

On the quotient representation for the essential manifold

Roberto Tron and Kostas Daniilidis*

Abstract

The essential matrix, which encodes the epipolar constraint between points in two projective views, is a cornerstone of modern computer vision. Previous works have proposed different characterizations of the space of essential matrices as a Riemannian manifold. However, they either do not consider the symmetric role played by the two views, or do not fully take into account the geometric peculiarities of the epipolar constraint. We address these limitations with a characterization as a quotient manifold which can be easily interpreted in terms of camera poses. While our main focus is on theoretical aspects, we include experiments in pose averaging, and show that the proposed formulation produces a meaningful distance between essential matrices.

1. Introduction

The essential matrix and the epipolar constraint, introduced in [10], have been a major mainstay of computer vision for the last thirty years, and are the basic building block in any Structure from Motion (SfM) system. Its robust estimation from image data is now textbook material [8, 11]. In practical terms, the space of essential matrices is a subset of $\mathbb{R}^{3 \times 3}$, but the algebraic relations imposed by the epipolar constraint render its geometry far from trivial. There have been a few attempts to interpret this space as a Riemannian manifold. The earliest works in this aspect are [12, 13], which use the relative pose between the two cameras (with a normalized translation) to parametrize the space of essential matrices (*i.e.*, each essential matrix is given as the product of a skew-symmetric matrix with a rotation). This implies a preferential treatment of one of the two cameras, whose local reference frame is chosen as a global reference frame, thus breaking the natural symmetry of the constraint. A different representation, based on the Singular Value Decomposition (SVD) of the essential matrix, was used in [5, 9, 15, 16]. While this representation has a natural symmetry, previous

works do not provide an intuitive geometric interpretation of its parameters. Also, they do not properly take into account the well-known twisted-pair ambiguity (the fact that four different pairs of poses correspond to the same essential matrix), and the algorithm used for the computation of the logarithm map (which is related to the notion of geodesics in this space) is neither efficient nor rigorously motivated.

In this paper, we propose a characterization related to the aforementioned SVD representation, and show how:

1. Our approach naturally arises from a particular choice of the global reference frame, and that the parameters have a clear geometric meaning.
2. The cheirality constraint (*i.e.*, the constraint that all the points lie in front of both cameras) impacts our representation and how it can be used to simplify the structure of the space.
3. To endow the space with a Riemannian manifold structure, and how to naturally obtain geodesics in this space from those in the space of camera rotations.
4. To efficiently compute the logarithm map and distance function, and how these can be used in a two-view SfM problem using the Weiszfeld algorithm.

Some material in this paper might appear quite basic for any reader versed in computer vision. However, it is necessary to revisit it and place it in the context of our work.

2. Definitions and notation

In this section we recall several notions from Riemannian geometry and group theory. We mention just the minimum necessary to follow the paper, and we refer the reader to the literature for the complete and rigorous definitions [3].

At a high level, a *manifold* \mathcal{M} is defined by a topological space together with a set of overlapping local coordinate charts, which allow to locally parametrize the space and smoothly pass from one chart to the other. The *tangent space* at a point $x \in \mathcal{M}$, denoted as $T_x\mathcal{M}$, can be defined as the linear space containing all the *tangent vectors* corresponding to the curves passing through x . We use the notation v^\vee to denote the vector of coordinates of v in some basis for $T_x\mathcal{M}$. A *vector field* X assigns a tangent vector to each x in \mathcal{M} or a subset of it. A Riemannian manifold $(\mathcal{M}, \langle \cdot, \cdot \rangle)$ is a manifold equipped with a metric, that is, a collection

*The authors are with the Department of Computer and Information Science, University of Pennsylvania, Philadelphia, PA 19104, {tron, kostas}@seas.upenn.edu. This work was supported by grants NSF-IIP-0742304, NSF-OIA-1028009, ARL MAST-CTA W911NF-08-2-0004, and ARL RCTA W911NF-10-2-0016, NSF-DGE-0966142, and NSF-IIS-1317788.

of inner products $\langle \cdot, \cdot \rangle_x$ over $T_x\mathcal{M}$ which varies smoothly with x . The metric is used to define the length of a curve $\gamma : \mathbb{R} \supset [a, b] \rightarrow \mathcal{M}$. A curve is a *geodesic* if the covariant derivative of its tangent is zero, *i.e.*, $\nabla_{\dot{\gamma}}\dot{\gamma} \equiv 0$ (where ∇ is the so-called *Levi-Civita connection*). The *exponential map* $\exp_x : T_x\mathcal{M} \rightarrow \mathcal{M}$ maps each tangent vector v to the endpoint of the unit-speed geodesic starting at x with tangent v . The *logarithm map* \log_x is the inverse of the \exp_x and is defined (in general) only on a neighborhood of x . We use the shorthand notation $\text{Log} = \log^\vee$. For any point x and any curve $\gamma(t)$ in \mathcal{M} sufficiently close together, the logarithm is related to the distance function by these two relations:

$$d(x, \gamma(t)) = \|\text{Log}_x(\gamma(t))\|, \quad (1)$$

$$\frac{d}{dt} \frac{1}{2} d^2(x, \gamma(t)) = -\text{Log}_x(\gamma(t))^T \dot{\gamma}(t)^\vee. \quad (2)$$

Given a map between manifolds $f : \mathcal{M} \rightarrow \tilde{\mathcal{M}}$, we define Df as the *differential* of the map, *i.e.*, the linear operator (the Jacobian, in local coordinates) that maps $T_x\mathcal{M}$ to $T_{f(x)}\tilde{\mathcal{M}}$ for any $x \in \mathcal{M}$ and satisfies, for any locally defined curve $\gamma(t) \in \mathcal{M}$, the expression

$$\left(\frac{d}{dt} f(\gamma(t)) \right)^\vee = Df(\gamma(t)) \dot{\gamma}^\vee. \quad (3)$$

