

Published in final edited form as:

SIAM J Imaging Sci. ; 5(1): 434–464. doi:10.1137/100801664.

## APPROXIMATING SYMMETRIC POSITIVE SEMIDEFINITE TENSORS OF EVEN ORDER\*

ANGELOS BARMPOUTIS, HO JEFFREY, and BABA C. VEMURI

### Abstract

Tensors of various orders can be used for modeling physical quantities such as strain and diffusion as well as curvature and other quantities of geometric origin. Depending on the physical properties of the modeled quantity, the estimated tensors are often required to satisfy the positivity constraint, which can be satisfied only with tensors of even order. Although the space  $\mathcal{P}_0^{2m}$  of  $2m^{\text{th}}$ -order symmetric positive semi-definite tensors is known to be a convex cone, enforcing positivity constraint directly on  $\mathcal{P}_0^{2m}$  is usually not straightforward computationally because there is no known analytic description of  $\mathcal{P}_0^{2m}$  for  $m > 1$ . In this paper, we propose a novel approach for enforcing the positivity constraint on even-order tensors by approximating the cone  $\mathcal{P}_0^{2m}$  for the cases  $0 < m < 3$ , and presenting an explicit characterization of the approximation  $\Sigma_{2m} \subset \Omega_{2m}$  for  $m \geq 1$ , using the subset  $\Omega_{2m} \subset \mathcal{P}_0^{2m}$  of semi-definite tensors that can be written as a sum of squares of tensors of order  $m$ . Furthermore, we show that this approximation leads to a non-negative linear least-squares (NNLS) optimization problem with the complexity that equals the number of generators in  $\Sigma_{2m}$ . Finally, we experimentally validate the proposed approach and we present an application for computing  $2m^{\text{th}}$ -order diffusion tensors from Diffusion Weighted Magnetic Resonance Images.

### Keywords

high-order tensors; sum of squares of polynomials; diffusion tensor imaging

## 1. Introduction

Multi-linear algebra is a generalization of linear algebra and tensors which are multi-linear forms are widely used for modeling various physical quantities commonly encountered in engineering and physics. Elasticity [34], stress, strain and diffusion [10] are some examples. In differential geometry, tensors are used to represent metrics, curvatures [40] and other geometric quantities. In image processing, structure tensors [46] have been used for texture analysis, trifocal tensors in multi-view geometry, etc. The tensors in most of these applications are required to satisfy certain properties. For example, the tensors that approximate the Bidirectional Reflectance Distribution Function (BRDF) [7] are anti-symmetric, while the diffusion [10] and the structure tensors [46] are antipodally symmetric. Furthermore, certain applications demand that the estimated tensors be positive-definite since they model positive-valued physical quantities such as the diffusivity function or the displacement probability of water molecules [8]. In this paper, we are interested in the case of fully symmetric positive-definite tensors of various orders and hence for sake of simplicity, every reference to the term *tensor* will imply this particular case of tensors unless otherwise stated.

\*This research was supported by the NIH grant EB007082 & NSF066340 to BCV.

Let  $\mathcal{P}^m$  denote the set of  $m^{\text{th}}$ -order symmetric positive-definite tensors in  $\mathbb{R}^3$ . As is well-known, positivity condition requires the order  $m$  to be even. Denote  $\mathcal{P}_0^m$  the closure of  $\mathcal{P}^m$  consisting of symmetric positive semi-definite tensors (PSD) in  $\mathbb{R}^3$ . As subsets of the space  $\mathbb{S}^m$  of  $m^{\text{th}}$ -order symmetric tensors,  $\mathcal{P}^m, \mathcal{P}_0^m$  are cones, convex subsets that are invariant under positive scaling [18]. In most applications, the main computational problem can be formulated as data interpolation problem with the domain being  $\mathcal{P}_0^{2m}$ . Specifically, the input data are often in the form  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_k, y_k)\}$  where  $\mathbf{x}_i$  are directions in  $\mathbb{R}^3$  represented as points on the unit sphere  $\mathbb{S}^2$ , and  $y_i$  are the values to be interpolated. The interpolation problem requires a non-negative tensor  $\mathbf{T} \in \mathcal{P}_0^{2m}$  that interpolates the input data. Formulated as a least-squares problem, it has the form

$$\mathbf{T} = \arg \min_{p \in \mathcal{P}_0^{2m}} \sum_{i=1}^k |y_i - p(\mathbf{x}_i)|^2.$$

We note that both the objective function and the domain  $\mathcal{P}_0^{2m}$  are convex, and therefore, the optimization problem above is in fact a convex optimization problem that, in principle, can be solved using existing techniques [12]. However, a formal and significant difficulty of applying these methods is that except for the  $m = 1$  case, there exists no known description of the cone  $\mathcal{P}_0^{2m}$  as it is well-known that the positivity test for polynomials of degree  $m > 2$  is a difficult problem. In the second-order case, the cone  $\mathcal{P}_0^2$  is known to be self-dual in the sense that there exists an inner product  $\langle \cdot, \cdot \rangle$  on  $\mathbb{S}^2$  such that  $\langle A, B \rangle \geq 0$  for any  $A, B \in \mathcal{P}_0^2$ . The inner product allows the extension of the usual duality theory using Lagrange multipliers to the cone  $\mathcal{P}_0^2$ , and there is a well-developed theory of semi-definite programming (SDP) [12] that deals with *linear* objective functions on  $\mathcal{P}_0^2$ .

While the difficulty of providing a complete description of  $\mathcal{P}_0^{2m}$  seems to be unsurmountable at this point, the main contribution of this paper is the realization of another formal difficulty that can be overcome relatively easily. A cone  $C$  in a vector space is said to be finitely-generated if there exists a finite number of elements  $v_1, \dots, v_n \in C$ , its generators, such that every element  $c \in C$  can be written as a non-negative linear combination of the generators

$$c = a_1 v_1 + \dots + a_n v_n, \quad a_1, \dots, a_n \geq 0.$$

If the cone  $\mathcal{P}_0^{2m}$  were finitely generated, the above optimization problem becomes a non-negative linear least-squares (NNLS) problem, with complexity (number of variables) equals to the number of generators. The advantage of solving an NNLS problem is that there are software packages that can efficiently solve NNLS problems containing thousands of variables [28]. While  $\mathcal{P}_0^{2m}$  is not finitely-generated, it follows naturally that we can try to approximate  $\mathcal{P}_0^{2m}$  with a finitely-generated subcone, and restrict the above optimization problem to the subcone. The restriction can be justified if the subcone can be shown to be a good approximation of  $\mathcal{P}_0^{2m}$ .

The second contribution of this paper is an explicit characterization of the approximations  $\Sigma_{2m} \subset \mathcal{P}_0^{2m}$  for  $0 < m < 3$ , and  $\Sigma_{2m} \subset \Omega_{2m}$  for  $m \geq 1$ , where  $\Sigma_{2m}$  is a finitely-generated

subcone in the respective spaces. More specifically, let  $\Omega_{2m}$  denote the subcone in  $\mathcal{P}_0^{2m}$  consisting of semi-definite tensors that can be written as a sum of squares of tensors of order  $m$ . We have the natural inclusions  $\sum_{2m} \subset \Omega_{2m} \subset \mathcal{P}_0^{2m}$ , and our result gives a detailed characterization of the approximation  $\sum_{2m} \subset \Omega_{2m}$  in terms of the geometry of the generators in  $\sum_{2m}$ . In particular, for  $m = 1, 2$ , it is known that  $\Omega_{2m} = \mathcal{P}_0^{2m}$ , and our result then gives a detailed characterization of the approximation  $\sum_{2m} \subset \mathcal{P}_0^{2m}$ . Our analysis have shown that, for the lower-order cases  $m = 1, 2, 3$ , which are of primary interest here, for a reasonable precision requirement,  $\Omega_{2m}$  can be approximated by  $\sum_{2m}$  containing a few hundreds or at most a few thousands of generators. It follows that the corresponding NNLS problems have the complexity that are well within the capability of currently available NNLS algorithms [28]. We quantitatively validate our method via several experiments, and we also present an application of the proposed technique for estimating the diffusivity function from diffusion-weighted MRI to demonstrate both the efficiency and accuracy of the proposed method.

The rest of this paper is organized as follows: In Sec. 2, we define the finitely-generated subcone  $\sum_{2m}$ . We also develop the theory that quantifies the approximation  $\sum_{2m} \subset \Omega_{2m}$ , and the main theorem proved in this section relates the approximation error with the geometry of the generators in  $\sum_{2m}$ . Using the theory developed in Sec. 2, in Sec. 3 we explicitly work out the formulas for the number of generators for  $\sum_{2m}$  required for a given accuracy requirement. The results show that, up to order-6 and depending on the order, it generally requires at most a few thousands of generators for  $\sum_{2m}$  in order to achieve a relative approximation error of less than 10%. Finally, in Sec. 4, we validate our theoretical findings using a set of experiments and we present an application of our method on diffusion-weighted MR datasets.

## Related Work

Symmetric positive-definite (SPD) tensors of order-2 have been used in modeling the diffusivity function in the so called Diffusion Tensor MR Imaging (DT-MRI) [10]. SPD matrices can be endowed with a Riemannian metric that is invariant under affine transforms. This metric or its approximations have been employed for estimating and processing diffusion tensor fields [48, 47, 29, 38, 18, 9]. Tensors of 3<sup>rd</sup> and 5<sup>th</sup> order can model reflectance distributions with specularities and cast shadows in facial images and have been used for re-lighting in [7]. In general, odd-order tensors are generalizations of the order-1 tensor, which have been commonly used in computer graphics for representing the Lambertian reflectance model. Similarly, 4<sup>th</sup>, 6<sup>th</sup> or higher even-order tensors generalize the 2<sup>nd</sup>-order tensors and have the ability to approximate multi-lobed functions [35, 30, 36] such as the kurtosis of diffusion [26]. In particular, some 4<sup>th</sup>-order tensors can be expressed as 2<sup>nd</sup>-order tensors in higher dimensions and their properties have been studied in detail by Moakher in [32, 33]. They however do not span the full space of the higher-order tensors as was shown in the case of order-4 tensors in [6, 5]. In [20], Ghosh et al. used the metric proposed by Moakher in [32, 33] to represent the space of 4<sup>th</sup>-order SPD tensors using the geometry of 2<sup>nd</sup>-order SPD tensors in higher dimensions. Recently, an algorithm for imposing positivity constraints on 4<sup>th</sup>-order tensors using their equivalent ternary quartic polynomial representation was proposed in [6] and this was further developed in [5] and [21, 49].

After estimating a field of high-order tensors, it can be processed using a Finsler metric by appropriately modifying the polynomial equivalent representation of the tensors that satisfy the properties of Finsler geometry [4]. This method can be used for neuronal fiber tracking from high angular resolution diffusion MRI data. Further processing of higher-order tensor fields can be achieved by using the eigenvalue decomposition of matrices which has been

extended for the case of high-order tensors in [23]. In this framework, the eigenvalues correspond to the extreme values (minima or maxima) of a tensor and they can be used to extract useful information from the kurtosis tensor [42] as well as the orientation of maximum diffusion [11, 22]. Another method for extracting the principal orientation of diffusion from a higher-order tensor was recently described in [44].

Although, high-order tensors have been employed in most of the aforementioned methods due to their simple polynomial form and their ability to model multi-lobed spherical functions, there are no existing methods for imposing positivity constraints in symmetric tensors of any order higher than two and four. The need to impose positivity constraints becomes essential especially in the case where the tensors approximate positive-valued physical quantities, and it has been shown that imposing the positivity constraint on the tensors approximating the diffusivity function being estimated reduces the approximation errors significantly [5]. Recently, Pasternak et al. [37] also emphasized the importance of enforcing positivity constraints in processing diffusion tensor MR images.

Finally, although Cartesian tensors basis have been widely used for modeling the diffusivity function in DW-MRI, we would like to mention that Spherical Harmonic basis have been employed in approximating other spherical functions involved in DW-MRI processing such as the diffusion propagator. A detailed review of several multi fiber reconstruction methods that employ spherical harmonic basis can be found in the recent article by Descoteaux et al. on Diffusion Propagator Imaging [17]. The orientation distribution function (ODF) is another example of a DW-MRI related spherical function, which can be reconstructed from Q-ball imaging data [16, 13, 3] and was recently done in [2] by using the mathematically correct definition of ODF and deriving a closed form expression for the same. In this article, however, our main focus is on the use of Cartesian tensor basis for parameterizing the diffusivity function in DW-MR datasets.

## 2. Theory

We will consider symmetric tensors of order  $m$  as functions defined on the unit sphere  $\mathbf{S}^2$  in  $\mathbb{R}^3$ . In particular, symmetric tensors of order  $m$  can be identified with homogeneous polynomials of degree  $m$ : for a symmetric tensor  $\mathbf{T}$  of order  $m$ , its associated homogeneous polynomial  $P(x, y, z)$  is given as

$$P(x, y, z) = \mathbf{T}(\underbrace{\mathbf{x}, \dots, \mathbf{x}}_m),$$

where  $\mathbf{x} = [x \ y \ z]^T$ . Under this identification,  $\mathcal{P}^m$  are homogeneous polynomials of degree  $m$  that do not vanish on  $\mathbf{S}^2$ , and similarly,  $\mathcal{P}_0^m$  are degree- $m$  homogeneous polynomials that do not take negative values in  $\mathbb{R}^3$ . Both are now considered as cones in  $\mathbf{H}_m$ , the set of homogeneous polynomials of degree  $m$ . For even degree  $2m$ , let  $\Omega_{2m}$  denote the subset of  $\mathcal{P}_0^{2m}$  consisting of polynomials that can be written as a sum of squares of polynomials of degree  $m$ .  $\Omega_{2m}$  is clearly a subcone of  $\mathcal{P}_0^{2m}$  for all  $m \geq 1$ , and for  $m = 1, 2$ , it is known that  $\Omega_{2m} = \mathcal{P}_0^{2m}$ : the  $m = 1$  case follows easily from linear algebra and  $m = 2$  case is the content of Hilbert's theorem on ternary quartics [24]. For  $m > 2$ , however, the inclusion is strict  $\Omega_{2m} \subsetneq \mathcal{P}_0^{2m}$ . In this section, we will describe a general method for approximating  $\Omega_{2m}$  using a finitely-generated subcone  $\Sigma_{2m}$  in  $\Omega_{2m}$ , and we will provide a characterization of the approximation error in terms of the geometry of the generators of  $\Sigma_{2m}$ . For the important

quadratic and quartic cases  $m = 1, 2$ , our result provides an approximation of the full PSD cone  $\mathcal{P}_0^{2m}$  using a finitely-generated subcone  $\Omega_{2m}$ .

The basic norm used in this paper is the  $L^1$ -norm over the sphere  $\mathbf{S}^2$ . More specifically, for any  $P \in \mathbf{H}_m$ , its  $L^1$ -norm  $\|P\|_1$  is the integral over  $\mathbf{S}^2$

$$\|P\|_1 = \int_{\mathbf{S}^2} |P(\mathbf{x})| d\mathbf{x}.$$

That it is indeed a norm follows from the fact that for two homogeneous polynomials  $P, Q, P = Q$  as polynomials if and only if  $\|P - Q\|_1 = 0$ . Note that the other norm properties are trivial to prove. For any  $P \in \mathcal{P}_0^{2m}$  and a subcone  $\Sigma_{2m} \subset \mathcal{P}_0^{2m}$ , we define the relative  $L^1$ -approximation error of  $P$  as

$$\mathbf{E}_{\Sigma_{2m}}(P) = \frac{\min_{p \in \Sigma_{2m}} \|P - p\|_1}{\|P\|_1}. \quad (2.1)$$

**Proposition 2.1:** Let  $\Sigma_{2m}$  be a closed subcone in  $\mathcal{P}_0^{2m}$  and  $P \in \mathcal{P}_0^{2m}$ .

1. The  $L^1$ -norm is convex: for any  $p \in \mathcal{P}_0^{2m}$ , the function  $g(q), q \in \mathcal{P}_0^{2m}$

$$g(q) = \|p - q\|_1$$

is a convex function on  $\mathcal{P}_0^{2m}$ .

