

A least-squares method for the inverse reflector problem in arbitrary orthogonal coordinates

Citation for published version (APA):

Beltman, R., ten Thije Boonkkamp, J., & IJzerman, W. (2018). A least-squares method for the inverse reflector problem in arbitrary orthogonal coordinates. *Journal of Computational Physics*, 367, 347-373.
<https://doi.org/10.1016/j.jcp.2018.04.041>

Document license:

TAVERNE

DOI:

[10.1016/j.jcp.2018.04.041](https://doi.org/10.1016/j.jcp.2018.04.041)

Document status and date:

Published: 15/08/2018

Document Version:

Publisher's PDF, also known as Version of Record (includes final page, issue and volume numbers)

Please check the document version of this publication:

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

www.tue.nl/taverne

Take down policy

If you believe that this document breaches copyright please contact us at:

openaccess@tue.nl

providing details and we will investigate your claim.

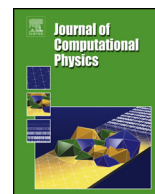


ELSEVIER

Contents lists available at ScienceDirect

Journal of Computational Physics

www.elsevier.com/locate/jcp



A least-squares method for the inverse reflector problem in arbitrary orthogonal coordinates

René Beltman^{a,*}, Jan ten Thije Boonkkamp^a, Wilbert IJzerman^b^a Dept. of Mathematics and Computer Science, Eindhoven University of Technology, PO Box 513 5600 MB Eindhoven, The Netherlands^b Philips Lighting, High Tech Campus 7, 5656 AE Eindhoven, The Netherlands

ARTICLE INFO

Article history:

Received 21 February 2018

Received in revised form 13 April 2018

Accepted 20 April 2018

Available online 27 April 2018

Keywords:

Monge–Ampère equation

Inverse reflector problem

Least-squares method

ABSTRACT

In this article we solve the inverse reflector problem for a light source emitting a parallel light bundle and a target in the far-field of the reflector by use of a least-squares method. We derive the Monge–Ampère equation, expressing conservation of energy, while assuming an arbitrary coordinate system. We generalize a Cartesian coordinate least-squares method presented earlier by Prins et al. [13] to arbitrary orthogonal coordinate systems. This generalized least-squares method provides us the freedom to choose a coordinate system suitable for the shape of the light source. This results in significantly increased numerical accuracy. Decrease of errors by factors up to 10^4 is reported. We present the generalized least-squares method and compare its numerical results with the Cartesian version for a disk-shaped light source.

© 2018 Elsevier Inc. All rights reserved.

1. Introduction

In the last decades LED lighting technology rapidly developed. The costs of LED lighting constantly decrease, as is expressed by Haitz' law which states that the cost per lumen (power perceived by the human eye) falls by a factor of 10 every decade [1]. Furthermore, LED lighting surpasses traditional lighting in efficacy (lumen per Watt) [2]. As a result, LED lighting systems, viz., LEDs integrated in optical systems for illumination, are used in illumination optics ever more frequently.

Two classes of methods are used to design these optical systems: *forward methods* and *inverse methods*. In forward methods the optimal optical system is determined through a process of trial and error. A given optical system is tested, the light output of the system is determined by Monte-Carlo ray tracing [3] and subsequent adjustments are made to improve the system. This process then iterates to a more or less satisfactory solution, of which the quality depends to a large extent on the skill of the designer. This method is widely applicable and straightforward, but time consuming. By contrast, in inverse methods, for given light source and desired light output, a partial differential equation can be derived relating these to the geometry of the optical system. The solution of this partial differential equation then gives the shape of the optical elements. Inverse methods are less straightforward to apply but lead to far more accurate results and are time efficient. Moreover, with inverse methods a diversity of new designs are possible that, due to their complexity, are completely unattainable by forward methods.

The rise of LED lighting has increased the interest in inverse methods because LED lighting operates at much lower temperatures than conventional lighting. This clears the path for the use of easy to mold transparent plastics instead of

* Corresponding author.

E-mail address: r.beltman@tue.nl (R. Beltman).

glass. The optimal shape of these plastic elements can be exactly determined by the inverse method. Moreover, due to active development in diamond turning techniques the arbitrarily shaped elements can be fabricated with increasingly high precision [4].

In this paper we consider an optical system consisting of an incoming parallel bundle of light and a reflecting surface. This optical system is relevant, because parallel bundles can be easily formed in LED lighting systems by placing a converging lens above divergently emitting LEDs. Given the intensity distribution of the incoming parallel bundle and a desired output distribution, a partial differential equation can be derived for the reflector surface. This partial differential equation is an equation of the Monge–Ampère type.

Monge–Ampère type equations also arise in the context of optimal mass transport (OMT). Inverse reflector problems and OMT problems are closely related [5]. OMT concerns, roughly speaking, the problem of filling a hole with a heap of sand from another location. The goal is to do this while minimizing a transportation cost. In inverse optical problems we do not consider a hole and heap of sand, but instead a light source with an emittance and a target with a desired light intensity distribution. It was shown that this problem can be viewed as an OMT problem [6].

Numerical methods for solving OMT problems and Monge–Ampère type equations have been scarce until recently. Benamou and Brenier introduced an augmented Lagrangian method to solve the OMT problem [7]. This approach was further developed by Haber et al. [8]. A numerical method for the Monge–Ampère equation using finite differences was introduced by Froese et al. [9,10]. This method is robust, but requires a convex target set. Brix et al. [11] solved the inverse reflector problem for a point source by using a collocation method with a tensor-product B-spline basis. For a comprehensive overview of the literature on numerical methods for the inverse reflector problem we refer to the thesis of Prins [12] and the aforementioned article by Brix et al.

In a recent publication, Prins et al. [13] introduced a least-squares method (LS method) to solve the OMT problem related to an inverse reflector problem. This method is based on a least-squares method for the Monge–Ampère equation with Dirichlet boundary conditions by Caboussat et al. [14]. The LS method solves the inverse reflector problem, i.e., the problem of finding the reflector surface that reflects a parallel bundle of light such that a prescribed luminous intensity pattern is achieved on a projection screen in the far-field of the reflector. The method can handle very complicated source and target intensities. It was used, for example, to determine the reflector surface that reflects a parallel bundle of light to form the luminous intensity pattern corresponding to a gray-scale image of a famous painting by the Dutch painter Johannes Vermeer.

The LS method determines the shape of the reflector surface by covering the light source with a rectangular grid and computing the height of the reflector in each grid point. This works fine for rectangular light sources, however, for differently shaped light sources the rectangular grid also contains grid points outside of the light source. For these grid points the emittance of the light source is taken to be zero. This approach to non-rectangular light sources is far from optimal and gives results much less satisfying than obtained for rectangular light sources. Most importantly, the boundary condition, which states that the boundary of the source must be mapped to the boundary of the target, is at places very badly satisfied and this makes the method troublesome for non-rectangular sources. This poses a severe restriction on the applicability of the method in illumination optics. The parallel bundles encountered in illumination optics often result from a converging lens and frequently have disk-shaped cross sections, therefore a numerical method that can handle disk-shaped light sources in a satisfactory way is highly desirable.

The goal of this paper is to present an improved generalized version of the LS method (GLS method) that is applicable to arbitrarily shaped light sources emitting a parallel bundle. We use some concepts from tensor calculus to derive the Monge–Ampère equation in coordinate-free form and generalize the LS method to arbitrary orthogonal coordinate systems. The GLS method like the LS method is an iterative method in which each iteration consists of three minimization steps. In one of these minimization steps a pair of boundary value problems is solved. For the LS method, these problems are decoupled because it uses Cartesian coordinates. However, in the GLS method they are coupled. We present how to deal with this issue. Furthermore, we compare the LS method from [13] with the GLS method presented in this paper. We show that for disk-shaped light sources the GLS method in polar coordinates outperforms the LS method significantly.

This paper is structured as follows. In Section 2 we derive the Monge–Ampère equation describing the reflector surface and formulate the reflector problem for an arbitrary coordinate system. In Section 3 we introduce the GLS method by generalizing the LS method to arbitrary orthogonal coordinate systems. We shed light on the different minimization steps in this method and show how they are different from the Cartesian version of the method. In Section 4 we compare the LS and GLS methods. We will consider three test cases. In all cases we will take a disk-shaped light source and therefore choose polar coordinates as the orthogonal coordinate system for the GLS method. In the first test the light source is mapped to a square gradient set, in the second test a Gaussian source emittance is used and in the third test we consider a target distribution corresponding to a lithograph by the artist M.C. Escher (Fig. 1). In Section 5 we summarize and discuss the results. This paper contains some appendices. In Appendix A we introduce some concepts from Tensor calculus needed in this paper, and give pointers to classical literature on these matters. In Appendix B one can find a proof of a result used in the main text. The reading of this proof should not be necessary for understanding the rest of the paper.

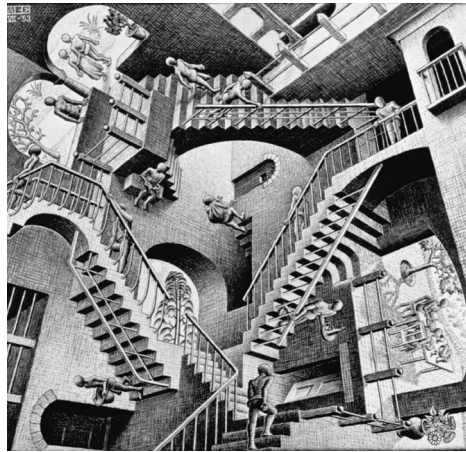


Fig. 1. Lithograph *Relativity* (1953) by the Dutch artist M.C. Escher who was frequently inspired by mathematics [15]. This lithograph, with its great detail, will serve as the ultimate test.

2. Monge–Ampère equation and inverse reflector problem

We introduce in this section the optical system and derive the corresponding Monge–Ampère equation. This formulation is not new (cf. the overview in [11] for earlier derivations), but it allows us to introduce our notation and reformulate the problem in the specific form that we will solve in Section 3.

The optical system consists of a light source and a reflector surface. We embed our optical system in three dimensional Euclidean space. We describe the light source by a set $\mathcal{E} \subset \mathbb{R}^2 \times \{-a\}$, a subset of a plane below and parallel to the x - y plane at a distance $a > 0$. For now, we assume a Cartesian coordinate system on \mathcal{E} with coordinates x and y , and, corresponding basis vectors \mathbf{e}_x and \mathbf{e}_y . Let (\cdot, \cdot) denote the Euclidean inner product on the ambient space \mathbb{R}^3 and let $\|\cdot\|$ be the corresponding norm. We assume that the light source emits a parallel bundle of light along the z -axis. The emittance of the light source at a point $\mathbf{x} = (x, y) \in \mathcal{E}$ is given in luminous flux per area by $E(x, y)$ [lm/m^2], where $E : \mathcal{E} \rightarrow (0, \infty)$ is the emittance function, which we assume to be continuous. $E(x, y)dx dy$ expresses the light flux through an infinitesimal area element on \mathcal{E} . For details on photometric quantities, see for example [16]. The light rays leaving the source will all hit upon the reflector surface. We describe the reflector surface by a function $u : \mathcal{E} \rightarrow (-a, \infty)$. A ray leaving from the point $\mathbf{x} \in \mathcal{E}$ will travel a distance $a + u(\mathbf{x})$ in the z -direction before hitting upon the reflector surface. The function u is the Monge parameterization of the reflector surface [17]. Note that by definition $u > -a$, because the reflector surface is situated above the source and not allowed to intersect with the source. In what follows we need the function u to be strictly convex and twice continuously differentiable. We will see, as a consequence of upcoming Lemma 1 and Lemma 2, that strict convexity of u implies that a pair of rays leaving \mathcal{E} from different points will be reflected in different directions. We assume the target to be positioned in the far-field of the reflector. Thus, we assume the rays after reflection to be all originating from one point and we discard the size of the reflector in this respect. In our embedding of the reflector system we let this point coincide with the origin of \mathbb{R}^3 .

The direction of reflection is given by the law of reflection, which in vector form is given by [3, p. 132]

$$\mathbf{r} = \mathbf{i} - 2(\mathbf{i}, \mathbf{n})\mathbf{n}, \tag{1}$$

where \mathbf{i} is the direction of the incoming ray, \mathbf{n} is the direction of the normal on the reflector surface and \mathbf{r} is the direction of the ray after reflection. These vectors all have unit length. The direction of an incoming ray will not depend on the point $\mathbf{x} \in \mathcal{E}$ at which it leaves the source, however, the normal \mathbf{n} on the reflector surface does depend on \mathbf{x} . The vector \mathbf{i} is the unit vector normal to the light source directed at the reflector. We denote this vector by the unit vector \mathbf{e}_z . This vector complements the two-dimensional bases on \mathcal{E} to a three-dimensional basis for \mathbb{R}^3 . The unit normal on the reflector surface pointing down towards the light source can be expressed in terms of the gradient of u and \mathbf{e}_z and when we substitute this in (1), we obtain

$$\mathbf{r}(\mathbf{x}) = \mathbf{e}_z + 2 \frac{\nabla u(\mathbf{x}) - \mathbf{e}_z}{\|\nabla u(\mathbf{x}) - \mathbf{e}_z\|^2}. \tag{2}$$

For all $\mathbf{x} \in \mathcal{E}$ the vector $\mathbf{r}(\mathbf{x})$ is of unit length and by the far-field approximation we may furthermore assume it to have its initial point at the origin. This implies that the vectors $\mathbf{r}(\mathbf{x})$ lie on the unit sphere, \mathcal{S}^2 . We can therefore interpret the map given by $\mathbf{x} \mapsto \mathbf{r}(\mathbf{x})$ to be mapping a point on the light source to a point on the unit sphere. We will denote this mapping by $r : \mathcal{E} \rightarrow \mathcal{S}^2$. The optical system is depicted in Fig. 2.

The reflected light will shine in a set of directions $\mathcal{G} \subset \mathcal{S}^2$. We assume a spherical coordinate system on \mathcal{G} with $\phi \in [0, \pi]$ the zenith angle between \mathbf{r} and \mathbf{e}_z , and, $\theta \in [0, 2\pi)$ the azimuth angle between \mathbf{e}_x and the projection of \mathbf{r} on

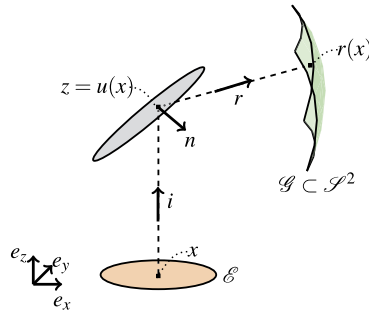


Fig. 2. The set-up of the reflector problem. A light ray from a point $x \in \mathcal{E}$ is emitted in the direction $i = \mathbf{e}_z$. Subsequently, the ray is reflected at the point $u(x)\mathbf{e}_z$ according to the law of reflection (1) in the direction \mathbf{r} . This direction corresponds to a point $r(x) \in \mathcal{G}$. The set \mathcal{G} (in green) is a subset of the unit sphere \mathcal{S}^2 . Note the unit vectors and unit sphere are not drawn in proportion. (For interpretation of the colors in the figure(s), the reader is referred to the web version of this article.)

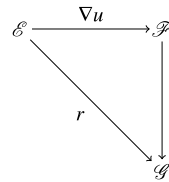


Fig. 3. The mappings and sets involved in the inverse reflector problem.

the x - y plane. Let us describe the luminous intensity in the directions \mathcal{G} by a continuous function $G : \mathcal{G} \rightarrow (0, \infty)$. The luminous intensity is the luminous flux per solid angle [lm/sr]. The luminous flux through an infinitesimal surface area element on \mathcal{G} is given by $G(\phi, \theta) \sin(\phi) d\phi d\theta$. In practice the couple \mathcal{G} and G will be such that a desired intensity pattern is projected on a screen in the far-field of the reflector. As long as \mathcal{G} is confined to one half of \mathcal{S}^2 there is a one-to-one correspondence between the couple \mathcal{G} and G on the one hand and the intensity pattern in the far-field on the other hand. Details can be found in [12]. We call \mathcal{G} the target set.

The problem we want to solve is informally stated as follows. *Given a light source \mathcal{E} with emittance function E , determine the shape of the reflector such that, after reflection, the intensity pattern in the far-field is given by the target set \mathcal{G} with luminous intensity function G .* This problem is known as the *inverse reflector problem*. Before we will state this problem more formally, we will first, under the assumption $u \in C^2(\mathcal{E})$, derive a partial differential equation from the principle of conservation of luminous flux. The luminous flux through $U \subset \mathcal{E}$ results in a luminous flux through the set $r(U) \subset \mathcal{S}^2$. By conservation of luminous flux these two fluxes must be equal and therefore we have

$$\int_U E(x, y) dx dy = \int_{r(U)} G(\phi, \theta) \sin(\phi) d\phi d\theta, \tag{3}$$

for every Lebesgue measurable set $U \subset \mathcal{E}$. We can use (3) to derive the partial differential equation. To see this we must closely examine the map $r : \mathcal{E} \rightarrow \mathcal{S}^2$. From (2) it can be seen that $\mathbf{r}(\mathbf{x})$ only depends on the gradient of u at the point \mathbf{x} . We can therefore interpret r as the composition $s \circ \nabla u$, i.e., the composition of ∇u and another map which we will denote by s . The relation between ∇u , s and r is depicted in Fig. 3. By the far-field approximation \mathbf{r} has its initial point at the origin. This implies that we should interpret ∇u also as a vector with its initial point at the origin. The vector ∇u is by definition parallel to \mathcal{E} and because it has its initial point in the origin it lies in the plane $\mathbb{R}^2 \times \{0\}$. From relation (2) we see that s maps a vector \mathbf{v} in this plane to the unit sphere according to

$$\mathbf{v} \mapsto \mathbf{e}_z + 2 \frac{\mathbf{v} - \mathbf{e}_z}{\|\mathbf{v} - \mathbf{e}_z\|^2}. \tag{4}$$

Closer inspection reveals that this map is the inverse of the stereographic projection pictured in Fig. 4 [18, p. 26]. It is the bijection between $\mathcal{S}^2 \setminus \mathbf{e}_z$, i.e., the unit-sphere without its north pole, and $\mathbb{R}^2 \times \{0\}$, the plane intersecting the equator of $\mathcal{S}^2 \setminus \mathbf{e}_z$. To substantiate this claim, consider a point on the unit sphere given by $\xi \mathbf{e}_x + \eta \mathbf{e}_y + \zeta \mathbf{e}_z$, $\xi^2 + \eta^2 + \zeta^2 = 1$. The stereographic projection maps $\xi \mathbf{e}_x + \eta \mathbf{e}_y + \zeta \mathbf{e}_z$ to the point $p(\xi, \eta, \zeta) \mathbf{e}_x + q(\xi, \eta, \zeta) \mathbf{e}_y$, with

$$p(\xi, \eta, \zeta) = \frac{\xi}{1 - \zeta}, \quad q(\xi, \eta, \zeta) = \frac{\eta}{1 - \zeta}. \tag{5}$$

We can easily verify that the inversion of these equations is given by

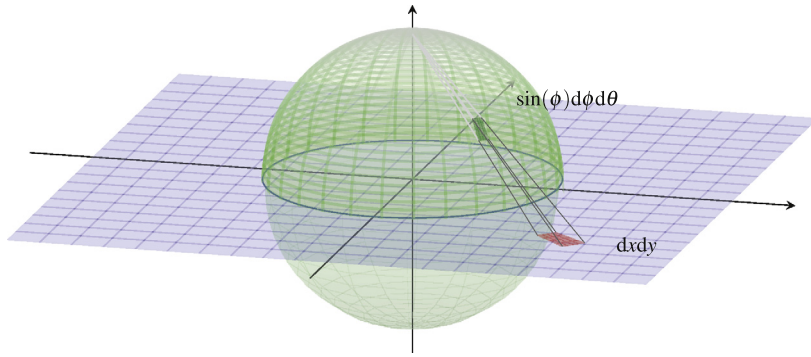


Fig. 4. The inverse stereographic projection maps the surface element $dxdy$ on $\mathbb{R}^2 \times \{0\}$ to the surface element $\sin(\phi)d\phi d\theta$ on $\mathcal{S}^2 \setminus \mathbf{e}_z$.

$$\xi(p, q) = \frac{2p}{p^2 + q^2 + 1}, \quad \eta(p, q) = \frac{2q}{p^2 + q^2 + 1}, \quad \zeta(p, q) = \frac{p^2 + q^2 - 1}{p^2 + q^2 + 1}.$$

We now see from (4) that $s : p\mathbf{e}_x + q\mathbf{e}_y \mapsto \xi(p, q)\mathbf{e}_x + \eta(p, q)\mathbf{e}_y + \zeta(p, q)\mathbf{e}_z$, which implies that s is indeed the inverse of the stereographic projection as defined in (5).