A *group* is a set G together with an operation $\circ : G \times G \rightarrow G$ which satisfy the four axioms of closure, associativity, identity element (denoted as $e \in G$) and inverse element. A *group action* “ \cdot ” on a set M is a mapping $\cdot : G \times M \rightarrow M$ which satisfies the properties $g \cdot (h \cdot x) = (g \circ h) \cdot x$ and $e \cdot x = x$ for all $g, h \in G, x \in M$. The group action induces an *equivalence relation* between the points in M , and we say that x is equivalent to y , *i.e.*, $x \sim y$ if there exist $g \in G$ such that $g \cdot x = y$. We denote all the elements equivalent to $x \in M$ as the *equivalence class* $[x]$. The *quotient space* M/G is the space of all equivalence classes. The *canonical projection* $\pi : G \rightarrow M/G$ maps each point $x \in \mathcal{M}$ to $[x]$.

In this paper, we will heavily use the space of 3-D rotations $SO(3) = \{R \in \mathbb{R}^{3 \times 3} : R^T R = I, \det(R) = 1\}$, and, to a less extent, the space of rigid body transformations $SE(3) = SO(3) \times \mathbb{R}^3$. We will also use $SO(3)^2$, the cartesian product of $SO(3)$ with itself. The space $SO(3)$ is a *Lie group*, *i.e.*, it is at the same time a group (with matrix multiplication as group operation) and a manifold. The tangent space at $R \in SO(3)$ is $T_R SO(3) = \{RV : V \in \mathfrak{so}(3)\}$, where $\mathfrak{so}(3)$ is the space of 3×3 skew-symmetric matrices. We can identify a tangent vector $v \in T_R SO(3)$ with a vector of local coordinates $w \in \mathbb{R}^3$ using the usual *hat* $(\cdot)^\wedge$ and *vee* $(\cdot)^\vee$ operators, given by the relations

$$w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} \begin{matrix} (\cdot)^\wedge \\ \rightleftarrows \\ (\cdot)^\vee \end{matrix} v = R \begin{bmatrix} 0 & -w_3 & w_2 \\ w_3 & 0 & -w_1 \\ -w_2 & w_1 & 0 \end{bmatrix}. \quad (4)$$

With this notation, the standard metric for $SO(3)$, with $v_1, v_2 \in T_R SO(3)$, is given by

$$\langle v_1, v_2 \rangle = (v_1^\vee)^T v_2^\vee. \quad (5)$$

The exponential and logarithm maps for $SO(3)$ can be computed in closed form by using the Rodrigues’ formula [11].

We denote as $R_x(\theta), R_y(\theta), R_z(\theta)$, the rotations around the x, y and z axes, respectively, with angle $\theta \in [-\pi, \pi]$, and as e_z the unit vector aligned with the z axis. We denote as $I \in \mathbb{R}^{3 \times 3}$ the identity matrix and as $P_z = \text{diag}(1, 1, 0)$ the standard projector on the xy -plane. We denote as $(A)_{ij}$ the element in row i , column j of the matrix A . For standard vectors $a \in \mathbb{R}^3$, $[a]_\times : \mathbb{R}^3 \rightarrow \mathfrak{so}(3)$ denotes the matrix representation of the cross product operator, *i.e.*, $[a]_\times b = a \times b$ for all $a, b \in \mathbb{R}^3$. We use $[a]_\times^{\text{inv}} : \mathfrak{so}(3) \rightarrow \mathbb{R}^3$ to denote the inverse of this linear mapping.

3. Derivation of the essential matrix

As customary, we model the pose of camera i as $g_i = (R'_i, T'_i) \in SE(3)$, where g_i represents the transformation from camera to world coordinates. Given an image x_i in homogeneous coordinates and the corresponding depth λ_i , the 3-D point in world coordinates is given by

$$X = \lambda_i R'_i x_i + T'_i. \quad (6)$$

Note that a change of world coordinates represented by $g = (R_0, T_0)$, *i.e.*, $X \mapsto R_0 X + T_0$ induces a transformation of the camera representation equivalent to multiplying g_i by g on the left, *i.e.*, $(R'_i, T'_i) \mapsto (R_0 R'_i, R_0 T'_i + T_0)$.

We now derive the essential matrix from two camera poses (R'_i, T'_i) and the two images $x_i, i = 1, 2$, of a same 3-D point X . We follow a general approach [2] as opposed to the traditional one which uses one camera as the global reference frame. From (6), and using the properties $[a]_\times a = 0$ and $b^T [a]_\times b = 0$ for all $a, b \in \mathbb{R}^3$, we have:

$$\lambda_1 R'_1 x_1 + T'_1 = \lambda_2 R'_2 x_2 + T'_2 \quad (7)$$

$$\lambda_1 R'_1 x_1 = \lambda_2 R'_2 x_2 + (T'_2 - T'_1) \quad (8)$$

$$\lambda_1 [T'_2 - T'_1]_\times R'_1 x_1 = \lambda_2 [T'_2 - T'_1]_\times R'_2 x_2 \quad (9)$$

$$x_1^T R_1^T [T'_2 - T'_1]_\times R'_2 x_2 = 0 \quad (10)$$

The essential matrix is then defined as

$$E = R_1^T [T'_2 - T'_1]_\times R'_2. \quad (11)$$

4. The normalized essential space

In this section, we define a canonical decomposition of the essential matrix in terms of two rotations by choosing a global reference frame aligned with the baseline between the two cameras. Then, we define the *normalized essential space*, and analyze its structure as a quotient space (which includes the twisted pair ambiguity). We also give an interpretation in terms of vector transformations.

4.1. The normalized canonical decomposition

Since (10) is a homogeneous equation, we cannot determine the scale of E from image data alone. Also, while E does not depend on the choice of global reference frame, this is not true for its decomposition (11). To remove most of the degrees of freedom, we use the following.