2. For  $P \neq 0, \mathbf{E}_{\Sigma_{2m}}(P) = 0$  if and only if  $P \in \Sigma_{2m}$ . For any  $s > 0$ ,

$$\mathbf{E}_{\Sigma_{2m}}(sP) = \mathbf{E}_{\Sigma_{2m}}(P).$$

**Proof:** For any  $q_1, q_2 \in \mathcal{P}_0^{2m}$ ,

$$g(tq_1 + (1-t)q_2) = \int_{\mathbf{S}^2} |tp - tq_1 - (1-t)q_2| d\mathbf{x} \leq \int_{\mathbf{S}^2} |tp - tq_1| d\mathbf{x} + \int_{\mathbf{S}^2} |(1-t)p - (1-t)q_2| d\mathbf{x},$$

and the convexity of the norm on  $\mathcal{P}_0^{2m}$  follows. (2) is clear because  $\Sigma_{2m}$  is closed. The invariance of  $\mathbf{E}_{\Sigma_{2m}}$  under positive scaling follows readily from the definition.

Let  $\mathbf{m}_1(\mathbf{x}), \dots, \mathbf{m}_{d(m)}(\mathbf{x})$  denote the  $d(m) = \frac{(m+2)(m+1)}{2}$  monomials in  $\mathbf{H}_m$ . (Note that  $d(m)$  also equals to the number of symmetric spherical harmonic basis elements, which can be mapped to the monomials in  $\mathbf{H}_m$  using an one-to-one transformation [35, 15].) The monomials form a basis in  $\mathbf{H}_m$  that identifies  $\mathbf{H}_m$  with  $\mathbb{R}^{d(m)}$ . We will denote  $\mathbf{HS}_m$  the unit sphere in  $\mathbf{H}_m$ , consisting of polynomials

\$watermark-text

\$watermark-text

\$watermark-text

$$p(\mathbf{x}) = \sum_{i=1}^{d(m)} a_i \mathbf{m}_i(\mathbf{x})$$

such that  $a_1^2 + \dots + a_{d(m)}^2 = 1$ . The subcone  $\Sigma_{2m}$  will be defined using polynomials in  $\mathbf{HS}_m$ , and this is accomplished through the square map  $\mathcal{F}_m^2: \mathbf{HS}_m \rightarrow \mathbf{H}_{2m}$ :

$$\mathcal{F}_m^2(p) = p^2.$$

Clearly  $\mathcal{F}_m^2$  is a smooth map, and  $\mathcal{F}_m^2(p) = \mathcal{F}_m^2(q)$  if and only if  $p = \pm q$ . While  $\mathcal{F}_m^2$  is not linear, it maps rays in  $\mathbf{HS}_m$  to rays in  $\mathbf{H}_{2m}$ :  $\mathcal{F}_m^2(tp) = t^2 \mathcal{F}_m^2(p)$ . The geometry of the map  $\mathcal{F}_m^2$  will play a crucial role in our analysis below, and it is quantified by its condition number  $\eta_m$ . First, we define two quantities.

$$\eta_m^{\max} = \max_{p \in \mathbf{HS}_m} \|\mathcal{F}_m^2(p)\|_1, \quad \eta_m^{\min} = \min_{p \in \mathbf{HS}_m} \|\mathcal{F}_m^2(p)\|_1.$$

Clearly we have  $\eta_m^{\min} > 0$  since  $\mathbf{HS}_m$  does not contain the zero polynomial. The two numbers measure the amount of stretching and shrinking  $\mathcal{F}_m^2$  does to the sphere  $\mathbf{HS}_m$ . Their ratio gives the condition number  $\eta_m$  for  $\mathcal{F}_m^2$

$$\eta_m = \frac{\eta_m^{\max}}{\eta_m^{\min}}.$$

In the following, we will often drop the subscript and denote the condition number simply as  $\eta$  when the degree  $m$  in the context is clear. Figure 2.1 illustrates the effect of  $\mathcal{F}_m^2$  and its condition number  $\eta$ .

**Proposition 2.2:**  $\eta_m^{\max}$ ,  $\eta_m^{\min}$  and hence  $\eta$  can be determined by evaluating  $\frac{d(m)^2 + d(m)}{2}$  trigonometric integrals.

**Proof:** Let  $\mathbf{m}_1, \dots, \mathbf{m}_{d(m)}$  denote the  $d(m)$  monomials in  $\mathbf{HS}_m$ . A polynomial  $p \in \mathbf{HS}_m$  is identified with the vector of coefficients  $\mathbf{a} = [a_1, \dots, a_{d(m)}]^\top$  as  $p = a_1 \mathbf{m}_1 + \dots + a_{d(m)} \mathbf{m}_{d(m)}$ . The  $L^1$ -norm  $\|\mathcal{F}_m^2(p)\|_1$  is the integral of  $p^2$  over  $\mathbf{S}^2$  that can be written as

$$\|\mathcal{F}_m^2(p)\|_1 = \sum_{i,j=1}^{d(m)} a_i a_j \int_{\mathbf{S}^2} \mathbf{m}_i(\mathbf{x}) \mathbf{m}_j(\mathbf{x}) \, d\mathbf{x}.$$

Let  $\Lambda^m$  denote the  $d(m) \times d(m)$  matrix whose components  $\Lambda_{ij}^m$  are the integrals  $\int_{\mathbf{S}^2} \mathbf{m}_i(\mathbf{x}) \mathbf{m}_j(\mathbf{x}) \, d\mathbf{x}$ , we have

$$\|\mathcal{F}_2(p)\|_1 = \mathbf{a}^\top \Lambda^m \mathbf{a}.$$

It follows that  $\eta_m^{\max}, \eta_m^{\min}$  can be determined as

$$\eta_m^{\max} = \min_{\mathbf{a}^\top \mathbf{a} = 1} \mathbf{a}^\top \Lambda^m \mathbf{a}, \quad \eta_m^{\min} = \min_{\mathbf{a}^\top \mathbf{a} = 1} \mathbf{a}^\top \Lambda^m \mathbf{a},$$

both of which can be solved once  $\Lambda^m$  is known using Singular Value Decomposition. The integrals  $\int_{\mathbf{S}^2} \mathbf{m}_i(\mathbf{x}) \mathbf{m}_j(\mathbf{x}) \, d\mathbf{x}$  can be computed in closed form since using spherical coordinates,  $x = \sin \psi \cos \theta$ ,  $y = \sin \psi \sin \theta$ ,  $z = \cos \psi$ , each integral is a product of two trigonometric integrals

$$\int_{\mathbf{S}^2} \mathbf{m}_i(\mathbf{x}) \mathbf{m}_j(\mathbf{x}) \, d\mathbf{x} = \left( \int_{\theta=0}^{2\pi} \cos^{b_1} \theta \sin^{b_2} \theta \, d\theta \right) \left( \int_{\psi=0}^{\pi} \cos^{b_3} \psi \sin^{b_4} \psi \, d\psi \right),$$

with exponents  $b_1, b_2, b_3, b_4$  depending on  $\mathbf{m}_i, \mathbf{m}_j$ .

In practice,  $\Lambda_{ij}^m$  can be numerically evaluated to any desired accuracy without appealing to the closed-form integral formulas. Next we prove a simple result that partially explains why the linear case  $m = 1$  is substantially easier than the nonlinear cases  $m > 1$ .

**Proposition 2.3:**  $\eta_m = 1$  if and only if  $m = 1$ . That is,  $\mathcal{F}_1^2$  is isotropic with respect to the  $L^1$ -norm in  $\mathbf{H}_2$ .

*Proof:* The ‘if’ part follows readily from the fact that

$$\int_{\mathbf{S}^2} xy \, d\mathbf{x} = \int_{\mathbf{S}^2} xz \, d\mathbf{x} = \int_{\mathbf{S}^2} yz \, d\mathbf{x} = 0,$$

and

$$\int_{\mathbf{S}^2} x^2 \, d\mathbf{x} = \int_{\mathbf{S}^2} y^2 \, d\mathbf{x} = \int_{\mathbf{S}^2} z^2 \, d\mathbf{x} = \frac{4\pi}{3}.$$

The matrix  $\Lambda^1$  is therefore diagonal with constant diagonal element  $\frac{4\pi}{3}$ , and  $\mathcal{F}_1^2$  is isotropic with respect to the  $L^1$ -norm in  $\mathbf{H}_2$ .

Conversely, for  $m > 1$ , let  $p = x^m$ ,  $q = x^{m-1}y$ . We show that  $\|\mathcal{F}_2(p)\|_1 \neq \|\mathcal{F}_2(q)\|_1$ :

$$\begin{aligned} \|\mathcal{F}_2(p)\|_1 &= \int_{\mathbf{S}^2} x^{2m} \, d\mathbf{x} = \int_{\psi=0}^{\pi} \int_{\theta=0}^{2\pi} \sin^{2m} \psi \cos^{2m} \theta \sin \psi \, d\theta d\psi, \\ \|\mathcal{F}_2(q)\|_1 &= \int_{\mathbf{S}^2} x^{2m-2} y^2 \, d\mathbf{x} = \int_{\psi=0}^{\pi} \int_{\theta=0}^{2\pi} \sin^{2m} \psi \cos^{2m-2} \theta \sin^2 \theta \sin \psi \, d\theta d\psi. \end{aligned}$$

Let  $c = \int_{\psi=0}^{\pi} \sin^{2m+1} \psi d\psi$ , we have

$$\begin{aligned} \|\mathcal{F}_2(p)\|_1 &= c \int_{\theta=0}^{2\pi} \cos^{2\pi} \theta d\theta, \\ \|\mathcal{F}_2(q)\|_1 &= c \int_{\theta=0}^{2\pi} \cos^{2m-2} \theta \sin^2 \theta d\theta. \end{aligned}$$

Therefore,

$$\begin{aligned} \|\mathcal{F}_2(p)\|_1 - \|\mathcal{F}_2(q)\|_1 &= c \int_{\theta=0}^{2\pi} \cos^{2m-2} \theta (\cos^2 \theta - \sin^2 \theta) d\theta \\ &= c \int_{\theta=0}^{2\pi} \cos^{2m-2} \theta (2\cos^2 \theta - 1) d\theta. \end{aligned}$$

Since  $\int_{\pi=0}^{2\pi} \cos^n \theta d\theta = \frac{n-1}{n} \int_{\pi=0}^{2\pi} \cos^{n-2} \theta d\theta$  for any  $n \geq 2$ , we have

$$\|\mathcal{F}_2(p)\|_1 - \|\mathcal{F}_2(q)\|_1 = c \left( \frac{2m-1}{m} - 1 \right) \int_{\theta=0}^{2\pi} \cos^{2m-2} \theta d\theta,$$

which shows that  $\|\mathcal{F}_2(p)\|_1 - \|\mathcal{F}_2(q)\|_1 \geq 0$  if  $m > 1$ . This implies  $\eta_m > 1$  if  $m > 1$ .

Using the square map  $\mathcal{F}_m^2$ , we will define the approximating subcone  $\sum_{2m}^{\mathcal{C}}$  by specifying its generators as polynomials in  $\mathbf{HS}_m$ . More specifically, let  $\mathcal{C} = \{p_1, \dots, p_k\}$  denote a finite set of  $k$  polynomials (points) in  $\mathbf{HS}_m$ . Its associated cone  $\sum_{2m}^{\mathcal{C}}$  in  $\mathbf{H}_{2m}$  is generated by the finite set of generators  $\mathcal{F}_m^2(\mathcal{C}) = \{p_1^2, \dots, p_k^2\}$ : elements in  $\sum_{2m}^{\mathcal{C}}$  are non-negative linear combinations of  $\mathcal{F}_m^2(p_i)$ :

$$p = a_1 p_1^2 + \dots + a_k p_k^2,$$

for some  $a_1, \dots, a_k \geq 0$ . It is immediately clear that  $\sum_{2m}^{\mathcal{C}} \subset \Omega_{2m} \subset \mathcal{P}_0^{2m}$  for any finite subset  $\mathcal{C} \subset \mathbf{HS}_m$ . Since  $\mathcal{F}_m^2(p) = \mathcal{F}_m^2(-p)$ , we can restrict points in  $\mathcal{C}$  to lie in one chosen hemisphere of  $\mathbf{HS}_m$ . For such  $\mathcal{C}$ , its completion  $\bar{\mathcal{C}} \subset \bar{\mathcal{C}}$  is obtained by joining all antipodal points of points in  $\mathcal{C}$ ,

$$\bar{\mathcal{C}} = \{p_1, \dots, p_k, -p_1, \dots, -p_k\}.$$

**Examples**—For  $m = 1$ ,  $\mathbf{H}_1$  is  $\mathbb{R}^3$  and  $\mathbf{HS}_1$  is  $\mathbb{S}^2$ . If  $\mathcal{C}$  consists of four points  $\{[1, 0, 0]^T, [0, 1, 0]^T, [0, 0, 1]^T, [\sqrt{1/3}, \sqrt{1/3}, \sqrt{1/3}]^T\}$ , the four polynomials  $p_1, p_2, p_3, p_4$  are  $x, y, z$  and  $\sqrt{1/3}x + \sqrt{1/3}y + \sqrt{1/3}z$ , respectively. Elements in  $\sum_2^{\mathcal{C}}$  are non-negative combinations of the four polynomials  $p_1^2, p_2^2, p_3^2, p_4^2$ . More precisely, any  $p \in \sum_2^{\mathcal{C}}$  is determined (in this case, uniquely) by four non-negative numbers  $a_1, a_2, a_3, a_4 \geq 0$  such that



$$p(x, y, z) = (a_1 + \frac{a_4}{3})x^2 + (a_2 + \frac{a_4}{3})y^2 + (a_3 + \frac{a_4}{3})z^2 + \frac{2a_4}{3}(xy + xz + yz).$$

For  $m = 2$ ,  $\mathbf{HS}_2$  can be identified with  $\mathbb{R}^6$  using the monomial basis  $\{x^2, y^2, z^2, xy, xz, yz\}$  and  $\mathbf{HS}_2$  is  $\mathbf{S}^5$ . If  $\mathcal{C}$  consists of three points

$$[\lambda, \lambda, 0, 0, 0, 0]^\top, [\lambda, 0, 0, -\lambda, 0, 0]^\top, [0, -\lambda, 0, 0, 0, \lambda]^\top,$$

where  $\lambda = \sqrt{1/2}$ , the three polynomials  $p_1, p_2, p_3$  are  $\lambda(x^2 + y^2), \lambda(x^2 - xy), \lambda(yz - y^2)$ . Any  $p \in \sum_4^{\mathcal{C}}$  can be written (again uniquely) as

$$p(x, y, z) = \frac{(a_1 + a_2)}{2}x^4 + \frac{(a_1 + a_3)}{2}y^4 + (a_1 + \frac{a_2}{2})x^2y^2 - a_2x^3y - a_3y^3z + \frac{a_3}{2}y^2z^2$$

for three non-negative  $a_1, a_2, a_3$ .

The inclusion  $\sum_{2m}^{\mathcal{C}} \subset \Omega_{2m}$  gives an approximation of  $\Omega_{2m}$  by  $\sum_{2m}^{\mathcal{C}}$  and it involves two main components: the square map  $\mathcal{F}_m^2$  and the chosen polynomials in  $\mathcal{C}$  that provide the generators in  $\sum_{2m}^{\mathcal{C}}$  through  $\mathcal{F}_m^2$ . The main result of our analysis on the approximation error of  $\sum_{2m}^{\mathcal{C}} \subset \Omega_{2m}$  is given in the next theorem, which asserts that the approximation error can be bounded by a product of contributions from both components: the **condition number**  $\eta_m$  of  $\mathcal{F}_m^2$  and the **condition number**  $\theta(\mathcal{C})$  of the set  $\mathcal{C}$  whose definition we now turn to.

