\mu_1$ , let  $(\mu, \nu)$  be the corresponding constant speed solution (cf. remark 2.18). If the gradient  $\nabla_w E$  at the point  $\mu_0$  exists, then it holds

$$v_0(x) = -\nabla_w E[\mu]. \quad (3.15)$$

*Proof.* By lemma 3.10, for the point  $\nu = \mu_1$  holds

$$\nabla_w H(\mu_0, \nu) = 0,$$

where the gradient is taken in the second argument. Equivalently,

$$\nabla_w E[\nu] + \nabla_w d_w(\mu_0, \nu) = 0.$$

To compute the Wasserstein gradient of the second summand, use remark 3.8: for any  $\nu \in \mathcal{P}_w(\mathbb{R}^d)$ , consider the constant speed solution  $(\mu, \nu)$  from  $\mu_0$  to  $\mu_1$ . Then,

$$\left. \frac{\partial}{\partial t} d_w(\mu_0, (Id + tv_0)_\# \mu_0) \right|_{t=0} = \left. \frac{\partial}{\partial t} \left( \epsilon \int_{\mathbb{R}^d} v_0^2(x) d\mu_0(x) \right) \right|_{t=0} = \langle v_0, v_0 \rangle_{L^2(\mu)},$$

from which formally follows

$$\nabla_w E[\mu] = -v_0(x). \quad (3.16)$$

Note that the gradient flow defined by this condition decreases the functional  $E$  by construction: indeed, for any small  $t \geq 0$  it holds

$$E((Id + tv_0)_\# \mu) = E(\mu) - t \int_{\mathbb{R}^d} v_0^2 d\mu + o\left(d_w((Id + tv_0)_\# \mu, \mu)\right). \quad (3.17)$$

**Example 3.12** (Steepest descent: differential equation approach). It is possible to search for curves  $(\mu_t)_{t \in \mathbb{R}_{\geq 0}}$  in the quadratic Wasserstein space defined by the equation

$$\frac{\partial \mu}{\partial t} = -\nabla_w E(\mu); \quad (3.18)$$

by construction, solutions of this equation would flow along the steepest descent of the functional  $E$ . This equation is valid in distributional sense; to interpret it in terms of a



density function  $g_t$  corresponding to  $\mu_t$ , consider a geodesic  $\mu_t = (1 - t\nabla_w E)_\# \mu_0$  along the steepest descent. Let  $u \in L^2(\mu)$  be a test function,  $t \in \mathbb{R}_{\geq 0}$ ,  $\varepsilon > 0$ . Then,

$$\begin{aligned} \int u(x) d\mu_{t+\varepsilon}(x) &= \int u(x - \varepsilon \nabla_w E(x)) d\mu_t(x) \\ &= \int u(x) d\mu_t(x) - \varepsilon \int \nabla u(x) \cdot \nabla_w E(x) d\mu_t(x) + o(\varepsilon). \end{aligned}$$

Hence, for regular measures  $\mu_t$  (3.18) is equivalent to

$$\partial_t \int u(x) g_t(x) dx = - \int \nabla u(x) \cdot \nabla E(x) g_t(x) dx,$$

or, integrating by parts and omitting test functions,

$$\partial_t g_t(x) = \nabla \cdot (\nabla E_w(x) \cdot g_t(x)) \quad (3.19)$$

with  $\nabla_w E_w(x)$  defined by eq. (3.10).

**Remark 3.13.** For infinitesimal times  $t$  the equation (3.18) can be rewritten in terms of the velocity. Consider the optimal map between  $\mu_0, \mu_t \in \mathcal{P}_{ac}(\mathbb{R}^d) \cap \mathcal{P}_w(\mathbb{R}^d)$  given by  $F_t(x) = x + tv_0(x)$ . Then (cf. eq. (2.4)),

$$\mu_t(x) = \frac{\mu_0(F_t^{-1}(x))}{|\det DF_t(x)|}. \quad (3.20)$$

We would like to make use of the Taylor expansion of  $\mu_t(x)$  at  $t = 0$ . To do so, observe that

$$|\det DF_t(x)| = 1 + t \nabla v_0(x) + O(t^2) \quad \text{for a.e. } x \in \mathbb{R}^d.$$