Proposition 4.1. *Any essential matrix E admits, up to scale, the following normalized canonical decomposition:*

$$E = R_1^T [e_z]_{\times} R_2. \quad (12)$$

Proof. Starting from (11), choose a global scale such that $\|T'_2 - T'_1\| = 1$ and let $R_0 \in SO(3)$ be such that $R_0(T'_2 - T'_1) = e_z$. There are infinite candidates for such rotation (we pick one using Householder transformations). Then, by applying the transformation $g_0 = (R_0, 0)$ and using the property $R[a]_{\times} R^T = [Ra]_{\times}$ for all $R \in SO(3)$, we have

$$E = (R_0 R'_1)^T [R_0(T'_2 - T'_1)]_{\times} R_0 R'_2 \quad (13)$$

which is of the form (12) with $R_i = R_0 R'_i$, $i = 1, 2$. \square

Intuitively, the change of world coordinates performed in the proof above aligns the vector $T'_2 - T'_1$ with the z -axis. In this way, the translation direction is known, and we are left with only the information about the two rotations.

Remark 1. *Notice that $[e_z]_{\times} R_z(\frac{\pi}{2}) = P_z = \text{diag}(1, 1, 0)$. Hence, $E = R_1^T P_z (R_z(\frac{\pi}{2}) R_2)$ is a valid SVD of E .*

The value of Remark 1 is twofold. First, it provides a practical way to compute the decomposition (12). Second, it relates our representation with the one of [15], giving a geometric meaning to the SVD of E .

We define the *normalized essential space* \mathcal{M}_E as the image of $SO(3)^2$ under the map (12). Since, according to Prop. 4.1, this map is surjective, \mathcal{M}_E corresponds to the space of all the essential matrices.

4.2. Ambiguities of the canonical form

While the map (12) is surjective, it cannot be also injective, because it is known that the space of essential matrices is five-dimensional, while $SO(3)^2$ is six-dimensional. The extra degree of freedom corresponds to a rotation of the global reference frame around the baseline (*i.e.*, to a particular choice of R_0 in the proof of Prop. 4.1). However, it turns out that this is not the only ambiguity. To be more precise, consider any two points $Q_a, Q_b \in SO(3)^2$ which, through (12), correspond to the essential matrices E_a, E_b . We define an equivalence relation “ \sim ” between points in $SO(3)^2$ as

$$Q_a \sim Q_b \iff E_a = E_b, \quad (14)$$

where, again, equality is intended up to scale (since E_a and E_b are normalized, this reduces to a “up to a sign flip”).

Proposition 4.2. *Define the groups*

$$H_z = \{(R_z(\theta), R_z(\theta)) : \theta \in [-\pi, \pi)\}, \quad (15)$$

$$H_\pi = \{(I, I), (R_x(\pi), R_x(\pi)), (I, R_z(\pi)), (R_x(\pi), R_y(\pi))\} \quad (16)$$

acting on the left on $SO(3)^2$ by simple component-wise left multiplication. Then, the equivalence class of a point Q with respect to “ \sim ” is exactly given by

$$[Q] = \{S_z S_\pi Q : S_z \in H_z, S_\pi \in H_\pi\}. \quad (17)$$

The proof involves first showing that H_z and H_π are subgroups of $SO(3)^2$, and then showing that the only matrices satisfying (17) are those in the equivalence class $[Q]$. The details can be found in the additional material. In the following we will use $S_z = (S_{z1}, S_{z2})$ and $S_\pi = (S_{\pi1}, S_{\pi2})$ to denote points in H_z and in H_π , respectively.

Intuitively, $[Q]$ has four components, each one isomorphic to $SO(2)$. In view of Prop. 4.2, the space \mathcal{M}_E can be identified with the quotient space

$$\mathcal{M}_E = (SO(3) \times SO(3)) / (H_z \times H_\pi), \quad (18)$$

where the actions of H_z and H_π are defined above.

Since $SO(3)^2$ has dimension six, and H_z has dimension one, we get the well known fact that the normalized essential space has dimension five (being discrete, H_π does not change the intrinsic dimension of the space).

4.3. Geometric interpretation

Using the geometric interpretation given by the proof of Prop. 4.1, we now show that also the epipolar constraint $x_1^T E x_2 = 0$ has a geometrical interpretation. Given an essential matrix $E = R_1^T [e_z]_{\times} R_2$, from Prop. 4.2 and the equivalence $[e_z]_{\times} = P_z^T R_z(\frac{\pi}{2}) P_z$, we have

$$x_1^T E x_2 = (P_z S_z S_{\pi1} R_1 x_1)^T R_z(\frac{\pi}{2}) (P_z S_z S_{\pi2} R_2 x_2) = 0. \quad (19)$$

This can be interpreted as the following procedure:

- Take the images x_i and rotate them as $R_i x_i$, $i = 1, 2$. This is equivalent to expressing in global coordinates the vectors corresponding to the images and centering them at the origin. Notice that, by construction, the transformed vectors and the z -axis e_z all lie in the same plane passing through the origin.
- Apply the action of an element of $S_\pi = H_\pi$ (see Figure 1). If $S_\pi = (I, I)$, no changes are made. If $S_\pi = (R_x(\pi), R_x(\pi))$ (and considering also H_z), the direction of the baseline is reversed. If $S_\pi = (I, R_z(\pi))$, one of the cameras is rotated from front-facing to rear-facing. Finally, if $S_\pi = (R_x(\pi), R_y(\pi))$, the last two cases are combined. Note that the coplanarity condition of the transformed vectors with e_z is preserved.