We proceed with examining the map $r : \mathcal{E} \rightarrow \mathcal{S}^2$ in order to derive the sought partial differential equation. When we identify the vector $\nabla u(\mathbf{x})$ with its endpoint, we can interpret $\nabla u(\mathcal{E})$ as a subset of $\mathbb{R}^2 \times \{0\}$. We will present bijectivity of ∇u in the following lemma.

Lemma 1. *Let $u \in C^2(\mathcal{E})$ be strictly convex and let us define*

$$p := \frac{\partial u}{\partial x} \quad \text{and} \quad q := \frac{\partial u}{\partial y}.$$

The map $\nabla u : \mathcal{E} \rightarrow \nabla u(\mathcal{E})$ is a continuously differentiable bijection, with Jacobian

$$\frac{\partial(p, q)}{\partial(x, y)} = \det \begin{pmatrix} \frac{\partial^2 u}{\partial x^2} & \frac{\partial^2 u}{\partial x \partial y} \\ \frac{\partial^2 u}{\partial x \partial y} & \frac{\partial^2 u}{\partial y^2} \end{pmatrix}, \tag{6}$$

i.e., the Jacobian of the map ∇u is the determinant of the Hessian tensor.

Proof. $u \in C^2(\mathcal{E})$ implies that ∇u is continuously differentiable. The bijectivity of ∇u follows from the strict convexity of u . ∇u is surjective by definition. To show injectivity, we argue by contradiction and will use a reasoning presented in [12, p. 93]. Suppose $\mathbf{x}, \mathbf{x}' \in \mathcal{E}$, such that $\mathbf{x} \neq \mathbf{x}'$ and $\nabla u(\mathbf{x}) = \nabla u(\mathbf{x}')$. Due to strict convexity u lies above its tangent planes, i.e., $u(\mathbf{x}') > u(\mathbf{x}) + (\nabla u(\mathbf{x}), \mathbf{x}' - \mathbf{x})$ and similarly $u(\mathbf{x}) > u(\mathbf{x}') + (\nabla u(\mathbf{x}'), \mathbf{x} - \mathbf{x}')$. Adding these two inequalities and subtracting $u(\mathbf{x}) + u(\mathbf{x}')$ from both sides we obtain $0 > (\nabla u(\mathbf{x}') - \nabla u(\mathbf{x}), \mathbf{x} - \mathbf{x}')$, which is contradicting the assumption $\nabla u(\mathbf{x}') = \nabla u(\mathbf{x})$. We have shown that ∇u is a continuously differentiable bijection. The fact that the Jacobian is the determinant of the Hessian follows directly from the definition of the Cartesian coordinates p and q on $\nabla u(\mathcal{E})$. \square

Note that the fact that u is convex implies that the Hessian is positive semi-definite and this again implies that the determinant of the Hessian is nonnegative. For the inverse of the stereographic projection, s , we have the following similar result.

Lemma 2. *The inverse of the stereographic projection $s : \mathbb{R}^2 \times \{0\} \rightarrow \mathcal{S}^2 \setminus \mathbf{e}_z$ is continuously differentiable and hence $s : \nabla u(\mathcal{E}) \rightarrow s(\nabla u(\mathcal{E}))$ is a continuously differentiable bijection. Let us denote by the Cartesian coordinates (p, q) the points in $\nabla u(\mathcal{E})$ and let $(\phi, \theta) = (\phi(p, q), \theta(p, q))$ be the image of (p, q) under s represented in spherical coordinates. For the Jacobian of s we have*

$$\sin(\phi) \frac{\partial(\phi, \theta)}{\partial(p, q)} = \frac{4}{(1 + p^2 + q^2)^2}. \tag{7}$$

Proof. We first prove injectivity. Suppose we have two distinct $\mathbf{v}, \mathbf{v}' \in \mathbb{R}^2 \times \{0\}$ such that $\mathbf{s}(\mathbf{v}) = \mathbf{s}(\mathbf{v}')$. This implies that

$$\frac{\mathbf{v} - \mathbf{e}_z}{\|\mathbf{v} - \mathbf{e}_z\|^2} = \frac{\mathbf{v}' - \mathbf{e}_z}{\|\mathbf{v}' - \mathbf{e}_z\|^2}.$$

Using the fact that \mathbf{e}_z is orthogonal to both \mathbf{v} and \mathbf{v}' we find that $\|\mathbf{v} - \mathbf{e}_z\|^2 = \|\mathbf{v}' - \mathbf{e}_z\|^2$ and this in turn implies that $\mathbf{v} = \mathbf{v}'$.

The map is also surjective from $\mathbb{R}^2 \times \{0\}$ to $\mathcal{S}^2 \setminus \mathbf{e}_z$. For $(\phi, \theta) \in (0, \pi] \times [0, 2\pi)$ we have [12, p. 77]

$$p(\phi, \theta) = \frac{\sin(\phi) \cos(\theta)}{1 - \cos(\phi)}, \quad q(\phi, \theta) = \frac{\sin(\phi) \sin(\theta)}{1 - \cos(\phi)},$$

indicating surjectivity of s . Conversely,

$$\theta(p, q) = \tan^{-1}(p, q), \quad \phi(p, q) = \arccos\left(\frac{p^2 + q^2 - 1}{p^2 + q^2 + 1}\right),$$

where $\theta(p, q) = \tan^{-1}(p, q)$ is the four-quadrant variant of $\arctan(q/p)$ [12, p. 77]. By direct calculation we find

$$\frac{\partial(\phi, \theta)}{\partial(p, q)} = \frac{1}{\sqrt{p^2 + q^2}} \frac{2}{p^2 + q^2 + 1}.$$

Furthermore, we have

$$\sin(\phi) = \sin\left(\arccos\left(\frac{p^2 + q^2 - 1}{p^2 + q^2 + 1}\right)\right) = \frac{2\sqrt{p^2 + q^2}}{p^2 + q^2 + 1},$$

where we have used the fact that $\sin(\phi) \geq 0$ and the identity $\sin(\arccos(x)) = \sqrt{1 - x^2}$. Combining the two relations we find (7). \square

Consequently, $s \circ \nabla u$ is a bijection implying that rays with different locations at the source are reflected in different directions.

We can use the results of the preceding two lemmas to derive the differential equation expressing conservation of energy. This is stated in the following theorem.

Theorem 1. Assume $\mathcal{E} \subset \mathbb{R}^2 \times \{-a\}$ is convex, closed and bounded, $u \in C^2(\mathcal{E})$ strictly convex and $r = s \circ \nabla u$, where s is the inverse of the stereographic projection. Let E and G be continuous, strictly positive and bounded functions. Furthermore, assume we have a coordinate system on \mathcal{E} with metric e_{ij} and $e = \det(e_{ij})$, and let $h_{ij}(u)$ be the coefficients of the Hessian tensor (given in Appendix A, equation (A.6)) in the basis of this coordinate system. Then the function u satisfies the Monge–Ampère type differential equation

$$\frac{4 \det(h_{ij}(u(\mathbf{x})))}{e(\mathbf{x})(1 + \|\nabla u(\mathbf{x})\|^2)^2} = \frac{E(\mathbf{x})}{G(r(\mathbf{x}))}, \tag{8}$$

for every $\mathbf{x} \in \mathcal{E}$.

Proof. We first derive an equation in the Cartesian coordinate system on \mathcal{E} and then generalize this equation to arbitrary coordinate systems on \mathcal{E} .

Both s and ∇u are continuously differentiable injections, therefore we can apply integration by substitution [19, Thm. 7.26]. For every Lebesgue measurable open subset $U \subset \mathcal{E}$ we have

$$\int_{r(U)} G(\phi, \theta) \sin(\phi) d\phi d\theta = \int_U G(\phi(x, y), \theta(x, y)) \sin(\phi) \frac{\partial(\phi, \theta)}{\partial(p, q)} \frac{\partial(p, q)}{\partial(x, y)} dx dy,$$

where no absolute values appear because both Jacobian determinants are nonnegative. Using identity (3) we find

$$\int_U E(x, y) dx dy = \int_U G(\phi(x, y), \theta(x, y)) \sin(\phi) \frac{\partial(\phi, \theta)}{\partial(p, q)} \frac{\partial(p, q)}{\partial(x, y)} dx dy,$$

for every Lebesgue measurable $U \subset \mathcal{E}$. The continuity of the functions E, G, r , the Jacobians and the sine function, Lemma 1 and Lemma 2 and the fact that the determinant of the Hessian of a convex function is positive imply that

$$\frac{4 \det(h_{ij}(u(x, y)))}{(1 + \|\nabla u(x, y)\|^2)^2} = \frac{E(x, y)}{G(\phi(x, y), \theta(x, y))} \quad \forall (x, y) \in \mathcal{E}.$$

This equation is the Cartesian coordinate expression of (8), because in the Cartesian coordinate system $e = 1$ and the Hessian is just the matrix with the second order partial derivatives.

The left hand side of (8) is a scalar and hence independent of the choice of coordinate system and basis. The term $4/(1 + \|\nabla u(\mathbf{x})\|^2)^2$ is also independent of the coordinate system and basis in use as it only involves the norm of a vector. To see that this also holds for $\det(h_{ij}(u(\mathbf{x}))) / (e(\mathbf{x}))$, we consider a basis transformation $a_j^i \mathbf{e}_i = \bar{\mathbf{e}}_j$, where $\mathbf{e}_1, \mathbf{e}_2$ are the old

basis vectors and \bar{e}_1, \bar{e}_2 are the new basis vectors. Let $B = (b_j^i)$ be the inverse of $A = (a_j^i)$. The Hessian and the metric are both tensors, as such they transform according to the tensor transformation law [20, p. 204], i.e., we have $\bar{h}_{ij} = b_i^k b_j^l h_{kl}$ and $\bar{e}_{ij} = b_i^k b_j^l e_{kl}$. From this it follows that $\det(\bar{h}_{ij}) = (\det(B))^2 \det(h_{ij})$ and $\bar{e} = (\det(B))^2 e$ for the determinants of the metrics, which implies $\det(\bar{h}_{ij})/\bar{e} = \det(h_{ij})/e$. Thus we see that (8) is indeed independent of the choice of coordinate system on \mathcal{E} . \square

We are now in the position to state the inverse reflector problem in more formal terms. *Given a convex, closed and bounded light source \mathcal{E} with strictly positive and bounded emittance $E \in C(\mathcal{E})$ and a closed target set $\mathcal{G} \subset \mathcal{S}^2$ with desired strictly positive and bounded luminous intensity $G \in C(\mathcal{G})$ such that*

$$\int_{\mathcal{E}} E(x, y) dx dy = \int_{\mathcal{G}} G(\phi, \theta) \sin(\phi) d\phi d\theta, \tag{9}$$

find a function $u \in C^2(\mathcal{E})$ that satisfies $r(\mathcal{E}) = \mathcal{G}$ and the Monge–Ampère type equation (8).

The condition $r(\mathcal{E}) = \mathcal{G}$ needs to be satisfied for equation (8) to have meaning. Alternatively, we can use the continuously differentiable bijection s to reformulate the problem in terms of a gradient set \mathcal{F} and a function F on this set instead of the target set \mathcal{G} and the luminous intensity G . Using the fact that s^{-1} exists, we define $\mathcal{F} := s^{-1}(\mathcal{G})$, and, using the differentiability of s , we define

$$F(\mathbf{y}) = \frac{4G(s(\mathbf{y}))}{(1 + \|\mathbf{y}\|^2)^2},$$

for all $\mathbf{y} \in \mathcal{F}$. Note that $F \in C(\mathcal{F})$. Furthermore, using integration by substitution we see that (9) implies

$$\int_{\mathcal{E}} E(x, y) dx dy = \int_{\mathcal{F}} F(p, q) dp dq. \tag{10}$$

The conditions $r(\mathcal{E}) = \mathcal{G}$ translates in the condition $\nabla u(\mathcal{E}) = \mathcal{F}$. These definitions allow us to reformulate the inverse reflector problem.

INVERSE REFLECTOR PROBLEM. *Given a convex, closed and bounded light source \mathcal{E} with strictly positive and bounded emittance $E \in C(\mathcal{E})$ and a closed gradient set \mathcal{F} with strictly positive, bounded and bounded away from zero, function $F \in C(\mathcal{F})$ that satisfy (10), find a function $u \in C^2(\mathcal{E})$ that satisfies $\nabla u(\partial\mathcal{E}) = \partial\mathcal{F}$, $e^{ij}h_{ij}(u) > 0$ and the Monge–Ampère type equation*

$$\frac{\det(h_{ij}(u(\mathbf{x})))}{e(\mathbf{x})} = \frac{E(\mathbf{x})}{F(\nabla u(\mathbf{x}))}. \tag{11}$$

Note that we replaced the implicit boundary condition $\nabla u(\mathcal{E}) = \mathcal{F}$ with the more explicit boundary condition $\nabla u(\partial\mathcal{E}) = \partial\mathcal{F}$. The explicit boundary condition is better manageable numerically. We will show in Appendix B that for strictly convex u these two conditions are equivalent. The reason that we demand F to be bounded away from zero, is to be able to show this equivalence. Furthermore, we added the earlier absent constraint $e^{ij}h_{ij}(u) = h_i^i(u) > 0$, demanding the trace (cf. Appendix (A.2)) of the Hessian to be strictly positive. The fraction E/F is by definition strictly positive, hence the determinant of the Hessian is strictly positive too. From this it follows that the Hessian matrix is strictly positive definite or strictly negative definite, corresponding to either a convex or concave solution, respectively. By demanding the trace of the Hessian to be positive we make sure that only a convex reflector surface is admitted. Thus, we conclude that a solution to the Inverse Reflector Problem as stated above and the same problem but with $\nabla u(\partial\mathcal{E}) = \partial\mathcal{F}$ replaced by $\nabla u(\mathcal{E}) = \mathcal{F}$ are truly equivalent as they both only admit strictly convex solutions and the boundary conditions are equivalent for strictly convex u .

In this paper we restrict ourselves to this convex solution, however, the algorithm can be easily adapted to find the concave solution instead. In [12, pp. 96–98] it is described how one can easily find the concave solution from the convex solution and vice versa.

A theorem by Brenier [5, p. 66] states that a weak formulation of the Inverse Reflector Problem admits a unique convex solution. It is, however, not clear that, for all pairs (\mathcal{E}, E) and (\mathcal{F}, F) , the unique weak solution of Brenier’s theorem is twice continuously differentiable. Regularity of the solution for the Monge–Ampère equation (11) is a complicated matter and beyond the scope of this paper. For u to be $C^2(\mathcal{E})$, certainly continuity of E and F seems to be required (see for example [5]). However, even this requirement will not be met in all the numerical experiments we present in Section 4. Though also in these cases good numerical approximations are found.

3. Least-squares method in arbitrary orthogonal coordinates

In [13] Prins et al. proposed the LS method to solve the Inverse Reflector Problem in Cartesian coordinates. We introduce in this section the GLS method, i.e., the generalization of the LS method by Prins et al. [13] to arbitrary *orthogonal* coordinate systems. We remind the reader that a brief explanation of the notions of tensor calculus used can be found in Appendix A.

We assume an arbitrary coordinate system on the source \mathcal{E} with orthogonal coordinates x^1, x^2 , local orthogonal basis vectors $\mathbf{e}_1, \mathbf{e}_2$ and a metric $e_{ij} = (\mathbf{e}_i, \mathbf{e}_j)$ with corresponding norm $\|\cdot\|$. The orthogonality of the basis vectors imply $e_{ij} = 0$ for $i \neq j$. We will not try to solve the Inverse Reflector Problem directly for u , but instead look for a mapping $\mathbf{m} : \mathcal{E} \rightarrow \mathcal{F}$ representing ∇u such that:

(i) \mathbf{m} solves the following boundary value problem

$$\frac{\det(\nabla \hat{\mathbf{m}}(\mathbf{x}))}{e(\mathbf{x})} = \frac{E(\mathbf{x})}{F(\mathbf{m}(\mathbf{x}))}, \quad \mathbf{x} \in \mathcal{E},$$

$$\mathbf{m}(\partial \mathcal{E}) = \partial \mathcal{F},$$

where $\hat{\mathbf{m}} = m_i \mathbf{e}^i = e_{ij} m^j \mathbf{e}^i$ and $\mathbf{m} = m^i \mathbf{e}_i$,

(ii) \mathbf{m} should be such that there exists a strictly convex $u \in C^2(\mathcal{E})$ such that $\mathbf{m} = \nabla u$.

From this mapping we will eventually find u . If \mathbf{m} satisfies (ii), then $\hat{\mathbf{m}} = du$ and hence the tensor

$$\nabla \hat{\mathbf{m}} := \nabla_{\mathbf{e}_j}(\hat{\mathbf{m}}) \otimes \mathbf{e}^j = (\nabla_{\mathbf{e}_j} m_i - \Gamma_{ij}^k m_k) \mathbf{e}^i \otimes \mathbf{e}^j$$

must be, by definition (Appendix A.3), the Hessian of some function and therefore needs to be symmetric (Appendix A.3). This condition is actually enough to ensure that \mathbf{m} equals the gradient of some function. The symmetry of $\nabla \hat{\mathbf{m}}$ implies $\nabla \times \mathbf{m} = 0$. To see this let us interpret \mathbf{m} as a vector in \mathbb{R}^3 . The component $m^3 = 0$ and hence the curl is given by [21, p. 170]

$$\nabla \times \mathbf{m} = \frac{1}{\sqrt{e}} \left((\nabla_{\mathbf{e}_2} m_1 - \Gamma_{12}^i m_i) - (\nabla_{\mathbf{e}_1} m_2 - \Gamma_{21}^i m_i) \right) \mathbf{e}_z.$$

From this we see that $\nabla \times \mathbf{m}$ vanishes if and only if $\nabla \hat{\mathbf{m}}$ is symmetric. A vector field with zero curl is called *conservative*. A conservative field on a simply connected domain always equals the gradient of some function, see for example [23, p. 551]. Thus we conclude that \mathbf{m} equals the gradient of some function $u \in C^2(\mathcal{E})$ if and only if $\nabla \hat{\mathbf{m}}$ is symmetric.