**Condition Number  $\theta(\mathcal{C})$  of  $\mathcal{C}$** —We use  $\theta(\mathcal{C})$  as the measure that quantifies the approximation of any  $q \in \mathbf{HS}_m$  considered as a point on the sphere, by the finite set  $\mathcal{C}$ . We will use the spherical distance  $\mathbf{d}_{\mathbf{HS}_m}(p, q)$  (arc-length in radians) to measure the distance between a pair of points  $p, q$  on the sphere  $\mathbf{HS}_m$ , and in particular,  $\mathbf{d}_{\mathbf{HS}_m}(p, q)$  is the angle between the two unit vectors  $p, q$  in  $\mathbf{HS}_m$ . A set  $\mathcal{C}$  is said to be good if there is a triangulation of  $\mathbf{HS}_m$  as a simplicial complex  $\mathcal{T}$  whose vertex set  $\mathcal{T}^0$  is the completion  $\bar{\mathcal{C}}$  of  $\mathcal{C}$ . Since  $\mathbf{HS}_m$  has dimension  $d(m) - 1$ , the top-dimensional simplices in  $\mathcal{T}$  have dimension  $d(m) - 1$  as well. Therefore, for any  $q \in \mathbf{HS}_m$ , there is a  $d(m) - 1$ -simplex  $\sigma \in \mathcal{T}^{d(m)-1}$  containing  $q$ . In particular, we will assume that  $q$  can be written as a non-negative linear combination of the vertices of  $\sigma$ :  $q = a_0p_0 + \dots + a_{d(m)-1}p_{d(m)-1}$  with  $a_0, \dots, a_{d(m)-1} \geq 0$ . While this is in general not true for an arbitrary triangulation  $\mathcal{T}$  of  $\mathbf{HS}_m$ , it is not difficult to show that  $\mathcal{T}$  can be modified (without changing its underlying abstract simplicial complex) to satisfy this property, e.g., by first defining a triangulation of the vertices in  $\mathcal{T}$  considered as points in the Euclidean space  $\mathbb{R}^{d(m)}$  using the same abstract simplicial complex as  $\mathcal{T}$  and radially projecting the simplices onto  $\mathbf{HS}_m$ . For  $0 \leq k \leq d(m) - 1$ ,  $\mathcal{T}^k$  will denote the set of  $k$ -simplices in  $\mathcal{T}$ , and for a  $k$ -simplex  $\sigma \in \mathcal{T}^k$ , its width  $\delta(\sigma)$  is defined as the maximal distance between its vertices,  $p_0, \dots, p_k$ .

$$\delta(\sigma) = \max_{0 \leq i, j \leq k} \mathbf{d}_{\mathbf{HS}_m}(p_i, p_j).$$

For a triangulation  $\mathcal{T}$ , we define its width to be the maximal width of its top-dimensional simplices:

$$\sigma(\mathcal{T}) = \max_{\sigma \in \mathcal{T}^{d(m)-1}} \delta(\sigma).$$

The condition number of  $\mathcal{C}$  is then defined as the minimal width of the triangulations  $\mathcal{T}$  that have  $\bar{c}$  as its vertex set:

$$\sigma(\mathcal{C}) = \min_{\mathcal{T}, \mathcal{T}^0 = \bar{c}} \delta(\mathcal{T}).$$

Since  $\mathcal{C}$  is finite, there exists a triangulation  $\Delta(\mathcal{C})$  whose width gives the condition number  $\theta(\mathcal{C})$ . We note that  $0 < \theta(\mathcal{C}) < \pi$ , and for a good set  $\mathcal{C}$ , the following conditions hold,

1. For each  $q \in \mathbf{HS}_m$ , there are  $d(m)$  elements,  $p_0, \dots, p_{d(m)-1}$ , in  $\bar{c}$  such that  $q = a_0 p_0 + \dots + a_{d(m)-1} p_{d(m)-1}$  for  $a_0, \dots, a_{d(m)-1} \geq 0$  and  $\mathbf{d}_{\mathbf{HS}_m}(p_i, p_j) < \theta(\mathcal{C})$  for any  $0 \leq i, j < d(m)$ .
2. For each  $q \in \mathbf{HS}_m$ , there exists  $p \in \bar{c}$  such that  $\mathbf{d}_{\mathbf{HS}_m}(q, p) < \theta(\mathcal{C})$ .

Property (1) follows immediately from the definition. Property (2) can be shown to follow from the requirement that if  $q \in \sigma \in \mathcal{T}^{d(m)-1}$ ,  $q$  is a non-negative linear combination of vertices in  $\sigma$ .

**Theorem 2.4:** Let  $\mathcal{C}$  denote a good finite subset in  $\mathbf{HS}_m$  and  $\sum_{2m}^{\mathcal{C}}$  its associated finitely-generated subcone in  $\mathbf{H}_{2m}$ . Let  $\theta = \theta(\mathcal{C})$  denote the condition number of  $\mathcal{C}$  as defined above and  $\eta_m$  the condition number of  $\mathcal{F}_m^2$ . Then, for any polynomial  $r \in \Omega_{2m}$  its  $L^1$ -relative approximation error  $\mathbf{E}_{\sum_{2m}^{\mathcal{C}}}(r)$  satisfies

$$\mathbf{E}_{\sum_{2m}^{\mathcal{C}}}(r) \leq 4 \tan \theta \sin^2 \frac{\theta}{2} \eta_m^2.$$

The bound above constitutes our quantitative characterization of the approximation  $\sum_{2m}^{\mathcal{C}} \subset \Omega_{2m}$ . Not surprisingly, the bound provided above depends on both the map  $\mathcal{F}_m^2$  as well as the set  $\mathcal{C}$  through  $\theta$  and  $\eta$ . The error measured by  $\mathbf{E}_{\sum_{2m}^{\mathcal{C}}}$  takes place in  $\mathbf{H}_{2m}$ , and the bound on the right factored into two components with contribution from  $\theta$  that essentially measures how well an arbitrary point  $q \in \mathbf{HS}_m$  can be approximated using  $\mathcal{C}$  and its associated triangulation  $\Delta(\mathcal{C})$ . In particular, as will be seen from the proof,  $\tan \theta$  arises from approximating  $q$  using its nearest neighbor in  $\bar{c}$  as in Property (2) above while  $\sin^2 \frac{\theta}{2}$  comes from approximating  $q$  using the simplex  $\sigma$  containing it as in Property (1).

We will prove the theorem through a sequence of lemmas given below. However, before delving into the proof, we remark that although using the triangulation  $\Delta(\mathcal{C})$  to define  $\theta(\mathcal{C})$  may seem unnecessary at first, it is in fact crucial to have Property (1) in order to produce a smaller bound on the error. For example, it is possible to define  $\theta(\mathcal{C})$  using only Property (2), i.e., each  $q \in \mathbf{HS}_m$  can be approximated by a  $p \in \bar{c}$  such that  $\mathbf{d}_{\mathbf{HS}_m}(q, p) < \theta(\mathcal{C})$ . However, this hypothesis itself is only strong enough to produce the bound given in Lemma 2.6 (Equation 2.4). Disregarding  $\eta_m$ , the bound given in Equation 2.4 is  $2 \sin \theta$ , which is

considerably inferior to the bound of  $4 \tan \theta \sin^2$  given in Theorem 2.4. In particular, for small  $\theta$ , the former is approximately  $2\theta$  while the latter is  $\theta^3$  (See Equation 3.1), two order of magnitude less. As will be clear in the proof, the main issue is to approximate the polynomial  $q^2$  for any  $q \in \mathbf{HS}_m$  with a sum of squares of polynomials in  $\mathcal{C}$ . Using only Property (2), it is difficult to determine what polynomials in  $\mathbf{HS}_m$  can be used to approximate  $q^2$  other than the polynomial  $p \in \mathcal{C}$  that is closest to  $q$ . With Property (1), we have more choices at our disposal as we can approximate  $q^2$  using the vertices  $p_i$  of the simplex  $\sigma$  that contains  $q$ , and more importantly, the remainder of this approximation (sum of  $(p_i - p_j)^2$ ) can be further approximated using polynomials in  $\mathcal{C}$ . This is the content of Lemma 2.8. In particular, when approximating  $q^2$ , Property (1) allows the access of not only the polynomials  $p_i \in \mathcal{C}$  that are neighbors of  $q$  but also polynomials in  $\mathcal{C}$  that are usually far away from  $q$ . See Figure 2.1. Furthermore, as will be detailed in Section 3, Property (1) allows us to formulate a simple method for estimating the minimal number of points (polynomials) in  $\mathcal{C}$  needed for a given precision requirement.

**Lemma 2.5:** Let  $p, q$  be two polynomials in  $\mathbf{HS}_m$  and  $\theta = \mathbf{d}_{\mathbf{HS}_m}(p, q)$  denote their geodesic distance considered as points on the sphere  $\mathbf{HS}_m$ . We have

$$\int_{\mathbb{S}^2} |p(\mathbf{x}) - q(\mathbf{x})|^2 d\mathbf{x} \leq 4 \sin^2 \frac{\theta}{2} \eta_m^{\max}.$$

**Proof:** Let  $r = p - q$ . As a vector in  $\mathbf{H}_m$ ,  $|r| = |p - q|$ . Using the law of cosines,

$$\gamma = |r| = |p - q| = \sqrt{2 - 2 \cos \theta} = 2 \sin \frac{\theta}{2}. \quad (2.2)$$

Therefore,  $r/\gamma \in \mathbf{HS}_m$  and we have

$$\int_{\mathbb{S}^2} r^2(\mathbf{x}) d\mathbf{x} = \gamma^2 \int_{\mathbb{S}^2} (r(\mathbf{x})/\gamma)^2 d\mathbf{x} \leq \gamma^2 \eta_m^{\max},$$

and the result follows.

Next we prove an important lemma which shows that for two nearby  $p, q$  in  $\mathbf{HS}_m$  we can approximate  $q^2$  using  $p^2$  such that the  $L^1$ -approximation error is a fraction (depending on the geodesic distance) of the  $L^1$ -norm of  $p^2$ .

**Lemma 2.6:** Let  $p, q$  be two polynomials in  $\mathbf{HS}_m$  and  $\theta = \mathbf{d}_{\mathbf{HS}_m}(p, q)$  denote their geodesic distance considered as points on the sphere  $\mathbf{HS}_m$ . Let  $\|\mathcal{F}_2(p) - \mathcal{F}_2(q)\|_1$  denote the  $L^1$ -difference between  $\mathcal{F}_2(p), \mathcal{F}_2(q)$

$$\|\mathcal{F}_2(p) - \mathcal{F}_2(q)\|_1 = \int_{\mathbb{S}^2} |p^2(\mathbf{x}) - q^2(\mathbf{x})| d\mathbf{x}. \quad (2.3)$$

We have

$$\|\mathcal{F}_2(p) - \mathcal{F}_2(q)\|_1 \leq 2 \sin \theta \eta_m \|\mathcal{F}_2(p)\|_1. \quad (2.4)$$

**Proof:** Using Hölder's inequality, we have

$$\begin{aligned} \int_{\mathbb{S}^2} |p^2(\mathbf{x}) - q^2(\mathbf{x})| d\mathbf{x} &= \int_{\mathbb{S}^2} |p(\mathbf{x}) - q(\mathbf{x})| |p(\mathbf{x}) + q(\mathbf{x})| d\mathbf{x} \\ &\leq \left( \int_{\mathbb{S}^2} |p(\mathbf{x}) - q(\mathbf{x})|^2 d\mathbf{x} \right)^{\frac{1}{2}} \left( \int_{\mathbb{S}^2} |p(\mathbf{x}) + q(\mathbf{x})|^2 d\mathbf{x} \right)^{\frac{1}{2}}. \end{aligned}$$

The proof will proceed to bound the two terms on the right. By the preceding lemma, we have

$$\left( \int_{\mathbb{S}^2} |p(\mathbf{x}) - q(\mathbf{x})|^2 d\mathbf{x} \right)^{\frac{1}{2}} \leq 2 \sin \frac{\theta}{2} \sqrt{\eta_m^{\max}}.$$

For the second term, we will consider the polynomial  $r = \gamma \frac{p+q}{2}$ , where  $\gamma > 1$  ensures that  $r \in \mathbf{HS}_m$ . A quick calculation shows that  $\gamma = 1/\cos \frac{\theta}{2}$ . Since, by definition,

$$\int_{\mathbb{S}^2} r^2(\mathbf{x}) d\mathbf{x} \leq \eta_m^{\max},$$

we have

$$\left( \int_{\mathbb{S}^2} |p(\mathbf{x}) + q(\mathbf{x})|^2 d\mathbf{x} \right)^{\frac{1}{2}} = \left( \int_{\mathbb{S}^2} \frac{4}{\gamma^2} r^2(\mathbf{x}) d\mathbf{x} \right)^{\frac{1}{2}} \leq \frac{2}{\gamma} \sqrt{\eta_m^{\max}}.$$

Combining the two inequalities, we have

$$\int_{\mathbb{S}^2} |p^2(\mathbf{x}) - q^2(\mathbf{x})| d\mathbf{x} \leq 4 \cos \frac{\theta}{2} \sin \frac{\theta}{2} \eta_m^{\max} = 2 \sin \theta \eta_m^{\max}.$$

Since

$$\|\mathcal{F}_2(p)\|_1 = \int_{\mathbb{S}^2} p^2(\mathbf{x}) d\mathbf{x} \leq \eta_m^{\min},$$

it follows that

$$\|\mathcal{F}_2(p) - \mathcal{F}_2(q)\|_1 \leq 2 \sin \theta \frac{\eta_m^{\max}}{\eta_m^{\min}} \|\mathcal{F}_2(p)\|_1.$$

This completes the proof.

We will use the preceding lemma to prove two basic error estimates. For any two points  $p, q$  in  $\mathcal{C}$ , the following lemma provides a bound on the approximation error for points that lie on the arc (geodesic path) joining  $p, q$ .