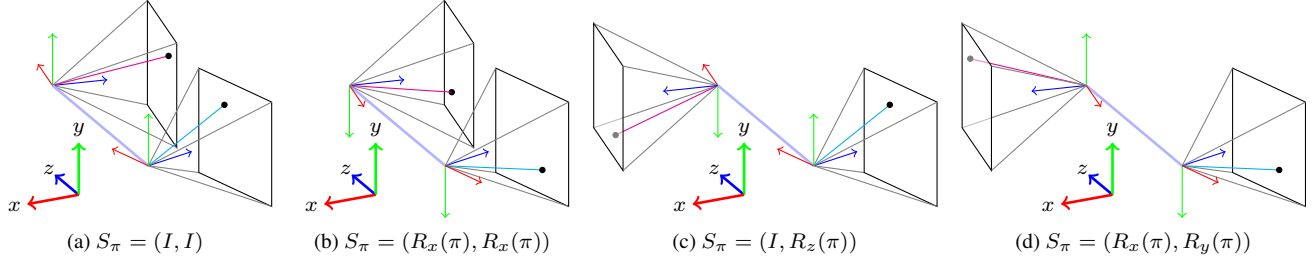


Figure 1: Diagram depicting the geometric twisted pair ambiguity given by the four elements of H_π

- Apply the action of an element of H_z , *i.e.*, rotate the two vectors around the z -axis by an arbitrary amount. This is equivalent to a rotation around the baseline, and does not change the coplanarity condition.
- Project the transformed vectors onto the xy -plane. In practice, this sets the last coordinate to zero. Since before the projection the vectors belonged to the same plane, after the projection they will have the same direction (but, in general, different lengths).
- Rotate one of the projected vectors by $R_z(\frac{\pi}{2})$, *e.g.*, $R_z(\frac{\pi}{2})(P_z R_2 x_2)$. Since the vectors were collinear, they are now orthogonal and the inner product is zero.

Although this interpretation of (19) is probably not new, in our context it shows that the action of H_π corresponds exactly to the well-known twisted-pair ambiguity in the decomposition of the essential matrix.

5. The signed normalized essential manifold

In this section we review how the cheirality constraint can be used to resolve the twisted pair ambiguity (*i.e.*, to choose an element of the group H_π), and how this simplifies the quotient structure of the normalized essential space into what we call the *signed normalized essential space*. We show that this space is a manifold, and that a metric and the corresponding geodesics can be naturally induced from $SO(3)^2$. Finally, we give a Newton-based algorithm for computing the logarithm map and the Riemannian distance.

5.1. Depth triangulation

We can use the simple geometrical interpretation of the canonical form to estimate the depths of the 3-D points enforce the cheirality constraint, *i.e.*, the fact that all these points need to be in front of both cameras.

From the discussion in Section 4.2, we have $T'_2 - T'_1 = e_z$ in the canonical form. Therefore, taking into account H_z and H_π , and assuming noiseless image points, (8) becomes

$$\lambda_1 S_{z1} S_{\pi1} R_1 x_1 = \lambda_2 S_{z2} S_{\pi2} R_2 x_2 + e_z. \quad (20)$$

Note that $e_z = S_{z1} e_z = S_{z2} e_z$, hence we can cancel S_z from (20). We then have the following proposition.

Proposition 5.1. *There is only one choice of S_π for which the solution of (20) is positive, *i.e.*, $\lambda_1, \lambda_2 > 0$.*

The proof, which can be found in the additional material, is similar to the one used to solve the twisted pair ambiguity, *i.e.*, to decide, given an essential matrix, which of the four possible pose estimates to use [11].

5.2. The signed normalized essential manifold

In our context, Prop. 5.1 allows us to pick one of four components in the equivalent class $[Q]$, *i.e.*, we can dispense with the group H_π and consider a new quotient space using H_z alone, which we call the signed essential space. Just for fun, we use the symbol $\mathcal{M}_{\mathcal{E}}$ (because it differs from “ \mathcal{M}_E ” by a 180 degrees rotation). Formally, we have

$$\mathcal{M}_{\mathcal{E}} = (SO(3) \times SO(3))/H_z. \quad (21)$$

In general, a quotient space of a Riemannian manifold is not a Riemannian manifold itself (because it does not satisfy some necessary topological conditions or the choice of a metric might not be an obvious). However, the action of H_z has some “nice” properties which lead to the following.

Proposition 5.2. *The space $\mathcal{M}_{\mathcal{E}}$ can be given a Riemannian manifold structure for which the natural projection $\pi_{\mathcal{M}_{\mathcal{E}}} : SO(3)^2 \rightarrow \mathcal{M}_{\mathcal{E}}$ introduces a local isometry between a subspace of $T_Q SO(3)^2$ and $T_{[Q]} \mathcal{M}_{\mathcal{E}}$.*

This proposition (for which the proof is somewhat technical, see the additional material) means that not only we can endow $\mathcal{M}_{\mathcal{E}}$ with a Riemannian structure, but also that the differential (*i.e.*, the Jacobian, in local coordinates) of the natural projection $\pi_{\mathcal{M}_{\mathcal{E}}}$ induces an orthogonal decomposition of the tangent space $T_Q SO(3)^2$ as:

$$T_Q SO(3)^2 = T_{VQ} SO(3)^2 \oplus T_{HQ} SO(3)^2, \quad (22)$$

where the *vertical space* $T_{VQ} SO(3)^2$ is the subspace of tangent vectors tangential to the equivalence class $[Q]$ and the *horizontal space* $T_{HQ} SO(3)^2$ is its orthogonal complement. Moreover, the same differential uniquely maps each vector $\tilde{v} \in T_{[Q]} \mathcal{M}_{\mathcal{E}}$ to a vector $v \in T_{HQ} SO(3)^2$, called the *horizontal lift* of \tilde{v} . The metric $\langle \cdot, \cdot \rangle_{[Q]}$ for $\mathcal{M}_{\mathcal{E}}$ is then defined

from the metric $\langle \cdot, \cdot \rangle_Q$ in $SO(3)^2$ as

$$\langle \tilde{u}, \tilde{v} \rangle_{[Q]} = \langle u, v \rangle_Q \quad (23)$$

In our case, the vertical space is one-dimensional, and $T_{VQ}SO(3)^2 = \text{span}(v_V)$, where

$$v_V = ([e_z]_{\times} R_1, [e_z]_{\times} R_2) = (R_1 [R_1^T e_z]_{\times}, R_2 [R_2^T e_z]_{\times}). \quad (24)$$

By definition, then, the horizontal space at Q includes all vectors v_H such that $v_H \perp v_V$, i.e.,

$$0 = \langle v_V, v_H \rangle = e_z^T (R_1 v_{H1} + R_2 v_{H2}), \quad (25)$$

where $(v_H)^{\vee} = \text{stack}(v_{H1}, v_{H2})$.