This equality is a well-known fact from the determinant calculus (see, e.g., [33]). Therefore,

$$\begin{aligned} \mu_t(x) &= \mu_0(x) + t \left( \nabla \mu_0(x) \frac{1}{\frac{d}{dt} F_t(x)|_{t=0}} - \mu_0(x) \nabla v_0(x) \right) + O(t^2) \\ &= \mu_0(x) + t \left( \nabla \mu_0(x) \frac{1}{v_0(x)} - \mu_0(x) \nabla v_0(x) \right) + O(t^2) \quad \text{as } t \rightarrow 0. \end{aligned}$$

We assume that  $\mu_0(\{v_0(x) = 0\}) = 0$ , so that the equation above is well-posed for  $\mu_0$ -a.e.  $x \in \mathbb{R}^d$ . Then,

$$\nabla_w E[\mu_t(x)] = \nabla_w E[\mu_0(x)] + O(t)$$

as  $t \rightarrow 0$ . Inserting these equations into (3.18), obtain the following equation:

$$\frac{\nabla \mu_0(x)}{v_0(x)} - \mu_0(x) \nabla v_0(x) = -\nabla_w E[\mu_0] \quad \text{for } \mu_0 \text{ - a.e. } x \in \mathbb{R}^d. \quad (3.21)$$

The Wasserstein gradient  $\nabla_w E[\mu_0]$  is defined by the equation (3.10). The equation above can be interpreted as an analogon of the ‘‘explicit Euler algorithm step’’: it allows us to calculate the velocity  $v_0(x)$  of the material distribution  $\mu_0$  at the time  $t = 0$ . Just as the explicit Euler algorithm, it can not be used for most practical purposes, since it is highly unstable: after a few steps, the approximation error becomes unacceptably large.

### 3.2 Phase retrieval on the quadratic Wasserstein space

This section applies the methods of the previous subsection to phase retrieval; specifically, to the modulus energy functional  $E_M$ . First, we observe the equation connected to the gradient flow of  $E_M$  in the space  $\mathcal{P}_w(\mathbb{R}^d)$ . Then, we discuss some issues regarding the practical applicability of the established connection between phase retrieval and energy minimization on Wasserstein space. We recall the definition of  $\lambda$ -geodesic convexity which is a central notion in the description of gradient flows on Wasserstein spaces (cf. [8]), and show an inequality indicating that the functional  $E_M$  is (with certain restrictions that are defined below) locally  $\lambda$ -geodesically convex. We conclude this section with the discussion of this result.

**Example 3.14** (Steepest descent in  $\mathcal{P}_w(\mathbb{R}^d)$  for the error modulus functional). As motivated in Section 1, we want to consider gradient flow for phase retrieval functional  $E_M$ . To do so, it is necessary to calculate  $\nabla_w E_M$ . Let  $\mu_0 \in \mathcal{P}_w(\mathbb{R}^d) \cap P_{ac}(\mathbb{R}^d)$  be the initial distribution with the density  $g \in \mathcal{S}(\mathbb{R}^d)$ . Assume that phase problem has a solution  $g_1(x) = f(x) \in \mathcal{S}(\mathbb{R}^d)$ ; let  $\mu_1 \in \mathcal{P}_w(\mathbb{R}^d)$  be the corresponding measure. Denote the time-dependent optimal map (cf. remark 2.18) from  $g_0$  to  $g_t$  by  $F_t(x) = x + t\nabla\varphi(x)$ ,  $\varphi \in C_c^\infty(\mathbb{R}^d)$ . Recall the modulus error functional as in (1.30):

$$E[g_t] = \frac{1}{2(2\pi)^d} \int_{\mathbb{R}^d} (|\hat{g}_t(k)| - A(k))^2 dk.$$

To determine  $\nabla_w E(x)$  from the equation (3.10), it is necessary to calculate  $\frac{d}{dt} E[g_t]$ . Notice that

$$\begin{aligned} \hat{g}_t(k) &= \int_{\mathbb{R}^d} e^{-ik \cdot x} g_t(x) dx = \int_{\mathbb{R}^d} e^{-ik \cdot x} d\mu_t(x) \\ &= \int_{\mathbb{R}^d} e^{-ik \cdot F_t(x)} d\mu_0(x) = \int_{\mathbb{R}^d} e^{-ik \cdot x - ikt \nabla\varphi(x)} d\mu_0(x), \end{aligned}$$

and

$$\left. \frac{d}{dt} \right|_{t=0} \hat{g}_t(k) = \int_{\mathbb{R}^d} (-ik) \nabla\varphi(x) e^{-ik \cdot x} g_0(x) dx =: \hat{g}'_0(k). \quad (3.22)$$

