



# A Geometric Approach to Study Aircraft Trajectories: the benefits of OpenSky Network ADS-B data

Rémi Perrichon, Thierry Klein, Xavier Gendre

## ► To cite this version:

Rémi Perrichon, Thierry Klein, Xavier Gendre. A Geometric Approach to Study Aircraft Trajectories: the benefits of OpenSky Network ADS-B data. 2022. hal-03848832v1

**HAL Id: hal-03848832**

**<https://enac.hal.science/hal-03848832v1>**

Preprint submitted on 11 Nov 2022 (v1), last revised 1 Jan 2023 (v2)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# A Geometric Approach to Study Aircraft Trajectories: the benefits of OpenSky Network ADS-B data

Rémi Perrichon <sup>1,\*</sup> , Thierry Klein <sup>1,2</sup> and Xavier Gendreau <sup>3,2</sup>

<sup>1</sup> ENAC - Ecole Nationale de l'Aviation Civile. Université de Toulouse, France

<sup>2</sup> Institut de Mathématiques de Toulouse (UMR5219). Université de Toulouse, France

<sup>3</sup> ISAE SUPAERO - Institut supérieur de l'aéronautique et de l'espace. Université de Toulouse, France

\* Correspondence: remi.perrichon@enac.fr

**Abstract:** To date, the statistical analysis of aircraft trajectories has been under-exploited in the Airspace Traffic Management (ATM) literature. One reason is the need for advanced methods to tackle the high sampling irregularity and temporal correlation that both characterize a trajectory. Differential geometry provides a relevant framework to study trajectories. Modeling trajectories as parametrized curves, shape analysis allows to answer operational questions. This work presents a geodesic distance that rigorously defines and quantifies shape differences between aircraft trajectories. The key idea is to compare how the shape of a given trajectory changes from one popular data set (the Eurocontrol R&D data archive) to another one (OpenSky Network ADS-B data). Distances as well as geodesic paths are computed for a sample of flights departing from Toulouse-Blagnac (LFBO) and landing at Paris-Orly (LFPO) in 2019. Its use for clustering purposes is illustrated and discussed.

**Keywords:** air traffic management; trajectory; geometry; clustering; elastic distance

## 1. Introduction

Geometry is widely used in aviation from the designing of airfoil to the one of procedures. For its part, statistics has been playing a key role for traffic forecast and predictive maintenance. Yet, contrary to geometry, the statistical analysis of aircraft trajectories has been under-exploited in the Airspace Traffic Management (ATM) literature. One possible reason is the irrelevance of usual statistical frameworks to model a trajectory. Indeed, techniques from multivariate statistics suffer from the very high correlation in time that exists between two consecutive points of a trajectory. More worryingly, when observation times are more numerous than the number of trajectories, the so-called curse of dimensionality occurs. The usual time series approach is not ideal either as trajectories can be seen as multidimensional time series with both irregular sampling and different durations.

The current institutional context evolves towards a greater availability of trajectory data. On the one hand, Eurocontrol has given access to a R&D data archive containing more than 17 million flights as of April 2022. On the other hand, the non-profit OpenSky Network has grown to 5,000 registered receivers all around the world, providing a large historical database of ADS-B data that is accessible to researchers [1].

Interestingly, past OpenSky Network symposia have demonstrated how statistics could provide some solutions to complex problems such as trajectory generation [2] or pattern identification [3], encouraging the development of new statistical methods to model trajectories.

An adapted statistical framework to study trajectories is the one of Functional Data Analysis (FDA), popularised by [4] and [5]. The promotion of FDA to study aircraft trajectories was early made by [6] and has been successfully used for Functional Principal Component Analysis (FPCA) carried out by [7] and applied to the detection of atypical energy behaviours by [8].

Recently, FDA has progressively made the most of differential geometry to study any kind of trajectories. As an example, [9] analyzed bird and hurricane trajectories on Riemannian manifolds. Incorporating geometry in statistics is particularly relevant for at least two reasons. First, parametric curves come as naturally good descriptors of

**Citation:** Perrichon, R.; Klein, T.; Gendreau, X. A Geometric Approach to Study Aircraft Trajectories. *Preprints* 2022, 1, 0. <https://doi.org/>

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Copyright:** © 2022 by the authors. Submitted to *Preprints* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

motion. It makes geometry a perfect framework to describe a dynamic system evolving over time. Second, geometry allows to carefully study the shape of a trajectory under arbitrary temporal evolution, taking into account translation, scale and rotational effects.

This work introduces a number of geometrical methods to quantify how the shape of a given trajectory changes from the R&D archive version to the ADS-B one. Additional to providing a proper shape-based distance between the two versions, we compute the optimal deformation required to go from the R&D archive version to the ADS-B one. In other words, a geodesic path allows to visualize all the bending and stretching transformations that are necessary to go from one version to the other.