**Lemma 2.7:** Let  $p, q$  be two neighboring points in  $\Delta(\mathcal{C})$ , i.e., there is a 1-simplex  $\sigma^1$  in  $\Delta(\mathcal{C})$  with  $p, q$  as its two vertices. Let  $r = ap + bq$  be a convex combination of  $p, q$  with  $a, b \geq 0$  and  $a + b = 1$ . If  $\theta = \theta(\mathcal{C})$  denotes the condition number of  $\mathcal{C}$ , then

$$\mathbf{E}_{\Sigma_{2m}^{\mathcal{C}}}(r^2) \leq 2\sin\theta \tan^2 \frac{\theta}{2} \eta_m^2.$$

**Proof:** By definition of  $\theta$ ,  $\mathbf{d}_{\mathbf{HS}_m}(p, q) \leq \theta$ . Let  $\phi = a^2 p^2 + b^2 q^2 + abp^2 + abq^2$  be an element in  $\Sigma_{2m}^{\mathcal{C}}$ . We have

$$\begin{aligned} \phi - r^2 &= (a^2 p^2 + b^2 q^2 + abp^2 + abq^2) - (ap + bq)^2 \\ &= ab(p - q)^2. \end{aligned}$$

Let  $\gamma = |p - q|$  and  $t = (p - q)/\gamma \in \mathbf{HS}_m$ . There exists  $s \in \bar{\mathcal{C}}$  such that the geodesic distance between  $t$  and  $s$  is less than  $\theta$ . By the preceding lemma,

$$\int_{s^2} |t^2(\mathbf{x}) - s^2(\mathbf{x})| d\mathbf{x} \leq 2\sin\theta \eta_m \|t^2(\mathbf{x})\|_1.$$

Now let  $\varphi = \phi + ab\gamma^2 s^2$  be another element in  $\Sigma_{2m}^{\mathcal{C}}$ . We have

$$\int_{s^2} |r^2(\mathbf{x}) - \varphi(\mathbf{x})| d\mathbf{x} = ab \int_{s^2} |(\gamma t)^2 - (\gamma s)^2| d\mathbf{x} \leq 2ab\sin\theta \eta_m \|(\gamma t)^2(\mathbf{x})\|_1.$$

By Lemma 2.5,  $\|(\gamma t)^2(\mathbf{x})\|_1 \leq 4\sin^2 \frac{\theta}{2} \eta_m^{\max}$ . This gives

$$\int_{s^2} |r^2(\mathbf{x}) - \varphi(\mathbf{x})| d\mathbf{x} \leq 2\sin\theta \sin^2 \frac{\theta}{2} \eta_m \eta_m^{\max}. \quad (2.5)$$

as  $ab \leq \frac{1}{4}$  for  $a, b \geq 0$  and  $a + b = 1$ . We next bound the  $L^1$ -norm of  $r^2(\mathbf{x})$ . Since  $r = ap + bq$ , there exists  $1 \leq \gamma \leq 1/\cos \frac{\theta}{2}$  such that  $\gamma r \in \mathbf{HS}_m$ . This implies that

$$\int_{w_S} \gamma^2 r^2(\mathbf{x}) d\mathbf{x} \geq \eta_m^{\min},$$

or

$$\int_{s^2} r^2(\mathbf{x})d\mathbf{x} \geq \cos^2 \frac{\theta}{2} \eta_m^{\min}. \quad (2.6)$$

Combining Equations 2.5 and 2.6 gives the desired result.

The preceding lemma can be generalized immediately to higher-order convex combinations.

**Lemma 2.8:** Let  $p_1, \dots, p_k$  denote the vertices of a  $k - 1$ -simplex  $\sigma^{k-1}$  in  $\Delta(\mathcal{C})$  as well as the corresponding homogeneous polynomials in  $\mathbf{HS}_{mr}$ . Let  $r = a_1 p_1 + \dots + a_k p_k$  be a convex combination of  $p_1, \dots, p_k$  with  $a_1, \dots, a_k \geq 0$  and  $a_1 + \dots + a_k = 1$ . If  $\theta$  denote the condition number of  $\mathcal{C}$ , Then

$$\mathbf{E}_{\sum_{\mathcal{C}}^{2m}}(r^2) \leq 4 \tan \theta \sin^2 \frac{\theta}{2} \eta_m^2.$$

**Proof:** Expanding  $r^2$ , we have

$$r^2 = \sum_{i=1}^k a_i^2 p_i^2 + 2 \sum_{i < j} a_i a_j p_i p_j.$$

The second sum contains  $C_2^k = \frac{k(k-1)}{2}$  terms. To approximate  $r^2$  using an element  $\phi \in \sum_{\mathcal{C}}^{2m}$ , we proceed similarly as before. We start with  $\phi$  equals the first sum above. For each cross-term  $2a_i a_j p_i p_j$  in the second sum, we add  $a_i a_j (p_i^2 + p_j^2)$  to  $\phi$ . This gives

$$\phi = \sum_{i=1}^k a_i^2 p_i^2 + \sum_{i < j} a_i a_j (p_i^2 + p_j^2).$$

It follows that

$$\phi - r^2 = \sum_{i < j} a_i a_j (p_i + p_j)^2.$$

Next, we will approximate the squares  $(p_i - p_j)^2$  using elements in  $\sum_{\mathcal{C}}^{2m}$  exactly as before. More specifically, let  $\gamma_{ij} = |p_i - p_j|$  and  $t_{ij} = (p_i - p_j) / \gamma_{ij}$ . There exists  $s_{ij} \in \mathcal{C}$  such that the geodesic distance between  $t_{ij}$  and  $s_{ij}$  is less than  $\theta$ . Now let  $\varphi = \phi + \sum_{i < j} a_i a_j (\gamma_{ij} s_{ij})^2$  be an element in  $\sum_{\mathcal{C}}^{2m}$ . We have

$$\int_{s^2} |r^2(\mathbf{x}) - \varphi(\mathbf{x})| d\mathbf{x} \leq \sum_{i < j} a_i a_j \int_{s^2} |(\gamma_{ij} t_{ij})^2(\mathbf{x}) - (\gamma_{ij} s_{ij})^2(\mathbf{x})| d\mathbf{x}.$$

It follows from Equations 2.2 and 2.4 that all the integrals on the right can be uniformly bounded

$$\int_{S^2} |(\gamma_{ij} t_{ij})^2(\mathbf{x}) - (\gamma_{ij} s_{ij})^2(\mathbf{x})| d\mathbf{x} \leq 8 \sin \theta \sin^2 \frac{\theta}{2} \eta_m^{\max} \eta_m,$$

and this gives

$$\int_{S^2} |r^2(\mathbf{x}) - \varphi(\mathbf{x})| d\mathbf{x} \leq 8 \sin \theta \sin^2 \frac{\theta}{2} \eta_m^{\max} \eta_m \sum_{i < j} a_i a_j.$$

Since  $a_1 + \dots + a_k = 1$ ,

$$\begin{aligned} \sum_{i < j} a_i a_j &= \frac{(a_1 + \dots + a_k)^2 - (a_1^2 + \dots + a_k^2)}{2} = \frac{1 - (a_1^2 + \dots + a_k^2)}{2} \\ &\leq \frac{1 - \frac{1}{k}}{2} = \frac{k-1}{2k} < \frac{1}{2} \end{aligned} \quad (2.7)$$

as  $a_1^2 + \dots + a_k^2 \geq \frac{1}{k}$  by Cauchy-Schwarz inequality. This yields the bound

$$\int_{S^2} |r^2(\mathbf{x}) - \varphi(\mathbf{x})| d\mathbf{x} \leq 4 \sin \theta \sin^2 \frac{\theta}{2} \eta_m^{\max} \eta_m. \quad (2.8)$$

We next bound the  $L^1$ -norm of  $r^2$ . Given that  $r = a_1 p_1 + \dots + a_k p_k$ , the following lemma shows that the  $L^2$ -magnitude  $|r|$  of the vector  $r$  satisfies

$$|r| \geq \sqrt{\cos \theta}.$$

Hence, there exists  $1 \leq \gamma \leq \frac{1}{\sqrt{\cos \theta}}$  such that  $\gamma r \in \mathbf{HS}_{m^r}$ . Exactly as before, we have

$$\int_{S^2} r^2(\mathbf{x}) d\mathbf{x} \geq \frac{1}{\gamma^2} \eta_m^{\min} \geq \cos \theta \eta_m^{\min}. \quad (2.9)$$

Equations 2.8 and 2.9 together complete the proof.

**Lemma 2.9:** Let  $\Delta$  denote a  $k$ -simplex in  $\mathbb{R}^{d(m)}$  whose vertices  $p_0, \dots, p_k$  are on the unit sphere, i.e.,  $\|p_0\|_2 = \dots = \|p_k\|_2 = 1$ . If there exists some  $\alpha$  such that  $1 > \alpha > 0$  and  $p_i^\top p_j \geq \alpha$  for all  $i < j$ , then for any  $\mathbf{x} \in \Delta$ ,

$$\|\mathbf{x}\|_2 > \sqrt{\alpha}.$$

**Proof:** Let  $\mathbf{x} = a_1 p_1 + \dots + a_k p_k$  with  $a_i \geq 0$  and  $a_1 + \dots + a_k = 1$ . It follows that

$$\mathbf{x}^\top \mathbf{x} \geq \sum_{i=0}^k a_i^2 + 2\alpha \sum_{i<j} a_i a_j.$$

Let  $s = 2\sum_{i<j} a_i a_j$ ; and the above inequality becomes  $\mathbf{x}^\top \mathbf{x} \geq 1 - (1 - \alpha)s$ . From Equation 2.7, we have  $0 \leq s \leq \frac{k-1}{k} < 1$ . It follows that

$$\mathbf{x}^\top \mathbf{x} > 1 - (1 - \alpha)\alpha.$$

We remark that when  $k = 2, \frac{k-1}{k} = \frac{1}{2}$  and the bound becomes tighter  $\mathbf{x}^\top \mathbf{x} \geq \frac{1}{2} + \frac{\alpha}{2}$ . This gives the  $\cos \theta$  term in Equation 2.9. Finally, we are ready to complete the proof of Theorem 2.4:

**Proof:** Since  $r(\mathbf{x})$  can be written as a sum of squares, by Proposition 2.10, it can be written as a sum of no more than  $d(m)$  terms with  $p_i \in \mathbf{HS}_m$ :

$$r(\mathbf{x}) = \sum_{i=1}^{d(m)} a_i p_i^2(\mathbf{x}).$$

Each  $p_i$  belongs to a  $(d(m) - 1)$ -dimensional simplex  $\sigma_i \in \Delta(\mathcal{C})$ . By the preceding lemma, each  $p_i^2$  can be approximated by an element  $\tilde{p}_i$  in  $\sum_{\mathcal{C}}^{2m}$  with uniformly bounded relative  $L^1$ -error

$$\|p_i^2(\mathbf{x}) - \tilde{p}_i(\mathbf{x})\|_1 \leq C \|p_i^2(\mathbf{x})\|_1,$$

where  $C = 4 \tan \theta \sin^2 \frac{\theta}{2} \eta_m^2$ . Define  $\tilde{r} \in \sum_{\mathcal{C}}^{2m}$  as

$$\tilde{r} = \sum_{i=1}^{d(m)} a_i \tilde{p}_i,$$

and we have

$$\|r(\mathbf{x}) - \tilde{r}(\mathbf{x})\|_1 \leq \sum_{i=1}^{d(m)} a_i \|p_i^2(\mathbf{x}) - \tilde{p}_i(\mathbf{x})\|_1 \leq C \sum_{i=1}^{d(m)} a_i \|p_i^2(\mathbf{x})\|_1.$$

On the other hand, we also have

$$\|r(\mathbf{x})\|_1 \leq \sum_{i=1}^{d(m)} a_i \int_{S^2} p_i^2(\mathbf{x}) d\mathbf{x} = \sum_{i=1}^{d(m)} a_i \|p_i^2(\mathbf{x})\|_1$$



Combining both inequalities yields the desired result.

In the proof above we made use of the following proposition.

**Proposition 2.10:** Let  $r$  denote a homogeneous polynomial of degree  $2m$  that can be written as a sum of squares of homogeneous polynomials of degree  $m$ . Then,  $r$  can be written as a sum of at most  $d(m)$  squares

$$r(\mathbf{x}) = \sum_{i=1}^{d(m)} a_i p_i^2(\mathbf{x}),$$

where  $a_1, \dots, a_{d(m)} \geq 0$  and  $p_1, \dots, p_{d(m)} \in \mathbf{HS}_m$ .

**Proof:** Suppose  $r$  is a sum of  $k$  squares of homogeneous polynomials  $\tilde{q}_1, \dots, \tilde{q}_k$  of degree  $m$

$$r(\mathbf{x}) = \tilde{q}_1^2(\mathbf{x}) + \dots + \tilde{q}_k^2(\mathbf{x}).$$

Denote  $\mathbf{m}_1, \dots, \mathbf{m}_{d(m)}$  the  $d(m)$  monomials of degree  $m$ , and  $\mathbf{X}$  the vector

$$\mathbf{X} = [\mathbf{m}_1(\mathbf{x}), \mathbf{m}_2(\mathbf{x}), \dots, \mathbf{m}_{d(m)}(\mathbf{x})]^\top$$

whose components are the monomials. It follows that  $\tilde{q}_i(\mathbf{x}) = \mathbf{a}_i^\top \mathbf{X}$  with  $\mathbf{a}_i$  the vector whose components are coefficients of  $\tilde{q}_i(\mathbf{x})$ , and

$$r(\mathbf{x}) = \mathbf{X}^\top (\mathbf{a}_1 \mathbf{a}_1^\top + \dots + \mathbf{a}_k \mathbf{a}_k^\top) \mathbf{X} = \mathbf{X}^\top \mathbf{S} \mathbf{X}.$$

The matrix  $\mathbf{S}$  is symmetric and positive semi-definite with non-negative eigenvalues. Let  $\lambda_1, \dots, \lambda_{d(m)}$  denote its complete set of eigenvalues and  $\mathbf{v}_1, \dots, \mathbf{v}_{d(m)}$  their associated unit eigenvectors,  $\|\mathbf{v}_j\|_2 = 1$ . It follows that

$$\mathbf{S} = \lambda_1 \mathbf{v}_1 \mathbf{v}_1^\top + \dots + \lambda_{d(m)} \mathbf{v}_{d(m)} \mathbf{v}_{d(m)}^\top,$$

and

$$\begin{aligned} r(\mathbf{x}) &= \lambda_1 \mathbf{X}^\top \mathbf{v}_1 \mathbf{v}_1^\top \mathbf{X} + \dots + \lambda_{d(m)} \mathbf{X}^\top \mathbf{v}_{d(m)} \mathbf{v}_{d(m)}^\top \mathbf{X} \\ &= \lambda_1 q_1^2(\mathbf{x}) + \dots + \lambda_{d(m)} q_{d(m)}^2(\mathbf{x}), \end{aligned}$$

where  $q_i(\mathbf{x}) = \mathbf{v}_i^\top \mathbf{X} \in \mathbf{HS}_m$  as  $\|\mathbf{v}_j\|_2 = 1$  for  $i = 1, \dots, d(m)$ .

### 3. Approximating PSD Tensors of Orders two, four and six

In this section, we apply Theorem 2.4 to derive formulas for the minimal number of generators in  $\sum_{2m}^{\mathcal{C}}$  needed to ensure that the approximation  $\sum_{2m}^{\mathcal{C}} \subset \Omega_{2m}$  is within a given accuracy requirement. Specifically, the accuracy requirement is presented in the form of the relative  $L^1$ -approximation error  $\mathbf{E}_{2m}^{\mathcal{C}}$  (cf. Equation 2.1): for  $0 < \varepsilon < 1$ , we derive a formula that gives the (approximated) minimal number  $\mathcal{N}(\varepsilon, m)$  of generators in  $\sum_{2m}^{\mathcal{C}}$  such that any  $r \in \Omega_{2m}$  can be approximated within  $\varepsilon$  using  $\sum_{2m}^{\mathcal{C}}$ , i.e.,

$$\mathbf{E}_{2m}^{\mathcal{C}}(r) < \varepsilon.$$

For PSD ternary tensors of orders two and four, it is known that they can be written as sums of squares of three tensors of order one and two, respectively. This follows from the well-known result that any ternary positive semi-definite homogeneous polynomial  $p(\mathbf{x})$  of degree two and four can be written as a sum of three squares of polynomials of degree one and two, respectively. The quadratic case follows easily from linear algebra while the quartic case follows from the celebrated theorem of Hilbert on ternary quartics [24]. We will first describe a general method for obtaining the formula  $\mathcal{N}(\varepsilon, m)$  for any order  $m$ , and we will then explicitly work out the three cases  $m = 1, 2, 3$  that are of most interest for various applications.

