



INSTITUT NATIONAL DE RECHERCHE EN INFORMATIQUE ET EN AUTOMATIQUE

*Probabilities and Statistics on Riemannian  
Manifolds:  
A Geometric approach*

Xavier Pennec

**N° 5093**

January 2004

THÈME 3



*R*apport  
de recherche



# Probabilities and Statistics on Riemannian Manifolds: A Geometric approach

Xavier Pennec

Thème 3 — Interaction homme-machine,  
images, données, connaissances  
Projet Epidaure

Rapport de recherche n° 5093 — January 2004 — 49 pages

**Abstract:** Measurements of geometric primitives are often noisy in real applications and we need to use statistics either to reduce the uncertainty (estimation), to compare measurements, or to test hypotheses. Unfortunately, geometric primitives often belong to manifolds that are not vector spaces. In previous works [Pennec, 1996, Pennec and Ayache, 1998], we used invariance requirements to develop some basic probability tools on transformation groups and homogeneous manifolds that avoids paradoxes.

In this paper, we consider the Riemannian metric as the basic structure for the manifold. Based on this metric, we develop the notions of mean value and covariance matrix of a random element, normal law, Mahalanobis distance and  $\chi^2$  test. We provide a simple (but highly non trivial) characterization of Karcher means and an original gradient descent algorithm to efficiently compute them. The notion of Normal law we propose is based on the the minimization of the information knowing the mean and covariance of the distribution. The resulting family of pdfs spans the whole range from uniform (on compact manifolds) to the point mass distribution. Moreover, we were able to provide tractable approximations (with their limits) for small variances which show that we can effectively implement and work with these definitions.

To come back to more practical cases, we then reconsider the case of connected Lie groups and homogeneous manifolds. In our Riemannian context, we investigate the use of invariance principles to choose the metric: we show that it can provide the stability of our statistical definitions w.r.t. geometric operations (composition, inversion and action of transformations). However, an invariant metric does not always exists for homogeneous manifolds, nor does a left and right invariant metric for non-compact Lie groups. In this case, we cannot guaranty the full consistency of geometric and statistical operations. Thus, future work will have to concentrate on constraints weaker than invariance.

**Key-words:** Riemannian geometry, statistics, probabilities, invariance.

# Probabilités et statistiques sur des variétés riemanniennes : une approche géométrique

**Résumé :** Les mesures de primitives géométriques sont souvent bruitées dans les applications réelles et nous devons utiliser les statistiques pour en réduire l'incertitude (estimation), pour les comparer ou pour tester des hypothèses. Malheureusement, ces primitives géométriques appartiennent souvent à des variétés qui ne sont pas des espaces vectoriels. Dans des travaux précédents [Pennec, 1996, Pennec and Ayache, 1998], nous nous sommes appuyés sur l'invariance pour développer des outils statistiques de base qui évitent les paradoxes dans des groupes de transformation et des variétés homogènes.

Dans ce rapport, nous considérons la métrique riemannienne comme la structure déterminant la variété. En se basant sur cette métrique, nous développons les notions de valeur moyenne et de matrice de covariance d'un élément aléatoire, de loi normale, de distance de Mahalanobis et de test du  $\chi^2$ . Nous présentons une caractérisation simple (mais hautement non triviale !) de la moyenne de Karcher ainsi qu'un algorithme de descente de gradient efficace pour l'obtenir. La notion de loi normale que nous proposons repose sur la minimisation de l'information de la distribution connaissant la moyenne et la covariance. La famille de densités qui en résulte va de la distribution uniforme (mesure ou densité pour le cas compacte) à la distribution ponctuelle (Dirac). Nous fournissons de plus des approximations simple (ainsi que leurs limites) pour de faibles variances qui montrent que l'on peut travailler efficacement avec ces définitions.

Pour en revenir à des cas plus pratiques, nous reconsidérons les groupes de Lie et les variétés homogènes. Dans le contexte riemannien, les principes d'invariance s'expriment sur la métrique. Nous montrons que cela apporte la stabilité de nos définitions statistiques par rapport aux opérations géométriques (composition, inversion et action de transformations). Cependant, une métrique invariante n'existe pas toujours pour une variété homogène, tout comme une métrique bi-invariante (à droite et à gauche) n'existe en général pas pour un groupe de Lie non compact. Dans ce cas, nous ne pouvons pas garantir la consistance de toutes les opérations statistiques et géométriques. Les travaux futurs devront donc se concentrer sur des contraintes plus faibles que l'invariance.