Far from being purely theoretical, the presented geometric framework has many operational applications. In this work, the geodesic distance enables to both visualize and quantify where and how much the flown trajectory has diverted from the flight plan. The geodesic distance is also used to solve a clustering problem based on shape.

This paper first introduces the geometric framework of shape analysis, departing from the early work of Kendall and co-authors to more recent methods. Second, the geodesic distance is defined. Third, a sample of flights is used to illustrate how the R&D archive version of a trajectory (Eurocontrol) changes from its ADS-B version (OpenSky Network).

## 2. What is the shape of a trajectory?

Whatever the data provider (Eurocontrol or OpenSky Network), a raw trajectory  $i$ , denoted  $\text{traj}_i$ , is a set of  $m_i$  pairs:

$$\text{traj}_i = \{(\mathbf{y}_{i,j}, t_{i,j}), j = \{1, \dots, m_i\}\}$$

$m_i$  being the number of observation times associated to flight  $i$ ,  $\mathbf{y}_{i,j}$  being a three-dimensional vector (longitude, latitude, altitude),  $t_{i,j}$  being a timestamp. In what follows, the components of  $\mathbf{y}_{i,j}$  are respectively denoted  $y_{i,j,x}$ ,  $y_{i,j,y}$ ,  $y_{i,j,z}$ . How to define the shape of such object? What do we mean by shape in the first place?

One of the earliest works in statistical analysis and modeling of shapes of objects came from Kendall and colleagues [10], [11]. The adopted definition of a shape is still commonly used in geometry and statistics. As defined by [12], shape is all the geometrical information that remains when location, scale and rotational effects are removed from an object. Two objects have the same shape if they can be translated, rescaled and rotated to each other so that they match exactly.

Whatever the nature of objects under study (images, sounds, curves, surfaces), shape was originally described by locating a finite number of points on each object. These so-called landmarks are points of correspondence and can be assigned by an expert (scientific landmark) or suggested by a property of the object such as a point of high curvature (mathematical landmark). The positions of these landmarks are points in  $\mathbb{R}^n$ . It is often assumed that the number of landmarks  $k$  is greater or equal than the dimension  $n$ . The configuration is the set of landmarks on a particular object. The configuration matrix  $X$  is the  $n \times k$  matrix of Cartesian coordinates of the  $k$  landmarks in  $n$  dimensions. The configuration space is the space of all landmark coordinates, usually the space of real  $n \times k$  matrices with possibly some special cases removed, such as coincident points. This space is denoted  $\mathcal{L}_{n,k}$ . Within a same shape class, many configurations are possible. Indeed, two objects have the same shape if their configurations are equivalent. These variations are formulated as actions of certain groups. In shape analysis, usual variations are all actions of Lie groups, that is to say mathematical objects that are both a group and a differential manifold. Lie groups of interest are the translation group  $\mathbb{R}^n$ , the scaling group  $\mathbb{R}^\times$  and the rotation group  $SO(n)$ . The three variations are combined together considering a product group denoted  $\mathbb{R}^n \rtimes (\mathbb{R}^\times \times SO(n))$  where  $\rtimes$  is the semi-direct product.

The action of  $\mathbb{R}^n \rtimes (\mathbb{R}^\times \times SO(n))$  on  $\mathcal{L}_{n,k}$  is  $(v, a, O) * X = aOX + v\mathbf{1}'_k$ , where  $\mathbf{1}_k$  is a column vector of length  $k$  containing all ones so that  $v\mathbf{1}'_k$  is an  $n \times k$  matrix with identical columns. More specifically, landmark space analysis relies on the study of spaces of equivalence classes that are called quotient spaces in differential geometry. The quotient

space of interest is  $\mathcal{L}_{n,k}/\mathbb{R}^n \rtimes (\mathbb{R}^\times \times SO(n))$  and is termed landmark shape space. In quotient space live orbits. In differential geometry, for any element  $p$  of a manifold  $M$ , the orbit of  $p$  under the action of a group  $G$  is defined as the set  $G \cdot p = \{g \cdot p : g \in G\}$  and is written  $[p]$ . The orbit of a point in a manifold refers to all possible points one can reach in the manifold using the action the group on that point.

Landmark shape analysis has led to many practical applications. In biology, it has been used to quantify the effects of selection for body weight on the shape of mouse vertebrae [13], and to highlight sexual dimorphism in hominoids using craniofacial shape differences [14]. In chemistry, [15] and [16] have analyzed a dataset of steroids to evidence how the shape ("steric") properties of the molecules are related to an activity class.

Despite these successful applications, the use of landmarks is not straightforward to study aircraft trajectories for at least three reasons.