#### 3.1. Preliminaries

Given a required precision  $\varepsilon > 0$ , the bound provided by Theorem 2.4 allows us to determine the condition number  $\theta = \theta(\mathcal{C})$  for the point set  $\mathcal{C}$  in  $\mathbf{HS}_m$  to ensure that the precision requirement is satisfied. The main result in this section is a simple estimate on the number  $\mathcal{N}(\varepsilon, m)$  of points in  $\mathcal{C}$  needed to achieve the desired  $\theta$  on the sphere  $\mathbf{HS}_m$ . Let

$C_\eta(\theta) = 4 \tan \theta \sin^2 \frac{\theta}{2} \eta^2$  denote the bound given in Theorem 2.4. Since

$$\tan \theta \sin^2 \frac{\theta}{2} = \frac{1}{2} (\tan \theta - \sin \theta),$$

$C_\eta(\theta)$  is a monotonically increasing function for  $0 \leq \theta \leq \frac{\pi}{2}$ , and we will denote its inverse by  $f_\eta(\varepsilon) = C_\eta^{-1}(\varepsilon)$ .  $f_\eta$  can be numerically evaluated and the plots for  $f_\eta$  over the range 0.01  $\leq \varepsilon \leq$  0.1 for several different  $\eta$ -values are shown in Figure 3.1. If  $\theta$  is assumed to be small,

$$\tan \theta \sin^2 \frac{\theta}{2} \approx \frac{\theta^3}{4}. \quad (3.1)$$

Therefore,  $4 \tan \theta \sin^2 \frac{\theta}{2} \eta^2 \approx \varepsilon$  implies that

$$\theta \approx \left( \frac{\varepsilon}{\eta^2} \right)^{\frac{1}{3}}. \quad (3.2)$$

The formula above gives an estimate on the condition number  $\theta = \theta(\mathcal{C})$  given  $\varepsilon$  and  $\eta$ . We next give an estimate on the size of  $\mathcal{C}$  for the given  $\theta(\mathcal{C})$ . Let  $n = d(m) - 1$  denote the dimension of the sphere  $\mathbf{HS}_m$  and  $\Delta(\mathcal{C})$  denote the triangulation associated with  $\mathcal{C}$ . A simplex in  $\Delta(\mathcal{C})$  is said to be  $\theta$ -regular if the distance between any pair of its vertices equals  $\theta$ , and the edge joining any pair of vertices is a geodesics on  $\mathbf{HS}_m$ . Due to the curvature on the sphere  $\mathbf{HS}_m$ , it is not possible to cover  $\mathbf{HS}_m$  with only  $\theta$ -regular simplices. Therefore, we assume that the  $n$ -simplices in  $\Delta(\mathcal{C})$  are approximately  $\theta$ -regular in the sense that the geodesic distance between any pair of vertices of a  $n$ -simplex in  $\mathbf{HS}_m$  is approximately  $\theta$  and the edge joining them is approximately a geodesic as well. For each vertex  $v$  in  $\Delta(\mathcal{C})$ , its degree is the number of  $n$ -dimensional simplices having it as a vertex. To estimate the number of points in  $\mathcal{C}$ , we will estimate two quantities: the number  $K$  of  $n$ -dimensional simplices in  $\Delta(\mathcal{C})$  and the *average* degree  $\nu$  of the vertices. The number of points in  $\mathcal{C}$  can then be estimated as

$$\mathcal{N}(\varepsilon, m) = \# \text{of points in } \mathcal{C} \simeq \frac{(n+1)K}{2\nu}.$$

The occurrence of 2 in the denominator accounts for the fact that points in  $\mathcal{C}$  are located only on a hemisphere.

**Estimate on  $K$** —Since  $\mathbf{HS}_m$  is covered by a collection of  $\theta$ -regular  $n$ -simplices,  $K$  can be estimated by taking the ratio between the volume of the sphere  $\mathbf{HS}_m$  and the volume of a  $\theta$ -regular  $n$ -simplex. Since  $\theta$  is in general assumed to be small, we will approximate the volume of a  $\theta$ -regular  $n$ -simplex on the sphere  $\mathbf{HS}_m$  with the volume  $\omega_n(\theta)$  of a corresponding  $\theta$ -regular  $n$ -simplex in the Euclidean space  $\mathbb{R}^n$ :

$$\omega_n(\theta) = \frac{\sqrt{n+1}}{n! \sqrt{2^n}} \theta^n. \quad (3.3)$$

It then follows that the number  $K$  of  $n$ -simplexes can be estimated as

$$K = \frac{\mathbf{V}_n}{\omega_n(\theta)}, \quad (3.4)$$

where the volume of the sphere  $\mathbf{V}_n$  is given by the formula [25]

$$\mathbf{V}_n = \begin{cases} \frac{(2\pi)^{(n+1)/2}}{2 \cdot 4 \cdots (n-1)} & \text{if } n \text{ is odd;} \\ \frac{2(2\pi)^{n/2}}{1 \cdot 3 \cdots (n-1)} & \text{if } n \text{ is even.} \end{cases}$$

**Estimate on  $\nu$** —For a typical vertex  $v$  in  $\Delta(\mathcal{C})$ , a small neighborhood  $\mathcal{U}$  around  $v$  in  $\mathbf{HS}_m$  is covered by the  $\theta$ -regular  $n$ -simplices having  $v$  as one of their vertices. Again, assuming  $\theta$  is small, we can approximate this using Euclidean geometry, by transforming the neighborhood  $\mathcal{U}$  onto the tangent space  $\mathbf{T}_v$  at  $v$  using the log map. The geodesic ball  $\mathbf{B}_\theta$  of radius  $\theta$  on  $\mathbf{HS}_m$  is mapped to the Euclidean ball of radius  $\theta$  and the image of each  $n$ -simplex under the log map can be approximated by a regular  $n$ -simplex in the Euclidean space with side length  $\theta$ . See Figure 3.2. It follows that the degree of  $v$  can be estimated as the ratio between the volume of the unit  $n$ -dimensional ball and the volume of regular  $n$ -

simplex in  $\mathbb{R}^n$  with side length  $\theta$ . The volume  $\mathbf{V}^n$  of an  $n$ -ball in  $\mathbb{R}^n$  with radius  $r=1$  is given by the formula [25]

$$\mathbf{V}^n = \begin{cases} \frac{(2\pi)^{n/2}}{2 \cdot 4 \cdots n} & \text{if } n \text{ is even;} \\ \frac{2(2\pi)^{(n-1)/2}}{1 \cdot 3 \cdots n} & \text{if } n \text{ is odd.} \end{cases}$$

The degree  $\nu$  is then estimated as

$$\nu = \frac{\mathbf{V}^n}{\omega_n(1)}. \quad (3.5)$$

Combining Equations 3.3, 3.4, 3.5, we have

$$\begin{aligned} \mathcal{N}(\varepsilon, m) = \# \text{of points in } \mathcal{C} &\approx \frac{1}{2} \frac{(n+1)\mathbf{V}_n}{\omega_n(\theta) \frac{\mathbf{V}_n}{\omega_n(1)}} = \frac{(n+1)\mathbf{V}_n}{2\mathbf{V}^n \theta^n} \\ &= \frac{(n+1)\mathbf{V}_n}{2\mathbf{V}^n} f_\eta(\varepsilon)^{-n}. \end{aligned} \quad (3.6)$$

In the remaining section, we will work out the implication of the above estimate for  $2^{nd}$ ,  $4^{th}$  and  $6^{th}$ -order tensors.

### 3.2. Second-Order Tensors

A quadratic homogeneous polynomial  $P(x, y, z)$  in  $\mathbb{R}^3$  has six coefficients  $P(x, y, z) = ax^2 + by^2 + cz^2 + dxy + exz + fyz$ . It can be written in a matrix form as,

$$P(x, y, z) = \begin{bmatrix} x & y & z \end{bmatrix} \begin{bmatrix} a & \frac{d}{2} & \frac{e}{2} \\ \frac{d}{2} & b & \frac{f}{2} \\ \frac{e}{2} & \frac{f}{2} & c \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \mathbf{x}^\top \mathbf{S} \mathbf{x}.$$

Positive semi-definiteness of the polynomial  $P(x, y, z)$  is equivalent to the positive semi-definiteness of the matrix  $\mathbf{S}$ . It follows that determining positive semi-definiteness of a homogeneous quadratic polynomial is straightforward by examining eigenvalues of  $\mathbf{S}$ :  $\mathbf{S}$  is positive semi-definite if and only its eigenvalues  $\lambda_1, \lambda_2, \lambda_3$  are all non-negative and  $\mathbf{S}$  can be written as

$$\mathbf{S} = \lambda_1 \mathbf{v}_1^\top \mathbf{v}_1 + \lambda_2 \mathbf{v}_2^\top \mathbf{v}_2 + \lambda_3 \mathbf{v}_3^\top \mathbf{v}_3,$$

where  $\mathbf{v}_i$  is the unit eigenvector with eigenvalue  $\lambda_i$  for  $i = 1, 2, 3$ . It follows that  $P(x, y, z)$  can be written as a sum of three linear polynomials  $p_1(\mathbf{x}), p_2(\mathbf{x}), p_3(\mathbf{x})$ ,

$$P(\mathbf{x}) = p_1(\mathbf{x})^2 + p_2(\mathbf{x})^2 + p_3(\mathbf{x})^2,$$

with  $p_i(\mathbf{x}) = \sqrt{\lambda_i} \mathbf{v}_i^\top \mathbf{x}$ .

With  $m = 1$ , the sphere  $\mathbf{HS}_m$  has dimension  $n = 2$ . According to Proposition 2.3, the map  $\mathcal{F}_1^2$  is isotropic with respect to the  $L^1$ -norm and  $\eta = 1$ . Equation 3.6 (together with Equation 3.2) then gives

$$N(\varepsilon, 1) \approx \frac{3V_2}{2V^2} \left(\frac{1}{\varepsilon}\right)^{\frac{2}{3}} = 6 \left(\frac{1}{\varepsilon}\right)^{\frac{2}{3}}. \quad (3.7)$$

**More Precise Estimate**—For the linear case  $m = 1$ , since  $\mathbf{HS}_m$  is the two-sphere  $\mathbf{S}^2$ , its geometry is well-known and a better estimate on  $N$  can be obtained. Given  $\theta$ ,  $\mathbf{S}^2$  is covered by geodesic triangles whose sides have lengths of approximately  $\theta$ . Approximating the areas of these geodesic triangles with the area of an Euclidean equilateral triangles with side  $\theta$  gives  $\theta^2 \sqrt{3}/4$ . Let  $F, E, V$  denote the number of triangles, edges and vertices in the triangulation  $\Delta(\mathcal{C})$ . According to Euler’s formula

$$F - E + V = \chi(\mathbf{S}^2) = 2$$

where  $\chi(\mathbf{S}^2)$  is the Euler characteristic of  $\mathbf{S}^2$ . Since  $E = 3F/2$ ,  $V = 2 + F/2 \approx F/2$ . This gives  $\nu = 6$  as the average degree of a vertex on  $\mathbf{S}^2$ . Our estimate on the degree  $\nu$  in Equation 3.5 in this case gives  $\nu = 4\pi / \sqrt{3} \approx 7.2$ , which gives a 20% overestimate.

The area  $A$  of a geodesic triangle on  $\mathbf{S}^2$  with three interior angles  $\alpha, \beta, \gamma$  is given as [1]

$$A = \alpha + \beta + \gamma - \pi.$$

In particular, for a geodesic equilateral triangle on  $\mathbf{S}^2$  with side length  $\theta$ , its angle  $\alpha$  is given as

$$\alpha = \cos^{-1} \left( \frac{\cos \theta - \cos^2 \theta}{\sin^2 \theta} \right),$$

and the estimate on the number of triangles is

$$K = \frac{4\pi}{3 \cos^{-1} \left( \frac{\cos \theta - \cos^2 \theta}{\sin^2 \theta} \right) - \pi}.$$

Let  $4 \tan \theta \sin^2 \frac{\theta}{2} = \varepsilon$  and  $\theta = f(\varepsilon)$  be the solution to the trigonometric equation. It then follows that

$$N(\varepsilon, 1) = \frac{\pi}{3 \cos^{-1} \left( \frac{\cos(f(\varepsilon)) - \cos^2(f(\varepsilon))}{\sin^2(f(\varepsilon))} \right) - \pi}. \quad (3.8)$$

In Figure 3.1, we compare the two estimates using Equations 3.8 and 3.7. For  $\varepsilon = 0.1$ , Equation 3.7 gives  $\mathcal{N} \approx 30$ . And for  $\varepsilon = 0.01$  and  $0.001$ , it gives  $\mathcal{N} \approx 130$  and  $600$ , respectively. As for Equation 3.8 it gives  $\mathcal{N} \approx 34, 156, 725$  for  $\varepsilon = 0.1, 0.01, 0.001$ , respectively.

### 3.3. Fourth-Order Tensors

In this case,  $m = 2$  and  $\mathbf{H}_2$  and  $\mathbf{HS}_2$  have dimensions six and five, respectively. The map  $\mathcal{F}_2^2$  is no longer isotropic with respect to  $L^1$ -norm in  $\mathbf{HS}_2$ . An analytic evaluation of the matrix  $\Lambda^2$  gives

$$\Lambda^2 = \frac{4\pi}{5} \begin{pmatrix} 1 & 1/3 & 1/3 & 0 & 0 & 0 \\ 1/3 & 1 & 1/3 & 0 & 0 & 0 \\ 1/3 & 1/3 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1/3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1/3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1/3 \end{pmatrix}.$$

The singular values of  $\Lambda^2$  arranged in the descending order are

$$\sigma(\Lambda^2) = \frac{4\pi}{3} [1, 2/5, 2/5, 1/5, 1/5, 1/5].$$

This gives  $\eta = 5$ , and Equation 3.6 gives

$$N \approx \frac{3\mathbf{V}_5}{\mathbf{V}^5} f_{\eta=5}(\varepsilon)^{-5} = \frac{90\pi}{16} f_{\eta=5}(\varepsilon)^{-5}.$$

For  $\varepsilon = 0.1$ , this yields  $N \approx 176790$ . However, in  $\mathbf{H}_2$ , the polynomial  $v(x, y, z) = x^2 + y^2 + z^2$  is the constant function 1 on  $\mathbf{S}^2$ . In particular,  $u(x, y, z) = v(x, y, z) / \sqrt{3} \in \mathbf{HS}_2$ , and  $\|\mathcal{F}_2^2(u)\|_1 = 4\pi/3$ . The map  $\mathcal{F}_2^2$  stretches the constant polynomial considerably more than any other quadratic polynomials, and this is the reason for the large condition number  $\eta$ . Let  $\mathbb{R}u$  denote the one-dimensional subspace in  $\mathbf{H}_2$  spanned by the constant polynomial  $u(x, y, z)$ , and  $\mathbf{W}$  its orthogonal complement,

$$\mathbf{H}_2 = \mathbb{R}u \oplus \mathbf{W}.$$

The intersection of the sphere  $\mathbf{HS}_2$  with the subspace  $\mathbf{W}$  is a four-sphere  $\mathbf{S}^4$ . If we specialize to this four-sphere, i.e., polynomials orthogonal to the constant polynomial  $x^2 + y^2 + z^2$ , the condition number  $\eta$  becomes 2 and the dimension of the sphere drops by one. Theorem 2.4 then provides the following estimate on the number of points

$$\mathcal{N} \approx \frac{5\mathbf{V}_4}{2\mathbf{V}^4} f_{\eta=2}(\varepsilon)^{-4} = \frac{40}{3} f_{\eta=2}(\varepsilon)^{-4}.$$

This number is considerably less than 176790. For example, for  $\varepsilon = 0.1$ , we have  $\mathcal{N} \approx 1800$  and for  $\varepsilon = 0.05, 0.01$ ,  $\mathcal{N} \approx 4670, 39620$ , respectively.