However, this condition alone will not suffice for our goals, because we also need u to be strictly convex. The function $u \in C^2(\mathcal{E})$ is convex if and only if \mathcal{E} is convex and the Hessian tensor $\mathbf{H}(u)$ is *positive semi-definite*, see for example [24, p. 71]. The Hessian tensor is positive semi-definite if and only if for every $\mathbf{x} = x^i \mathbf{e}_i$ we have $H(u)(\mathbf{x}, \mathbf{x}) \geq 0$, where

$$\mathbf{H}(u)(\mathbf{x}, \mathbf{x}) = h_{ij} x^i x^j = x_k e^{ki} h_{ij} x^j = \mathbf{x}^T (e^{ki} h_{ij}) \mathbf{x}.$$

From this we see that $\mathbf{H}(u)$ is positive semi-definite if and only if the matrix $(e^{ki} h_{ij})$ is positive semi-definite. For our orthogonal basis the metric is diagonal and therefore

$$(e^{ki} h_{ij}) = \begin{pmatrix} e^{11} h_{11} & e^{11} h_{12} \\ e^{22} h_{21} & e^{22} h_{22} \end{pmatrix}.$$

Unfortunately, we can not demand positive definiteness, because, although every $u \in C^2(\mathcal{E})$ with positive definite Hessian tensor is strictly convex, not every strictly convex $u \in C^2(\mathcal{E})$ has a positive definite Hessian tensor.¹ Thus asking for more than $\nabla \hat{\mathbf{m}}$ to be positive semi-definite would be too restrictive. The numerical method that we will introduce solves the following boundary value problem (BVP):

TRANSPORT BVP. Find a continuously differentiable map $\mathbf{m} : \mathcal{E} \rightarrow \mathcal{F}$ that satisfies

$$\frac{\det(\nabla \hat{\mathbf{m}}(\mathbf{x}))}{e(\mathbf{x})} = \frac{E(\mathbf{x})}{F(\mathbf{m}(\mathbf{x}))}, \quad \mathbf{x} \in \mathcal{E}, \tag{12a}$$

$$\mathbf{m}(\partial \mathcal{E}) = \partial \mathcal{F}, \tag{12b}$$

and for which $\nabla \hat{\mathbf{m}}$ is a symmetric positive semi-definite tensor. In this problem the functions E and F are strictly positive and bounded on \mathcal{E} and \mathcal{F} , respectively, such that (10) is satisfied and F is bounded away from zero.

¹ Consider for example the strictly convex function $f(x) = x^4$ on the real line. Although f is strictly convex, the Hessian tensor, i.e. f'' , is zero for $x = 0$ and hence not positive definite.

If u is a solution to the Inverse Reflector Problem, then u is strictly convex and $\mathbf{m} = \nabla u$ will be a solution to the Transport BVP. The reverse statement is not necessarily true because a solution \mathbf{m} of the Transport BVP may be such that u satisfying $\mathbf{m} = \nabla u$ is convex but not strictly convex. The Transport BVP thus admits for a slightly larger solution class.

We numerically solve the Transport BVP by starting with an initial guess \mathbf{m}^0 and improve this initial guess in an iterative manner. One iteration consists of three stages. These three stages together give an improved mapping \mathbf{m}^{n+1} from the current mapping \mathbf{m}^n . We now explain these three stages.

First, we approximate $\nabla \hat{\mathbf{m}}$ by a symmetric positive semi-definite tensor \mathbf{P} by minimizing the functional

$$J_1(\mathbf{m}, \mathbf{P}) := \frac{1}{2} \int_{\mathcal{E}} \|\nabla \hat{\mathbf{m}} - \mathbf{P}\|^2 \, dA, \tag{13a}$$

over the space

$$\mathcal{P}(\mathbf{m}) := \left\{ \mathbf{P} \in \mathbf{T}_2^0(\mathcal{E})_{C^1} \mid \det(p_{ij}(\mathbf{x})) = \frac{e(\mathbf{x})E(\mathbf{x})}{F(\mathbf{m}(\mathbf{x}))}, \mathbf{P}(\mathbf{x}) \text{ is spsd} \right\}, \tag{13b}$$

where “spsd” stands for symmetric positive semi-definite and $\mathbf{P} = p_{ij} \mathbf{e}^i \otimes \mathbf{e}^j$. Furthermore, we use $\mathbf{T}_2^0(\mathcal{E})_{C^1}$ to denote the space of continuously differentiable tensor fields of contravariant rank 0 and covariant rank 2 (Appendix A.2). It seems as if we demand more smoothness than necessary, because Transport BVP and (13a) suggest that \mathbf{P} only needs to be continuous for continuous E and F . However, in one of the minimization procedures that follows we need $\nabla \hat{\mathbf{m}}$ to be continuously differentiable and therefore we also need \mathbf{P} to be continuously differentiable.

The norm in functional (13a) is defined in the following way. Let $\mathbf{A}, \mathbf{B} \in \mathbf{T}_2^0(\mathcal{E})$, where $\mathbf{A} = a_{ij} \mathbf{e}^i \otimes \mathbf{e}^j$ and $\mathbf{B} = b_{ij} \mathbf{e}^i \otimes \mathbf{e}^j$, then

$$\mathbf{A} : \mathbf{B} := e^{ik} e^{jl} a_{ij} b_{kl} = a^{ij} b_{ij}, \tag{14}$$

defines an inner product on $\mathbf{T}_2^0(\mathcal{E})$. This inner product on $\mathbf{T}_2^0(\mathcal{E})$ is induced by the metric. The fact that this is indeed an inner product follows from the symmetry, linearity and positivity of the metric \mathbf{e} . Let $\|\cdot\|$ be the norm associated with this inner product.² It is clear that if $J_1 = 0$, \mathbf{m} will satisfy (12a) and $\nabla \hat{\mathbf{m}}$ will be symmetric positive semi-definite.

Secondly, we approximate the boundary value of \mathbf{m} at $\partial \mathcal{E}$ by a vector field \mathbf{b} that exactly satisfies the boundary condition (12b). This is done by minimizing

$$J_B(\mathbf{m}, \mathbf{b}) := \frac{1}{2} \int_{\partial \mathcal{E}} \|\mathbf{m} - \mathbf{b}\|^2 \, ds, \tag{15a}$$

over the space

$$\mathcal{B} := \left\{ \mathbf{b} = \bar{\mathbf{b}}|_{\partial \mathcal{E}} \mid \bar{\mathbf{b}} \in \mathbf{T}_0^1(\mathbb{R}^2)_C, \mathbf{b}(\mathbf{x}) \in \partial \mathcal{F} \right\}, \tag{15b}$$

for an arc-length parameterization of the boundary. Analogously to the functional J_1 we notice that if $J_B = 0$, \mathbf{m} satisfies (12b).

Finally, to find the improved mapping \mathbf{m}^{n+1} we minimize J_1 and J_B simultaneously. To do so we define a third functional:

$$J(\mathbf{m}, \mathbf{P}, \mathbf{b}) := \alpha J_1(\mathbf{m}, \mathbf{P}) + (1 - \alpha) J_B(\mathbf{m}, \mathbf{b}) \tag{16}$$

with $\alpha \in (0, 1)$. We minimize this functional for \mathbf{m} over the space $\mathcal{M} := \mathbf{T}_0^1(\mathcal{E})_{C^2}$. Thus, to determine \mathbf{m}^{n+1} three stages are performed subsequently:

$$\mathbf{b}^{n+1} = \operatorname{argmin}_{\mathbf{b} \in \mathcal{B}} J_B(\mathbf{m}^n, \mathbf{b}), \tag{17a}$$

$$\mathbf{P}^{n+1} = \operatorname{argmin}_{\mathbf{P} \in \mathcal{P}(\mathbf{m}^n)} J_1(\mathbf{m}^n, \mathbf{P}), \tag{17b}$$

$$\mathbf{m}^{n+1} = \operatorname{argmin}_{\mathbf{m} \in \mathcal{M}} J(\mathbf{m}, \mathbf{P}^{n+1}, \mathbf{b}^{n+1}). \tag{17c}$$

To solve these minimization problems we will cover our light source with a grid. The grid will be an orthogonal curvilinear grid with as grid lines a finite set of the coordinate lines. The minimization problems (17c) will then be translated to discrete problems on this grid.

The first minimization step, step (17a), can be performed in an efficient point-wise way as discussed in [13]. No changes are made to this minimization step and therefore we do not further discuss it here. The minimization step (17b) is discussed

² We use the same notation as for the vector norm, but this is not very likely to cause confusion because it will be clear from the argument which norm we mean.

quite extensively. A new geometrical interpretation of this minimization is presented, which provides increased insight and clarifies the intricate expressions of [13]. This allows us to algebraically determine the minimizer for this problem. Also minimization problem (17c) is covered in great detail, because this minimization problem becomes substantially more involved for arbitrary coordinate systems. We start with minimization problem (17b).

3.1. Minimization of J_1

The integrand of J_1 does not contain derivatives of \mathbf{P} , therefore we can carry out the minimization for each grid point $\mathbf{x} \in \mathcal{E}$ individually. For each grid point $\mathbf{x} \in \mathcal{E}$ we want to minimize $\|\nabla \hat{\mathbf{m}}(\mathbf{x}) - \mathbf{P}(\mathbf{x})\|^2/2$. Let us denote by $\delta_{\mathbf{e}_i} m_j$ the central difference approximation of $\nabla_{\mathbf{e}_i} m_j$, i.e., the difference of the value in two neighboring grid points in the \mathbf{e}_i direction divided by the distance between those two points. The tensor $\nabla \hat{\mathbf{m}}$ will then be approximated by $\mathbf{D} = d_{ij} \mathbf{e}^i \otimes \mathbf{e}^j$, where $d_{ij} := \delta_{\mathbf{e}_j} m_i - \Gamma_{ij}^k m_k$. Assuming this approximation of $\nabla \hat{\mathbf{m}}$, we will minimize

$$\begin{aligned} \frac{1}{2} \|\mathbf{D} - \mathbf{P}\|^2 &= \frac{1}{2} e^{ik} e^{jl} (d_{ij} - p_{ij})(d_{kl} - p_{kl}) \\ &= \frac{1}{2e} \left(e^{11} e_{22} (d_{11} - p_{11})^2 + (d_{12} - p_{12})^2 + (d_{21} - p_{12})^2 + e^{22} e_{11} (d_{22} - p_{22})^2 \right), \end{aligned}$$

where we used the fact that the basis $\{\mathbf{e}_1, \mathbf{e}_2\}$ is orthogonal and hence (e_{ij}) is diagonal. The tensor $\mathbf{P}(\mathbf{x})$ is positive semi-definite if and only if the matrix $(e^{ij} p_{jk})$ is positive semi-definite. Recall that symmetric 2×2 matrices are positive semi-definite if and only if their trace and determinant are both nonnegative. However, the matrix $(e^{ij} p_{jk})$ is not symmetric, because

$$(e^{ij} p_{jk}) = \begin{pmatrix} e^{11} p_{11} & e^{11} p_{12} \\ e^{22} p_{12} & e^{22} p_{22} \end{pmatrix}, \tag{18}$$

where we used that $p_{21} = p_{12}$. Let the transformation matrix T be given by $T = \text{diag}(\sqrt{e_{11}}, \sqrt{e_{22}})$. We use this transformation to make $(e^{ij} p_{jk})$ symmetric:

$$T(e^{ij} p_{jk})T^{-1} = \begin{pmatrix} e^{11} p_{11} & p_{12}/\sqrt{e} \\ p_{12}/\sqrt{e} & e^{22} p_{22} \end{pmatrix}. \tag{19}$$

A quick calculation shows that the eigenvalues of the matrix $(e^{ij} p_{jk})$, and hence also of the matrix $T(e^{ij} p_{jk})T^{-1}$, are given by

$$\mu_{\pm} = \frac{1}{2e} \left(e_{22} p_{11} + e_{11} p_{22} \pm \sqrt{(e_{22} p_{11} + e_{11} p_{22})^2 - 4e \det(p_{ij})} \right), \tag{20}$$

which are both real since the matrix $T(e^{ij} p_{jk})T^{-1}$ is symmetric. It is a familiar result that a matrix is positive semi-definite if and only if its eigenvalues are nonnegative. The matrix $(e^{ij} p_{jk})$ is positive semi-definite if and only if the matrix in (19) is positive semi-definite. The matrix in (19) is symmetric, hence we can conclude that $\mathbf{P}(\mathbf{x})$ is positive semi-definite if and only if the trace and determinant of the matrix in (19) are nonnegative, i.e., if and only if $e^{11} p_{11} + e^{22} p_{22} \geq 0$ and $(p_{11} p_{22} - p_{12}^2)/e \geq 0$. The metric e_{ij} is derived from an ordinary Pythagorean inner product hence we have $e > 0$ and therefore we can simplify the last requirement to $\det(p_{ij}) \geq 0$.

The determinant of (p_{ij}) needs to equal eE/F . This quotient is positive by definition and hence $\det(p_{ij}) > 0$ is always satisfied. Let us now, to get rid of the metric altogether, introduce the variables

$$\bar{p}_{11} := e^{11} p_{11}, \quad \bar{p}_{12} := p_{12}/\sqrt{e}, \quad \bar{p}_{22} := e^{22} p_{22}, \tag{21a}$$

$$\bar{d}_{11} := e^{11} d_{11}, \quad \bar{d}_{22} := e^{22} d_{22}, \quad \bar{d}_{12} := (d_{12} + d_{21})/(2\sqrt{e}). \tag{21b}$$

We can give a more convenient reformulation of the minimization problem in terms of these variables. Moreover, we also drop the constant term $(d_{12} - d_{21})^2/(4e)$ from the function to minimize. We may do this as it does not effect the minimizers. The reformulated problem then reads as follows.

MINIMIZATION PROBLEM. Given the symmetric matrix

$$\bar{\mathbf{D}} = \begin{pmatrix} \bar{d}_{11} & \bar{d}_{12} \\ \bar{d}_{12} & \bar{d}_{22} \end{pmatrix},$$

with \bar{d}_{11} , \bar{d}_{12} and \bar{d}_{22} as defined in (21b), find the symmetric matrix

$$\bar{P} = \begin{pmatrix} \bar{p}_{11} & \bar{p}_{12} \\ \bar{p}_{12} & \bar{p}_{22} \end{pmatrix},$$

that minimizes the function

$$H(\bar{P}) := \frac{1}{2} \|\bar{D} - \bar{P}\|_F^2, \tag{22}$$

under the constraints $\det(\bar{P}) = E/F$ and $\text{tr}(\bar{P}) \geq 0$, where the norm used in (22) is the Frobenius norm for matrices, defined as $\|A\|_F = \sqrt{\sum_{i,j} a_{ij}^2}$ for a matrix $A = (a_{ij})$.

From the relations (21a) the minimizer (p_{11}, p_{12}, p_{22}) can be found once the minimizer $(\bar{p}_{11}, \bar{p}_{12}, \bar{p}_{22})$ of Minimization Problem has been found. Furthermore, we have $H(\bar{P}) = \|\mathbf{D} - \mathbf{P}\|^2/2 - (d_{12} - d_{21})^2/(4e)$. We solve Minimization Problem algebraically by using the method of Lagrange multipliers. Besides this we give a graphical representation of this problem. This serves to get more intuition for the problem and also provides a convenient way to verify the algebraically found solutions.

3.1.1. Lagrange minimizers and their geometric representation

We find the minimizers of Minimization Problem with the help of the Lagrange function

$$\Lambda(\bar{P}; \lambda) = H(\bar{P}) + \lambda \left(\det \bar{P} - \frac{E}{F} \right). \tag{23}$$

In a minimum of this function all the partial derivatives have to equal zero, hence we find the following set of equations,

$$\bar{p}_{11} + \lambda \bar{p}_{22} = \bar{d}_{11}, \tag{24a}$$

$$\lambda \bar{p}_{11} + \bar{p}_{22} = \bar{d}_{22}, \tag{24b}$$

$$(1 - \lambda) \bar{p}_{12} = \bar{d}_{12}, \tag{24c}$$

$$\bar{p}_{11} \bar{p}_{22} - \bar{p}_{12}^2 = E/F. \tag{24d}$$

In the Lagrange function (23) the condition $\text{tr}(\bar{P}) \geq 0$ has not been taken into account, hence a solution of (24a)–(24d) might have $\text{tr}(\bar{P}) < 0$. In what follows, we will show that there always exists a solution to (24a)–(24d) such that $\text{tr}(\bar{P}) \geq 0$.

Let us now give a geometric interpretation of the minimizers of the Lagrange function. The minimizers correspond to a joint tangent plane of a hyperboloid and an ellipsoid. We introduce the iso-surfaces $H(\bar{P}) = C_H$. By definition $C_H \geq 0$. Every value of C_H corresponds to an iso-surface of the function H . By definition of H we have

$$\left(\frac{\bar{p}_{11} - \bar{d}_{11}}{\sqrt{2C_H}} \right)^2 + \left(\frac{\bar{p}_{12} - \bar{d}_{12}}{\sqrt{C_H}} \right)^2 + \left(\frac{\bar{p}_{22} - \bar{d}_{22}}{\sqrt{2C_H}} \right)^2 = 1. \tag{25}$$

Equation (25) describes an ellipsoid in \mathbb{R}^3 with center $(\bar{d}_{11}, \bar{d}_{12}, \bar{d}_{22})$ and semi-axes $\sqrt{2C_H}$, $\sqrt{C_H}$ and $\sqrt{2C_H}$. Thus the iso-surfaces of H can be interpreted as ellipsoids in \mathbb{R}^3 .

The constraint $\det(\bar{P}) = E/F$ describes an hyperboloid in \mathbb{R}^3 with symmetry axes given by $\bar{p}_{11} = \bar{p}_{22}$ and $\bar{p}_{12} = 0$. To see this we will rotate our coordinate system to align the symmetry axes with the new coordinate axes. We perform the rotation given by

$$\begin{aligned} p_1 &:= (\bar{p}_{11} - \bar{p}_{22})/\sqrt{2}, & p_2 &:= \text{tr}(\bar{P})/\sqrt{2}, & p_3 &:= \bar{p}_{12}, \\ d_1 &:= (\bar{d}_{11} - \bar{d}_{22})/\sqrt{2}, & d_2 &:= \text{tr}(\bar{D})/\sqrt{2}, & d_3 &:= \bar{d}_{12}. \end{aligned}$$

Using this transformation, the constraint $\det(\bar{P}) = E/F$ can be rewritten as

$$\left(\frac{p_1}{\sqrt{2E/F}} \right)^2 - \left(\frac{p_2}{\sqrt{2E/F}} \right)^2 + \left(\frac{p_3}{\sqrt{E/F}} \right)^2 = -1. \tag{26}$$

This equation describes a hyperboloid of two separate sheets. One sheet is located in the half-space $p_2 > 0$ and the other one is located in the half-space $p_2 < 0$. The distance from the origin to the extremum of the sheet with $\text{tr}(\bar{P}) > 0$ and the extremum of the sheet with $\text{tr}(\bar{P}) < 0$ is both $\sqrt{2E/F}$. Equation (25) transforms to

$$\left(\frac{p_1 - d_1}{\sqrt{2C_H}} \right)^2 + \left(\frac{p_2 - d_2}{\sqrt{2C_H}} \right)^2 + \left(\frac{p_3 - d_3}{\sqrt{C_H}} \right)^2 = 1.$$

We see (Fig. 5) that the principal axes of both the ellipsoids and the hyperboloids are such that the p_1 - and p_2 -principal axis are equally long and $\sqrt{2}$ times the length of the p_3 -principal axis. This fact will play a role in the minimization problem.