1. The choice of landmark is very subjective. How many landmarks are relevant to summarize the shape of an aircraft trajectory? Should this number depend on the origin-destination we consider? Would landmarks based on flight phases be enough to capture the shape of a trajectory?
2. Scientific landmarks are not available in raw data. Both Eurocontrol and OpenSky Network data sources do not include expert knowledge.
3. Mathematical landmarks are likely to be imprecise for Eurocontrol data. Indeed, recovering flight phases thanks to existing algorithms ([17] or [18]) will perform differently depending on the number of observation points.

These limits are not new. They were exemplified in medical imaging problems for which landmarks are not obvious to find (think of soft tissues where boundaries have no sharp edges). To account for these problems but in the same spirit as Kendall's formulation, [19] was one of the first to propose a convenient shape representation of curves in  $\mathbb{R}^n$ . This is the framework we consider in the sequel.

### 3. Trajectories as parameterized curves in $\mathbb{R}^3$

As said, raw data are stored in a usual matrix format. A few preliminary remarks should be made. First, the duration changes from one flight to another. As a consequence, the number of observation points is not the same from a trajectory to another. Second, for a given trajectory, observation times are not equally spaced. Third, between two trajectories, the spacing is irregular but is never the same. For a given trajectory, observation points are the same for all components of  $y_{i,j}$ .

Given that most methods in time series and multivariate statistics are inappropriate in this context, a relevant statistical framework is to study trajectories from a Functional Data Analysis (FDA) perspective. The main assumption in FDA is the existence of an underlying function that has given birth to a set of values the statistician observes on the aforementioned irregular grid. When the number of observation times is high enough, each underlying function is retrieved thanks to interpolation or smoothing. It is then effortless to resample trajectories on a regular grid.

In FDA, a trajectory is viewed as the realisation of a functional random vector. A functional random vector of dimension  $p$  is denoted  $X \equiv (X^1, X^2, \dots, X^p)'$ . Note that  $\forall k \in \{1, \dots, p\}, X^k \equiv (X_t^k)_{t \in [0,1]}$ . More precisely, each component  $X^k$  is a functional random variable taking its values in an infinite dimensional vector space, often chosen to be  $\mathbb{L}^2([0,1], \mathbb{R})$  associated to the scalar product  $\langle f, g \rangle_{\mathbb{L}^2} = \int_0^1 f(t)g(t)dt$  for  $f, g \in \mathbb{L}^2([0,1], \mathbb{R})$ .

As the three dimensions of a trajectory (longitude, latitude, altitude), are the one of an underlying coordinate system, once reconstructed, a trajectory can be viewed as an open parameterised curve  $\beta$  of  $\mathbb{R}^3$ .

As each trajectory has a different duration, a scaling is made so that each curve  $\beta_i$  goes from  $[0,1]$  to  $\mathbb{R}^3$ .

### 3.1. Interpolation: from raw data to curves

Going from raw data to smooth curves calls for a smoothing or interpolation step. Smoothing refers to the case in which observations are assumed to be noisy whereas interpolation hypotheses that observations are recorded with negligible errors. Many smoothing and interpolation methods are available to reconstruct individual trajectories. A review of four popular non-parametric techniques is given by [20]. Whatever the method, the construction of functional observations using the discrete data takes place separately for each dimension of a flight: the longitude, the latitude, the altitude.

Linear interpolation is chosen in this work because Eurocontrol R&D data are based on updated filled flight plans, that is to say a series of straight lines.

### 3.2. The Square-Root Velocity Function (SRVF) captures shape

Once trajectories are interpolated, their shape can be studied. To this end, the Square-Root Velocity framework has been developed in [19]. The main assumption is that curves are differentiable with a first derivative in  $\mathbb{L}^2([0, 1], \mathbb{R}^n)$ . In this framework, the shape of a curve is no more captured by landmarks but by the Square-Root Velocity Function (SRVF). Let  $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be a mapping given by:

$$F(v) = \begin{cases} \frac{v}{\sqrt{|v|}} & \text{if } |v| \neq 0 \\ 0 & \text{if } |v| = 0 \end{cases}$$

where  $|\cdot|$  denotes the usual Euclidean norm of  $\mathbb{R}^n$ . The SRVF of  $\beta$  is defined to be  $q(t) \equiv F(\dot{\beta}(t))$  where  $\dot{\beta}$  is the time derivative of  $\beta$ . The curve  $\beta$  can be reconstructed from  $q$  up to a translation. Luckily, the SRVF naturally removes translation effects as it is based on the curve derivative.

Scale effects are easy to remove as all curves can be rescaled to be of unit length in a pre-processing step.