**Mots-clés :** Géométrie riemannienne, statistiques, probabilités, invariance.

## Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
1.1	Riemannian metric, distance and geodesics . . . . .	6
1.2	Exponential map and cut locus . . . . .	7
1.3	Riemannian measure or volume form . . . . .	9
1.4	Gradient, Hessian and Taylor expansion of a real function . . . . .	9
1.5	Curvature and regular geodesic balls . . . . .	10
<b>2</b>	<b>Probabilities and statistics on a Riemannian Manifold</b>	<b>11</b>
2.1	Probability density function . . . . .	11
2.2	Expectation and Mean value . . . . .	12
2.3	Characterizing a Karcher mean . . . . .	14
2.4	A gradient descent algorithm to obtain the mean . . . . .	16
2.5	Covariance matrix . . . . .	17
2.6	Several random primitives . . . . .	18
2.7	A generalization of the Normal distribution . . . . .	19
2.8	Mahalanobis distance and $\chi^2$ law . . . . .	23
<b>3</b>	<b>The case of connected Lie groups</b>	<b>24</b>
3.1	Left and Right Invariant Metrics . . . . .	24
3.2	Principal chart . . . . .	25
3.3	Propagation of the pdfs for some simple group operations . . . . .	26
3.4	Obtaining the Karcher mean values . . . . .	28
3.5	Properties of the Karcher expectation for the group operations . . . . .	28
3.6	A word on the propagation of the Normal law parameters . . . . .	31
<b>4</b>	<b>The case of connected homogeneous manifolds</b>	<b>32</b>
4.1	Invariant Riemannian metric . . . . .	32
4.2	Principal chart . . . . .	33
4.3	Propagation of pdfs . . . . .	33
4.4	Obtaining the Karcher mean values . . . . .	34
4.5	Stability of the Fréchet expectation and Normal law . . . . .	35
<b>5</b>	<b>Discussion</b>	<b>36</b>
	<b>References</b>	<b>38</b>
<b>A</b>	<b>Gradient of the variance</b>	<b>40</b>
<b>B</b>	<b>Approximation of the generalized Normal density</b>	<b>43</b>
B.1	Manifolds of non positive curvature at the mean point . . . . .	43
B.2	Manifolds with a cut locus at the mean point . . . . .	47
<b>C</b>	<b>Approximated generalized <math>\chi^2</math> law</b>	<b>48</b>



## 1 Introduction

To represent the results of a random experiment, one usually construct a probabilized space  $(\Omega, \mathcal{A}, \text{Pr})$  where  $\Omega$  is the space of elementary events (the possible outcomes of the experiment),  $\mathcal{A}$  is a tribe on  $\Omega$ , usually the Borelian one<sup>1</sup>, and  $\text{Pr}$  is a measure of probability (a normalized  $\sigma$ -additive function from  $\mathcal{A}$  to  $\mathbb{R}_+$ ).

Although this probabilized space contains all the information about the random experiment, we are often only interested in some measurements depending of the outcome of the experiment. The mathematical way to formalize that is to investigate *random variables* or *observables* which are maps<sup>2</sup>  $\mathbf{x} = x(\omega)$  from  $\Omega$  to  $\mathbb{R}$ . This formalism allows to “forget” the original probabilized space  $(\Omega, \mathcal{A}, \text{Pr})$  and work directly in  $\mathbb{R}$  by associating to each random variable  $\mathbf{x}$  or  $\mathbf{y}$  different probabilities.