Hence,

$$\begin{aligned} \left. \frac{d}{dt} \right|_{t=0} E[g_t] &= \frac{1}{2(2\pi)^d} \int_{\mathbb{R}^d} (|\hat{g}_0(k)| - A(k)) \frac{1}{|\hat{g}_0(k)|} (\hat{g}_0^*(k) \hat{g}'_0(k) + \hat{g}'_0^*(k) \hat{g}_0(k)) dk \\ &= \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} (|\hat{g}_0(k)| - A(k)) \frac{\hat{g}_0(k)}{|\hat{g}_0(k)|} \hat{g}'_0^*(k) dk \end{aligned} \quad (3.23)$$

(since  $\hat{g}_0(-k) = \hat{g}_0^*(k)$  and  $\hat{g}'_0(-k) = \hat{g}'_0^*(k)$ )

$$\begin{aligned} &= \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} \left( 1 - A(k) \frac{\hat{g}_0(k)}{|\hat{g}_0(k)|} \right) \int_{\mathbb{R}^d} (ik) \nabla\varphi(x) e^{ik \cdot x} g_0(x) dx dk \\ &= \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} \nabla\varphi(x) \nabla_x \left( \int_{\mathbb{R}^d} e^{ik \cdot x} \left( 1 - A(k) \frac{\hat{g}_0(k)}{|\hat{g}_0(k)|} \right) dk \right) g_0(x) dx \\ &= \int_{\mathbb{R}^d} \nabla\varphi(x) \nabla((1 - P_M)g_0)(x) g_0(x) dx, \end{aligned} \quad (3.24)$$

recalling the modulus projection operator  $P_M(g) = \mathcal{F}^{-1} \left( A(k) \frac{\hat{g}_0(k)}{|\hat{g}_0(k)|} \right)$ . Note that here and in further similar calculations it would be more precise to use the expression  $e^{-i \arg \hat{g}}$  instead of the expression  $\hat{g}(k)/|\hat{g}(k)|$ , since the former is well-defined for  $\hat{g}_0(k) = 0$ . A precise treatment of the integration in (3.23) would include splitting it into two integrations over the domains  $A = \{k \in \mathbb{R}^d \mid \hat{g}(k) \neq 0\}$  and  $\mathbb{R}^d/A$ . This calculation does not alter the result (3.24) as long as we use the convention  $\hat{g}(k)/|\hat{g}(k)| = e^{-i \arg \hat{g}}$  (with other words,  $\frac{\hat{g}(k)}{|\hat{g}(k)|} = 1$  for  $k \in A$ ). Therefore,

$$\nabla E[g_t](x) = \nabla \psi(x), \quad (3.25)$$

and, by (3.19),

$$\partial_t g_t(x) = \nabla \left( g_t(x) \cdot \nabla ((1 - P_M)g_0)(x) \right). \quad (3.26)$$

The equation presented in the example above shows that phase retrieval can be expressed as a differential equation. This equation can be interpreted as the gradient flow of the material towards a minimum of the energy  $E_M$ . Analysis of this equation may lead to better understanding of the existence and uniqueness of phase problem solutions.

Gradient flows generated by the method used above are related to functionals of the form

$$E[\mu_t] + \frac{d_w^2(\mu_0, \mu_t)}{2T}, \quad (3.27)$$

with a time parameter  $T \in \mathbb{R} [1, 7, 8]$ . From the applied point of view, it is easier to search for minimizers of (3.27) rather than to solve such equation as (3.26). This is still a cumbersome task, since a search for  $\mu_t$  minimizing (3.27) requires search for a solution of the Kantorovich problem in each step (due to the implicit definition of  $d_w(\mu_0, \mu_1)$ ).

One can try investigate the gradient flow (3.26) using convergence results from the optimal transport theory. The behaviour of the gradient flow corresponding to a general energy functional  $E$  is well understood if the functional is well-behaved along the geodesics in the Wasserstein space [8]; the precise description requires the following notion.

**Definition 3.15** ( $\lambda$ -geodesically convex). A functional  $E: \mathcal{P}_w(\mathbb{R}^d) \rightarrow \mathbb{R}$  is called  $\lambda$ -geodesically convex, if for any  $\mu, \nu \in \mathcal{P}_w(\mathbb{R}^d), t \in (0, 1)$ ,

$$E[t\mu + (1-t)\nu] \leq tE[\mu] + (1-t)E[\nu] + \lambda t(1-t)d_w(\mu, \nu) \quad (3.28)$$

for some  $\lambda \in \mathbb{R}$ .