The space of interest is then

$$\mathcal{C} \equiv \left\{ q \in \mathbb{L}^2([0, 1], \mathbb{R}^n), \int_{[0,1]} |q(t)|^2 dt = 1 \right\}.$$

It is the set of all Square-Root Velocity Functions with unit  $\mathbb{L}^2$ -norm. As recalled in [21], this corresponds to a sphere in  $\mathbb{L}^2([0, 1], \mathbb{R}^n)$ . A lot is known about the geometry of such a sphere, including geodesics. Shape analysis of curves under this representation (without additional constraints) is relatively simple. Yet, we still need to take care of different re-parameterizations of curves as well as the effect of rotations.

### 3.3. How to deal with rotation and re-parameterization?

The action of the rotation group  $SO(n)$  is usual and similar in landmark-based shape analysis:  $SO(n) \times \mathcal{C} \rightarrow \mathcal{C}, (O, q(t)) = Oq(t)$ .

Yet, parameterization effects are specific to the Square-Root Velocity (SRV) framework because there is no parametrization in the landmark-based approach. Note that it is natural that shape should be invariant by re-parametrization: two curves may have the same shape but be traveled along at different speeds.

In the SRV framework, a re-parameterization is an element of  $\Gamma_{[0,1]}$ , the set of all orientation-perserving diffeomorphisms of  $[0, 1]$ . Orientation-perserving means that the transformed passage of time is still coherent (time is not moving back, 0 maps to 0 and 1 maps to 1). Working with diffeomorphisms allows to actually re-parameterize curves. For a  $\gamma \in \Gamma_{[0,1]}$ , the composition  $\beta \circ \gamma$  denotes its re-parameterization. It can be shown that the associated group action is  $\Gamma_{[0,1]} \times \mathcal{C} \rightarrow \mathcal{C}, (q, \gamma) = (q \circ \gamma) \sqrt{\dot{\gamma}}$ .

The quotient space of interest is  $\mathcal{S} \equiv \mathcal{C} / (\Gamma_{[0,1]} \times SO(n))$ . Crucially, [19] have shown that if the Riemannian metric is well-chosen to be an elastic Riemannian metric, the rotation group  $SO(n)$  and the re-parameterization group  $\Gamma_{[0,1]}$  both act by isometries. This property

is key to define distances between orbits, that is to say distances between elements of the quotient space.  $\mathcal{S}$  is called the shape space and is a metric space with the distance inherited from  $\mathcal{C}$ . It is a collection of orbits (individual shapes):  $\mathcal{S} \equiv \{[q] : q \in \mathcal{C}\}$ .

### 3.4. Distance between two shapes, geodesic path

The distance between two orbits  $[q_1]$  and  $[q_2]$  is given by:

$$d_{\mathcal{S}}([q_1], [q_2]) = \inf_{(\gamma \times O) \in \Gamma_{[0,1]} \times SO(n)} d_{\mathcal{C}}(q_1, \sqrt{\gamma} O(q_2 \circ \gamma))$$

with  $d_{\mathcal{C}}(a_1, a_2) = \cos^{-1} \left( \int_0^1 \langle a_1(t), a_2(t) \rangle dt \right)$ , where  $\langle \cdot, \cdot \rangle$  is the usual inner product between vectors in  $\mathbb{R}^n$ .

The actual geodesic between  $[q_1]$  and  $[q_2]$  in  $\mathcal{S}$  is given by  $[\alpha(\tau)]$ , where  $\alpha(\tau)$  is the geodesic in  $\mathcal{C}$  between  $q_1$  and  $\sqrt{\gamma^*} O^*(q_2 \circ \gamma^*)$ . Computing  $\alpha(\tau)$  is very easy because  $\mathcal{C}$  is a sphere.  $O^*$  and  $\gamma^*$  minimize the right side of the above distance. This is a joint optimization problem with a well-known solution as pinpointed in [21].

## 4. Application

### 4.1. Data

#### 4.1.1. R&D Eurocontrol

Eurocontrol is an international organisation working to achieve safe and seamless air traffic management across Europe. Since 2020, Eurocontrol has given access to a R&D data archive containing more than 17 million flights as of April 2022. The data are collected from all commercial flights operating in and over Europe. To be more specific, Eurocontrol receives flight plans for all Instrument Flight Rules (IFR) flights. These flight plans are activated and updated based on live data from air navigation service providers. Data are available for 4 months each year: March, June, September and December. Only two data subsets are used in this work: flight metadata and actual flight points. We are interested in studying the 15,628 flights that have been departing from Toulouse-Blagnac (LFBO) and landing at Paris-Orly (LFPO) from 2015 to March 2020.

In the R&D data archive, all flights are identified by a unique code coming from the PRISME Data Warehouse (DWH). The registration number is also available. There are 386 unique registration numbers for the 15,628 flights.

The actual off-block time as well as the actual arrival time are respectively corresponding to the first and the last point of a trajectory.