We can take (25) as the condition defining horizontal vectors at Q . Given a vector $v \in T_Q SO(3)^2$, let

$$p_Q(v) = e_z^T (R_1 v_1 + R_2 v_2). \quad (26)$$

We define the orthogonal projection of v onto $T_{HQ}SO(3)^2$ as

$$\Pi_H(v) = v - \frac{p_Q(v)}{2} \begin{bmatrix} R_1^T e_z \\ R_2^T e_z \end{bmatrix}. \quad (27)$$

5.3. Geodesics and the exponential map

The goal of this section is to show that the natural projection of geodesics in $SO(3)^2$ are geodesics in $\mathcal{M}_{\mathcal{A}}$. The key insight we will use is the following (see the additional material for a proof).

Proposition 5.3. *Let $Q(t) : \mathbb{R} \rightarrow SO(3)^2$ be a geodesic curve such that $\dot{Q}(t) \in T_{HQ}SO(3)$ for all t . Then, $\tilde{Q} = \pi_{\mathcal{M}_{\mathcal{A}}}(Q)$ is a geodesic curve in $\mathcal{M}_{\mathcal{A}}$.*

This result tells us that to find geodesics in $\mathcal{M}_{\mathcal{A}}$, we can focus on finding geodesics in $SO(3)^2$ for which the tangent vector is always horizontal. This idea is repeatedly used in [4] to give expressions for the geodesics in the Stiefel and Grassmann manifolds. Here we now show that if a geodesic $Q(t) \in SO(3)^2$ has a horizontal initial tangent vector $\dot{Q}(0)$, then the tangent is horizontal for every t .

Proposition 5.4. *Let $V \in TSO(3)^2$ be a vector field of the form*

$$W(t)^{\vee} = \text{stack}(R_1(t)^T e_z, R_2(t)^T e_z) \quad (28)$$

defined along a geodesic $Q(t) \in SO(3)^2$. Then we have

$$\langle \dot{Q}(t), W(t) \rangle = \langle \dot{Q}(0), W(0) \rangle. \quad (29)$$

Proof. Denote the tangent to the geodesic $Q(t)$ as $\dot{Q}(t)^{\vee} = \text{stack}(v_1, v_2)$, and let

$$m(t) = \langle \dot{Q}(t), W(t) \rangle = v_1^T R_1^T e_z + v_2^T R_2^T e_z. \quad (30)$$

Taking the derivative we have

$$\dot{m}(t) = v_1^T [v_1]_{\times}^T R_1^T w_1 + v_2^T [v_2]_{\times}^T R_2^T w_2 \equiv 0. \quad (31)$$

Since the first derivative of $m(t)$ is identically zero, $m(t)$ must be constant, which implies (29). \square

Combining Propositions 5.3 and 5.4, we get that the exponential map in $\mathcal{M}_{\mathcal{A}}$, i.e.,

$$[Q_b] = \exp_{[Q_a]}(v_a), \quad [Q_a] \in \mathcal{M}_{\mathcal{A}}, \quad v_a \in T_{[Q_a]}\mathcal{M}_{\mathcal{A}}, \quad (32)$$

is obtained by computing

$$Q_b = \exp_{Q_a}(\tilde{v}_a), \quad Q_a \in SO(3)^2 \quad (33)$$

where \tilde{v}_a is the horizontal lift of v_a .

5.4. The distance and the logarithm map

Let $Q_a = (R_{a1}, R_{a2})$ and $Q_b = (R_{b1}, R_{b2})$ be two points in $SO(3)^2$. We would like to find the distance between $[Q_a]$ and $[Q_b]$ and the logarithm map $\log_{[Q_a]}[Q_b]$. In general, we cannot directly use the distance and logarithm map in $SO(3)^2$, because the tangent of the corresponding geodesic is not horizontal. However, we can “move” Q_b to another representative of the equivalence class $[Q_b]$, so that the geodesic between Q_a and Q_b corresponds to a geodesic between $[Q_a]$ and $[Q_b]$. This is formalized in the following:

Proposition 5.5. *Define the cost*

$$f(t) = \sum_{i=1,2} f_i, \quad f_i = \frac{1}{2} \theta_i^2(t), \quad \theta_i(t) = d(R_{ai}, R_z(t)R_{bi}), \quad (34)$$

and let $t_{\text{opt}} = \text{argmin}_t f(t)$. Then, the logarithm

$$\log_{Q_a}(S_z(t_{\text{opt}})Q_b) = \text{stack}(\{\text{Log}(R_{ai}^T R_z(t_{\text{opt}})R_{bi})\}_{i=1,2}) \quad (35)$$

is an horizontal vector in $T_{HQ}SO(3)^2$.

Using (1) and the isometry given by horizontal lifts, the distance between the two elements in $\mathcal{M}_{\mathcal{A}}$ is then given by

$$d([Q_a], [Q_b]) = \|\log_{[Q_a]}[Q_b]\| = \|\log_{Q_a}(S_z(t_{\text{opt}})Q_b)\|. \quad (36)$$

Intuitively, this distance is the least amount of rotation needed to align two camera pairs (corresponding to two essential matrices) with a common baseline.