We can not expect the functional  $E_M$  to be  $\lambda$ -geodesically convex. (This is suggested by the fact that non-uniqueness is common for phase retrieval solutions and by the fact that  $E_M$  is not convex.) However, it is possible to show that  $E_M$  is locally  $\lambda$ -geodesically convex in a certain sense described in the following proposition. Before presenting it, we recall some useful identities and fix the notation in the lemma 3.16. Beware that in this proposition the function  $f$  denotes not the solution of the phase problem, but just an approximation to the solution (and is on equal footing with the notation  $g$ ).

**Lemma 3.16.** Let  $f, g \in \mathcal{S}(\mathbb{R}^d)$  be the densities of two measures in quadratic Wasserstein space. Assume that there is a constant speed geodesic between these measures in  $\mathcal{P}_w(\mathbb{R}^d) \cap \mathcal{P}_{ac}(\mathbb{R}^d)$ . Let

$$F_h(x) = x + hv(x) \quad \text{for a.e. } x \in \mathbb{R}^d, \quad (3.29)$$

be the optimal map from  $f$  to  $g$ , where  $h \geq 0, v = \nabla\phi$  for some  $\phi \in C^\infty(\mathbb{R}^d)$ . Recall

$$|\det DF_h(x)|g(F_h(x)) = f(x) \quad \text{for a.e. } x \in \mathbb{R}^d. \quad (3.30)$$

Assume that the set  $\{v(x) = 0\}$  has the Lebesgue measure 0. Further, let  $F_h(x)$  be bijective, with  $v_b$  defined through

$$F_h^{-1}(x) = x + hv_b(x). \quad (3.31)$$

Note that with this definition holds

$$g(x) = |\det DF_h^{-1}(x)|f(F_h^{-1}(x)) \quad \text{for a.e. } x \in \mathbb{R}^d. \quad (3.32)$$

Then:

$$v_b(x) = \frac{1}{v(x)} + O(h) \quad \text{for a.e. } x \in \mathbb{R}^d, \quad (3.33)$$

as  $h \rightarrow 0$ . Also, the derivative of  $g$  by  $h$  from equation (3.32) is

$$\left. \frac{dg}{dh} \right|_{h=0} = \nabla f(x) \cdot \frac{1}{v(x)} - f(x)\nabla v(x) \quad \text{for a.e. } x \in \mathbb{R}^d, \quad (3.34)$$

and the derivative of  $f$  by  $h$  from equation (3.30) is

$$\left. \frac{df}{dh} \right|_{h=0} = \nabla g(x) \cdot v(x) + g(x)\nabla v(x) \quad \text{for a.e. } x \in \mathbb{R}^d, \quad (3.35)$$

*Proof.* Proof of the equation (3.33). By inverse function theorem, we obtain using Taylor expansion:

$$\begin{aligned} F_h^{-1}(x) &= F_h^{-1}(x)|_{h=0} + h \left. \frac{d}{dh} \right|_{h=0} F_h^{-1}(x) + O(h^2) \\ &= x + h \frac{1}{\left. \frac{d}{dh} \right|_{h=0} F_h(x)} + O(h^2) = x + h \cdot \frac{1}{v(x)} + O(h^2) \quad \text{for a.e. } x \in \mathbb{R}^d, \end{aligned}$$

as  $h \rightarrow 0$ . From (3.31),

$$x + hv_b(x) = x + h \cdot \frac{1}{v(x)} + O(h^2) \quad \text{for a.e. } x \in \mathbb{R}^d,$$

from which follows (3.33). This also means that

$$\left. \frac{d}{dh} \right|_{h=0} F_h^{-1}(x) = \frac{1}{v(x)} \quad \text{for a.e. } x \in \mathbb{R}^d,$$

which will be useful in the proof of (3.34).

Proof of equations (3.34), (3.35). As mentioned before (cf. [33]),

$$|\det DF_h(x)| = x + h\nabla \cdot v(x) + O(h^2), \quad \text{as } h \rightarrow 0.$$

From that follows,

$$|\det DF_h(x)|^{-1} = x - h\nabla \cdot v(x) + O(h^2), \quad \text{as } h \rightarrow 0.$$