#### 4.1.2. Retrieving ADS-B flights thanks to OpenSky Network

To query a flight in OpenSky Network database (Impala Shell), one should use an ICAO24 identifier and a date. The ICAO24 code is a 24-bit unique number that is assigned to each vehicle or object that can transmit ADS-B messages. The correspondence between the registration number and the ICAO24 code is made thanks to the aircraft database of OpenSky Network. Among 15,628 flights, only 999 cannot be matched.

To find a flight in the OpenSky Network flight table, the query is made using the actual off-block time and the actual arrival time additional to a time buffer. A query outputs several candidate flights.

The proportion of flights for which there is at least one candidate goes from 0% in 2015 (it was expected as ADS-B data were not collected by OpenSky Network in 2015) to 85% in 2018. For each Eurocontrol flight, the ADS-B flight that is most likely (same ICAO code and closest off-block time) is kept. At this step, 6,053 flights are remaining. State vectors are then queried with a time buffer. Thus, 6,031 flights are remaining.

A great challenge when matching ADS-B data and Eurocontrol data is the problem of the beginning/end of a flight. The two first points of any flight in the R&D archive are located at the same geographical coordinates (the ones of the departure airport), with zero altitude. The time between the first two points is most likely the actual taxi-out time, that



is to say the period between the actual off-block time and the actual take off time. As the taxi-out time is not visible with ADS-B data, the first point of any flight from the R&D archive was removed. Duration is computed as being the time difference between the first point and the last point of the flight.

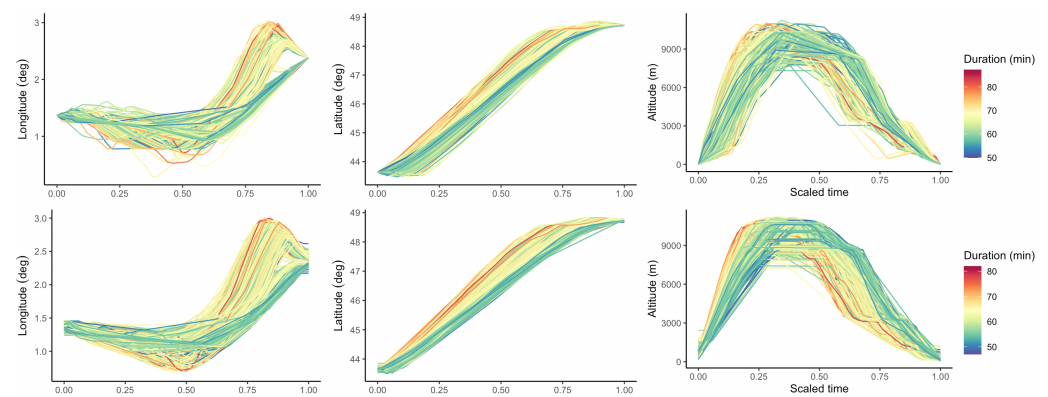
ADS-B data can be very noisy. Some basic filters are used to withdraw outliers. For instance, ADS-B flights with less than 20 points are withdrawn. ADS-B flights for which there is at least one time gap between two points that exceeds 20 minutes are withdrawn. 5,226 flights are remaining after the cleaning step.

In the spirit of the ADS-B cleaning step, Eurocontrol flights with less than 10 points are removed from the analysis as well as flights having at least a time gap between two consecutive points that is above 20 minutes. 5,180 trajectories are remaining after this step.

#### 4.2. Interpolation

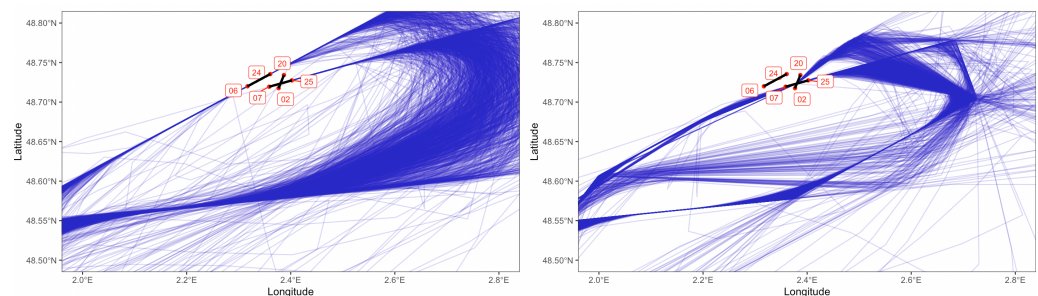
In total, there are 5,180 flights for which we have both a clean ADS-B version and a clean corresponding Eurocontrol version. For example in 2019, there are 1,746 valid trajectories.