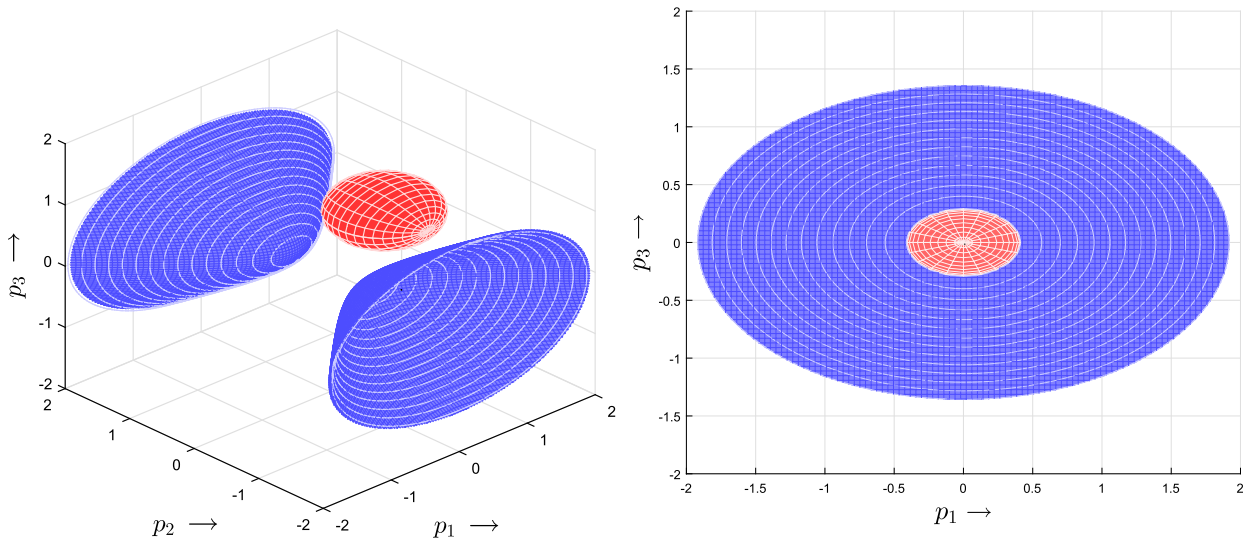


Fig. 5. An example of an ellipsoidal iso-surface of H and a hyperboloid are shown from two different perspectives for two instances of \bar{D} . We see that the principal p_1 - and p_3 -axis have the same proportion for the hyperboloid and the ellipsoid.

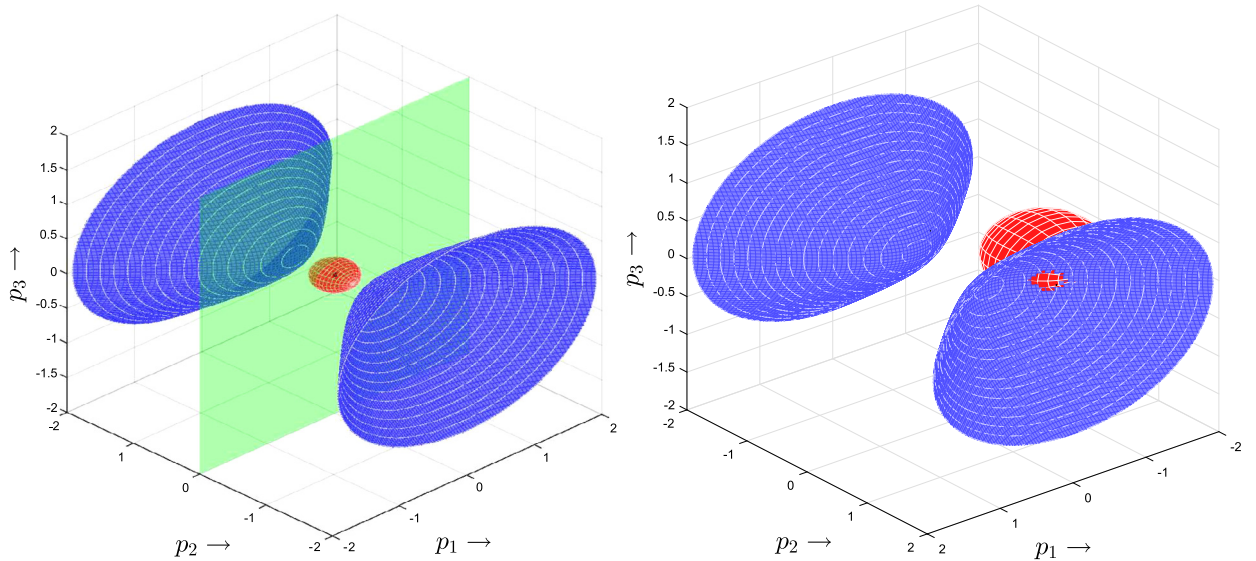


Fig. 6. On the left side the two sheets of the hyperboloid and the dividing plane $p_2 = \text{tr}(\bar{P})/\sqrt{2} = 0$ are shown. On the right side an example of an ellipsoid which is tangent to the hyperboloid with $\text{tr}(\bar{P}) > 0$ is shown. Some red of the ellipsoid can be seen through the hyperboloid. This point is the minimizer.

The local minimizers of the Lagrange function (23) are exactly the points where an iso-surface of H is tangent to the hyperboloid. The plane $p_2 = \text{tr}(\bar{P})/\sqrt{2} = 0$ lies precisely between the two sheets of the hyperboloid. Thus, only the points where an iso-surface of H is tangent to the sheet of the hyperboloid with $\text{tr}(\bar{P}) > 0$ are actual minimizers of Minimization Problem. In Fig. 6 this is illustrated. The global minimizer corresponds to the smallest ellipsoid that is tangent to the upper sheet of the hyperboloid.

In the remaining part of this section we will algebraically solve the system (24). We will verify the algebraic solutions that we find by these graphical representations. This allows us to get more intuition for the problem and visualize symmetries that are not directly apparent from (24a)–(24d).

3.1.2. Determining the minimizers

We will show that for each given \bar{D} we can find \bar{P} that is the solution of Minimization Problem. If $\lambda \neq \pm 1$, we can invert (24a)–(24c). Doing this we obtain

$$\bar{p}_{11} = \frac{\lambda \bar{d}_{22} - \bar{d}_{11}}{\lambda^2 - 1}, \quad \bar{p}_{12} = \frac{\bar{d}_{12}}{1 - \lambda}, \quad \bar{p}_{22} = \frac{\lambda \bar{d}_{11} - \bar{d}_{22}}{\lambda^2 - 1}. \tag{27}$$

However, these equations only hold if $\lambda \neq \pm 1$. From (24a)–(24c) we have the following immediate logical implications:

$$\lambda = 1 \implies (\bar{d}_{11} = \bar{d}_{22} \wedge \bar{d}_{12} = 0), \quad \lambda = -1 \implies (\bar{d}_{11} = -\bar{d}_{22}).$$

From these implications we see there are only two situations that have to be dealt with separately, namely the cases $(\bar{d}_{11} = \bar{d}_{22} \wedge \bar{d}_{12} = 0)$ and $(\bar{d}_{11} = -\bar{d}_{22})$. When we are not in one of these two cases, the solution (27) holds. We will now treat the three different cases in turn, starting out with the general case.

Lemma 3. *If $(\bar{d}_{11} \neq \bar{d}_{22} \vee \bar{d}_{12} \neq 0)$ and $(\bar{d}_{11} \neq -\bar{d}_{22})$, the global minimizer to Minimization Problem is given by (27), where λ is given by one of the following expressions:*

$$\begin{aligned} \lambda_i &= -\sqrt{\frac{y}{2}} + (-1)^i \sqrt{-\frac{y}{2} - \frac{a_2}{2a_4} + \frac{a_1}{2a_4\sqrt{2y}}}, & i = 1, 2, \\ \lambda_i &= \sqrt{\frac{y}{2}} + (-1)^i \sqrt{-\frac{y}{2} - \frac{a_2}{2a_4} - \frac{a_1}{2a_4\sqrt{2y}}}, & i = 3, 4. \end{aligned} \tag{28}$$

In (28) y is given by the following two sets of equations:

$$\begin{aligned} y &= A + \frac{Q}{A} - \frac{b_2}{3}, & A &= -\operatorname{sgn}(R)(|A| + \sqrt{R^2 - Q^3})^{1/3}, \\ R &= \frac{2b_2^3 - 9b_1b_2 + 27b_0}{54}, & Q &= \frac{b_2^3 - 3b_1}{9}, \end{aligned} \tag{29}$$

and

$$\begin{aligned} a_4 &= \frac{E}{F}, & a_2 &= -2a_4 - \det(\bar{D}), & a_1 &= \|\bar{D}\|^2, & a_0 &= a_4 - \det(\bar{D}), \\ b_0 &= -\frac{a_1^2}{8a_4^2}, & b_1 &= \frac{a_2^2 - 4a_0a_4}{4a_4^2}, & b_2 &= \frac{a_2}{a_4}. \end{aligned} \tag{30}$$

At least one of the four choices for λ is such that the requirement $\operatorname{tr}(\bar{P}) > 0$ is satisfied by (27).

Proof. Substituting the expressions (27) in (24d) we obtain the following quartic equation $\Pi(\lambda) := a_4\lambda^4 + a_2\lambda^2 + a_1\lambda + a_0 = 0$, where the coefficients are as given in (30). In [13] it is shown that this polynomial admits the four solutions (28). Since $a_4 = E/F > 0$ we have $\lim_{\lambda \rightarrow \pm\infty} \Pi(\lambda) = \infty$. Furthermore, we can rewrite $\Pi(\lambda)$ as

$$\Pi(\lambda) = a_4(\lambda^2 - 1)^2 - (a_0 - a_4)(\lambda^2 + 1) + a_1\lambda.$$

From this we see that $\Pi(-1) = -(\bar{d}_{11} + \bar{d}_{22})^2$. By assumption $\bar{d}_{11} \neq -\bar{d}_{22}$, hence $\Pi(-1) < 0$. From this inequality combined with the fact that $\Pi(\lambda) \rightarrow +\infty$ for $\lambda \rightarrow \pm\infty$ it follows by the Intermediate Value Theorem that Π must have at least two real roots, one smaller than -1 and one larger than -1 . From (24) it follows that $\operatorname{tr}(\bar{P}) = \operatorname{tr}(\bar{D})/(1 + \lambda)$. This shows that for one of the two real roots $\operatorname{tr}(\bar{P}) > 0$, while for the other real root $\operatorname{tr}(\bar{P}) < 0$.

We now have established the fact that one of the four roots λ in (28) gives the minimum of the Lagrange function such that $\operatorname{tr}(\bar{P}) > 0$, thereby it follows that a global minimizer exists. Moreover, the minimizer is given by (27), with λ given by one of the real roots of (28). The global minimizer will be found by checking for which of the four λ_i ($i = 1, \dots, 4$) the function H is minimal. \square

Now that we have dealt with the general case we will turn our attention to the cases $(\bar{d}_{11} = \bar{d}_{22} \wedge \bar{d}_{12} = 0)$ and $(\bar{d}_{11} = -\bar{d}_{22})$. We first handle $(\bar{d}_{11} = -\bar{d}_{22})$.

Lemma 4. *When $\bar{d}_{11} = -\bar{d}_{22}$, the global minimizer to Minimization Problem is given by*

$$\bar{p}_{11} = \frac{1}{2} \left(\bar{d}_{11} + \sqrt{\bar{d}_{11}^2 + 4E/F + \bar{d}_{12}^2} \right), \quad \bar{p}_{12} = \frac{\bar{d}_{12}}{2}, \quad \bar{p}_{22} = \bar{p}_{11} - \bar{d}_{11}. \tag{31}$$

Proof. When $\bar{d}_{11} = -\bar{d}_{22}$, the Lagrange conditions (24a) and (24b) imply that $(\lambda + 1)(\bar{p}_{11} + \bar{p}_{22}) = 0$. From this it follows that we have either $\lambda = -1$ or $\bar{p}_{11} = -\bar{p}_{22}$, or both. When $\bar{p}_{11} = -\bar{p}_{22}$, it follows from (24d) that $-\bar{p}_{11}^2 - \bar{p}_{12}^2 = E/F$. However, this situation cannot occur because $E/F > 0$. We conclude that $\lambda = -1$ must hold. The Lagrange conditions (24a)–(24d) now simplify to

$$\bar{p}_{11} - \bar{p}_{22} = \bar{d}_{11}, \quad 2\bar{p}_{12} = \bar{d}_{12}, \quad \bar{p}_{11}\bar{p}_{22} = \frac{E}{F} + \frac{\bar{d}_{12}^2}{4}.$$

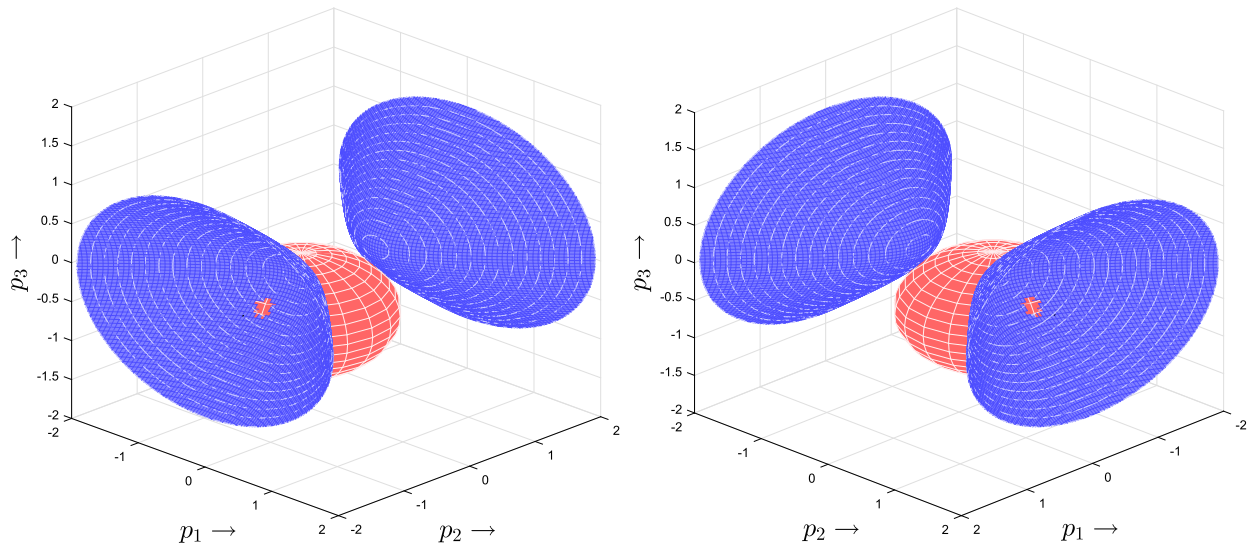


Fig. 7. These figures correspond to Lemma 4. The ellipsoid is centered around the same point ($d_1 = (\bar{d}_{11} - \bar{d}_{22})/\sqrt{2} = \sqrt{2}\bar{d}_{11}$, $d_2 = (\bar{d}_{11} + \bar{d}_{22})/\sqrt{2} = 0$, $d_3 = \bar{d}_{12}$) in both figures. This results in two local minima with the same function value for H . In the figure on the left we see the minimum on the hyperboloid sheet with $\text{tr}(\bar{P}) < 0$, which has to be discarded, and in the figure on the right we see the minimum on the sheet with $\text{tr}(\bar{P}) > 0$.

Combining the first and third of these equations gives us

$$\bar{p}_{11}^2 - \bar{d}_{11}\bar{p}_{11} - \frac{E}{F} - \frac{\bar{d}_{12}^2}{4} = 0.$$

This polynomial has for any \bar{d}_{11} and \bar{d}_{12} always two real solutions, which are given by $\bar{p}_{11} = (\bar{d}_{11} \pm \sqrt{\bar{d}_{11}^2 + 4E/F + \bar{d}_{12}^2})/2$. However, if the minus sign holds we see that $\text{tr}(\bar{P}) = -\sqrt{\bar{d}_{11}^2 + 4E/F + \bar{d}_{12}^2} < 0$. Thus, when $\bar{d}_{11} = -\bar{d}_{22}$, the global minimizer to Minimization Problem is (31). In Fig. 7 these findings are illustrated. \square

Now we still have to deal with the case $(\bar{d}_{11} = \bar{d}_{22} \wedge \bar{d}_{12} = 0)$.

Lemma 5. Suppose $\bar{d}_{11} = \bar{d}_{22}$ and $\bar{d}_{12} = 0$. When $\bar{d}_{11} < 2\sqrt{E/F}$, the solution to Minimization Problem is the global minimum given by

$$\bar{p}_{11} = \sqrt{E/F}, \quad \bar{p}_{12} = 0, \quad \bar{p}_{22} = \sqrt{E/F}, \tag{32}$$

otherwise, when $\bar{d}_{11} \geq 2\sqrt{E/F}$, the solution is a continuum of global minimizers given by

$$\bar{p}_{11} \in \left[\frac{\bar{d}_{11} - l}{2}, \frac{\bar{d}_{11} + l}{2} \right], \quad \bar{p}_{12} = \pm \sqrt{\bar{d}_{11}\bar{p}_{11} - \bar{p}_{11}^2 - \frac{E}{F}}, \quad \bar{p}_{22} = \bar{d}_{11} - \bar{p}_{11}, \tag{33}$$

where $l = \sqrt{\bar{d}_{11}^2 - 4E/F}$.

Proof. In the case that $\bar{d}_{11} = \bar{d}_{22}$ and $\bar{d}_{12} = 0$, Lagrange conditions (24a) and (24b) imply that $(1 - \lambda)(\bar{p}_{11} - \bar{p}_{22}) = 0$. From this it follows that we must either have $\lambda = 1$, or, $\lambda \neq 1$ and $\bar{p}_{11} = \bar{p}_{22}$. Let us first deal with the case $\lambda \neq 1$. When $\lambda \neq 1$, the Lagrange conditions (24c) and (24d) read

$$(1 - \lambda)\bar{p}_{12} = \bar{d}_{12} = 0, \quad \bar{p}_{11}^2 - \bar{p}_{12}^2 = E/F.$$

As $\lambda \neq 1$, the first of these equations implies that $\bar{p}_{12} = 0$. This fact combined with the second equation implies that $\bar{p}_{11} = \bar{p}_{22} = \pm\sqrt{E/F}$. The condition $\text{tr}(\bar{P}) > 0$ is only satisfied when the plus sign holds, hence we find the minimizer given by (32).