Proof of Proposition 5.5. We will need the following result

$$\frac{d}{dt} R_z(t)R_{bi} = [e_z]_{\times} R_{bi} = R_{bi} [R_{bi}^T e_z]_{\times}, \quad i = 1, 2. \quad (37)$$

Taking the derivative of each term f_i we have

$$\begin{aligned} \dot{f}_i(t) &= -\langle \log_{R_{ai}}(R_z(t)R_{bi}), R_{bi}(R_{bi}^T e_z)^{\wedge} \rangle \\ &= -\text{Log}(R_{ai}^T R_z(t)R_{bi})^T R_{bi}^T e_z \\ &= -\text{Log}(R_{ai}^T R_z(t)R_{bi})^T R_{ai}^T R_z(t)R_{bi} R_{bi}^T e_z \\ &= -e_z^T R_{ai} \text{Log}(R_{ai}^T R_z(t)R_{bi}), \end{aligned} \quad (38)$$

where we used the fact that $R^T \text{Log}(R) = \text{Log}(R)$ and, similarly, $R_z(t)e_z = e_z$. For $t = t_{\text{opt}}$ we have $\dot{f}_1(t_{\text{opt}}) + \dot{f}_2(t_{\text{opt}}) = 0$, which, together with (25), implies that the vector is in the horizontal space at Q_a . \square

The problem now is to find t_{opt} , the minimizer of f . In general, this is a nonlinear optimization problem with multiple local minima (see Figure 2 for an example). However, we can exploit its special structure to reliably and efficiently find the global minimizer t_{opt} . First, consider each function f_i separately. The derivative of f_i can be computed as

$$\dot{f}_i(t) = e_z^T R_{ai} \text{Log}(R_{ai}^T R_z(t) R_{bi}) = \theta_i(t) e_z^T R_{ai} u_i \quad (39)$$

where (using the closed form expression of Log from [11])

$$u_i = \frac{1}{2 \sin \theta_i(t)} [(R_{ai}^T R_z(t) R_{bi}) - (R_{ai}^T R_z(t) R_{bi})^T]_{\times}^{\text{inv}}, \quad (40)$$

Notice that the derivative of f exists everywhere except at a point t_{di} for which $\sin(\theta_i(t_{di})) = 0$. The following proposition gives a way to compute the location of this point.

Proposition 5.6. *Let θ_i be defined as in (34), and define*

$$c_{1i} = (R_{bi} R_{ai}^T)_{1,1} + (R_{bi} R_{ai}^T)_{2,2} \quad (41)$$

$$c_{2i} = (R_{bi} R_{ai}^T)_{1,2} - (R_{bi} R_{ai}^T)_{2,1} \quad (42)$$

$$\phi_i = \arctan_2(c_{1i}, c_{2i}). \quad (43)$$

Then, the function $\theta_i(t)$ is continuous, 2π -periodic and

$$\sin(\theta_i(t_{di})) = 0 \text{ for } t_{di} = \frac{3}{2}\pi - \phi_i. \quad (44)$$

Using the definition of DLog and its closed-form expression from [17], the second derivative of f_i is given by:

$$\begin{aligned} \ddot{f}_i(t) &= e_z^T R_{ai} \text{DLog}(R_{ai}^T R_z(t) R_{bi}) R_{ai}^T e_z \\ &= (e_z^T R_{ai} u_i)^2 + \frac{\theta}{2} \cot\left(\frac{\theta}{2}\right) (1 - (e_z^T R_{ai} u_i)^2) \end{aligned} \quad (45)$$

Note that (as a simple plot can confirm)

$$0 \leq \frac{\theta}{2} \cot\left(\frac{\theta}{2}\right) \leq 1 \text{ for } \theta \in [-\pi, \pi]. \quad (46)$$

This implies that $\ddot{f} \geq 0$, and that f is convex between discontinuity points.

In summary, from the results above, we have that the function f is continuous, 2π -periodic and with positive second derivative except at $\{t_{di} + 2k\pi\}$, $k \in \mathbb{Z}$. Assuming (without loss of generality) the ordering $-\frac{\pi}{2} \leq t_{d1} \leq t_{d2} \leq \frac{\pi}{2}$, this suggests an algorithm to find all the global minimizers of f which considering separately the two intervals $[t_{d1}, t_{d2}]$ and $[t_{d2}, t_{d1} + 2\pi]$ (on which the function is convex and differentiable). Since we have a closed form expression for \dot{f} , we can use Newton's method (with the additional projection of the iterates to the interval). In addition, one can easily show (using the intermediate value's theorem on \dot{f}) that if \dot{f} has the same sign at the two extremum points of an interval, then

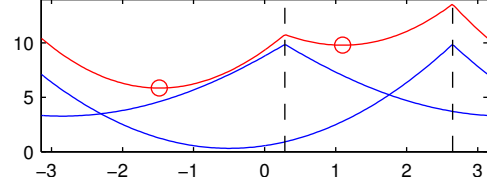


Figure 2: An example realization of the cost $f(t)$ from (34). Blue and red lines: value of each term f_i and of f , respectively. Black dashed line: location of the discontinuity points $\{t_{di}\}$ computed using Prop. 5.6. Red circles: local minimizers $\{t_{\text{opt},i}\}$ computed in Algorithm 1.

that interval does not contain a local minimizer, and it can be skipped. These steps are summarized in Algorithm 1 (see also Figure 2). We use the notation \dot{f}^+ and \dot{f}^- to denote right and left derivatives, respectively. Note that Algorithm 1 is only a basic version. A complete version would also consider degenerate cases, where $m_i = 0$ for some $i \in \{1, 2\}$ or where $t_{d1} = t_{d2}$. In our experiments, we saw that an interval could be skipped about 25% of the time, and that the Newton's iteration took about 5 to 8 iterations to converge to the machine's precision of $2 \cdot 10^{-16}$ (as a comparison, the method suggested in [15] only achieves a precision of 10^{-4} after about 5 iterations).