The claims (3.34), (3.35) follow by the chain rule.  $\square$

**Proposition 3.17.** Consider two regular measures in quadratic Wasserstein space with densities  $f, g$  and the notation and assumptions fixed in the lemma above. Then, for all  $t \in [0, 1]$ :

$$\begin{aligned} E[tf + (1-t)g] &\leq tE[f] + (1-t)E[g] \\ &\quad + t(t-1)(C_f + C_g)d_w(g, f) + \mathcal{O}(d_w^2(g, f)), \end{aligned} \quad (3.36)$$

as  $h \rightarrow 0$ , with the notation

$$C_f = \|(\nabla f) \cdot (1 - P_M)f\|_{L^2} + \|\nabla(f(1 - P_M)f)\|_{L^2}, \quad (3.37)$$

and  $C_g$  mutatis mutandis.  $\square$

*Proof.* To show the claim, use the Taylor expansion of  $E[tf + (1-t)g]$  at  $f$  in  $h$  (and at  $g$  in  $h$ ) using (3.35) (resp. (3.34)) and apply Hölder's inequality. First, consider the Taylor-expansion of  $E[tf + (1-t)g]$  at  $f$  in  $h$ :

$$\begin{aligned} E[tf + (1-t)g] &= \frac{1}{2(2\pi)^d} \int_{\mathbb{R}^d} \left( |\mathcal{F}(tf + (1-t)g)(k)| - A(k) \right)^2 dk = \\ &= \frac{1}{2(2\pi)^d} \int_{\mathbb{R}^d} \left( |\mathcal{F}(tf + (1-t)g)(k)| - A(k) \right)^2 dk \Big|_{h=0} \\ &\quad + \frac{h}{(2\pi)^d} \int_{\mathbb{R}^d} \left( |\mathcal{F}(tf + (1-t)g)(k)| - A(k) \right) \frac{\mathcal{F}(tf + (1-t)g)(k)}{|\mathcal{F}(tf + (1-t)g)(k)|} \\ &\quad \quad \times (1-t) \int_{\mathbb{R}^d} e^{-ikx} \frac{d}{dh} \Big|_{h=0} g(x) dx dk \Big|_{h=0} + O(h^2) \\ &= E[f] + h \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} (|\hat{f}(k)| - A(k)) \frac{\hat{f}(k)}{|\hat{f}(k)|} \\ &\quad \quad \times (1-t) \int_{\mathbb{R}^d} e^{-ikx} (\nabla f(x) \frac{1}{v(x)} - f(x) \nabla v(x)) dx dk + O(h^2) \\ &= E[f] + h(1-t) \int_{\mathbb{R}^d} (1 - P_M)f(x) (\nabla f(x) \frac{1}{v(x)} - f(x) \nabla v(x)) dx dk + O(h^2) \\ &= E[f] + h(1-t) \int_{\mathbb{R}^d} (1 - P_M)f(x) (\nabla f(x) v_b(x) - f(x) \nabla v(x)) dx dk + O(h^2) \\ &\leq E[f] + h(1-t) \|(1 - P_M)f(x) (\nabla f(x))\|_{L^2} \|v_b\|_{L^2} \\ &\quad \quad + h(1-t) \|\nabla((1 - P_M)f(x) f(x))\|_{L^2} \|v\|_{L^2} + O(h^2) \\ &= E[f] + (1-t)C_f d_w(g, f) + O(h^2). \end{aligned} \quad (3.38)$$

We used the definition of the projection operator

$$(1 - P_M)f(x) = \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} (|\hat{f}(k)| - A(k)) \frac{\hat{f}(k)}{|\hat{f}(k)|} e^{-ikx} dk,$$

and the fact that for constant speed geodesics by Brenier's theorem holds

$$h\|v\|_{L^2} = d_w(g, f) = d_w(f, g) = h\|v_b\|_{L^2}.$$

In the second to last row of (3.38) we used integration; the boundary terms vanish due to the assumption  $f, g \in \mathcal{S}(\mathbb{R}^d)$ .

Analogously, expanding  $E[tf + (1-t)g]$  at  $g$  in  $h$ :

$$\begin{aligned}
E[tf + (1-t)g] &= E[g] + h \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} (|\hat{g}(k)| - A(k)) \frac{\hat{g}(k)}{|\hat{g}(k)|} \\
&\quad \times t \int_{\mathbb{R}^d} e^{-ikx} (\nabla g(x)v(x) + g(x)\nabla v(x)) dx dk + O(h^2) \\
&= E[g] + ht \int_{\mathbb{R}^d} (1 - P_M)g(x) (\nabla g(x)v(x) + g(x)\nabla v(x)) dx dk + O(h^2) \\
&\leq E[f] + ht \|(1 - P_M)g(x) (\nabla g(x))\|_{L^2} \|v\|_{L^2} \\
&\quad + ht \|\nabla((1 - P_M)g(x)g(x))\|_{L^2} \|v\|_{L^2} + O(h^2) \\
&= E[f] + tC_g d_w(g, f) + O(h^2). \tag{3.39}
\end{aligned}$$