Often, working directly with probabilities on events is difficult and one can restrict to random variables that have a *probability density function*, i.e. a normalized function  $p_{\mathbf{x}}$  from  $\mathbb{R}$  to  $\mathbb{R}_+$  such that for all interval  $]a; b[$ :

$$\text{Pr}(\mathbf{x} \in ]a, b[) = \int_a^b p_{\mathbf{x}}(y).dy$$

However, from a computational point of view, the pdf is too informative and we have to restrict the measurements to a few numeric characteristics of a random variable. Thus, one usually approximate a unimodal pdf by a central value and a dispersion value around it. The most used central value is the *mean value* or *expectation* of the random variable:

$$\bar{x} = \mathbf{E}[\mathbf{x}] = \int \mathbf{x}.d\text{Pr} = \int_{\mathbb{R}} y.p_{\mathbf{x}}(y).dy$$

The corresponding dispersion value is the *variance*  $\sigma_{\mathbf{x}}^2 = \mathbf{E}[(\mathbf{x} - \bar{x})^2]$ .

In real problems, we can have several simultaneous measurements of the same random experiment. If we arrange these  $n$  random variables  $\mathbf{x}_i$  into a vector  $\mathbf{x} = (\mathbf{x}_1 \dots \mathbf{x}_n)$ , we obtain a *random vector* which is a map  $\mathbf{x} = x(\omega)$  from  $\Omega$  to  $\mathbb{R}^n$ . One can easily generalize the probability density function to random vector (using open sets or “paves” instead of intervals). Also, as the expectation is a linear operator, it is easily generalized to vectorial or matricial functions in order to define the mean value and the covariance matrix of a random vector:

$$\begin{aligned} \bar{\mathbf{x}} &= \mathbf{E}[\mathbf{x}] = (\mathbf{E}[\mathbf{x}_1], \dots, \mathbf{E}[\mathbf{x}_n])^T = (\bar{x}_1, \dots, \bar{x}_n)^T \\ \Sigma_{\mathbf{x}\mathbf{x}} &= \mathbf{E}[(\mathbf{x} - \bar{\mathbf{x}}).(\mathbf{x} - \bar{\mathbf{x}})^T] = \int (y - \bar{\mathbf{x}}).(y - \bar{\mathbf{x}})^T.p_{\mathbf{x}}(y).dy \end{aligned}$$

However, for many statistical problems, one have to assume a probability distribution. The Gaussian distribution is especially well adapted, as it is completely determined by the mean and the covariance. It is moreover the distribution that minimizes the information knowing only these moments. Then, one can use standard statistical tests such as the  $\chi^2$  test and use a probabilistic distance between distributions such as the Mahalanobis distance. For more details about all these notions, one can refer to [Neveu, 1990, Papoulis, 1991, Pelat, 1992].

<sup>1</sup>A tribe of  $\Omega$  is a family of parts of  $\Omega$  containing  $\emptyset$  which is stable by complementation, countable union and intersection. The Borelian tribe of a topological space is the tribe generated by the class of open sets of this space.

<sup>2</sup>The random variables should however map an element of the tribe  $\mathcal{A}$  to an element of the considered tribe of  $\mathbb{R}$ . One usually use the Borelian tribes on each space and restrict the random variables to functions that respect these structures (Borelian functions).

The problem we investigate in this article is to generalize this framework to measurements in Riemannian manifolds instead of measurements in a vector space. We call them *random primitives*. Examples of manifolds we routinely use are 3D rotations and 3D rigid transformations as transformation groups and frames (a 3D point and an orthonormal trihedron) semi- or non-oriented frames (where 2 (resp. 3) of the trihedron unit vectors are given up to their sign), oriented or directed points. We have already shown in [Pennec and Thirion, 1997, Pennec and Ayache, 1998] that this is not an easy problem and that some paradoxes can arise. In particular, we cannot generalize the expectation to give a mean value since it would be an integral with value in the manifold.

We review in the remainder of this section some basic notions of differential and Riemannian geometry that will be used afterward. This synthesis was inspired from [Spivak, 1979, chap. 9], [Klingenberg, 1982] and [Carmo, 1992], and the reader can refer to these books to find more details. In section 2, we develop the general theory for (connected and geodesically complete) Riemannian manifolds. In section 3 we show how to apply and simplify this framework in the case of transformation (Lie) groups, by using the left invariant Riemannian metric. In section 4, we focus on homogeneous manifolds and emphasize the very simple results obtained using the invariant metric if it exists.