5.5. Comparison with previous formulations

Among the papers that use the relative pose between cameras to parametrize the essential space, the definition of normalized essential space used in [12] is compatible with \mathcal{M}_E , while the definition used in [13] (which includes the cheirality constraints explicitly) is compatible \mathcal{M}_G . For the papers using the parametrization derived from the SVD [5, 9, 15, 16], the definition used is the same as \mathcal{M}_G . However, these papers (mistakenly) do not consider the action of the group H_π and the cheirality constraint. In particular, Prop. 4.2 shows that the claim made in [15] that an essential matrix E corresponds uniquely to a point in \mathcal{M}_G is false.

Algorithm 1 Global minimization of $f(t)$

- 1: Compute the points t_{di} , $i = 1, 2$ (assume $t_{d1} < t_{d2}$).
 - 2: Define intervals $S_1 = [t_{d1}, t_{d2}]$ and $S_2 = [t_{d2}, t_{d1} + 2\pi]$.
 - 3: **for** $i \in 1, 2$ **do**
 - 4: **if** $\text{sgn}(\dot{f}^+(\min(S_i))) \neq \text{sgn}(\dot{f}^-(\max(S_i)))$ **then**
 - 5: Compute $t_{\text{opt},i} = \text{argmin}_{t \in S_i} f(t)$ using the projected Newton's method.
 - 6: **end if**
 - 7: **end for**
 - 8: Select t_{opt} as the point $t_{\text{opt},i}$ for which f is minimum.
-

6. The Weiszfeld algorithm and pose averaging

In principle, the Riemannian framework established in this paper could be used in any optimization problem involving essential matrices (e.g., through any of the algorithms given in [1]). Here, however, we show that the distance obtained in §5.4 can be naturally and meaningfully used in a proof-of-concept application to the two-view structure from motion problem.

In a standard pipeline, the relative pose (R, T) between two calibrated views is computed using RANSAC (see [8]):

- Extract pairs of matching image points $\{x_1^i, x_2^j\} \in \mathbb{R}^2$.
- For $i \in \{1, \dots, N\}$, select a random subset S_i of point pairs $\{x_1^i, x_2^j\}_{j \in S_i}$, estimate the essential matrix E_i and compute its support (i.e., the number of points which approximately satisfy the epipolar constraint).
- Compute the pose (R, T) from the matrix E_i with the largest support.

In [6] an alternative approach is suggested, where instead of using RANSAC, each sample E_i is decomposed into a pose estimate (R_i, T_i) and then all the rotations $\{R_i\}$ are averaged. Toward this, they propose to minimize the cost

$$\varphi(R) = \sum_i d(R, R_i)^p, \quad (47)$$

where $p = 1$ (L1 averaging) or $p = 2$ (L2 averaging), by using the Weiszfeld algorithm, which we report in Algorithm 2 for points lying in a general Riemannian manifold \mathcal{M} . Strictly speaking, the traditional Weiszfeld algorithm refers only to the version $p = 1$, but here we give a generalized version for ease of exposition. The set I in (49) is used to take into account the fact that w_i becomes ill-defined when $p = 1$ and the iterate x falls on one of the input points. Intuitively, each iteration of the algorithm maps the input points to the tangent space of the current iterate $x(t)$, take the average (with weights given by the relative distances) and use the resulting vector to obtain the next iterate $x(t+1)$.

In this section, we follow the same approach proposed by [6], but we average essential matrices instead of rotations.

Algorithm 2 The Weiszfeld algorithm

Input: Points $x_i \in \mathcal{M}$, $i \in \{1, \dots, N\}$.

- 1: Initialize $x(0)$
- 2: **for** $t \in \{0, \dots, N_t\}$ **do**
- 3: Update x using:

$$w_i(t) = d(x(t), x_i)^{p-2} \quad (48)$$

$$I(t) = \{i \in \{1, \dots, N\} : x(t) \neq x_i\} \quad (49)$$

$$x(t+1) = \exp_x \left(\frac{\sum_{i \in I} w_i(t) \log_x(x_i)}{\sum_{i \in I} w_i(t)} \right) \quad (50)$$

- 4: **end for**
-

In practice, the only difference is the use of the definition of \exp , \log and Riemannian distance for $\mathcal{M}_{\mathcal{G}}$ in Algorithm 2. Note that the approach proposed here has the immediate advantage of naturally considering both rotation and translation components together, while the approach of [6] considers only rotations. We compare the two approaches against standard RANSAC on the `fountain-P11` dataset from [14], which includes the ground-truth pose for the cameras.

We used SIFT features extraction and matching [18] to find corresponding points between every possible pair of cameras. We excluded image pairs with less than thirty good matches (as determined using the essential matrix from the ground truth pose). We then use the five point algorithm [7] to generate the RANSAC samples E_i . We compare four versions of the Weiszfeld algorithm corresponding to the four possible combinations of $p = 1, 2$ and $\mathcal{M} = \mathcal{M}_{\mathcal{G}}, SO(3)$ by using between 1 and 50 RANSAC samples. To initialize the algorithm, we evaluate the cost at every input sample, and use the half-way point between the two samples with lower costs. Also, we set the number of iterations N_t to 30 (although, in our preliminary tests, the algorithms usually converged in less than 15 iterations). As baseline, we use the errors of the RANSAC solution after the same number of samples and after 2000 samples. As the error measure, we consider the geodesic distance between estimated and ground-truth rotations. For our approach and the RANSAC-based solutions, we also consider the angle between the estimated translation direction and the ground truth. All the results are averaged across all the image pairs and 30 independent sampling realizations.