Multiplying (3.38) by  $t$  and (3.39) by  $1-t$  and adding them together, obtain the claim.  $\square$

This proposition indicates that the behaviour of the functional is connected to the gradient  $\nabla g$  of an approximation  $g$ . The inequality shown above indicates that in the regions with high values of the gradient  $\nabla$  the modulus energy functional may be not  $\lambda$ -geodesically convex. It is still possible to employ algorithms from the optimal transportation, for example, such as described in [34].

From the analytical point of view, it may be beneficial to consider a space of approximations  $\mathcal{X}$  containing only approximations with the appropriately bounded gradient. This requires an investigation of properties of the modulus energy functional on restrictions of the Wasserstein space. More specifically: is it possible to show convergence properties of gradient flows (cf., e.g., [8]) in the restriction of the Wasserstein space? From the applied point of view, such investigation includes design of energy minimizing algorithm behaviour on the boundary of the subset of the Wasserstein space. With other words: is it possible (and if so, how) to consider a Wasserstein flow while simultaneously restricting the gradient norm of the approximations.

Another point of interest concerns the generalization of the energy modulus functional to the space of distributions:  $E_{M,general}: \mathcal{S}'(\mathbb{R}^d) \rightarrow \mathbb{R}_{\geq 0}$ , which would allow consideration of general (and not only absolutely continuous measures) in  $\mathcal{P}_w(\mathbb{R}^d)$ . For example, one could consider the functional

$$E[g] := \frac{1}{2(2\pi)^d} \int_{\mathbb{R}^d} (|\hat{g}(k)|^2 - A^2(k))^2 dk,$$

whose domain can be easily generalized to  $\mathcal{S}'(\mathbb{R}^d)$  using the autocorrelation measure (cf. remark 1.10).

Summing up, it is possible to consider phase retrieval as a gradient flow in a Wasserstein space. The connection of the phase problem to the PDE theory may be used to gain deeper insight for the existence and uniqueness of phase problem solutions. Another point of investigation is the question whether the optimal transport algorithms can be applied to the case of the modulus energy  $E_M$  or related functionals (convergence of such algorithms often may be generalized more directly to gradient flows without convexity properties [34]). Also, it may be profitable to discuss gradient flows on the subset of the

Wasserstein space on which the phase retrieval energy functionals exhibit geodesically convex behaviour. The latter consideration includes investigating the behaviour of the gradient flows on a boundary of a subset of the quadratic Wasserstein space.

## 4 Outlook

To conclude this thesis, we briefly summarize some of its results.

### Phase retrieval and quadratic Wasserstein space

It can be beneficial to discuss phase retrieval (and related questions from X-ray crystallography, such as phase problem) on the quadratic Wasserstein space. This setting allows a separation of the non-negativity and the Fourier modulus constraints that are necessary for phase retrieval. The non-negativity and mass conservation properties are thus satisfied in a natural way. In this setting, phase retrieval can be thought of as a transportation problem from some random initial density to the density minimizing an appropriate energy functional.

### Phase retrieval as a gradient flow

In [3] it was shown that the error reduction algorithm can be accessed via the method of steepest descent. The method of steepest descent was applied to an appropriate functional (the energy modulus functional  $E_M[g] = \|(1 - P_M)g\|_{L^2}$ ) using the metric  $\|\cdot\|_{L^2}$ . (The modulus projection operator is defined by  $g \mapsto P_M g := \mathcal{F}^{-1}(e^{i \arg \hat{g}} \cdot A)$  for any  $g \in \mathcal{S}(\mathbb{R}^d)$  and for the scattering amplitude  $A \in \mathcal{S}(\mathbb{R}^d)$ .) In that case, the non-negativity constraint must be imposed explicitly. The example 3.14 shows that it may be beneficial to explore this approach with a different metric (namely, the quadratic Wasserstein metric  $d_w$ ). This approach leads to the evolution equation

$$\partial_t g_t = \nabla(g_t \cdot \nabla(1 - P_M)g_t), \quad (4.1)$$

where  $g_t \in \mathcal{S}(\mathbb{R}^d)$  is the density of an absolutely continuous measure in the quadratic Wasserstein space. This equation connects phase retrieval to the PDE theory, which may be used to gain a deeper insight on the phase problem.