Now suppose that $\lambda = 1$. Lagrange condition (24a) implies $\bar{p}_{22} = \bar{d}_{11} - \bar{p}_{11}$ and from Lagrange condition (24d) we obtain $\bar{p}_{12} = \pm\sqrt{\bar{p}_{11}\bar{p}_{22} - E/F}$. Substituting the former expression in the latter gives us $\bar{p}_{12} = \pm\sqrt{\bar{d}_{11}\bar{p}_{11} - \bar{p}_{11}^2 - E/F}$, which is only real if $\bar{p}_{11}^2 - \bar{d}_{11}\bar{p}_{11} + E/F \leq 0$, that is, when $\bar{p}_{11} \in [(\bar{d}_{11} - l)/2, (\bar{d}_{11} + l)/2]$, where $l = \sqrt{\bar{d}_{11}^2 - 4E/F}$. This gives us

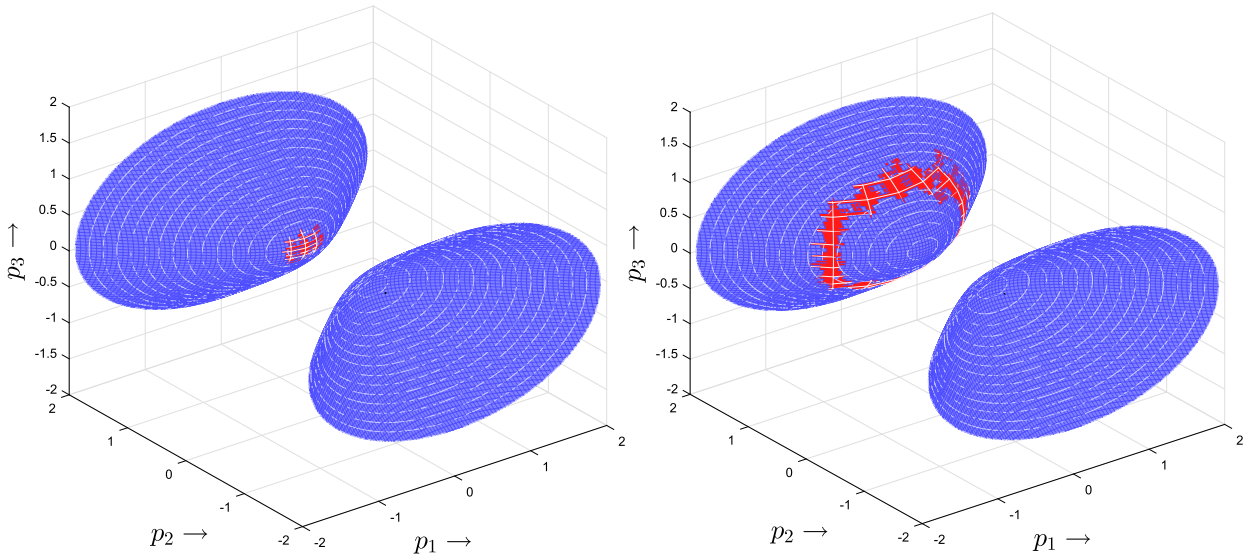


Fig. 8. These figures correspond to the minimizers in Lemma 5. The ellipsoid is located behind the hyperboloid. We see the sheet of the hyperboloid with $\text{tr}(\bar{P}) > 0$ on the left. The ellipsoid is centered around the point $(d_1 = (\bar{d}_{11} - \bar{d}_{22})/\sqrt{2} = 0, d_2 = (\bar{d}_{11} + \bar{d}_{22})/\sqrt{2} = \sqrt{2}\bar{d}_{11}, d_3 = \bar{d}_{12} = 0)$. In the figure on the left $d_2 < 2\sqrt{2E/F}$ and we find the extremum of the hyperboloid as minimizer. In the figure on the right $d_2 \geq 2\sqrt{2E/F}$ and we find an elliptical continuum of minimizers.

the continuum of minimizers (33). However, \bar{p}_{11} is only real if $|\bar{d}_{11}| \geq 2\sqrt{E/F}$. Moreover, because $\text{tr}(\bar{P}) = \bar{d}_{11}$, we see that $\text{tr}(\bar{P}) > 0$ is only satisfied when $\bar{d}_{11} > 0$. From this it follows that the continuum of minimizers can only be a solution to Minimization Problem when $\bar{d}_{11} \geq 2\sqrt{E/F}$. Thus, when $\bar{d}_{11} < 2\sqrt{E/F}$, the global minimizer is given by (32). To decide for $\bar{d}_{11} \geq 2\sqrt{E/F}$ whether the global minimizer is given by (32) or by an element of the continuum (33), we must compare the values of the function being minimized, i.e. H , for the local minimizers.

$H(\bar{P})$ has the same value for every element of the continuum of minimizers, because otherwise not all the elements of the continuum would have been local minima. For the value of $H(\bar{P})$ in the continuum we have

$$H_{\text{cont}} = \frac{1}{2} \left\| \begin{pmatrix} \bar{d}_{11} - \bar{p}_{11} & \sqrt{\bar{d}_{11}\bar{p}_{11} - \bar{p}_{11}^2 - E/F} \\ \sqrt{\bar{d}_{11}\bar{p}_{11} - \bar{p}_{11}^2 - E/F} & \bar{p}_{11} \end{pmatrix} \right\|^2 = \frac{\bar{d}_{11}^2}{2} - \frac{E}{F}.$$

On the other hand, for the local minimizer at the extremum of the hyperboloid, given by (32), we have

$$H_{\text{ext}} = \frac{1}{2} \left\| \begin{pmatrix} \bar{d}_{11} - \sqrt{E/F} & 0 \\ 0 & \bar{d}_{11} - \sqrt{E/F} \end{pmatrix} \right\|^2 = \bar{d}_{11}^2 - 2\bar{d}_{11}\sqrt{\frac{E}{F}} + \frac{E}{F}.$$

This implies that $H_{\text{cont}} - H_{\text{ext}} = -\bar{d}_{11}^2/2 + 2\bar{d}_{11}\sqrt{E/F} - 2E/F$. This polynomial in \bar{d}_{11} has its maximal value in $\bar{d}_{11} = 2\sqrt{E/F}$ where it equals 0, therefore it is negative for every $\bar{d}_{11} > 2\sqrt{E/F}$. This implies that if $\bar{d}_{11} \geq 2\sqrt{E/F}$, the solution to Minimization Problem is given by the continuum of minimizers (33). □

In Fig. 8 examples of the results from Lemma 5 are shown. Recall that the extrema of the two sheets of the hyperboloid are located at

$$(p_1, p_2, p_3) = \pm((\bar{p}_{11} - \bar{p}_{22})/\sqrt{2}, (\bar{p}_{11} + \bar{p}_{22})/\sqrt{2}, \bar{p}_{12}) = \pm(0, \sqrt{2E/F}, 0).$$

Thus Lemma 5 implies the following.

- The point $(p_1, p_2, p_3) = (0, \sqrt{2E/F}, 0)$ is the global minimizer when $\bar{d}_2 < 2\sqrt{2E/F}$, $\bar{d}_{11} = \bar{d}_{22}$ and $\bar{d}_{12} = 0$, i.e., when the center of the ellipsoid is located in $(0, p_2, 0)$, where $\bar{p}_2 = \sqrt{2}\bar{d}_{11} < 2\sqrt{2E/F}$, in other words, if the distance from the center of the ellipsoid to the origin is less than two times the distance to the extremum of the sheet with $\text{tr}(\bar{P}) > 0$, or if the center of the ellipsoid is located beneath the plane $p_2 = \text{tr}(\bar{P})/\sqrt{2} = 0$, then the global minimizer is given by the extremum of the upper sheet of the hyperboloid. This is depicted in the graph on the left in Fig. 8.
- If $\bar{d}_{11} = \bar{d}_{22}$, $\bar{d}_{12} = 0$, the center of the ellipsoid is located above the plane given by $p_2 = \text{tr}(\bar{P})/\sqrt{2} = 0$ and its distance to the origin is more than twice the distance from the extremum to the origin, then we have the continuum of global minimizers. This case is depicted in the graph on the right in Fig. 8.

Summarizing, we have proved the following theorem.

Theorem 2. *Minimization Problem, can be solved algebraically. In the general case, when $(\bar{d}_{11} \neq \bar{d}_{22} \vee \bar{d}_{12} \neq 0)$ and $(\bar{d}_{11} \neq -\bar{d}_{22})$, the solution to Minimization Problem is given by (27), with λ given by one of the four possibilities in (28). At least two of the λ 's in (28) are real. Explicit calculation of the function value $H(\bar{P})$ shows which of the real λ 's gives the global minimizer. In the case that $(\bar{d}_{11} = -\bar{d}_{22})$, there is a unique solution to Minimization Problem given by (31). Finally, in the case that $(\bar{d}_{11} = \bar{d}_{22} \wedge \bar{d}_{12} = 0)$, there is unique solution to Minimization Problem if $\bar{d}_{11} < 2\sqrt{E/F}$ and it is given by (32). If $\bar{d}_{11} \geq 2\sqrt{E/F}$, there is a whole continuum of solutions to Minimization Problem given by (33).*

3.2. Minimization of J

In this section we focus on the last step of the least-squares method, i.e. (17c). We will minimize the functional J , defined in (16), for $\mathbf{m} \in \mathcal{M}$, while keeping \mathbf{P} and \mathbf{b} constant. Again we do this for an arbitrary coordinate system on \mathcal{E} with basis vectors $\mathbf{e}_1, \mathbf{e}_2$ and corresponding metric $\mathbf{e} = e_{ij}\mathbf{e}^i \otimes \mathbf{e}^j$. We derive a coordinate-independent boundary value problem for the mapping \mathbf{m} and subsequently derive from this the boundary value problem in Cartesian and polar coordinates. We will see that in the Cartesian case we end up with the same boundary value problem for \mathbf{m} as derived in [12, pp. 142–144].

3.2.1. Derivation of the boundary value problem for the mapping

We will use Calculus of Variations to determine the minimizer \mathbf{m} for J . For a minimum to be attained the Gâteaux derivative of the J must be 0 in every direction, i.e.,

$$\delta J(\mathbf{m}, \mathbf{P}, \mathbf{b}; \boldsymbol{\eta}) := \lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} (J(\mathbf{m} + \varepsilon \boldsymbol{\eta}, \mathbf{P}, \mathbf{b}) - J(\mathbf{m}, \mathbf{P}, \mathbf{b})) = 0,$$

for every direction $\boldsymbol{\eta} \in \mathcal{M}$. δJ_I and δJ_B are defined analogously. By linearity of the Gâteaux derivative we have

$$\delta J(\mathbf{m}, \mathbf{P}, \mathbf{b}; \boldsymbol{\eta}) = \alpha \delta J_I(\mathbf{m}, \mathbf{P}, \mathbf{b}; \boldsymbol{\eta}) + (1 - \alpha) \delta J_B(\mathbf{m}, \mathbf{P}, \mathbf{b}; \boldsymbol{\eta}).$$

We first determine $\delta J_I(\mathbf{m}, \mathbf{P}, \mathbf{b}; \boldsymbol{\eta})$. By linearity of the derivative we find

$$\delta J_I(\mathbf{m}, \mathbf{P}, \mathbf{b}; \boldsymbol{\eta}) = \lim_{\varepsilon \rightarrow 0} \frac{1}{2\varepsilon} \int_{\mathcal{E}} \left(\|\varepsilon \nabla \hat{\boldsymbol{\eta}} + \nabla \hat{\mathbf{m}} - \mathbf{P}\|^2 - \|\nabla \hat{\mathbf{m}} - \mathbf{P}\|^2 \right) dA.$$

We will now need the following convenient property of inner product on $\mathbf{T}_2^0(\mathcal{E})$ as defined in (14). Let $\mathbf{A}, \mathbf{B} \in \mathbf{T}_2^0(\mathcal{E})$, then we have

$$\|\mathbf{A} + \mathbf{B}\|^2 = e^{ik}e^{jl}(a_{ij} + b_{ij})(a_{kl} + b_{kl}) = \|\mathbf{A}\|^2 + 2\mathbf{A} : \mathbf{B} + \|\mathbf{B}\|^2.$$

Applying this property to $\|\varepsilon \nabla \hat{\boldsymbol{\eta}} + \nabla \hat{\mathbf{m}} - \mathbf{P}\|^2$, with $\mathbf{A} = \varepsilon \nabla \hat{\boldsymbol{\eta}}$ and $\mathbf{B} = \nabla \hat{\mathbf{m}} - \mathbf{P}$, we obtain

$$\delta J_I(\mathbf{m}, \mathbf{P}, \mathbf{b}; \boldsymbol{\eta}) = \lim_{\varepsilon \rightarrow 0} \frac{1}{2\varepsilon} \int_{\mathcal{E}} \left(\varepsilon^2 \|\nabla \hat{\boldsymbol{\eta}}\|^2 + 2\varepsilon \nabla \hat{\boldsymbol{\eta}} : (\nabla \hat{\mathbf{m}} - \mathbf{P}) \right) dA = \int_{\mathcal{E}} \nabla \hat{\boldsymbol{\eta}} : (\nabla \hat{\mathbf{m}} - \mathbf{P}) dA.$$

In the same fashion, using the fact that

$$\|\varepsilon \boldsymbol{\eta} + \mathbf{m} - \mathbf{b}\|^2 - \|\mathbf{m} - \mathbf{b}\|^2 = \varepsilon^2 \|\boldsymbol{\eta}\|^2 + 2\varepsilon(\boldsymbol{\eta}, \mathbf{m} - \mathbf{b}), \tag{34}$$

we find the Gâteaux derivative of J_B to be

$$\delta J_B(\mathbf{m}, \mathbf{P}, \mathbf{b}; \boldsymbol{\eta}) = \int_{\partial \mathcal{E}} (\boldsymbol{\eta}, \mathbf{m} - \mathbf{b}) ds.$$

Combining the results for J_I and J_B we find that

$$\forall \boldsymbol{\eta} \in \mathcal{M} : \alpha \int_{\mathcal{E}} \nabla \hat{\boldsymbol{\eta}} : (\nabla \hat{\mathbf{m}} - \mathbf{P}) dA + (1 - \alpha) \int_{\partial \mathcal{E}} (\boldsymbol{\eta}, \mathbf{m} - \mathbf{b}) ds = 0. \tag{35}$$

In order to proceed we will rewrite the integrands in terms of their components. For the first integrand in (35) we have by (14)

$$\nabla \hat{\boldsymbol{\eta}} : (\nabla \hat{\mathbf{m}} - \mathbf{P}) = D_j \eta_i (D^j m^i - p^{ij}),$$

where $D_j \eta_i$ are the components of the covariant derivative of $\hat{\boldsymbol{\eta}}$ (Appendix A.3), $D^j = e^{ij} D_i$ and we used bijection (A.1). The product rule implies

$$D_j \eta_i (D^j m^i - p^{ij}) = D_j (\eta_i (D^j m^i - p^{ij})) - \eta_i D_j (D^j m^i - p^{ij}).$$

If we integrate the first term and apply the *Green's Theorem* [22, p. 135] we find

$$\int_{\mathcal{E}} D_j (\eta_i (D^j m^i - p^{ij})) dA = \int_{\partial \mathcal{E}} (D^j m^i - p^{ij}) \eta_i n_j ds,$$

where n_j are the covariant components of the outward unit normal vector on the boundary $\partial \mathcal{E}$ and the orientation on $\partial \mathcal{E}$ is the one induced by \mathcal{E} [22, p. 119]. It follows that

$$\int_{\mathcal{E}} \nabla \hat{\boldsymbol{\eta}} : (\nabla \hat{\boldsymbol{m}} - \mathbf{P}) dA = \int_{\partial \mathcal{E}} (D^j m^i - p^{ij}) \eta_i n_j ds - \int_{\mathcal{E}} D_j (D^j m^i - p^{ij}) \eta_i dA.$$

Combining this result with identity (35) we obtain

$$0 = \int_{\partial \mathcal{E}} [\alpha (D^j m^i - p^{ij}) n_j + (1 - \alpha) (m^i - b^i)] \eta_i ds - \alpha \int_{\mathcal{E}} D_j (D^j m^i - p^{ij}) \eta_i dA,$$

for all $\boldsymbol{\eta} \in \mathcal{M}$. Invoking the *Fundamental Lemma of Calculus of Variations* [25, p. 185] we find from this the boundary value problem

$$D_j D^j m^i = D_j p^{ij} \quad \text{in } \mathcal{E}, \tag{36a}$$

$$\alpha (D^j m^i) n_j + (1 - \alpha) m^i = \alpha p^{ij} n_j + (1 - \alpha) b^i \quad \text{on } \partial \mathcal{E}. \tag{36b}$$

The solution of boundary value problem (36b) will minimize J for given \mathbf{P} and \mathbf{b} . Note that (36b) are vector equations. The term $D_j D^j m^i$ is the so-called *vector Laplacian* [26, p. 91], which can be written in terms of vector calculus operators as $\nabla(\nabla \cdot \mathbf{m}) - \nabla \times (\nabla \times \mathbf{m})$. In Cartesian coordinates $D_j D^j m^i = (\partial_1^2 + \partial_2^2) m^i$, thus, in Cartesian coordinates the Laplacian of a vector amounts to just taking the Laplacian component-wise. However, in different coordinate systems this is not true, because nonzero Christoffel symbols imply that $[D_j D^j m^i]_{i=1,2}$ depend both on both m^1 and m^2 . This results for an arbitrary coordinate system in two coupled equations, while for a Cartesian coordinate system these two decouple. This will become more apparent when we derive the coordinate specific boundary value problem for Cartesian and polar coordinates.

3.2.2. The boundary value problem in specific coordinate systems

In Cartesian coordinates the partial differential equations in (36b) decouple. Let us denote the standard Cartesian basis vectors by \mathbf{e}_x and \mathbf{e}_y , define

$$\mathbf{p}_x = \begin{pmatrix} p^{xx} \\ p^{xy} \end{pmatrix} = \begin{pmatrix} p^{11} \\ p^{12} \end{pmatrix} \quad \text{and} \quad \mathbf{p}_y = \begin{pmatrix} p^{xy} \\ p^{yy} \end{pmatrix} = \begin{pmatrix} p^{12} \\ p^{22} \end{pmatrix},$$

and write $\mathbf{m} = m^x \mathbf{e}_x + m^y \mathbf{e}_y$. With the use of this definition we can rewrite $(D_j p^{ij})_{i=1}$ as $\nabla \cdot \mathbf{p}_x$ and $(D_j p^{ij})_{i=2}$ as $\nabla \cdot \mathbf{p}_y$. From this we see that in Cartesian coordinates (36b) reduces to the decoupled set of equations

$$\begin{aligned} \Delta m^x &= \nabla \cdot \mathbf{p}_x && \text{in } \mathcal{E}, \\ \alpha (\nabla m^x, \mathbf{n}) + (1 - \alpha) m^x &= \alpha (\mathbf{p}_x, \mathbf{n}) + (1 - \alpha) b^x && \text{on } \partial \mathcal{E}, \end{aligned} \tag{37a}$$

$$\begin{aligned} \Delta m^y &= \nabla \cdot \mathbf{p}_y && \text{in } \mathcal{E}, \\ \alpha (\nabla m^y, \mathbf{n}) + (1 - \alpha) m^y &= \alpha (\mathbf{p}_y, \mathbf{n}) + (1 - \alpha) b^y && \text{on } \partial \mathcal{E}. \end{aligned} \tag{37b}$$

The boundary value problems (37a) and (37b) are exactly the same as in [12, p. 143].

In polar coordinates the equations do not decouple. Notice that the coordinate specific boundary value problem that we deduce from (36b) does depend on the choice of basis for polar coordinates, because (36b) is a vector equation. Thus, we find for an anholonomic basis different expressions than for a holonomic basis (Appendix A.2).