Interpolation is done year by year. For both Eurocontrol and ADS-B data, the linear interpolation resamples all trajectories on a regular grid ranging from 0 to 1 with an incrementation of 0.02.



**Figure 1.** Linear interpolation of Eurocontrol trajectories [top] and ADS-B versions [bottom] for the flights departing from Toulouse-Blagnac (LFBO) and landing at Paris-Orly (LFPO) in 2019.

Contrary to the altitude profile of Eurocontrol trajectories, the take-off and landing of ADS-B trajectories are not always associated to a null altitude. This phenomenon and the small difference between the two color scales come from the aforementioned problem of finding the true beginning/end of an ADS-B flight.

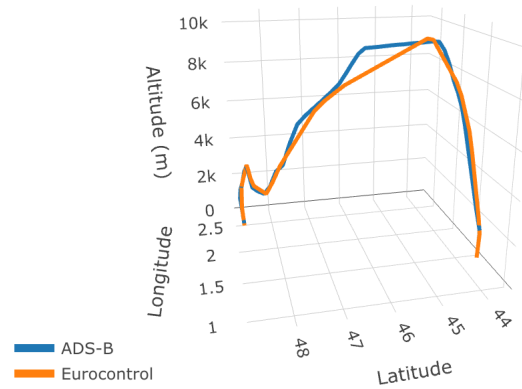


**Figure 2.** Zoom on Paris-Orly (LFPO) where runway displaced thresholds (DTHR) are indicated by red dots. Linear interpolation of ADS-B trajectories [left] and Eurocontrol versions [right] for the flights departing from Toulouse-Blagnac (LFBO) and landing at Paris-Orly (LFPO) in 2019.

As expected, Figure 2 shows a striking difference between the shape of trajectories from their ADS-B version to their Eurocontrol ones.

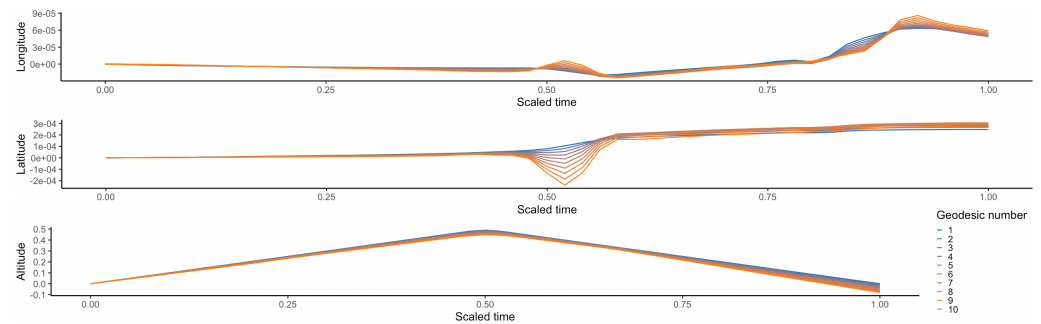
#### 4.3. How to compute the geodesic path in practice?

Most geometric concepts developed by [19] and [22] have been implemented in R, in the `fdasrvf` package maintained by J. Derek Tucker. The computation of the geodesic path is illustrated for a given flight that is shown on Figure 3.



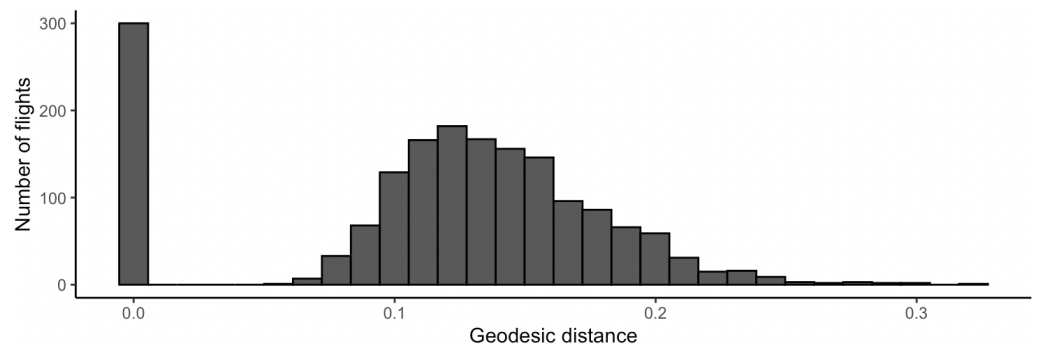
**Figure 3.** Linear interpolation of a given flight between Toulouse-Blagnac (LFBO) and Paris-Orly (LFPO) in 2019. Both its ADS-B and Eurocontrol versions are represented.