### 3.4. Sixth-Order Tensors

In this case,  $m = 3$  and  $\mathbf{H}_3, \mathbf{HS}_m$  have dimensions 10, 9, respectively. The map  $\mathcal{F}_3^2$  is again non-isotropic with respect to  $L^1$ -norm in  $\mathbf{H}_6$ . The singular values of  $\Lambda^3$  arranged in the descending order are

$$\sigma(\Lambda^3) = 4\pi \left[ \frac{19 + \sqrt{193}}{210}, \frac{19 + \sqrt{193}}{210}, \frac{19 + \sqrt{193}}{210}, \frac{19 - \sqrt{193}}{210}, \frac{19 - \sqrt{193}}{210}, \frac{19 - \sqrt{193}}{210}, 2/105, 2/105, 2/105, 1/105 \right].$$

The condition number  $\eta = 16.44$ , which is quite substantial. However, similar analysis as above can be applied to eliminate polynomials in  $\mathbf{HS}_3$  coming from polynomials of lower degree to substantially decrease the condition number. First, the three linear polynomials  $x, y, z$  are now embedded in  $\mathbf{H}_3$  as  $x(x^2 + y^2 + z^2), y(x^2 + y^2 + z^2), z(x^2 + y^2 + z^2)$ . Let  $\hat{r}(\mathbf{x}), \hat{s}(\mathbf{x}), \hat{t}(\mathbf{x}), r(\mathbf{x}), s(\mathbf{x}), t(\mathbf{x})$  be the following polynomials

$$\begin{aligned} \hat{r}(\mathbf{x}) &= x(x^2 + y^2 + z^2) / \sqrt{3}, & r(\mathbf{x}) &= 0.7184x^3 + 0.3951\hat{r}(\mathbf{x}), \\ \hat{s}(\mathbf{x}) &= y(x^2 + y^2 + z^2) / \sqrt{3}, & s(\mathbf{x}) &= 0.7184y^3 + 0.3951\hat{s}(\mathbf{x}), \\ \hat{t}(\mathbf{x}) &= z(x^2 + y^2 + z^2) / \sqrt{3}, & t(\mathbf{x}) &= 0.7184z^3 + 0.3951\hat{t}(\mathbf{x}). \end{aligned}$$

The three polynomials  $r(\mathbf{x}), s(\mathbf{x}), t(\mathbf{x})$  are responsible for the three largest singular values of  $\Lambda^3$ . The smallest singular value of  $4\pi/105$  comes from the polynomial  $q(\mathbf{x}) = xyz$ . Let  $\mathbf{W}$  denote the six-dimensional subspace in  $\mathbf{H}_3$  that is the orthogonal complement of the subspace spanned by  $r(\mathbf{x}), s(\mathbf{x}), t(\mathbf{x})$  and  $q(\mathbf{x})$ ,

$$\mathbf{H}_3 = \mathbb{R}r \oplus \mathbb{R}s \oplus \mathbb{R}t \oplus \mathbb{R}q \oplus \mathbf{W}.$$

The sphere in  $\mathbf{W}$  is five-dimensional, and the condition number of  $\mathcal{F}_3^2$  on  $\mathbf{S}^5$  is  $\eta = 1.2769$ .

$$\mathcal{N} \approx \frac{3\mathbf{V}_5}{\mathbf{V}^5} f_{\eta=1.27}(\varepsilon)^{-5} = \frac{90\pi}{16} f_{\eta=1.27}(\varepsilon)^{-5}.$$

For  $\varepsilon = 0.1, 0.05, 0.01$ , the result above gives  $\mathcal{N} \approx 1943, 6021, 85495$ , respectively.

## 4. Experimental Results

In this section we experimentally validate the proposed theory and at the end of this section we present an application to Diffusion-Weighted MRI. In all the experiments we use tensors in  $\mathbb{R}^3$ , which can be visualized by plotting the corresponding homogeneous polynomial  $P(x, y, z)$  as a spherical function (see Fig. 4.1). Such tensor glyphs can be generated by scaling the radius of a unit sphere at orientation  $\mathbf{x} = [x \ y \ z]^T$  with the value of  $P(x, y, z)$ . Additionally, we assign a color to each tensor glyph by using the following coloring scheme: we use the method in [11, 22] to compute the unit vector  $[x \ y \ z]^T$  that maximizes  $P(x, y, z)$  and then we assign to the R, G, B color channels the squares of the three components in the

vector  $\mathbf{x}$  (i.e.  $R = x^2$ ,  $G = y^2$ ,  $B = z^2$ ). This color map produces smooth color transitions when visualizing fields of tensors such as the diffusion tensor fields.

First, we construct a dataset with samples from  $\Omega_{2m}$  as follows: we first generate random vectors in  $\mathbb{R}^{d(m)}$  using the normal distribution  $\mathcal{N}(\mu = 0, \sigma^2 = 1)$  in  $d(m) = 3, 6, \text{ and } 10$  dimensions, and we use them as coefficients of linear, quadratic and cubic homogeneous polynomials  $p \in \mathbf{HS}_1, \mathbf{HS}_2, \mathbf{HS}_3$  in three variables, respectively. Then we construct  $2^{nd}, 4^{th}$  and  $6^{th}$ -order positive semi-definite tensors that belong to  $\Omega_{2m}$  by taking sums of squares of the polynomials in  $\mathbf{HS}_1, \mathbf{HS}_2, \mathbf{HS}_3$ , respectively. This process is repeated for 5000 times for each order, producing a dataset of 15000 tensors in total. Several of the generated tensors are shown in Fig. 4.1(right). The primary goal of the aforementioned process is to generate samples from  $\Omega_{2m}$  in order to test the error analysis presented in Section 3, and it should not be perceived as a DW-MRI simulation as in this section we do not discuss any application of the proposed method to DW-MR imaging.

In order to investigate how many generators in the finitely-generated cone  $\sum_{2m}^{\mathcal{C}}$  are necessary for our algorithm to approximate accurately a set of given tensors, we apply our framework to the previously described synthetic dataset using finite subsets  $\mathcal{C} \in \mathbf{HS}_m$  of various sizes  $N$ . The sets  $\mathcal{C}$  are constructed as the vertices computed by triangulating the unit  $n$ -sphere. The triangulation is based on a variation of the algorithm for mesh generation presented in [39], which extends to any dimension  $n$  of the  $n$ -sphere. This method is an iterative force-based technique that uses a force displacement function to move the nodes of the mesh and the Delaunay triangulation [14], which is a fundamental and widely used triangulation process, to adjust the topology (i.e. the edges). Obviously, in our particular case we discard the edge information since we only need the finite set of nodes. This algorithm produces at the end the finite subsets  $\mathcal{C} \in \mathbf{HS}_m$  for different predefined sizes  $N$ .

We first use the constructed finite sets  $\mathcal{C}$  in a numerical framework for approximating the

error rate  $\varepsilon$  achieved by the finitely-generated cone  $\sum_{2m}^{\mathcal{C}}$  for  $m = 1, 2, 3$ . The numerical calculations were performed by randomly generating points in the  $n$ -sphere and testing if

each point lies inside or outside the cone  $\sum_{2m}^{\mathcal{C}}$ . The error rate  $\varepsilon$  is the ratio of the points outside the cone over the total number of generated points. For each numerical computation we used 100k points. The numerical approximations are shown as circles in Fig. 4.2. By observing the figures we can see that in most of the cases the numerical approximations are close to the proposed formulas for computing  $N$ . We should note that the results are based on the computed sets  $\mathcal{C}$  using the method in [39]. One may expect that the results will be slightly different if another method is employed for triangulating the  $n$ -sphere.

We also use the sets  $\mathcal{C}$  in a non-negative least squares (NNLS) optimization framework [28]

in order to estimate tensors from the finitely-generated cone  $\sum_{2m}^{\mathcal{C}}$  that approximate the given 15000 tensors. For each order of tensors, the NNLS system is formulated as  $\mathbf{A}\mathbf{w} = \mathbf{b}$ , where  $\mathbf{A}$  a matrix constructed from  $\mathcal{C}$ ,  $\mathbf{w}$  the unknown solution vector and  $\mathbf{b}$  contained the values of the given positive-semidefinite homogeneous polynomial at  $K = 81$  three-dimensional unit vectors  $\mathbf{x}_1 \cdots \mathbf{x}_{81}$  (producing 81 components of  $\mathbf{b}$  as  $b_1 = P(\mathbf{x}_1) \cdots b_{81} = P(\mathbf{x}_{81})$ ) for each tensor in the dataset. Although this problem seems extremely unconstrained in general, in our particular case the NNLS algorithm by definition constrains the number of non-zero elements in the solution vector to be at most  $d(2m)$ , which is significantly smaller than the number of data points  $K$  in all of our experiments. In order to estimate such a constrained solution the NNLS algorithm implements a basis selection mechanism that starts with a set of possible basis vectors in  $\mathcal{C}$ , computes the associated dual vector, and then



reselects the basis in the solution by iteratively performing swaps in order to minimize the entries in the dual vector until they are all non positive. In our particular case of  $m = 1, 2, 3$  the estimated unknown non-zero entries are 6, 15, 28 respectively which are all significantly smaller than the number of given samples  $K = 81$ . For a detailed description of the NNLS algorithm the reader is referred to [28].

The solutions  $\mathbf{w}$  provide tensors in  $\sum_{2m}^{\mathcal{C}}$  that approximate the given tensors in  $\Omega_{2m}$ , for  $m = 1, 2$ , and 3. The computed tensors are compared to the ground truth (given) tensors using the relative  $L^1$ -error (fitting error):

$$\frac{\int_{\mathbf{s}^2} |P_{\text{given}}(\mathbf{x}) - P(\mathbf{x})| d\mathbf{x}}{\int_{\mathbf{s}^2} |P_{\text{given}}(\mathbf{x})| d\mathbf{x}}. \quad (4.1)$$

The histograms of the errors found in the experiments (measured by Eq. 4.1) are plotted in Fig. 4.3 for the case of  $2^{nd}$ ,  $4^{th}$ , and  $6^{th}$ -order tensors, respectively. Obviously, by increasing

$N$ , i.e. the number of generators in the finitely-generated subcone  $\sum_{2m}^{\mathcal{C}}$ , the error decreases correspondingly. The table in Fig. 4.3 reports the mean errors for various difference sizes  $N$  of the generator set.

The experimental results presented in Fig. 4.3 and Fig. 4.4 validate empirically our method as the results corroborate well with our previous analysis on the number of generators required for a given relative error bound. For  $2^{nd}$ -order tensors, the analysis in Section 3 shows that for the error to be less than  $\epsilon = 10\%$ ,  $1\%$ ,  $0.1\%$ , it requires approximately  $N \approx 30$ ,  $130$  and  $600$  generators, respectively. The first plot in Fig. 4.3 shows that with  $N = 45$ , there are no occurrences of error greater than  $10\%$ , and with  $N = 150$ , there are no occurrences of error greater than  $1\%$ . With  $321$  generators, the error becomes negligible. For  $4^{th}$ -order tensors, our analysis shows that for the error to be less than  $\epsilon = 10\%$ ,  $5\%$ , it requires approximately  $N \approx 1800$ ,  $4670$  generators, respectively. This can be seen from the second plot in Fig. 4.3. With  $N < 1500$  generators, there are occurrences of  $10\%$  error, and with  $N = 1500$ , there are no occurrences of error greater than  $10\%$ . To decrease the error under  $5\%$  level, the plot shows that we need at least  $N = 3000$  generators. Finally, for  $6^{th}$ -order tensors, our analysis shows that for the error to be less than  $\epsilon = 10\%$  and  $5\%$ , it requires approximately  $N \approx 1943$  and  $6021$  generators, respectively. The third plot in the figure show that at  $N = 3000$ , there is only a small percentage of errors greater than  $10\%$ , and with  $N = 6000$ , there is an even smaller percentage (less than  $1\%$ ) of errors greater than  $5\%$ . In most cases, our earlier analysis underestimate the required numbers of generators, and this is not surprising as these analysis are themselves based on several approximations. Nevertheless, the experimental results do agree in general with the predictions made in Section 3.

Figure 4.4 shows the running time of the optimization method for fitting one tensor versus the approximation error for various orders and number of generators  $N$  in the set  $\mathcal{C}$ . The running times are measured using an Intel Pentium Dual CPU at  $1.60$  GHz and  $1$ GB RAM. The plots demonstrate that the proposed technique can efficiently estimate positive tensors of various orders. More specifically,  $2^{nd}$ ,  $4^{th}$ , and  $6^{th}$ -order tensors can be estimated using finitely-generated subcones of size  $N = 45$ ,  $N = 900$ , and  $N = 6000$  at  $0.5ms$ ,  $12ms$ , and  $243ms$ , respectively.

#### 4.1. Application: Diffusion-Weighted MRI

Finally, we present an application of the proposed tensor approximation theory to Diffusion-Weighted MRI (DW-MRI). In several DW-MRI processing methods, a diffusion tensor is computed from the acquired diffusion-weighted signals. Negative diffusion values are non-physical; therefore, appropriate methods such as our proposed framework are necessary to ensure positive semi-definiteness of the estimated Diffusion tensors.

In order to demonstrate the necessity for estimating tensors with the positivity constraints, we compare our method with an existing one that computes tensors without the constraints [35]. In this experiment, we use the aforementioned synthetic dataset of 6<sup>th</sup>-order tensors, and we sample the corresponding homogeneous polynomials using  $K = 81$  3-dimensional unit vectors  $\mathbf{x}_1 \cdots \mathbf{x}_{81}$  in the Stejskal-Tanner model [45], producing 81 DW-MRI samples for each tensor in the dataset. Various levels of Rician noise are added to the samples with standard deviations ranging from  $\sigma = 0.04$  up to  $\sigma = 0.12$ . The noisy datasets are given as inputs to: a) the proposed algorithm (using  $N = 6000$ ), and b) the method proposed in in [35], which is one of the several existing methods in the literature [15, 19] that estimate 6<sup>th</sup>-order tensors. For both, the computed 6<sup>th</sup>-order tensors  $P(\mathbf{x})$  are compared to the ground truth tensors using the error defined in Eq. 4.1.

Figure 4.5 shows the comparison of the fitting errors between the two methods for various levels of noise in the data. The results conclusively demonstrate that tensors estimated using positivity constraints approximate the data significantly better than the ones without. We also note that this result agrees with similar comparisons reported earlier for tensors of lower orders (e.g. 4<sup>th</sup>-order comparison in [5]), showing that the errors incurred in approximating positive-valued functions are significantly smaller when positivity constraints are enforced in the process. Our current results have provided further evidence that supports the importance of imposing positivity constraints in this context.

In order to illustrate the performance of our framework on real data sets, we applied the method to a DW-MRI data set of an excised rat hippocampus (shown in Fig. 4.6). The data set contains 46 images acquired using a pulsed gradient spin echo pulse sequence, with 45 different diffusion gradients and approximate  $b$  value of  $1250 \text{ s/mm}^2$ . Figure 4.6 shows the computed 6<sup>th</sup>-order diffusion tensor field. The highlighted regions of interest demonstrate the variability of the estimated structures. At each voxel, the fiber orientations can be estimated from the peaks of the displacement probability, which can be computed from the diffusion tensors as was shown in [5].

Finally, Fig. 4.7 presents the results obtained by applying our method to a DW-MRI dataset from an excised rat optic chiasm. The data acquisition protocol was the same as in the rat hippocampus dataset. The computed field of 4<sup>th</sup>-order diffusion tensors is shown in the center. Using the estimated diffusion tensors, we can compute the underlying fiber orientations by computing the orientations that correspond to the maxima of the water molecule displacement probabilities. The computed fiber orientations are shown on the right and they agree with the known fiber orientations in the optic chiasm. Further quantitative validations of these orientations with respect to those from histology will be performed as part of our future work.