To derive the boundary value problem in polar coordinates, let us first elaborate the components of the covariant derivatives appearing in (36b). We start with the vector Laplacian. By the fact that $D_i (e^{jk}) = 0$ (Appendix A.3) it follows by (A.5c) that $D_j D^j m^i = e^{jk} (\nabla_{\mathbf{e}_j} (D_k m^i) - \Gamma_{kj}^l D_l m^i + \Gamma_{lj}^i D_k m^l)$ and by (A.4a) that $D_k m^i = \nabla_{\mathbf{e}_k} m^i + \Gamma_{lk}^i m^l$, hence

$$D_j D^j m^i = e^{jk} (\nabla_{\mathbf{e}_j} \nabla_{\mathbf{e}_k} m^i + \nabla_{\mathbf{e}_j} (\Gamma_{lk}^i) m^l + \Gamma_{lk}^i \nabla_{\mathbf{e}_j} m^l - \Gamma_{kj}^l \nabla_{\mathbf{e}_l} m^i - \Gamma_{kj}^l \Gamma_{sl}^i m^s + \Gamma_{lj}^i \nabla_{\mathbf{e}_k} m^l + \Gamma_{lj}^i \Gamma_{sk}^l m^s). \tag{38}$$

Doing the same derivation for the divergence of \mathbf{P} we obtain³

³ Note, that due to the symmetry of \mathbf{P} it is clear what we mean when we speak of the divergence of \mathbf{P} . It does not matter if we contract D_k with the first or second index of p^{ij} , the result is the same.

$$D_j p^{ij} = \nabla_{\mathbf{e}_j} p^{ij} + \Gamma_{lj}^i p^{lj} + \Gamma_{lj}^j p^{il}. \tag{39}$$

Similarly, we find for the normal derivative of \mathbf{m} in (36b)

$$(D^j m^i) n_j = e^{jk} (D_k m^i) n_j = e^{jk} (\nabla_{\mathbf{e}_k} m^i + \Gamma_{lk}^i m^l) n_j. \tag{40}$$

We use (38)–(40) to expand the boundary value problem (36b) in polar coordinates. We consider the anholonomic orthonormal basis ($e_{ij} = \delta_{ij}$), because this is the basis we used in the implementation. The only nonzero Christoffel symbols in the anholonomic basis are $\Gamma_{\theta\theta}^r = -r^{-1}$ and $\Gamma_{r\theta}^\theta = r^{-1}$ [20, p. 218]. After doing the tedious calculations of determining the coordinate system specific expressions of the various terms in (38) we find

$$(D_j D^j m^i)_{i=r} = \frac{\partial^2 m^r}{\partial r^2} + \frac{1}{r} \frac{\partial m^r}{\partial r} + \frac{1}{r^2} \frac{\partial^2 m^r}{\partial \theta^2} - \frac{m^r}{r^2} - \frac{2}{r^2} \frac{\partial m^\theta}{\partial \theta},$$

$$(D_j D^j m^i)_{i=\theta} = \frac{\partial^2 m^\theta}{\partial r^2} + \frac{1}{r} \frac{\partial m^\theta}{\partial r} + \frac{1}{r^2} \frac{\partial^2 m^\theta}{\partial \theta^2} - \frac{m^\theta}{r^2} + \frac{2}{r^2} \frac{\partial m^r}{\partial \theta}.$$

In the same way we calculate the expressions for the divergence of \mathbf{P} and obtain

$$(D_j p^{ij})_{i=r} = \frac{\partial p^{rr}}{\partial r} + \frac{1}{r} \frac{\partial p^{r\theta}}{\partial \theta} + \frac{p^{rr} - p^{\theta\theta}}{r},$$

$$(D_j p^{ij})_{i=\theta} = \frac{\partial p^{r\theta}}{\partial r} + \frac{1}{r} \frac{\partial p^{\theta\theta}}{\partial \theta} + \frac{2p^{r\theta}}{r}.$$

Finally, we determine the expressions for the normal derivative of \mathbf{m} from (40) and find

$$((D^j m^i) n_j)_{i=r} = \frac{\partial m^r}{\partial r} n^r + \frac{\partial m^r}{\partial \theta} \frac{n^\theta}{r} - \frac{m^\theta n^\theta}{r},$$

$$((D^j m^i) n_j)_{i=\theta} = \frac{\partial m^\theta}{\partial r} n^r + \frac{\partial m^\theta}{\partial \theta} \frac{n^\theta}{r} + \frac{m^\theta n^\theta}{r}.$$

We define

$$\mathbf{p}_r = \begin{pmatrix} p^{rr} \\ p^{r\theta} \end{pmatrix} \quad \text{and} \quad \mathbf{p}_\theta = \begin{pmatrix} p^{r\theta} \\ p^{\theta\theta} \end{pmatrix}, \tag{41}$$

and collect all the different terms and find that for polar coordinates with an orthonormal basis (36b) is given by

$$\Delta m^r - \frac{1}{r^2} \left(m^r + 2 \frac{\partial m^\theta}{\partial \theta} \right) = \nabla \cdot \mathbf{p}_r - \frac{1}{r} p^{\theta\theta} \quad \text{in } \mathcal{E}, \tag{42}$$

$$\alpha (\nabla m^r, \mathbf{n}) - \alpha \frac{m^\theta n^\theta}{r} + (1 - \alpha) m^r = \alpha (\mathbf{p}_r, \mathbf{n}) + (1 - \alpha) b^r \quad \text{on } \partial \mathcal{E},$$

and

$$\Delta m^\theta - \frac{1}{r^2} \left(m^\theta - 2 \frac{\partial m^r}{\partial \theta} \right) = \nabla \cdot \mathbf{p}_\theta + \frac{1}{r} p^{r\theta} \quad \text{in } \mathcal{E}, \tag{43}$$

$$\alpha (\nabla m^\theta, \mathbf{n}) + \alpha \frac{m^r n^\theta}{r} + (1 - \alpha) m^\theta = \alpha (\mathbf{p}_\theta, \mathbf{n}) + (1 - \alpha) b^\theta \quad \text{on } \partial \mathcal{E},$$

where Δ , $\nabla \cdot$ and ∇ are the familiar Laplace, divergence and gradient operator in polar coordinates with orthonormal basis [23, pp. 542–543]. The equations (42) and (43) are coupled.

In the implementation of the GLS method we solve these two boundary value problems by using the standard second order central difference discretization. This provides us with a linear system, the matrix of which does not change during an entire run of the GLS method. This allows us to determine the LU-decomposition before the start of the algorithm and we can use this to solve the linear system when necessary.

To deal with the fact that, in polar coordinates, (42) and (43) are coupled, we will iterate between the two. Starting with (42) we keep m^θ fixed and solve for m^r . Next we keep m^r fixed and solve (43) for m^θ . We stop this iterative procedure when $J^{(n+1,i)} < c J^{(n)}$, where n is the outer iteration count of (17c) and i is the inner iteration count, or, when the number of inner iterations is larger than a specified value d , i.e., $i > d$.

The optimal choice for these values are problem specific and if the number of inner iterations is increased by demanding more precision in (17c), the outer iterative procedure might converge faster. However, demanding far more precision in (17c) than is achieved by the outer iterative procedure up to that point is a waste of time. A maximum on the number of iterations is introduced to make sure that the method does not stall when $J^{(n+1,i)} < c J^{(n)}$ is a too severe requirement. This will come into play in the final part of the iteration sequence. We have no proof of convergence, but in practice this procedure always converged in a few iterations, because the mapping \mathbf{m}^n provides a very good initial guess for (17c). In practice we took $c = 0.9$ and $d = 5$ and these values seem to be a good choice for the problems tested so far. Alternatively, the complete coupled system can be solved at once. However, this is likely to be more expensive.

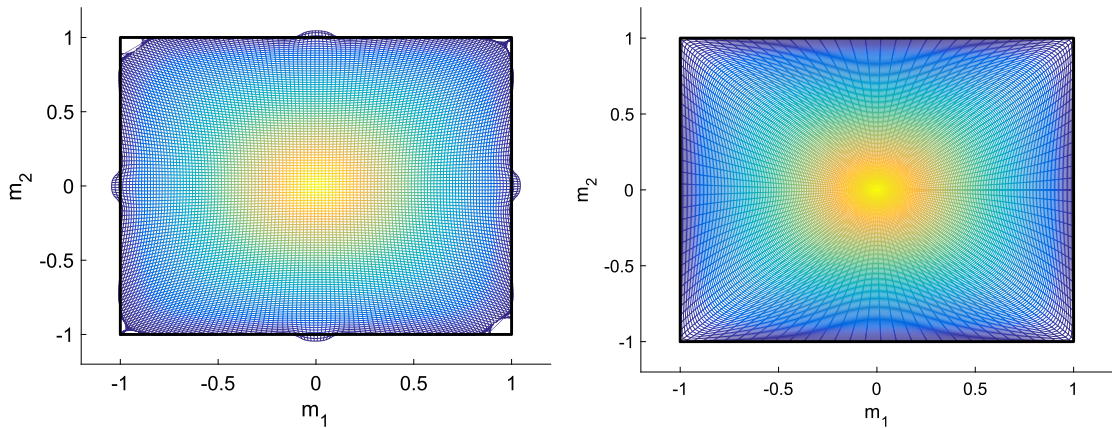


Fig. 9. The resulting mapping: on the left for Cartesian coordinates and on the right for polar coordinates. For both a 500×500 grid is used and $\alpha = 0.2$. We see the grid as it is mapped on \mathcal{F}_S . Grid points that initially had the same distance to the center of \mathcal{E} have the same color. Bright yellow corresponds to points in the center of \mathcal{E} and dark blue corresponds to points on $\partial\mathcal{E}$.

3.3. Determining the reflector surface from the mapping

To determine the reflector surface from the mapping \mathbf{m} we generalize the derivation given by Prins et al. [13] to arbitrary coordinate systems. We remarked earlier that \mathbf{m} equals the gradient of u if and only if $\nabla\hat{\mathbf{m}}$ is symmetric. However, in the GLS method J_1 is minimized in the L^2 -norm and hence $\nabla\hat{\mathbf{m}}$ is not exactly symmetric. We can therefore only search for a function $u : \mathcal{E} \rightarrow (0, \infty)$ with gradient equal to \mathbf{m} in an L^2 -sense, hence we will search for u that minimizes

$$I(u) := \int_{\mathcal{E}} \|\nabla u - \mathbf{m}\|^2 dA.$$

After a derivation very similar to the one by which we arrived at boundary value problem (36b), which we leave out for brevity, we obtain the Poisson problem

$$\begin{aligned} \Delta u &= \nabla \cdot \mathbf{m} && \text{in } \mathcal{E}, \\ (\nabla u, \mathbf{n}) &= (\mathbf{m}, \mathbf{n}) && \text{on } \partial\mathcal{E}. \end{aligned}$$

It is this problem that we solve to find the reflector surface for the problems presented in the next section. We discretize this Poisson problem using second order central differences, giving us a linear system. The solution of this linear system gives us the reflector surface.

4. Numerical results

We show the performance of the least-squares method in polar coordinates on the basis of three test cases. In the first test case we compare the method in polar coordinates with the method in Cartesian coordinates as presented in [13], in the second test case we test the performance of the method for a source with non-uniform emittance and in the third test case we investigate the performance of the method in polar coordinates for a complex problem with a desired light output with a lot of contrast.

For all test cases we take as source domain \mathcal{E} the unit disk. In the first and last test case we use a source with uniform emittance E_U and in the second test case we use a Gaussian emittance E_G . We choose these sources, because they frequently occur in lighting systems. For the first and second test case we take as the target $\mathcal{F}_S = [-1, 1] \times [-1, 1]$ with a uniform intensity function F_S . For the third test case we have determined the pair (\mathcal{F}_E, F_E) such that an intensity pattern corresponding to the lithograph by M. C. Escher (Fig. 1) is projected on a screen in the far-field. We normalize the source emittances and target intensities such that (10) holds.

4.1. From a circle to a square: uniform emittance

For the first test the mapping is presented in Fig. 9. In this figure we see that near the boundary $\partial\mathcal{F}_S$ the method in Cartesian coordinates has great difficulties. This results from the implementation where as actual source the smallest bounding box of \mathcal{E} is used with emittance zero in the points outside \mathcal{E} (see [13] for details). In the polar coordinate method the grid perfectly aligns with $\partial\mathcal{E}$. In Fig. 9 it can be seen that now all the difficulties at $\partial\mathcal{F}_S$ are resolved.

Fig. 10 shows the convergence history of the method for different values of α . The value of α determines approximately the ratio between J_1 and J_B . In general, for values of α close to 1 the method finds a reflector that closely satisfies (11),

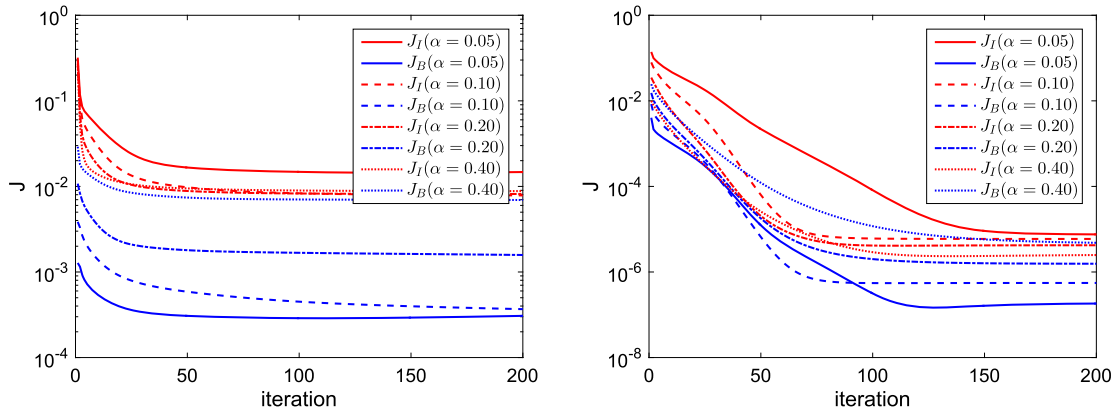


Fig. 10. For the first test case the interior functional J_I and the boundary functional J_B are shown as function of the number of iterations for different α 's, on the left for the method in Cartesian coordinates and on the right for the method in polar coordinates. In both cases a 100×100 grid was used.

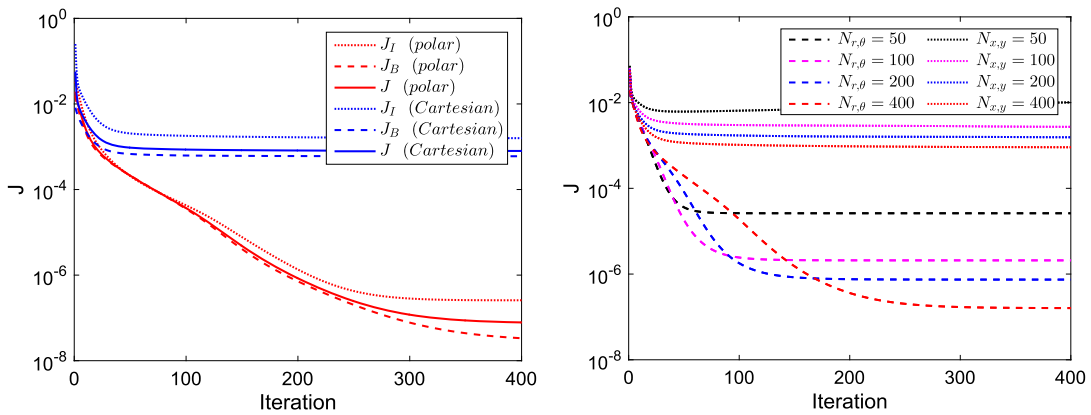


Fig. 11. The convergence history for Cartesian and polar coordinates is compared for the first test case. In the left plot a 500×500 grid is used and the different error components are shown. In the right plot $N_{r,\theta}$ refers tot the number of grid points that is used in both the radial and angular directions and $N_{x,y}$ has an analogous meaning. An improvement by a factor 10^4 is observed when using polar instead of Cartesian coordinates. On the right J is shown for different grids. For both plots we used $\alpha = 0.2$.

but might possibly be less accurate concerning the boundary condition of the Inverse Reflector Problem, and vice versa for α close to 0. For smoother (\mathcal{E}, E) and (\mathcal{F}, F) the solution found by the method seems less dependable on the choice of α . However, in cases where for example $\partial\mathcal{F}$ is not differentiable, as in the current test case, the boundary condition and (11) seem to be conflicting goals. For such problems they cannot be satisfied exactly and simultaneously. In order to clarify and quantify this alleged conflict further study has to be done. Nonetheless, the freedom in α provides the user of the GLS method with an opportunity to choose the best balance of relative weight given to the boundary condition and (11) for the specific application at hand. Moreover, it can be seen that due to better handling of the boundary $\partial\mathcal{E}_U$ the gap between J_I and J_B is far smaller for the method in polar coordinates for all choices of α .

Fig. 11 shows that the method in polar coordinates significantly outperforms its Cartesian counterpart. In the figure on the left it is seen that for a 500×500 grid the convergence of the Cartesian method stalls after approximately 75 iterations. The convergence of the polar method proceeds for another 300 iterations and this eventually leads to a value for J_I that is 10^4 times as small as the J_I found with the Cartesian method. In the figure on the right it is seen that the use of increasingly finer grids has more effect for the polar method. However, even for the polar method, the final value for J_I seems not to convergence to zero when ever finer grids are used.

In Fig. 12 the convergence of the method as a function of the number of variables is analyzed. It is seen that the convergence rate is approximately second order for the method in polar coordinates. In Cartesian coordinates only a convergence rate of approximately first order is achieved. This is probably a result of the boundary issue alluded to above.

4.2. From a circle to a square: Gaussian emittance

As a second test we analyze the performance of the method for a source with non-uniform intensity. We again take \mathcal{E} to be the unit disk, but now take the intensity function to be a (radial-symmetric) Gaussian:

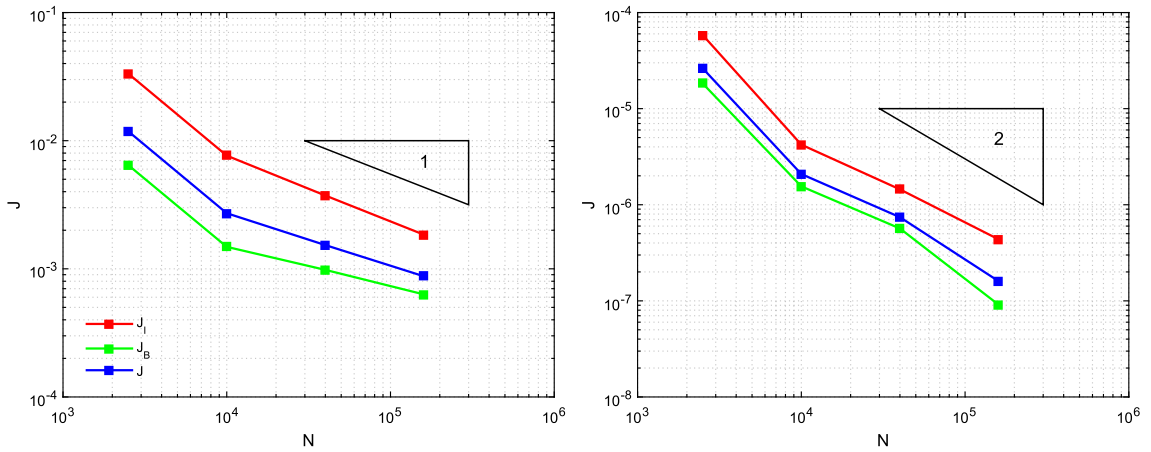


Fig. 12. The final errors J_I , J_B and J are shown as a function of the total number of variables N for the method applied to first test case. Left: Cartesian coordinates with $\sqrt{N} = N_x = N_y$. Right: polar coordinates with $\sqrt{N} = N_r = N_\theta$.

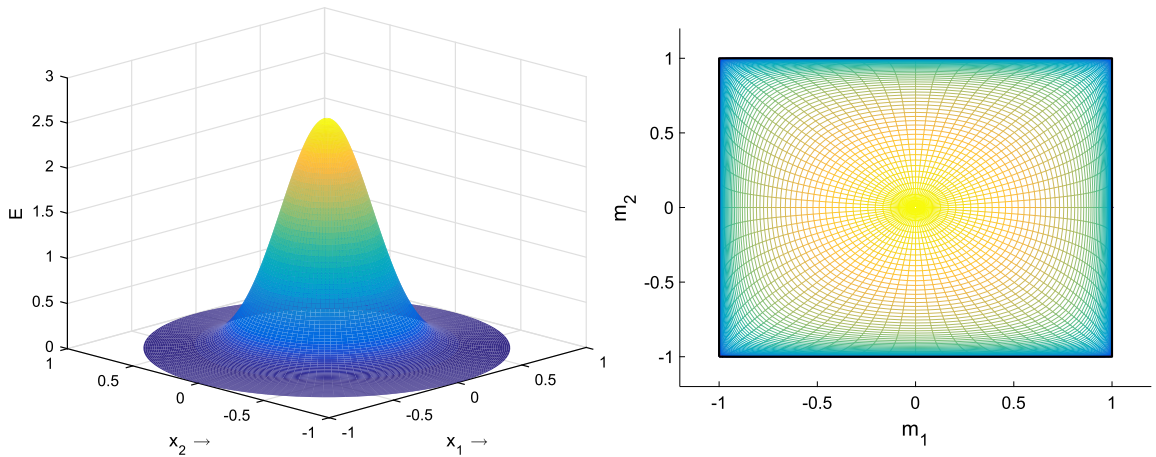


Fig. 13. Left: Emittance function E_G . Right: mapping for $N_r = N_\theta = 400$ and $\alpha = 0.2$.

$$E_G(r) = a \exp\left(-\frac{r^2}{2\sigma^2}\right),$$

where a is chosen such that the integral of E_G over \mathcal{E} is 1 and we taken the standard deviation $\sigma = 0.25$. The source intensity at $\partial\mathcal{E}$ is equal to

$$E_G(1) = \frac{1}{2\pi\sigma^2} \left(\exp\left(\frac{1}{2\sigma^2}\right) - 1 \right)^{-1} \approx 8.5454 \cdot 10^{-4}.$$

As a result the ratio between the maximum and minimum intensity approximately equals $E_G(0)/E_G(1) \approx 2.9786 \cdot 10^3$.