We report the results in Figure 3. As one can see, the Weiszfeld algorithm using the proposed distance on $\mathcal{M}_{\mathcal{G}}$ outperforms the corresponding version using the distance on $SO(3)$ for both $p = 1$ and $p = 2$. Moreover, the use of the cost with $p = 1$ produces better results than those using $p = 2$, likely due to the fact that the first cost is more robust to outliers in the samples. This dataset also shows that, while our approach (which does not require setting a threshold) gives reasonably good results, the efficiency of RANSAC with a well-tuned threshold is quite hard to beat.

7. Conclusion

In this paper we considered a Riemannian structure for the essential manifold, and introduced a novel, geometrical interpretation which shed light on the limitations of previous approaches and on the connections with traditional concepts in computer vision. We also proposed efficient algorithms for computing the distance and logarithm map, and considered an application to the problem of two-view pose estimation using averages. In our future work we will investigate relations between three views, and determine if similar ideas can be applied to the space of trifocal tensors and other similar objects.

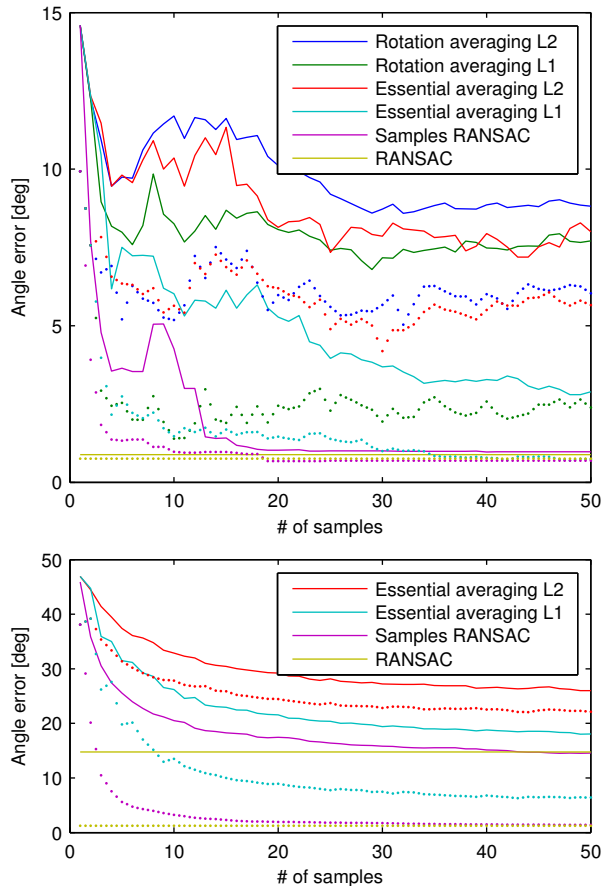


Figure 3: Results for two-view pose estimation. Rotation (top) and translation angle errors (bottom) for the different methods on the fountain-P11 dataset. Solid and dotted lines represent the mean and median errors, respectively.

References

- [1] P.-A. Absil, R. Mahony, and R. Sepulchre. *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, Princeton, NJ, 2008. 7
- [2] M. Arie-Nachimson, S. Kovalsky, I. Kemelmacher-Shlizerman, A. Singer, and R. Basri. Global motion estimation from point matches. In *International Conference on 3D Imaging, Modeling, Processing, Visualization and Transmission*, pages 81–88, 2012. 2
- [3] M. P. do Carmo. *Riemannian geometry*. Birkhäuser, Boston, MA, 1992. 1
- [4] A. Edelman, T. A. Arias, and S. T. Smith. The geometry of algorithms with orthogonality constraints. *SIAM Journal on Matrix Analysis and Applications*, 20(2):303–353, 1998. 5
- [5] C. Geyer and K. Daniilidis. Mirrors in motion: Epipolar geometry and motion estimation. In *IEEE International Conference on Computer Vision*, pages 766–773, 2003. 1, 6
- [6] R. Hartley, K. Aftab, and J. Trunpf. L1 rotation averaging using the Weiszfeld algorithm. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2011. 7
- [7] R. Hartley and H. Li. An efficient hidden variable approach to minimal-case camera motion estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2012. 7
- [8] R. I. Hartley and A. Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, second edition, 2004. 1, 7
- [9] U. Helmke, K. Hüper, P. Y. Lee, and J. Moore. Essential matrix estimation using gauss-newton iterations on a manifold. *International Journal of Computer Vision*, 74(2):117–136, 2007. 1, 6
- [10] H. C. Longuet-Higgins. A computer algorithm for reconstructing a scene from two projections. *Nature*, 293:133–135, 1981. 1
- [11] Y. Ma. *An invitation to 3-D vision: from images to geometric models*. Springer, 2004. 1, 2, 4, 6
- [12] Y. Ma, J. Košecská, and S. Sastry. Optimization criteria and geometric algorithms for motion and structure estimation. *International Journal of Computer Vision*, 44(3):219–249, 2001. 1, 6
- [13] S. Soatto, R. Frezza, and P. Perona. Motion estimation via dynamic vision. *IEEE Transactions on Automatic Control*, 41(3):393–413, 1996. 1, 6
- [14] C. Strecha, W. von Hansen, L. V. Gool, P. Fua, and U. Thoennessen. On benchmarking camera calibration and multi-view stereo for high resolution imagery. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, 2008. 7
- [15] R. Subbarao, Y. Genc, and P. Meer. Robust unambiguous parametrization of the essential manifold. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2008. 1, 3, 6
- [16] R. Subbarao and P. Meer. Nonlinear mean shift over riemannian manifolds. *International Journal of Computer Vision*, 84(1):1–20, 2009. 1, 6
- [17] R. Tron. *Distributed optimization on manifolds for consensus algorithms and camera network localization*. PhD thesis, The Johns Hopkins University, 2012. 6
- [18] A. Vedaldi and B. Fulkerson. VLFeat: An open and portable library of computer vision algorithms. <http://www.vlfeat.org/>, 2008. 7