## 5. Discussion and conclusions

Symmetric positive semi-definite tensors have been used in many applications. Although there are existing methods for imposing positivity constraints on the estimated tensors of order two and four, none of these techniques can be easily extended to higher orders. In this paper, we presented a framework for estimating PSD tensors of any order by approximating

the space (cone) of PSD tensors with a finitely-generated subcone  $\Sigma_{2m}$ . We discussed in detail the geometry of the higher-order tensors, and we presented an explicit characterization of the approximation, using the subset of semi-definite tensors that can be written as a sum of squares of tensors of order  $m$ . This approximation leads to a non-negative linear least-squares (NNLS) optimization problem, which can be efficiently solved, as it was demonstrated using synthetic datasets and real diffusion-weighted MR images.

An interesting property of the NNLS optimization algorithm is that it produces sparse solution vectors. In our particular case, although the problem seems significantly unconstrained, the solution vector contains at most  $d(2m)$  non-zero weights, which corresponds to the rank of the basis matrix. Therefore if the finitely-generated set  $\mathcal{C}$  contains a few thousands bases, the algorithm will select only 6, 15, 28 for tensors of order 2, 4, and 6 respectively. Note that the number of non-zero weights in the solution vector equals to the number of the unique unknown parameters of the symmetric tensor in each case. The sparsity of NNLS in comparison with other optimization techniques for modeling the diffusion-weighted MR signal has also been studied in [27].

In our experiments the sets  $\mathcal{C}$  were generated by tessellating the unit  $n$ -sphere using the iterative force-based technique in [39]. The vertices produced by this algorithm form the finite subset  $\mathcal{C} \in \mathbf{HS}_m$  for different predefined sizes  $N$ . An alternative approach could involve constructing  $\mathcal{C}$  as a finite dictionary of elements in  $\mathbf{HS}_m$  by running a training algorithm on a control dataset [31]. A finite set of diffusion basis for multi-fiber reconstruction is also employed by the method in [43].

One of the advantages of the proposed algorithm is that it enforces positive semi-definite constraints to the estimated tensors. The need for positivity constraints in DW-MRI has been demonstrated in [6] and [5]. It has been shown that unconstrained methods may yield negative diffusivities in real datasets, especially in voxels with high anisotropy or in the presence of noise in the data.

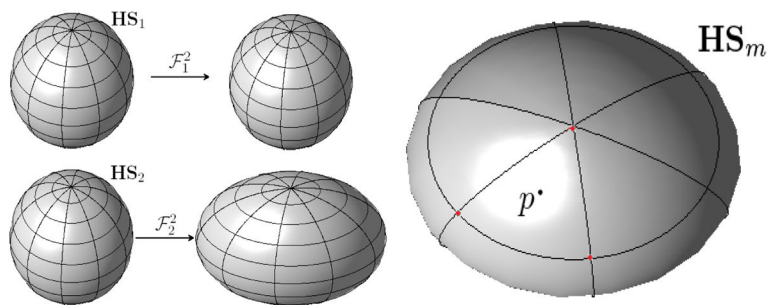
Finally, although high order tensors can approximate several distinct fiber orientations, in the current standard clinical settings for DW-MRI acquisition most of the multi-fiber reconstruction techniques cannot estimate more than two fiber orientations [41], due to the low diffusion weighting (b-value) and the small number of gradient orientations. However, theoretically or in experimental settings with higher b-values and larger sets of diffusion gradient orientations, the proposed technique can estimate up to 2 and 3 distinct fiber orientations using tensors of order 4 and 6 respectively, which also agrees with the results presented in [35].

## References

1. Abramowitz, M.; Stegun, IA. Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables. New York: Dover; 1972.
2. Aganj, Iman; Lenglet, Christophe; Sapiro, Guillermo. Odf reconstruction in q-ball imaging with solid angle consideration. ISBI; 2009. p. 1398-1401.
3. Alexander, Daniel C. Maximum entropy spherical deconvolution for diffusion MRI. IPMI; 2005. p. 76-87.
4. Astola L, Florack L. Finsler geometry on higher order tensor fields and applications to high angular resolution diffusion imaging. Scale Space and Variational Methods in Computer Vision. 2009:224–234.
5. Barmpoutis A, Hwang MS, Howland D, Forder JR, Vemuri BC. Regularized positive-definite fourth-order tensor field estimation from DW-MRI. NeuroImage. 2009; 45:153–162.
6. Barmpoutis, Angelos; Jian, Bing; Vemuri, Baba C.; Shepherd, Timothy M. Symmetric positive 4th order tensors and their estimation from diffusion weighted MRI. IPMI. 2007; 4584:308–319.

7. Barmpoutis, A.; Kumar, R.; Vemuri, BC.; Banerjee, A. Beyond the Lambertian assumption: A generative model for apparent BRDF fields of faces using anti-symmetric tensor splines. Proceedings of CVPR08: IEEE Conference on Computer Vision and Pattern Recognition; 2008. p. 1-6.
8. Barmpoutis, A.; Vemuri, BC.; Forder, JR. Fast displacement probability profile approximation from hardi using 4th-order tensors. Proceedings of ISBI08: IEEE International Symposium on Biomedical Imaging; 2008. p. 911-914.
9. Barmpoutis A, Vemuri BC, Shepherd TM, Forder JR. Tensor splines for interpolation and approximation of DT-MRI with applications to segmentation of isolated rat hippocampi. TMI: Transactions on Medical Imaging. 2007; 26:1537–1546.
10. Bassar PJ, Mattiello J, Lebihan D. Estimation of the Effective Self-Diffusion Tensor from the NMR Spin Echo. J Magn Reson B. 1994; 103:247–254. [PubMed: 8019776]
11. Bloy, L.; Verma, R. On computing the underlying fiber directions from the diffusion orientation distribution function. In the proceedings of MICCAI; 2008. p. 1-8.
12. Boyd, SP.; Vandenberghe, L. Convex optimization. Cambridge University Press; 2004.
13. Cho, Kuan-Hung; Yeh, Chun-Hung; Tournier, Jacques-Donald; Chao, Yi-Ping; Chen, Jyh-Horng; Lin, Ching-Po. Evaluation of the accuracy and angular resolution of q-ball imaging. NeuroImage. 2008; 42:262–271. [PubMed: 18502152]
14. Delaunay B. Sur la sphre vide. Izvestia Akademii Nauk SSSR, Otdelenie Matematicheskikh i Estestvennykh Nauk. 1934; 7:793800.
15. Descoteaux, Maxime; Angelino, Elaine; Fitzgibbons, Shaun; Deriche, Rachid. Apparent diffusion coefficients from high angular resolution diffusion imaging: Estimation and applications. Magnetic Resonance in Medicine. 2006; 56:395–410. [PubMed: 16802316]
16. Descoteaux, Maxime; Angelino, Elaine; Fitzgibbons, Shaun; Deriche, Rachid. Regularized, fast and robust analytical q-ball imaging. MRM. 2007; 58:497–510.
17. Descoteaux, Maxime; Deriche, Rachid; Le Bihan, Denis; Mangin, Jean-Francois; Poupon, Cyril. Diffusion propagator imaging: Using laplace’s equation and multiple shell acquisitions to reconstruct the diffusion propagator. IPMI; 2009. p. 1-13.
18. Fletcher PT, Lu Conglin, Pizer SM, Joshi Sarang. Principal geodesic analysis for the study of nonlinear statistics of shape. IEEE Transactions on Medical Imaging. 2004; 23:995–1005. [PubMed: 15338733]
19. Florack, LMJ.; Balmachnov Sizykh, EG. Two canonical representations for regularized high angular resolution diffusion imaging. MICCAI Workshop on Computational Diffusion MRI; 2008. p. 94-105.
20. Ghosh, A.; Descoteaux, M.; Deriche, R. Riemannian framework for estimating symmetric positive definite 4th order diffusion tensors. Proceedings of MICCAI; 2008. p. 858-865.
21. Ghosh, A.; Moakher, M.; Deriche, R. Ternary quartic approach for positive 4th-order diffusion tensors revisited. Proceedings of ISBI; 2009. p. 618-621.
22. Ghosh, A.; Tsigaridas, E.; Descoteaux, M.; Comon, P.; Mourrain, B.; Deriche, R. A polynomial based approach to extract the maxima of an antipodally symmetric spherical function and its application to extract directions from the orientation distribution function in diffusion MRI. Workshop on Computational Diffusion MRI; MICCAI. 2008.
23. Han D, Qi L, Wu EX. Extreme diffusion values for non-gaussian diffusions. Optimization Methods and Software. 2008; 23:703–716.
24. Hilbert D. Über die darstellung definiter formen als summe von formenquadraten. Math Ann. 1888; 32:342–350.
25. Huber, Greg. Gamma function derivation of n-sphere volumes. Am Math Monthly. 1982; 89:301–302.
26. Jensen JH, Helpers JA, Ramani A, Lu H, Kaczynski K. Diffusional kurtosis imaging: The quantification of non-gaussian water diffusion by means of magnetic resonance imaging. MRM. 2005; 53:1432–1440.
27. Jian, B.; Vemuri, BC. Multi-fiber reconstruction from diffusion mri using mixture of wisharts and sparse deconvolution. In the proceedings of IPMI; 2007. p. 384-395.
28. Lawson, CL.; Hanson, RJ. Solving Least Squares Problems. Prentice-Hall; 1974.

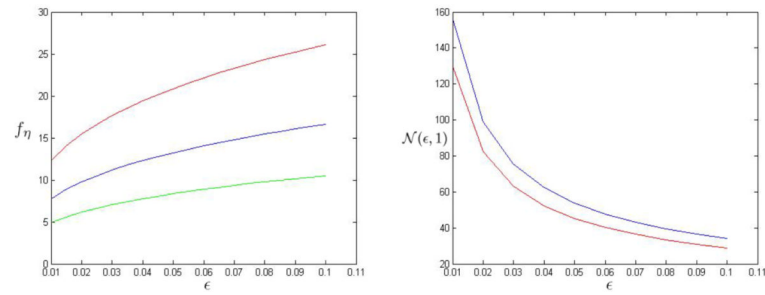
29. Lenglet C, Rousson M, Deriche R. DTI segmentation by statistical surface evolution. *IEEE Trans Med Imaging*. 2006; 25:685–700. [PubMed: 16768234]
30. Liu C, Acar B, Moseley ME. Characterizing non-gaussian diffusion by using generalized diffusion tensors. *Magnetic Resonance in Medicine*. 2004; 51:924–937. [PubMed: 15122674]
31. Mallat S, Zhang Z. Matching pursuits with time-frequency dictionaries. *Trans on Signal Processing*. 1993; 41:3397–2415.
32. Moakher M. Fourth-order cartesian tensors old and new facts, notions applications. *Quarterly Journal of Mechanics and Applied Mathematics*. 2008; 61:181–203.
33. Moakher, M. The algebra of fourth-order tensors with applications to diffusion MRI. In: Laidlaw, D.; Weickert, J., editors. *Visualization and Processing of Tensor Fields*. 2009. p. 57-80.
34. Moakher M, Norris AN. The closest elastic tensor of arbitrary symmetry to an elasticity tensor of lower symmetry. *Journal of Elasticity*. 2006; 85(3):215–263.
35. Ozarslan E, Mareci TH. Generalized diffusion tensor imaging and analytical relationships between DTI and HARDI. *MRM*. 2003; 50:955–965.
36. Ozarslan E, Vemuri BC, Mareci TH. Generalized scalar measures for diffusion MRI using trace, variance, and entropy. *Magn Reson Med*. 2005; 53:866–76. [PubMed: 15799039]
37. Pasternak O, Sochen N, Basser PJ. The effect of metric selection on the analysis of diffusion tensor mri data. *NeuroImage*. 2010; 49:2190–2204. [PubMed: 19879947]
38. Pennec X, Fillard P, Ayache N. A Riemannian framework for tensor computing. *International Journal of Computer Vision*. 2005; 65
39. Persson PO, Strang G. A simple mesh generator in matlab. *SIAM Review*. 2004; 46:329–345.
40. Petrovic V. Concircular curvature tensor. *Pub Inst Math*. 1979; 25:131–137.
41. Prckovska, V., et al. Optimal acquisition schemes in high angular resolution diffusion weighted imaging. In the proceedings of MICCAI; 2008. p. 9-17.
42. Qi L, Han D, Wu EX. Principal invariants and inherent parameters of diffusion kurtosis tensors. *Journal of Mathematical Analysis and Applications*. 2009; 349:165–180.
43. Ramirez-Manzanares A, et al. Diffusion basis functions decomposition for estimating white matter intravoxel fiber geometry. *IEEE Transactions on Medical Imaging*. 2007; 26:1091–1102. [PubMed: 17695129]
44. Schultz T, Seidel HP. Estimating crossing fibers: A tensor decomposition approach. *IEEE Trans Vis Comput Graph*. 2008; 14:1635–1642. [PubMed: 18989020]
45. Stejskal EO, Tanner JE. Spin diffusion measurements: Spin echoes in the presence of a time-dependent field gradient. *Journal of Chemical Physics*. 1965; 42:288–292.
46. Wang, Wei; Gao, Jinghuai; Li, Kang. Structure-adaptive anisotropic filter with local structure tensors. *Intelligent Information Technology Applications, 2007 Workshop on*; 2008. p. 1005-1010.
47. Wang Z, Vemuri BC. DTI segmentation using an information theoretic tensor dissimilarity measure. *IEEE Transactions on Medical Imaging*. 2005; 24:1267–1277. [PubMed: 16229414]
48. Wang, Zhizhou; Vemuri, Baba C.; Chen, Yunmei; Mareci, Thomas H. A constrained variational principle for direct estimation and smoothing of the diffusion tensor field from complex dwi. *IEEE Trans Med Imaging*. 2004; 23:930–939. [PubMed: 15338727]
49. Yassine, I.; McGraw, T. 4th order diffusion tensor interpolation with divergence and curl constrained bezier patches. In *Proceedings of ISBI*; 2009. p. 634-637.



**Fig. 2.1.**

**Left:** Comparison between  $\mathcal{F}_1^2$  and  $\mathcal{F}_2^2$ . Let  $S_r^{2m}$ ,  $r > 0$  denote the circle with radius  $r$  centered at origin in  $\mathbf{H}_{2m}$ .  $\mathcal{F}_1^2$  is isotropic in the sense that  $\|\mathcal{F}_1^2(p)\|_1 \in S_{4\pi/3}^2$  for all  $p \in \mathbf{HS}_1$ .  $\mathcal{F}_2^2$ , on the other hand, is not isotropic.  $\mathbf{HS}_2$  is the five-dimensional sphere  $\mathbf{S}^5$ . Its equator can be identified with  $\mathbf{S}^4$ , and the two polynomials  $\pm 1/\sqrt{3}(x^2+y^2+z^2)$  form the two poles. Inside the equator, are embedded  $\mathbf{S}^1$  and  $\mathbf{S}^2$ .  $\mathcal{F}_2^2$  maps the poles  $\pm(x^2+y^2+z^2)$  to  $S_{4\pi/3}^4$ , and it maps the embedded  $\mathbf{S}^1$  and  $\mathbf{S}^2$  to  $S_{8\pi/15}^4, S_{4\pi/15}^4$ , respectively. The condition number  $\eta$  for  $\mathcal{F}_2^2$  on  $\mathbf{HS}_2$  is 5. Restricting  $\mathcal{F}_2^2$  to the equator  $\mathbf{S}^4$ , the condition number improves to 2.

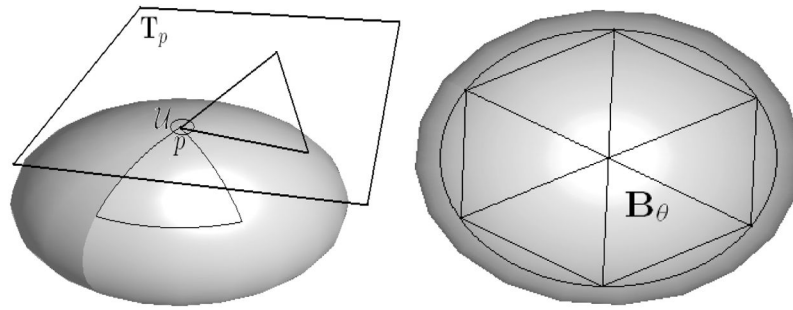
**Right:** Local and non-local approximations. For each  $p \in \mathbf{HS}_m$ , Lemma 2.8 approximates  $p$  first with the vertices of the simplex containing  $p$ . This local approximation is further improved using non-local approximations as the polynomials  $(p_i - p_j)^2$  are approximated by polynomials in  $\mathbf{HS}_m$  that are generally far from  $p$ .