In Fig. 13 the emittance function E_G and the resulting mapping are shown for the method in polar coordinates. For this test case the method in Cartesian coordinates does not even converge. We again expect this to be a result of the boundary issue discussed before in the first paragraph of Section 4.1. Moreover, for the current test E_G is close to 0 near $\partial\mathcal{E}$ which the Cartesian method cannot deal with.

The convergence of the spatial discretization for this test case as a function of the number of variables is shown in Fig. 14. The order of convergence for this problem is at least 4. It is not entirely clear to us why the convergence for this test case is more than fourth order while it was only approximately second order for the test case with uniform emittance.

4.3. From a circle to an Escher lithograph in the far-field

In the final test case we calculate the reflector for the target intensity (\mathcal{F}_E, F_E) corresponding to the lithograph by Escher. We take the projection screen at a distance 100 times the radius of the source and we take (\mathcal{F}_E, F_E) such that width and height of the projection are 4.3 times the radius of the source. To avoid division by zero in (12a), we increased the minimum

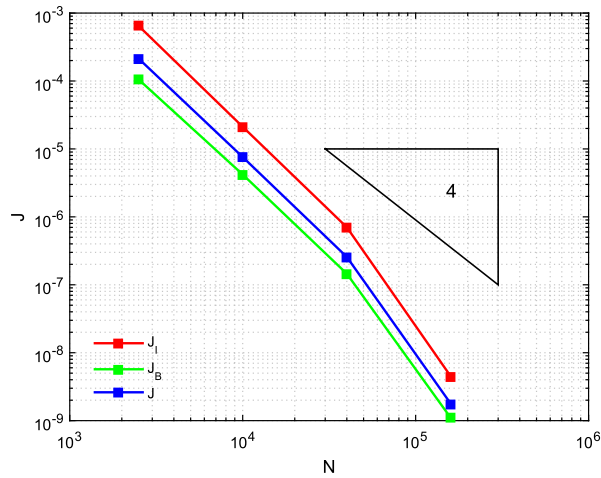


Fig. 14. The final errors J_1 , J_B and J are shown as a function of the total number of variables N for the method applied to second test case with $\sqrt{N} = N_r = N_\theta$.

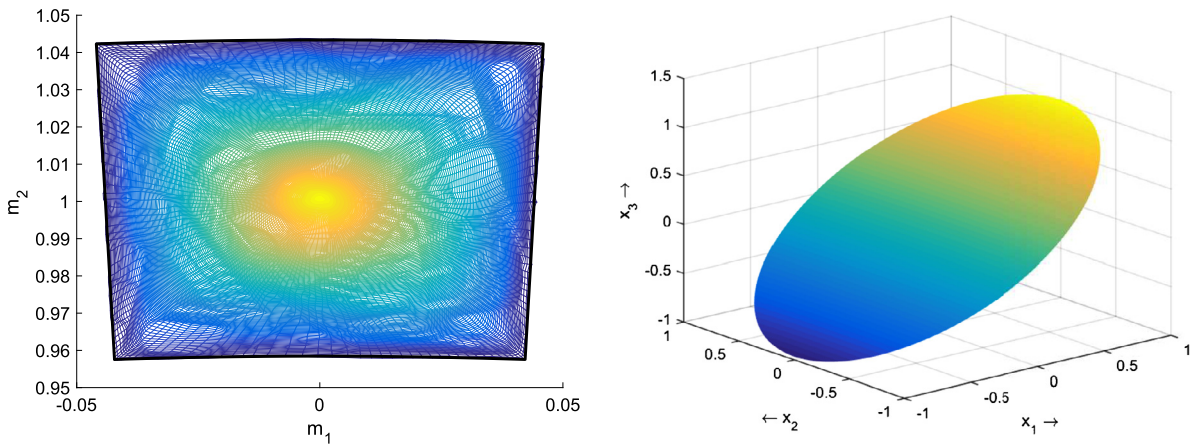


Fig. 15. The mapping and corresponding reflector are shown for the Escher test case.

intensity of the lithograph to be 5% of the maximum intensity. For this very demanding test case a 1400×1400 grid was used and $\alpha = 0.1$.

Fig. 15 shows the results of applying the method in polar coordinates. The elliptic shaped reflector is globally close to flat but locally contains great detail. The integrands of the final errors J_1 and J_B are shown in Fig. 16. It is seen that the errors are small for this fine mesh.

Subsequently, the reflector was simulated by using a ray trace method [3]. The result can be seen in Fig. 17. The ray trace result closely resembles the original picture, although there is some decrease in contrast. In the algorithm the reflector surface is in the class of twice continuously differentiable functions. This naturally results in smoothing of the, often discontinuous, intensity function of the original picture. Nonetheless the resolution obtained is high enough to carry over all the minute details of the original picture.

5. Summary and concluding remarks

In Section 2 we derived the Monge–Ampère equation, describing the reflector surface, for an arbitrary coordinate system. We found the map r , which maps a point on the source \mathcal{E} to the direction of the reflected ray, to be the composition of the gradient of the reflector surface, ∇u , and the inverse of the stereographic projection, s . Furthermore, we formulated the Inverse Reflector Problem in terms of the source and emittance (\mathcal{E}, E) and the gradient set and intensity function (\mathcal{F}, F) .

In Section 3 we introduced the GLS method by generalizing the LS method, earlier introduced in [13], to general coordinate systems. Moreover, we gave a new geometric interpretation to the minimization problem for the functional J_1 and found that the minimization problem for the total functional J consists of two Poisson problems which, contrary to the Cartesian case, are coupled in general coordinate systems.

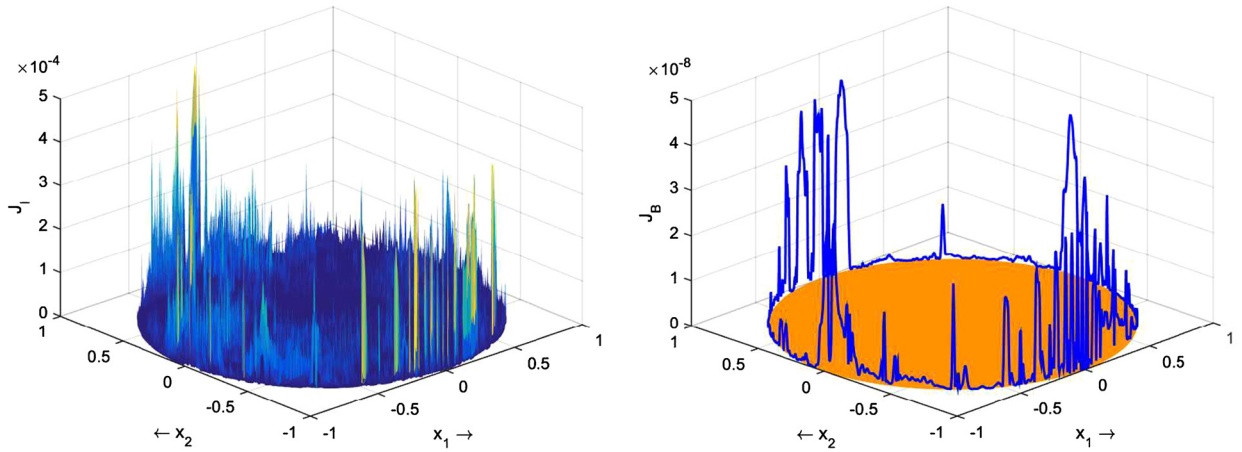


Fig. 16. Left: The interior error $\|\nabla \hat{m} - \mathbf{P}\|^2/2$ as a function of $x \in \mathcal{E}$. Right: The boundary error $\|\mathbf{m} - \mathbf{b}\|^2/2$ as a function of $x \in \partial \mathcal{E}$.

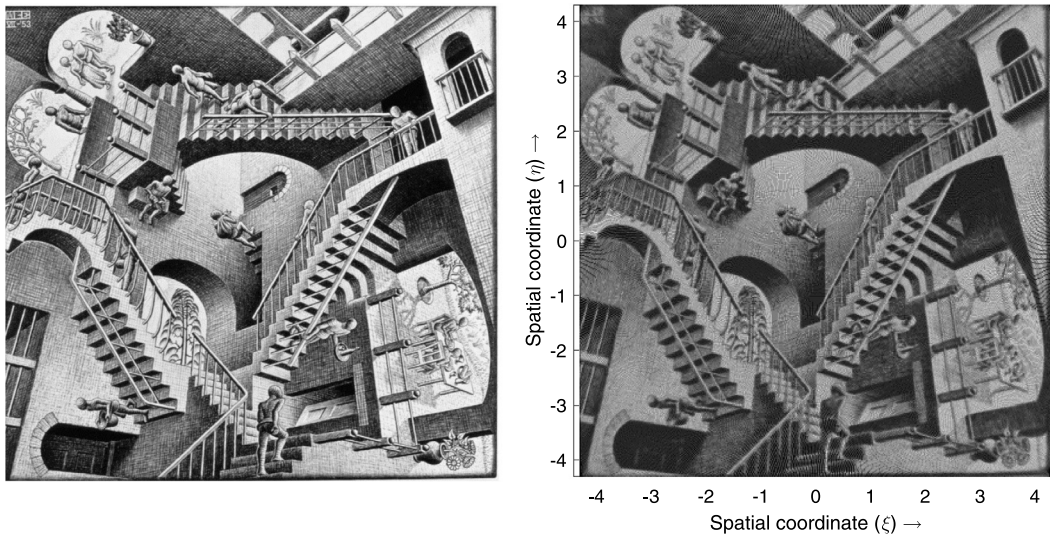


Fig. 17. The image projected on the screen by the reflector of Fig. 15 is determined by ray tracing about 4 million rays, with uniform fixed spacing, that leave \mathcal{E} . The original is shown on the left and the ray trace result is shown on the right.

In Section 4 we showed that the GLS method has far wider applicability than the LS method. We showed that for a disk-shaped source the GLS method in polar coordinates gave a significant improvement over the LS method, decreasing the error by four orders of magnitude. It was seen that for problems with non-smooth desired output intensity the final ratio between J_I and J_B depends on the value of α in (16). Further research and literature study into this relation should be done. It would be for example important to know for which combination of source pair (\mathcal{E}, E) and target pair (\mathcal{F}, F) the solution of the method depends on α and to quantify to what extent.

Lastly, the method was applied to a very challenging problem concerning a detailed piece of art and still the method obtained a high resolution preserving the details of the original picture. This gives confidence in the wide applicability of the method in an industrial context.

Appendix A. Tensor Calculus on \mathbb{R}^n in a nutshell

A.1. Tensors

Consider the space of linear functionals from the n -dimensional vector space V to the real numbers: $\mathcal{L}(V; \mathbb{R})$. This space is called the dual space of V and we will denote it by $V^* := \mathcal{L}(V; \mathbb{R})$. Suppose $\mathbf{v} \in V$ and $\mathbf{v} = v^i \mathbf{e}_i$, where v^i are the components with respect to a basis $(\mathbf{e}_i)_{i=1}^n$ of the vector space V and we employ the Einstein summation convention. Furthermore, let $\hat{\mathbf{w}} \in V^*$, where we use a hat on \mathbf{w} to indicate that its from V^* and not from V . By linearity of $\hat{\mathbf{w}}$ it follows that $\hat{\mathbf{w}}(\mathbf{v}) = v^i \hat{\mathbf{w}}(\mathbf{e}_i)$. This shows that a linear mapping $\hat{\mathbf{v}}$ is completely determined by specifying its action on the

basis elements of V . We can define in this way a set of elements of V^* , denoted by $(\mathbf{e}^i)_{i=1}^n$, according to $\mathbf{e}^i(\mathbf{e}_j) = \delta_j^i$. The elements $(\mathbf{e}^i)_{i=1}^n$ form a basis for V^* . They span V^* , because we have

$$\hat{\mathbf{w}}(\mathbf{v}) = v^i \hat{\mathbf{w}}(\mathbf{e}_i) = \mathbf{e}^i(\mathbf{v}) \hat{\mathbf{w}}(\mathbf{e}_i) = w_i \mathbf{e}^i(\mathbf{v}),$$

where $w_i := \hat{\mathbf{w}}(\mathbf{e}_i)$. Furthermore, suppose $\hat{\mathbf{w}} = w_i \mathbf{e}^i = 0$ then

$$0 = \hat{\mathbf{w}}(\mathbf{e}_j) = w_i \mathbf{e}^i(\mathbf{e}_j) = w_i \delta_j^i = w_j,$$

which indicates that $w_j = 0$ for $1 \leq j \leq n$, implying that $(\mathbf{e}^i)_{i=1}^n$ are linearly independent. The basis $(\mathbf{e}^i)_{i=1}^n$ is called the dual basis and elements of V^* are also called covectors.

The double dual $V^{**} := \mathcal{L}(V^*; \mathbb{R})$ can be identified with V . This follows because when we define the double dual basis $(\mathbf{f}_i)_{i=1}^n$ according to $\mathbf{f}_j(\mathbf{e}^i) = \delta_j^i$, the identification of \mathbf{f}_i with \mathbf{e}_i gives an isomorphism between V and V^{**} . This shows that we can equally well consider V to be the dual to V^* and write $\mathbf{v}(\hat{\mathbf{w}}) := \hat{\mathbf{w}}(\mathbf{v})$. Similar to $w_i := \hat{\mathbf{w}}(\mathbf{e}_i) = \mathbf{e}_i(\hat{\mathbf{w}})$, we define $v^i := \mathbf{e}^i(\mathbf{v}) = \mathbf{v}(\mathbf{e}^i)$.

The space $\mathcal{L}(V, V^*; \mathbb{R})$ is the space of multilinear functionals (mappings linear in each argument) from $V \times V^*$ to \mathbb{R} . We will denote this space by $\mathbf{T}_1^1(V)$. An example of a functional in $V \times V^* \rightarrow \mathbb{R}$ is $\mathbf{e}^i \otimes \mathbf{e}_j$ which is defined as

$$\mathbf{e}^i \otimes \mathbf{e}_j(\mathbf{v}, \hat{\mathbf{w}}) = \mathbf{e}^i(\mathbf{v}) \mathbf{e}_j(\hat{\mathbf{w}}) = v^i w_j.$$

This element of $\mathbf{T}_1^1(V)$ is called the tensor product of \mathbf{e}^i and \mathbf{e}_j . It can be shown that $(\mathbf{e}^i \otimes \mathbf{e}_j)_{i,j=1}^n$ is actually a basis for $\mathbf{T}_1^1(V)$ [27, p. 25]. For an element $\mathbf{T} \in \mathbf{T}_1^1(V)$ we define its coefficients as $T_j^i = \mathbf{T}(\mathbf{e}_j, \mathbf{e}^i)$ and have $\mathbf{T} = T_j^i \mathbf{e}^j \otimes \mathbf{e}_i$. By linearity it follows that for a vector \mathbf{v} and covector $\hat{\mathbf{w}}$ we have

$$\mathbf{T}(\mathbf{v}, \hat{\mathbf{w}}) = T_j^i v^j w_i.$$

Analogously we can define the spaces $\mathbf{T}_0^2(V)$ and $\mathbf{T}_2^0(V)$, with bases $(\mathbf{e}_i \otimes \mathbf{e}_j)_{i,j=1}^n$ and $(\mathbf{e}^i \otimes \mathbf{e}^j)_{i,j=1}^n$, respectively. Or, more generally, we can define the space $\mathbf{T}_q^p(V)$ and its corresponding basis, formed by taking consecutive tensor products between basis elements of V and V^* . The elements of the space $\mathbf{T}_q^p(V)$ are called tensors of contravariant rank p and covariant rank q . The spaces V and V^* are denoted by $\mathbf{T}_0^1(V)$ and $\mathbf{T}_1^0(V)$, respectively.

An example of a tensor of covariant rank 2 is given by the inner product $(\cdot, \cdot) : V \times V \rightarrow \mathbb{R}$. This tensor is known as the metric tensor or simply the metric. We denote the coefficients of the metric by $e_{ij} = (\mathbf{e}_i, \mathbf{e}_j)$. It follows that the metric can also be written as $\mathbf{e} = e_{ij} \mathbf{e}^i \otimes \mathbf{e}^j$. The inner product also provides us with a bijection between V and V^* . This bijection is given by the map

$$\mathbf{v} \mapsto (\mathbf{v}, \cdot). \tag{A.1}$$

In terms of the components of \mathbf{v} this is given by $v^i \mathbf{e}_i \mapsto v^i e_{ij} \mathbf{e}^j$. The inverse of this bijection is given by the inverse of the matrix (e_{ij}) , whose components we denote by (e^{ij}) . Using this matrix we can express the inverse of (A.1) as $v_i \mathbf{e}^i \mapsto v_j e^{ij} \mathbf{e}_i$. It is conventional to write $v_j = v^i e_{ij}$ and $v^i = v_j e^{ij}$. We will follow this convention in this paper. The symmetry of the inner product implies that the matrices (e_{ij}) and (e^{ij}) are symmetric as well.

For tensors $\mathbf{T} \in \mathbf{T}_0^2(V)$, $\mathbf{S} \in \mathbf{T}_1^1(V)$ and $\mathbf{R} \in \mathbf{T}_2^0(V)$ we can define the trace as

$$\text{tr}(\mathbf{T}) := T^{ij} e_{ij}, \quad \text{tr}(\mathbf{S}) := S_i^i, \quad \text{tr}(\mathbf{R}) := R_{ij} e^{ij}. \tag{A.2}$$

This definition as a contraction involves only the tensor itself and possibly the metric tensor and therefore independent of the choice of coordinate system [20, p. 82]. For an orthonormal coordinate system the traces of all three are just the ordinary traces of their coefficient matrix.

A.2. Tensor fields and coordinate systems on \mathbb{R}^n

We restrict the discussion now to $n = 2$, the dimension of the source \mathcal{E} . Just as we have vector fields on \mathbb{R}^2 , which assign to each point $\mathbf{x} \in \mathbb{R}^2$ a vector in $V = \mathbb{R}^2$, we can define a tensor field, which assigns to each point $\mathbf{x} \in \mathbb{R}^n$ a tensor in $\mathbf{T}_q^p(V) = \mathbf{T}_q^p(\mathbb{R}^2)$. It is important to note that the basis, with respect to which a tensor in this tensor field can be expressed in a certain point $\mathbf{x} \in \mathbb{R}^2$, in general depends on this point. For the basis vectors of a Cartesian coordinate system this is not the case, but for polar coordinates the basis vectors do depend on position. We will denote the space of tensor fields of contravariant rank p and covariant rank q on some subset \mathcal{U} of \mathbb{R}^2 by $\mathbf{T}_q^p(\mathcal{U})$.

In principle the basis in each point \mathbf{x} can be chosen independently of the coordinate system used to represent the points in \mathbb{R}^2 . However, a logical choice is to use the vectors found by taking the tangent vectors to the coordinate lines. For Cartesian coordinates this gives the usual Cartesian basis $\{\mathbf{e}_x, \mathbf{e}_y\}$ and for polar coordinates we find

$$\mathbf{e}_r := \frac{\partial}{\partial r}(r \cos(\theta)\mathbf{e}_x + r \sin(\theta)\mathbf{e}_y) = \cos(\theta)\mathbf{e}_x + \sin(\theta)\mathbf{e}_y, \tag{A.3a}$$

$$\mathbf{e}_\theta := \frac{\partial}{\partial \theta}(r \cos(\theta)\mathbf{e}_x + r \sin(\theta)\mathbf{e}_y) = -r \sin(\theta)\mathbf{e}_x + r \cos(\theta)\mathbf{e}_y. \tag{A.3b}$$

A basis defined by taking the tangent vectors to the coordinate lines is called a coordinate basis or holonomic basis. It is clear that not for every choice of basis vector fields there exists a coordinate system to which these basis vector fields are the tangent vector fields. An often used example of such a choice of basis is the orthonormal polar coordinate basis, which is found by rescaling (A.3b) to unit length. Such a type of basis is called anholonomic.