Once re-parametrization, location, scale and rotational effects have been taken into account, it is clear most bending and/or stretching is needed to match the ends of the highest flight level. Figure 4 shows that to optimally go from the Eurocontrol version to the ADS-B one, altitude must be slowly increased. A rise in latitude is needed in the middle of the flight which can be seen from Figure 3.



**Figure 4.** One possible representation of the shortest path from the Eurocontrol version to the ADS-B version in the quotient space  $\mathcal{S}$  thanks to the computation of  $\hat{\alpha}(\tau)$  for 10 values of  $\tau$ .

#### 4.4. Geodesic distance



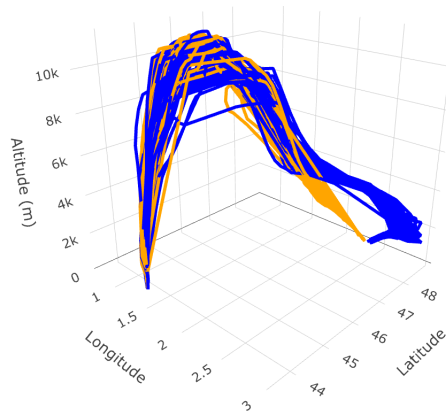
**Figure 5.** Histogram of the geodesic distances between Eurocontrol and ADS-B versions of the flights departing from Toulouse-Blagnac (LFBO) and landing at Paris-Orly (LFPO) in 2019.



The histogram in Figure 5 clearly shows two groups of flights. When the geodesic distance is null, no bending/stretching is needed to match trajectories once location, scale, rotational and re-parametrization effects are taken into account. In this case, Eurocontrol and OpenSky versions are carrying the same shape information. Yet, it is not the case when the geodesic distance is not null. In most cases, one should benefit from OpenSky Network ADS-B data, as there are more detailed.

#### 4.5. Hierarchical clustering

The geodesic distance may be used for the clustering and classification of trajectory shapes. [21] and [19] successfully cluster shapes of helices in  $\mathbb{R}^3$  by matching and deforming one into another, the underlying motivation being the analysis of protein structures. In a similar spirit, a clustering of ADS-B trajectories based on their shape can be performed. In this section, a random sample of 100 trajectories departing from Toulouse-Blagnac (LFBO) and landing at Paris-Orly (LFPO) in 2019 is drawn. A hierarchical clustering is implemented. The geodesic distance matrix is computed and the Ward's method is chosen to perform the clustering. When the quotient space we consider is taken to be  $\mathcal{S}$  (translation, scale, rotational and re-parametrization effects are removed), shape clusters are mostly based on the arrival direction at Paris-Orly airport (LFPO) as one may note from Figure 6.



**Figure 6.** Two clusters of trajectories (orange and blue), based on the geodesic distance. The sample is of size 100 and is randomly drawn from flights going from Toulouse-Blagnac (LFBO) to Paris-Orly (LFPO) in 2019.

## 5. Discussion

The growing availability of ADS-B data is likely to have a huge impact in statistical studies. As compared to the Eurocontrol R&D data archive that is based on updated filled flight plans, ADS-B data allow for a reliable study of trajectory shapes. Tools from differential geometry, namely the geodesic distance, help to quantify and clarify the empirical assessment that flown trajectory are *far* from the flight plan. This distance can be used in clustering problems.

This paper focused on a sample of flights that is not representative of all flights in Europe. Future works may quantify the benefits of ADS-B data for each flight phase and for a larger scope. The parsimonious sampling of Eurocontrol R&D data is likely to be satisfactory for shape analysis focusing on the en-route phase.

In this work, the ADS-B cleaning process is basic and should be refined in future works. Matching Eurocontrol flights with their ADS-B versions is far from being trivial and require an in-depth study.

Depending on operational needs, the quotient space defining what is meant by *shape* should be modified. An interesting extension could be to investigate for the relevant geometrical framework needed to model go-around patterns.

**Author Contributions:** Conceptualization, all; methodology, all; software, R.P.; validation, all; formal analysis, all; investigation, all; resources, all; data curation, R.P.; writing—original draft preparation, R.P.; writing—review and editing, all; visualization, R.P.; supervision, T.K. and X.G.; project administration, T.K. and X.G.; funding acquisition, T.K. and X.G. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research received no external funding.

**Data Availability Statement:** All data is freely available through the OpenSky Network <https://opensky-network.org> and the Eurocontrol website: <https://www.eurocontrol.int/dashboard/rnd-data-archive>.