### Convexity properties of $E_M$

An algorithmical implementation of (4.1) can be explored via the related energy minimization problem for the functional

$$E_M[g] + \frac{d_w^2(g_0, g)}{2T}, \quad (4.2)$$

with an initial density  $g_0 \in \mathcal{S}(\mathbb{R}^d)$ , the density  $g \in \mathcal{S}(\mathbb{R}^d)$  that minimizes the functional, and the parameter  $T \in \mathbb{R}$ . It may be beneficial to investigate the numerical behaviour of the gradient flow defined by (4.2), and to check whether the convergence results for such algorithms are transferrable to the phase retrieval case (cf., e.g., [34] and references therein).

The proposition 3.17 indicates that  $E_M$  shows properties of  $\lambda$ -geodesic convexity under certain restrictions; namely, conditions

$$\begin{aligned} \|(\nabla g) \cdot (1 - P_M)g\|_{L^2} &< C_1, \\ \|\nabla(g \cdot (1 - P_M)g)\|_{L^2} &< C_2, \end{aligned}$$



$C1, C2 \in \mathbb{R}$ , would be sufficient to show local  $\lambda$ -geodesic convexity for in the vicinity of the absolutely continuous measure with density  $g$ .

### Open questions

These results indicate that the discussion of phase retrieval in the context of optimal transport may lead to new results. The gradient flow formulation allows an investigation of phase retrieval from the PDE point of view.

The challenges presented by the Fourier modulus constraint of phase retrieval do not fall directly into the  $\lambda$ -geodesically convex setting that is preferable for optimal transport. However, they exhibit similar properties if the underlying quadratic Wasserstein space is restricted to a subset. One of the possible choices would be the subset

$$\{\mu \in \mathcal{P}_w(\mathbb{R}^d) \mid \mu \in \mathcal{P}_{ac}(\mathbb{R}^d) \text{ with density } g \in \mathcal{S}(\mathbb{R}^d), \text{ so that } \|\nabla g\|_{L^2} < R\} \quad (4.3)$$

for some  $R \in \mathbb{R}$  (or an appropriate closure of this subset). The possible line of investigation would include the following questions: is  $E_M$  (or any similar functional representing the modulus constraint)  $\lambda$ -geodesically convex on that set? Which properties of gradient flows are valid under this restriction? Another possibility is to try a direct transfer of the existing gradient flow convergence proofs to the phase retrieval case.

To construct a more general theory of phase retrieval on the quadratic Wasserstein space, it may be beneficial to consider other energy functionals implementing the modulus constraint, e.g.

$$E[g] = \frac{1}{2(2\pi)^d} \int_{\mathbb{R}^d} (|\hat{g}(k)|^2 - A^2(k))^2 dk, \quad (4.4)$$

that may be easier to generalize to the domain of tempered distributions.

Another promising direction could be the phase retrieval on the quadratic Wasserstein space with the additional atomicity constraint. With other words, one assumes that the number of atoms (peaks) in the solution density is known. In the simplest case, the transportation problem is then reduced to the discrete rearrangement problem on the Wasserstein space. The atomicity constraint is relevant for the X-ray crystallography and can be used to greatly simplify the phase problem description on the Wasserstein space.