A.3. Covariant derivative of tensor fields

We now introduce the covariant derivative on \mathbb{R}^2 . On Euclidean spaces the covariant derivative is nothing more than the directional derivative. Suppose $f \in C^1(\mathbb{R}^2)$ then the directional derivative of f in the direction of the vector $\mathbf{v} \in T_0^1(\mathbb{R}^2)$ at $\mathbf{x} \in \mathbb{R}^2$ is defined as

$$\nabla_{\mathbf{v}} f(\mathbf{x}) := \left. \frac{d}{dt} (f \circ \gamma_{\mathbf{v}})(t) \right|_{t=0},$$

where $\gamma_{\mathbf{v}}$ is a curve such that $\gamma_{\mathbf{v}}(0) = \mathbf{x}$ and $d\gamma_{\mathbf{v}}(t)/dt|_{t=0} = \mathbf{v}$. When \mathbf{e}_i is a coordinate basis vector then we have $\nabla_{\mathbf{e}_i} f = \frac{\partial f}{\partial x^i}$, where x^i is the coordinate corresponding to \mathbf{e}_i . The differential of f is the covector given by $df := (\nabla_{\mathbf{e}_i} f)\mathbf{e}^i$. The vector corresponding to this covector by bijection (A.1) is $\nabla f = e^{ij}(\nabla_{\mathbf{e}_i} f)\mathbf{e}_j$ and called the gradient of f .

When we take the directional derivatives of vectors, covectors and general tensors we have to take care of the fact that the basis vectors might also change in the direction of the derivative. The directional derivatives of vectors in \mathbb{R}^2 are itself vectors in \mathbb{R}^2 and therefore there exist coefficients Γ_{ji}^k such that $\nabla_{\mathbf{e}_i}(\mathbf{e}_j) = \Gamma_{ji}^k \mathbf{e}_k$. These coefficients are called Christoffel symbols. Using these we can determine the directional derivative of contravariant tensors. For $\mathbf{v} \in T_0^1(\mathbb{R}^2)$, $\mathbf{T} \in T_0^2(\mathbb{R}^2)$, we have by the product rule

$$\nabla_{\mathbf{e}_i} \mathbf{v} = \nabla_{\mathbf{e}_i}(v^j)\mathbf{e}_j + v^j \nabla_{\mathbf{e}_i}(\mathbf{e}_j) = (\nabla_{\mathbf{e}_i}(v^j) + v^k \Gamma_{ki}^j)\mathbf{e}_j, \tag{A.4a}$$

$$\nabla_{\mathbf{e}_i} \mathbf{T} = \nabla_{\mathbf{e}_i}(T^{jk})\mathbf{e}_j \otimes \mathbf{e}_k + T^{jk}(\nabla_{\mathbf{e}_i}\mathbf{e}_j) \otimes \mathbf{e}_k + T^{jk}\mathbf{e}_j \otimes (\nabla_{\mathbf{e}_i}\mathbf{e}_k) = (\nabla_{\mathbf{e}_i}(T^{jk}) + T^{lk}\Gamma_{li}^j + T^{jl}\Gamma_{li}^k)\mathbf{e}_j \otimes \mathbf{e}_k. \tag{A.4b}$$

We will sometimes use the notation $D_i v^j = \nabla_{\mathbf{e}_i}(v^j) + v^k \Gamma_{ki}^j$ and $D_i T^{jk} = \nabla_{\mathbf{e}_i}(T^{jk}) + T^{lk}\Gamma_{li}^j + T^{jl}\Gamma_{li}^k$. Furthermore, we will use the notational convention $D^i = e^{ij}D_j$. Using the fact that the inner product and hence the corresponding metric tensor \mathbf{e} does not change with position [20, p. 215] we have $D_k(e^{ij}) = 0$ from which it follows that $\nabla_{\mathbf{e}_k}(e^{ij}) = -\Gamma_{lk}^j e^{il} - \Gamma_{lk}^i e^{lj}$ by (A.4b). Using this it follows from calculating $D_k(v^i) = D_k(e^{ij}v_j)$ and $D_m(S^{ij}) = D_m(e^{ik}e^{jl}S_{kl})$ that we have for $\hat{\mathbf{v}} \in T_1^0(\mathbb{R}^2)$, $\mathbf{S} \in T_2^0(\mathbb{R}^2)$, $\mathbf{R} \in T_1^1(\mathbb{R}^2)$ that

$$\nabla_{\mathbf{e}_i} \hat{\mathbf{v}} = D_i(v_j)\mathbf{e}^j = (\nabla_{\mathbf{e}_i}(v_j) - v_k \Gamma_{ji}^k)\mathbf{e}^j, \tag{A.5a}$$

$$\nabla_{\mathbf{e}_i} \mathbf{S} = D_i(S_{jk})\mathbf{e}^j \otimes \mathbf{e}^k = (\nabla_{\mathbf{e}_i}(S_{jk}) - S_{lk}\Gamma_{ji}^l - S_{jl}\Gamma_{ki}^l)\mathbf{e}^j \otimes \mathbf{e}^k. \tag{A.5b}$$

$$\nabla_{\mathbf{e}_i} \mathbf{R} = D_i(R_k^j)\mathbf{e}_j \otimes \mathbf{e}^k = (\nabla_{\mathbf{e}_i}(R_k^j) + R_k^l \Gamma_{li}^j - R_l^j \Gamma_{ki}^l)\mathbf{e}_j \otimes \mathbf{e}^k. \tag{A.5c}$$

A tensor of special interest in this paper is found by again differentiating the differential of a function. This tensor is the Hessian tensor and defined as [26, p. 172]

$$\mathbf{H}(v) := \nabla_{\mathbf{e}_j}(dv) \otimes \mathbf{e}^j = D_j(\nabla_{\mathbf{e}_i} v)\mathbf{e}^i \otimes \mathbf{e}^j = (\nabla_{\mathbf{e}_j}(\nabla_{\mathbf{e}_i} v) - \Gamma_{ij}^k \nabla_{\mathbf{e}_k} v)\mathbf{e}^i \otimes \mathbf{e}^j. \tag{A.6}$$

The covariant directional derivative as introduced above can be generalized to Riemannian manifolds and is in that context called a Levi-Civita connection [26, p. 160]. For a Levi-Civita connection the Hessian matrix is symmetric [28, p. 4], hence the Hessian matrix will always be symmetric in this paper. Note that in Cartesian coordinates $(h_{ij}(v))$, the coefficient matrix of $\mathbf{H}(v)$, is the matrix with second derivatives of v .

Appendix B. Equivalence of boundary conditions for a strictly convex reflector surface

In this appendix we show that the Inverse Reflector Problem as stated on page 353 is equivalent to this same problem but with the boundary condition $\nabla u(\partial \mathcal{E}) = \partial \mathcal{F}$ replaced by $\nabla u(\mathcal{E}) = \mathcal{F}$. When doing this we need to make use of the fact that ∇u is an open map, i.e., a map that maps open sets to open sets. This we show in the following lemma.

Lemma 6. Suppose that $u \in C^2(\mathcal{E})$ is the strictly convex solution to the Inverse Reflector Problem with either $\nabla u(\mathcal{E}) = \mathcal{F}$ or $\nabla u(\partial \mathcal{E}) = \partial \mathcal{F}$ as boundary condition. Then the map ∇u is open, i.e., for each open subset $A \subset \mathcal{E}$ the image $\nabla u(A)$ is an open subset of \mathcal{F} .

Proof. In Lemma 1 we saw that for the strictly convex solution $u \in C^2(\mathcal{E})$, ∇u is a bijection. Moreover, because u is twice continuously differentiable, the mapping ∇u is a continuously differentiable mapping. In Cartesian coordinates, the matrix (h_{ij}) is also the Jacobian matrix of ∇u . The fact that $\det(h_{ij}) > 0$ therefore implies that the Jacobian of ∇u is always strictly positive. This implies that the conditions for the inverse function theorem [29] are satisfied. The inverse function theorem states, among other things, that for every open subset A of \mathcal{E} and $\mathbf{x} \in A$, there exists an open set $U_{\mathbf{x}}$ in A containing \mathbf{x} , and an open set $V_{\mathbf{x}}$ in \mathcal{F} containing $\nabla u(\mathbf{x})$ such that ∇u is a bijection from $U_{\mathbf{x}}$ to $V_{\mathbf{x}}$ and the inverse $(\nabla u)^{-1}$ is continuously differentiable on $V_{\mathbf{x}}$.

From this it follows that ∇u is open. To see this, suppose A is some open set in \mathcal{E} . By the inverse function theorem there exists for every $\mathbf{x} \in \mathcal{E}$ open sets $U_{\mathbf{x}}$ and $V_{\mathbf{x}}$ such that $\mathbf{x} \in U_{\mathbf{x}}$, $\nabla u(\mathbf{x}) \in V_{\mathbf{x}}$ and $U_{\mathbf{x}} \subset A$. $\nabla u(U_{\mathbf{x}}) = V_{\mathbf{x}}$ is open for every $\mathbf{x} \in A$. Notice that $\cup_{\mathbf{x} \in A} U_{\mathbf{x}} = A$ and that $\nabla u(A) = \cup_{\mathbf{x} \in A} \nabla u(U_{\mathbf{x}}) = \cup_{\mathbf{x} \in A} V_{\mathbf{x}}$. Thus ∇u is an open map. \square

The map ∇u is a homeomorphism from \mathcal{E} to \mathcal{F} , because it is a continuous bijection which is open and hence also has continuous inverse. We will use this convenient property of ∇u in the following lemma.

Lemma 7. Let $u \in C^2(\mathcal{E})$ be the strictly convex solution to the Inverse Reflector Problem with $\nabla u(\mathcal{E}) = \mathcal{F}$ instead of $\nabla u(\partial \mathcal{E}) = \partial \mathcal{F}$. Then also $\nabla u(\partial \mathcal{E}) = \partial \mathcal{F}$.

Proof. The map ∇u is a homeomorphism and therefore it links every open set in \mathcal{E} with an open set in \mathcal{F} and vice versa. Let us denote by $\text{int}(A)$ the interior of a set A . Suppose $A \subset \mathcal{E}$. We have $\nabla u(\text{int}(A)) \subset \nabla u(A)$ and because ∇u is an open map $\nabla u(\text{int}(A))$ is also open. The largest open subset of $\nabla u(A)$ is the interior $\text{int}(\nabla u(A))$, therefore $\nabla u(\text{int}(A)) \subset \text{int}(\nabla u(A))$. If $\nabla u : \mathcal{E} \rightarrow \mathcal{F}$ is a homeomorphism, then $(\nabla u)^{-1} : \mathcal{F} \rightarrow \mathcal{E}$ is a homeomorphism also, hence $(\nabla u)^{-1}(\text{int}(B)) \subset \text{int}((\nabla u)^{-1}(B))$ for all $B \subset \mathcal{F}$. From this it follows that we have both $\nabla u(\text{int}(\mathcal{E})) \subset \text{int}(\nabla u(\mathcal{E})) = \text{int}(\mathcal{F})$ and $(\nabla u)^{-1}(\text{int}(\mathcal{F})) \subset \text{int}((\nabla u)^{-1}(\mathcal{F})) = \text{int}(\mathcal{E})$. Using this we see that

$$\text{int}(\mathcal{F}) = \nabla u((\nabla u)^{-1}(\text{int}(\mathcal{F}))) \subset \nabla u(\text{int}(\mathcal{E})) \subset \text{int}(\mathcal{F}).$$

Thus, we see that $\nabla u(\text{int}(\mathcal{E})) = \text{int}(\mathcal{F})$. Now, because ∇u is a bijection this implies that we must have $\nabla u(\partial \mathcal{E}) = \partial \mathcal{F}$. \square

Thus the strictly convex solution of the Inverse Reflector Problem with boundary condition $\nabla u(\mathcal{E}) = \mathcal{F}$ is also a solution to the Inverse Reflector Problem with boundary condition $\nabla u(\partial \mathcal{E}) = \partial \mathcal{F}$. Now the following lemma states the converse.

Lemma 8. Let $u \in C^2(\mathcal{E})$ be a strictly convex solution to Inverse Reflector Problem. Then $\nabla u(\mathcal{E}) = \mathcal{F}$.

Proof. The map ∇u is a homeomorphism from \mathcal{E} to $\nabla u(\mathcal{E}) \subset \mathbb{R}^2$. The set \mathcal{E} is convex and hence simply connected. The set $\partial \mathcal{E}$ is a simple and closed curve, i.e., a Jordan curve. The map ∇u is continuous and injective and hence $\nabla u(\partial \mathcal{E}) = \partial \mathcal{F}$ is also a Jordan curve. Now the Jordan curve theorem [30, p. 198] states that the complement $\mathbb{R}^2 \setminus \partial \mathcal{F}$ has two connected components one of which is bounded and one of which is not, namely the interior and the exterior of the curve, and the boundary of both these sets is $\partial \mathcal{F}$. The set \mathcal{E} is simply connected, therefore $\nabla u(\mathcal{E})$ is simply connected also. The interior and exterior to the curve $\nabla u(\partial \mathcal{E}) = \partial \mathcal{F}$ are the only two subsets of \mathbb{R}^2 with $\partial \mathcal{F}$ as boundary. The fact that ∇u is a homeomorphism implies $\nabla u(\partial \mathcal{E}) = \partial(\nabla u(\mathcal{E}))$, because $\nabla u(\text{int}(\mathcal{E})) = \text{int}(\nabla u(\mathcal{E}))$ as we showed in the proof of Lemma 7. The fact that $\partial \mathcal{F} = \nabla u(\partial \mathcal{E}) = \partial(\nabla u(\mathcal{E}))$ implies that $\text{int}(\nabla u(\mathcal{E}))$ is one of two sets of the Jordan curve theorem. The exterior set is clearly not simply connected, while $\nabla u(\mathcal{E})$ is, therefore $\text{int}(\nabla u(\mathcal{E}))$ is the interior set in the Jordan curve theorem. The fact that \mathcal{E} is bounded, that the functions E and F are strictly positive and bounded and that F is bounded away from zero implies by (10) that the set \mathcal{F} is bounded also. This implies that $\text{int}(\mathcal{F})$ needs to be the interior set also and hence we find that $\nabla u(\mathcal{E}) = \mathcal{F}$. \square

We have established that u is a strictly convex solution to the Inverse Reflector Problem with boundary conditions $\nabla u(\mathcal{E}) = \mathcal{F}$ if and only if it is a strictly convex solution to the Inverse Reflector Problem with boundary condition $\nabla u(\partial \mathcal{E}) = \partial \mathcal{F}$. Thus the two boundary conditions are equivalent.

References

- [1] R. Haitz, Y.T. Tsao, Solid-state lighting: 'The Case' 10 years after and future prospects, *Phys. Status Solidi A* 208 (1) (2011) 17–29.
- [2] T. Taguchi, Present status of energy saving technologies and future prospects in white LED lighting, *IEEJ Trans. Electr. Electron. Eng.* 3 (2008) 21–26.
- [3] A.S. Glassner, *An introduction to Ray Tracing*, Academic Press, 1991.
- [4] F.Z. Fang, X.D. Zhang, A. Weckenmann, G.X. Zhang, C. Evans, Manufacturing and measurement of freeform optics, *CIRP Ann. – Manuf. Technol.* 62 (2) (2013) 823–846.
- [5] C. Villani, *Topics in Optimal Transportation*, American Mathematical Society, Providence, 2003.
- [6] V. Oliker, Designing freeform lenses for intensity and phase control of coherent light with help from geometry and mass transport, *Arch. Ration. Mech. Anal.* 201 (3) (2011).
- [7] J.D. Benamou, Y. Brenier, A computational fluid mechanics solution to the Monge–Kantorovich mass transfer problem, *Numer. Math.* 84 (2000) 375–393.

- [8] E. Haber, T. Rehman, A. Tannenbaum, An efficient numerical method for the solution of the L_2 optimal mass transfer problem, *SIAM J. Sci. Comput.* 32 (2010) 197–211.
- [9] B.D. Froese, A.M. Oberman, Fast finite difference solvers for singular solutions of the elliptic Monge–Ampère equation, *J. Comput. Phys.* 230 (3) (2011) 818–834.
- [10] J.D. Benamou, B.D. Froese, A.M. Oberman, Numerical solution of the optimal transportation problem using the Monge–Ampère equation, *J. Comput. Phys.* 260 (2014) 107–126.
- [11] K. Brix, Y. Hafizogullari, A. Platen, Solving the Monge–Ampère equations for the inverse reflector problem, *Math. Models Methods Appl. Sci.* 25 (2015) 803–837.
- [12] C.R. Prins, Inverse Methods for Illumination Optics, Ph.D. Thesis, Eindhoven University of Technology, 2014.
- [13] C.R. Prins, R. Beltman, J.H.M. ten Thije Boonkkamp, W.L. IJzerman, T.W. Tukker, A least-squares method for optimal transport using the Monge–Ampère equation, *SIAM J. Sci. Comput.* 37 (6) (2015) B937–B961.
- [14] A. Caboussat, R. Glowinski, D.C. Sorensen, A least-squares method for the numerical solution of the Dirichlet problem for the elliptic Monge–Ampère equation in dimension two, *ESAIM Control Optim. Calc. Var.* 19 (2013) 780–810.
- [15] Official website on Escher: <http://www.mcescher.com/about/>.
- [16] M. Bass (Ed.), *Handbook of Optics Volume II – Devices, Measurements and Properties*, 2nd ed, McGraw-Hill, 1995.
- [17] G. Monge, *Application de l'analyse a la geometrie, a l'usage de l'Ecole imperiale polytechnique*, Bernard, Paris, 1807.
- [18] H. Cohn, *Conformal Mapping on Riemann Surfaces*, McGraw-Hill, 1967.
- [19] W. Rudin, *Real and Complex Analysis*, 3rd ed., McGraw-Hill, London, 1987.
- [20] C.W. Misner, K.S. Thorne, J.A. Wheeler, *Gravitation*, W. H. Freeman, San Francisco, 1973.
- [21] R. Aris, *Vectors, Tensors, and the Basic Equations of Fluid Mechanics*, Dover Publications, 1989.
- [22] M. Spivak, *Calculus on Manifolds: a Modern Approach to Classical Theorems of Advanced Calculus*, Benjamin, Amsterdam, 1965.
- [23] J.E. Marsden, A.J. Tromba, *Vector Calculus*, 5th ed., W.H. Freeman and Company, 2003.
- [24] S. Boyd, L. Vandenberghe, *Convex Optimization*, Cambridge University Press, 2004.
- [25] R. Courant, D. Hilbert, *Methods of Mathematical Physics*, vol. 1, 1989 ed., John Wiley & Sons, 1937.
- [26] J. Jost, *Riemannian Geometry and Geometric Analysis*, 6th ed., Springer-Verlag, Berlin Heidelberg, 2011.
- [27] S. Kobayashi, K. Nomizu, *Foundations of Differential Geometry*, Interscience, S.I., 1963.
- [28] J. Morgan, G. Tian, *Ricci Flow and the Poincaré Conjecture*, American Mathematical Society, Providence RI, 2007.
- [29] T. Tao, *Analysis II*, 2nd ed., Hindustan Book Agency, 2009.
- [30] E.H. Spanier, *Algebraic Topology*, McGraw-Hill, London, 1966.