**Fig. 3.1.**

**Left:** Plots of  $f_\eta$  for  $\eta = 1, 2, 4$  in red, blue and green, respectively.  $\epsilon$  varies from 0.01 to 0.1 and  $\theta$  is given in degree. **Right:** Comparison plot of  $\mathcal{N}(\epsilon, 1)$  according to Equations 3.7 (in red) and 3.8 (in blue). The estimate using Equation 3.7 is between 17% and 20% less than the estimate using Equation 3.8.

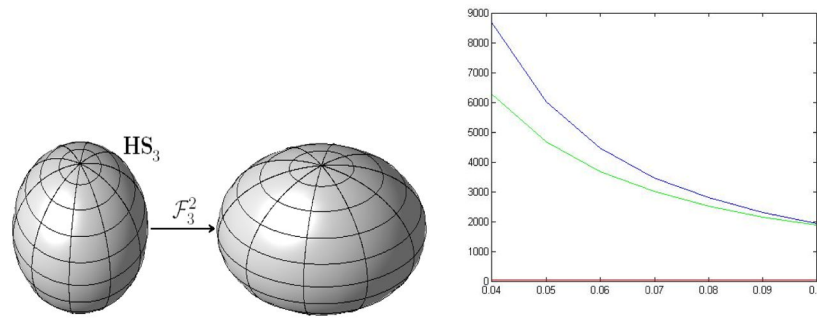




**Fig. 3.2.**

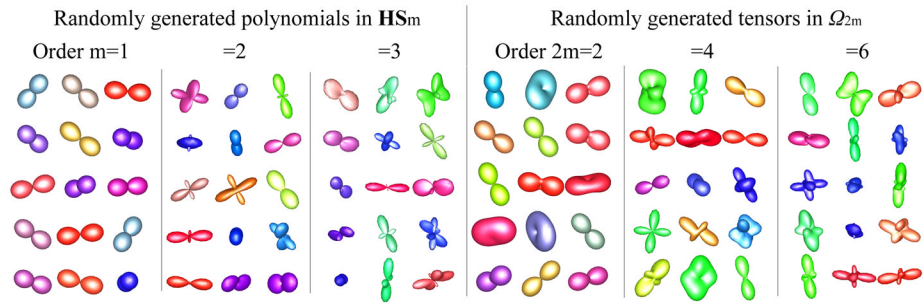
**Left:** For small  $\theta$ , we can approximate the volume of a  $\theta$ -regular spherical simplex by the volume of a  $\theta$ -regular Euclidean simplex. The exponential map  $\mathbf{Exp}_p$  maps a neighborhood of the origin in the tangent space  $\mathbf{T}_p$  diffeomorphically onto a neighborhood  $\mathcal{U}$  at  $p$ . Since the derivative of  $\mathbf{Exp}_p$  at  $p$  is the identity, for small enough  $\theta$ ,  $\mathbf{Exp}_p$  is close to an isometry in  $\mathbf{B}_\theta$ . **Right:** The average degree of a vertex,  $\nu$ , can be approximated by the number of  $\theta$ -regular simplexes contained in the ball of radius  $\theta$ .



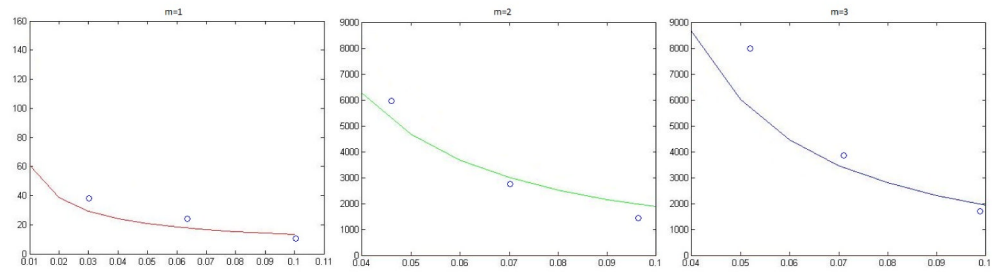


**Fig. 3.3.**

The geometry of the map  $\mathcal{F}_3^2$ . **Left:**  $HS_3$  is the nine-dimensional sphere  $S^9$ . The decomposition of  $H_3$  into four subspaces of dimensions of 3, 3, 3, 1 respectively implies that  $HS_3$  contains separate copies of sphere  $S^2$ ,  $S^2$ ,  $S^2$  and  $S^0$ .  $\mathcal{F}_3^2$  maps these spheres to spheres of radii  $52\pi/83$ ,  $68\pi/699$ ,  $8\pi/105$  and  $4\pi/105$ , respectively. **Right:** The number of generators in  $\Sigma_2$ ,  $\Sigma_4$  and  $\Sigma_6$  that can ensure the given accuracy requirement. The plots for  $m = 1, 2, 3$  are in red, blue and green, respectively.



**Fig. 4.1.** Examples of randomly computed symmetric positive semi-definite tensors in  $\Omega_2$ ,  $\Omega_4$ ,  $\Omega_6$ . The tensor glyphs are shown.

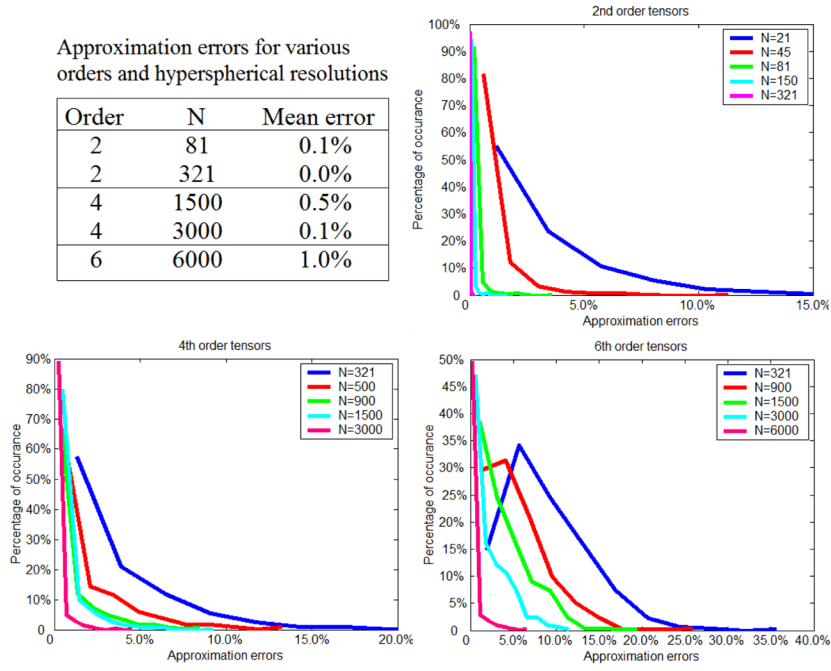


**Fig. 4.2.**

Comparison of the proposed formulas for computing  $N$  for  $m=1,2,3$  with results produced using a numerical approximation algorithm. The horizontal axis show the the accuracy achieved by  $N$  finite generators (vertical axis) in the unit  $n$ -sphere. The circles show the numerical results produced for specific sets  $\mathcal{C}$  of various sizes  $N$ .

Approximation errors for various orders and hyperspherical resolutions

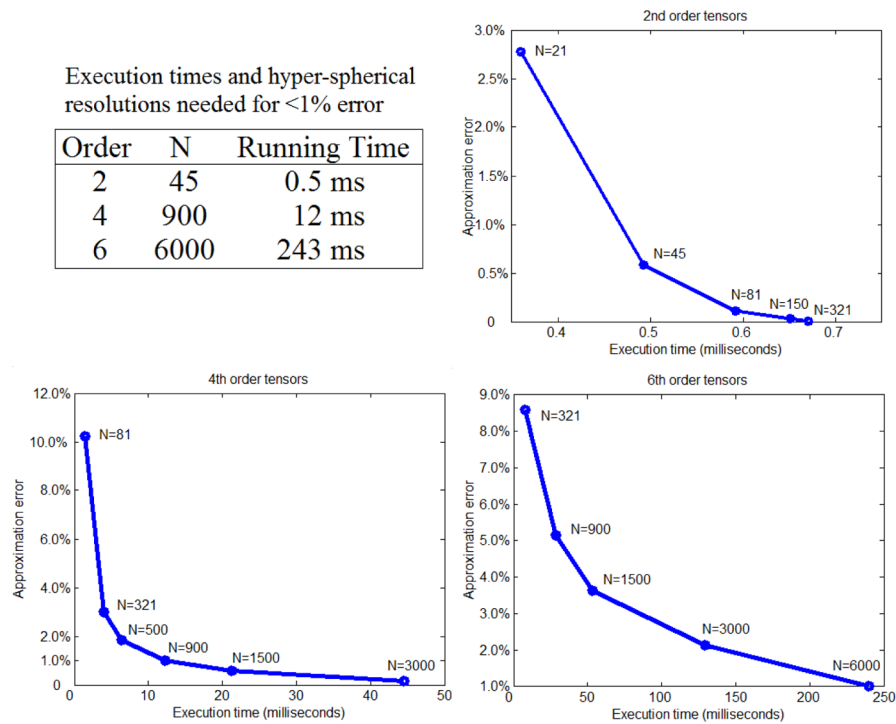
Order	N	Mean error
2	81	0.1%
2	321	0.0%
4	1500	0.5%
4	3000	0.1%
6	6000	1.0%



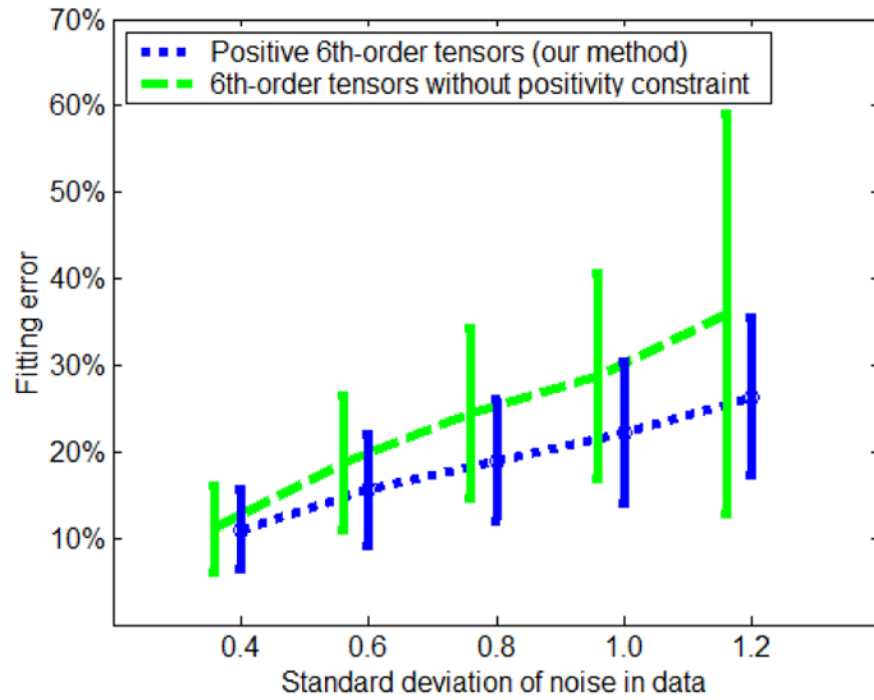
**Fig. 4.3.** Histograms of tensor fitting errors obtained by our method for the case of 2<sup>nd</sup>, 4<sup>th</sup>, and 6<sup>th</sup>-order tensors respectively, using various sizes N of the set  $\mathcal{C}$ . The vertical axis corresponds to the percentage of the tensors in the dataset (i.e. number of occurrences), and the horizontal axis corresponds to the given fitting error.

Execution times and hyper-spherical resolutions needed for <1% error

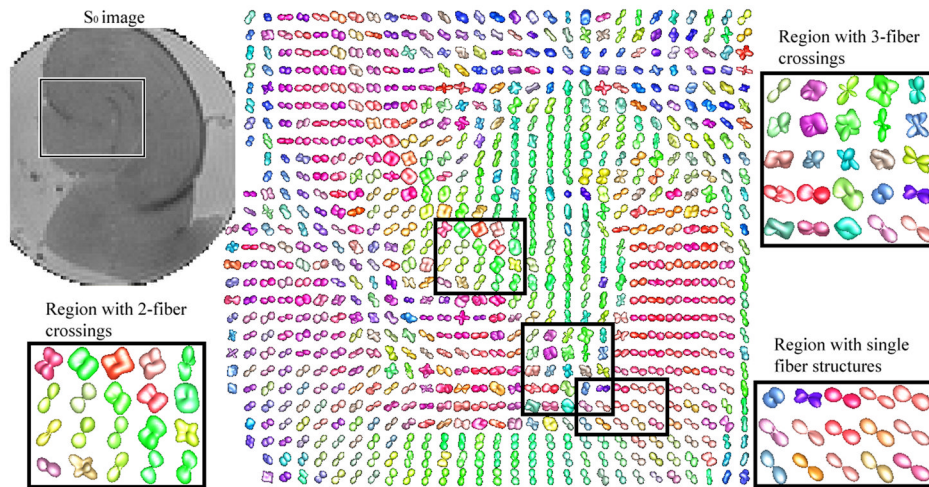
Order	N	Running Time
2	45	0.5 ms
4	900	12 ms
6	6000	243 ms



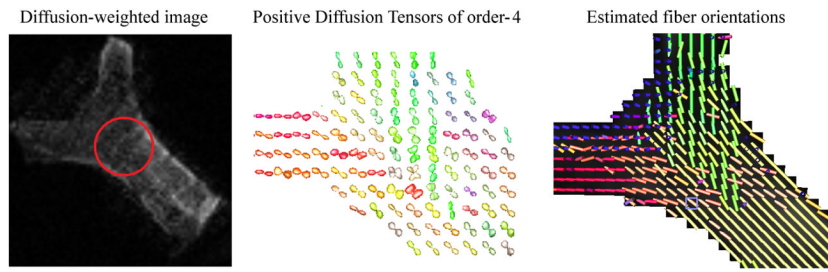
**Fig. 4.4.** Plots of the running time of our method for fitting one tensor versus the approximation error for the case of 2<sup>nd</sup>, 4<sup>th</sup>, and 6<sup>th</sup>-order tensors, using various sizes N of the set  $\mathcal{C}$ . The vertical axis corresponds to the obtained mean fitting error, and the horizontal axis corresponds to the execution time.



**Fig. 4.5.** Comparison of the 6<sup>th</sup>-order tensor fitting errors obtained by the proposed method and the technique in [35] for various Rician noise levels in the data.



**Fig. 4.6.** DW-MRI dataset from an isolated rat hippocampus. The image without diffusion weighting ( $S_0$ ) is shown on the top left. The 6<sup>th</sup>-order diffusion tensors estimated by the proposed method are shown as a field of spherical functions. The three regions of interest depict 6<sup>th</sup>-order diffusion tensors that model one, two, and three fiber structures.



**Fig. 4.7.** DW-MRI dataset from an isolated rat optic chiasm. A field of 4<sup>th</sup>-order diffusion tensors computed by the proposed method is shown in the central plate. The corresponding estimated fiber orientations are shown on the right.