## References

- [1] R. Jordan, D. Kinderlehrer, and F. Otto, “The Variational Formulation of the Fokker-Planck Equation,” *Society for Industrial and Applied Mathematics*, vol. 29, pp. 1–17, January 1998.
- [2] D. Jüstel, *Radiation for the Analysis of Molecular Structures with Non-Crystalline Symmetry: Modelling and Representation Theoretic Design*. PhD thesis, Technische Universität München, 2014.
- [3] J. R. Fienup, “Phase Retrieval Algorithms: a Comparison,” *Applied Optics*, vol. 21, pp. 2758–2769, 1982.
- [4] V. Elser, “Phase Retrieval by Iterated Projections,” *J. Opt. Soc. Am. A.*, vol. 20, pp. 40–55, 2003.
- [5] J.-D. Benamou and Y. Brenier, “A Computational Fluid Mechanics Solution to the Monge-Kantorovich Mass Transfer Problem,” *Numer. Math.*, vol. 84, pp. 375–393, 2000.
- [6] C. Villani, *Topics in Optimal Transportation*. AMS, 2003.
- [7] N. Gigli, *On the Geometry of the Space of Probability Measures in  $R_n$  Endowed with the Quadratic Optimal Transport Distance*. PhD thesis, Scuola Normale Superiore di Pisa, 2004.
- [8] L. Ambrosio, N. Gigli, and G. Savare, *Gradient Flows in Metric Spaces and in the Space of Probability Measures*. Birkhäuser, 2008.
- [9] H. Stark (Editor), *Image Recovery: Theory and Application*. Academic Press, 1987.
- [10] J. Miao, R. Charalambous, J. Kirz, and D. Sayre, “Extending the Methodology of X-ray Crystallography to Allow Imaging of Micrometer-Sized Non-Crystalline Specimens,” *Nature*, pp. 342–422, 1999.
- [11] P. Thibault, M. Dierolf, A. Menzel, O. Bunk, C. David, and F. Pfeiffer, “High-Resolution Scanning X-ray Diffraction Microscopy,” *Science*, vol. 321, no. 5887, pp. 379–382, 2008.
- [12] W. Friedrich, P. Knipping, and M. von Laue, “Interferenzerscheinungen bei Röntgenstrahlen,” *Sitzungsberichte der Mathematisch-Physikalischen Classe der Königlich-Bayerischen Akademie der Wissenschaften zu München*, 1912.
- [13] W. L. Bragg, “The Structure of Some Crystals as Indicated by Their Diffraction of X-rays,” *Proc. R. Cos. Lond. A*, vol. 89, pp. 248–277, 1913.
- [14] P. P. Ewald (Editor), *50 Years of X-Ray Diffraction*.
- [15] H. Hauptman and J. Karle, “Solution of the Phase Problem. I. The Centrosymmetric Crystal,” *A.C.A. Monograph No. 3*, 1953.

- [16] D. Shechtman, I. Blech, D. Gradias, and J. Cahn, “Metallic Phase with Long-Range Orientational Order and No Translational Symmetry,” *Phys. Rev. Lett.*, vol. 53, pp. 1951–1954, 1984.
- [17] S. Amelinckx (Editor), *Diffraction and Imaging Techniques in Material Science*, vol. 2. Elsevier, 1978.
- [18] “International Union of Crystallography.” <http://reference.iucr.org/>. Accessed: 2014-11-11.
- [19] “Worldwide Protein Data Bank.” <http://www wwptdb.org/>. Accessed: 2014-08-27.
- [20] G. Friesecke, “Lectures on Fourier Analysis.” Lecture Notes, 2007.
- [21] G. Friesecke, R. D. James, and D. Jüstel, “Design of Radiation for Nanostructure Determination via Mathematical Modelling of Diffraction Patterns,”
- [22] R. P. Feynman, R. B. Leighton, and M. Sands, *The Feynman Lectures on Physics*, vol. 1. Addison-Wesley, 1963.
- [23] A. Hof, “On Diffraction by Aperiodic Structures,” *Comm. Math. Phys.*, vol. 169, pp. 25–43, 1995.
- [24] W. Schempp and B. Dreseler, *Einführung in die Harmonische Analyse*. Teubner, 1980.
- [25] G. Folland, *A Course in Abstract Harmonic Analysis*. CRC Press, 1995.
- [26] U. Grimm and M. Baake, “Recent Progress in Mathematical Diffraction,” *Acta Physica Polonica A*, vol. 126, pp. 474–478, 2014.
- [27] M. Senechal and J. Taylor, “Quasicrystals: The View from Les Houches,” *Math. Intell.*, vol. 12, pp. 54–64, 1990.
- [28] H. H. Bauschke, P. L. Combettes, and R. D. Luke, “Phase Retrieval, Error Reduction Algorithm, and Fienup Variants: a View from Convex Optimization,” *J. Opt. Soc. Am. A.*, vol. 19, pp. 1334–1345, 2002.
- [29] C. Villani, *Optimal Transport: Old and New*. Springer, 2008.
- [30] R. C. Givens and R. M. Shortt, “A Class of Wasserstein Metrics for Probability Distributions,” *Michigan Math. J.*, vol. 31, pp. 231–240, 1984.
- [31] Y. Brenier, “Polar Factorization and Monotone Rearrangement of Vector-Valued Functions,” *Comm. Pure Appl. Math.*, vol. 44, pp. 375–417, 1991.
- [32] V. Arnold, *Mathematical Methods of Classical Mechanics*. Springer, 1989.
- [33] T. Tao, “Matrix Identities as Derivatives of Determinant Identities.”
- [34] D. Matthes and H. Osberger, “Convergence of a Variational Lagrangial Scheme for a Nonlinear Drift Diffusion Equation,” *ESAIM: M2AN*, vol. 48.