**Conflicts of Interest:** The authors declare no conflict of interest.

## References

- Sun, J.; Basora, L.; Olive, X.; Strohmeier, M.; Schafer, M.; Martinovic, I.; Lenders, V. OpenSky Report 2022: Evaluating Aviation Emissions Using Crowdsourced Open Flight Data. 2022, p. 8.
- Krauth, T.; Morio, J.; Olive, X.; Figuet, B.; Monstein, R. Synthetic Aircraft Trajectories Generated with Multivariate Density Models. *Engineering Proceedings* **2021**, *13*, 7. <https://doi.org/10.3390/engproc2021013007>.
- Olive, X.; Sun, J.; Lafage, A.; Basora, L. Detecting Events in Aircraft Trajectories: Rule-Based and Data-Driven Approaches. *Proceedings* **2020**, *59*, 8. <https://doi.org/10.3390/proceedings2020059008>.
- Ramsay, J.O.; Silverman, B.W. *Functional data analysis*, 2nd ed.; Springer series in statistics, Springer-Verlag New York Inc., 2005.
- Ferraty, F.; Vieu, P. Nonparametric Functional Data Analysis: Theory and Practice **2006**. <https://doi.org/10.1007/0-387-36620-2>.
- Puechmorel, S.; Delahaye, D. 4D trajectories : a functional data perspective. *IEEE*, 2007, p. pp 1.C.6. <https://doi.org/10.1109/DASC.2007.4391832>.
- Nicol, F. Statistical Analysis of Aircraft Trajectories: a Functional Data Analysis Approach. Alldata 2017, The Third International Conference on Big Data, Small Data, Linked Data and Open Data, 2017, p. pp.51.
- Jarry, G.; Delahaye, D.; Nicol, F.; Feron, E. Aircraft atypical approach detection using functional principal component analysis. *Journal of Air Transport Management* **2020**, *84*, 101787. <https://doi.org/10.1016/j.jairtraman.2020.101787>.
- Su, J.; Kurtek, S.; Klassen, E.; Srivastava, A. Statistical analysis of trajectories on Riemannian manifolds: Bird migration, hurricane tracking and video surveillance. *The Annals of Applied Statistics* **2014**, *8*. <https://doi.org/10.1214/13-AOAS701>.
- Kendall, D.G. Shape Manifolds, Procrustean Metrics, and Complex Projective Spaces. *Bulletin of the London Mathematical Society* **1984**, *16*, 81–121. <https://doi.org/10.1112/blms/16.2.81>.
- Dryden, I.L.; Mardia, K.V. *Statistical Shape Analysis*; Wiley, 1998.
- Kendall, D.G. The diffusion of shape. *Advances in Applied Probability* **1977**, *9*, 428–430. <https://doi.org/10.2307/1426091>.
- Mardia, K.V.; Dryden, I.L. The Statistical Analysis of Shape Data. *Biometrika* **1989**, *76*, 271–281. <https://doi.org/10.2307/2336660>.
- O'Higgins, P.; Dryden, I.L. Sexual dimorphism in hominoids: further studies of craniofacial shape differences in Pan, Gorilla and Pongo. *Journal of Human Evolution* **1993**, *24*, 183–205. <https://doi.org/10.1006/jhev.1993.1014>.
- Dryden, I.L.; Hirst, J.D.; Melville, J.L. Statistical Analysis of Unlabeled Point Sets: Comparing Molecules in Chemoinformatics. *Biometrics* **2007**, *63*, 237–251.
- Czogiel, I.; Dryden, I.L.; Brignell, C.J. Bayesian matching of unlabeled marked point sets using random fields, with an application to molecular alignment. *The Annals of Applied Statistics* **2011**, *5*, 2603–2629. <https://doi.org/10.1214/11-AOAS486>.
- Sun, J.; Ellerbroek, J.; Hoekstra, J.M. Large-Scale Flight Phase Identification from ADS-B Data Using Machine Learning Methods. *7th International Conference on Research in Air Transportation* **2016**.
- Liu, D.; Xiao, N.; Zhang, Y.; Peng, X. Unsupervised Flight Phase Recognition with Flight Data Clustering based on GMM. In *Proceedings of the 2020 IEEE International Instrumentation and Measurement Technology Conference (I2MTC)*, 2020, pp. 1–6. <https://doi.org/10.1109/I2MTC43012.2020.9128596>.
- Srivastava, A.; Klassen, E.; Joshi, S.H.; Jermyn, I.H. Shape Analysis of Elastic Curves in Euclidean Spaces. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **2011**, *33*, 1415–1428. <https://doi.org/10.1109/TPAMI.2010.184>.
- Zhang, J.T. *Analysis of Variance for Functional Data*, 1st ed.; Chapman and Hall/CRC, 2013.
- Srivastava, A.; Klassen, E.P. *Functional and Shape Data Analysis*, 1st ed.; Springer-Verlag New York Inc., 2016.
- Tucker, J.D.; Wu, W.; Srivastava, A. Generative models for functional data using phase and amplitude separation. *Computational Statistics & Data Analysis* **2013**, *61*, 50–66. <https://doi.org/10.1016/j.csda.2012.12.001>.