### 1.1 Riemannian metric, distance and geodesics

In the geometric framework, one specifies the structure of a manifold  $\mathcal{M}$  by a *Riemannian metric*. This is a continuous collection of dot products  $\langle \cdot | \cdot \rangle_x$  on the tangent space  $T_x\mathcal{M}$  at each point  $x$  of the manifold. A local coordinate system  $x = (x^1, \dots, x^n)$  induces a basis  $\frac{\partial}{\partial x} = (\partial_1, \dots, \partial_n)$  of the tangent spaces ( $\partial_i$  is a shorter notation for  $\partial/\partial x^i$ ). Thus, we can express the metric in this basis by a symmetric positive definite matrix  $G(x) = [g_{ij}(x)]$  where each element is given by the dot product of the tangent vector to the coordinate curves:  $g_{ij}(x) = \langle \partial_i | \partial_j \rangle$ . This matrix is called the *local representation of the Riemannian metric* in the chart  $x$  and the dot products of two vectors  $v$  and  $w$  in  $T_x\mathcal{M}$  is now

$$\langle v | w \rangle_x = v^T \cdot G(x) \cdot w$$

The matrix  $G(x)$  is called the *local representation of the Riemannian metric* in the chart  $x$ .

If we consider a curve  $\gamma(t)$  on the manifold, we can compute at each point its instantaneous speed vector  $\dot{\gamma}(t)$  and its norm, the instantaneous speed. To compute the length of the curve, we can proceed as usual by integrating this value along the curve:

$$\mathcal{L}_a^b(\gamma) = \int_a^b \|\dot{\gamma}(t)\| dt = \int_a^b \left( \langle \dot{\gamma}(t) | \dot{\gamma}(t) \rangle_{\gamma(t)} \right)^{\frac{1}{2}} dt \quad (1)$$

The Riemannian metric is the intrinsic way of measuring length on a manifold. The extrinsic way is to consider the manifold as embedded in a larger vector space  $E$  (think for instance to the sphere  $\mathcal{S}_2$  in  $\mathbb{R}^3$ ) and compute the length of a curve in  $\mathcal{M}$  as for any curve in  $E$ . In this case, the corresponding Riemannian metric is the restriction of the dot product of  $E$  onto the tangent space at each point of the manifold. By Whitney theorem, there always exists such an embedding for a large enough vector space  $E$  ( $\dim(E) \leq 2\dim(\mathcal{M}) + 1$ ).

To obtain a distance between two points of a connected Riemannian manifold, we simply have to take the minimum length among the smooth curves joining these points:

$$\text{dist}(x, y) = \min_{\gamma} \mathcal{L}(\gamma) \quad \text{with} \quad \gamma(0) = x \quad \text{and} \quad \gamma(1) = y \quad (2)$$

The curves realizing this minimum for any two points of the manifold are called geodesics<sup>3</sup>. Let  $[g^{ij}] = [g_{ij}]^{(-1)}$  be the inverse of the metric matrix (in a given coordinate system  $x$ ) and  $\Gamma_{jk}^i$  the Christoffel symbols (using Einstein summation convention<sup>4</sup>):

$$\Gamma_{jk}^i = \frac{1}{2}g^{im} (\partial_k g_{mj} + \partial_j g_{mk} - \partial_m g_{jk}) \tag{3}$$

The calculus of variations shows the geodesics are the curves satisfying the following second order differential system (in the chart  $x = (x_1, \dots, x_n)$ ):

$$\ddot{\gamma}_i + \Gamma_{jk}^i \dot{\gamma}^j \dot{\gamma}^k = 0 \tag{4}$$

The manifold is said to be *geodesically complete* if the definition domain of all geodesics can be extended to  $\mathbb{R}$ . This means that the manifold has no boundary nor any singular point that we can reach in a finite time (for instance,  $\mathbb{R}^n - \{0\}$  with the usual metric is not geodesically complete, but  $\mathbb{R}^n$  or  $\mathcal{S}_n$  are). As an important consequence, the Hopf-Rinow-De Rham theorem state that such a manifold is complete for the induced distance (equation 2), and that there always exist at least one minimizing geodesic between any two points of the manifold (i.e. which length is the distance between the two points). From now on, we will assume that the manifold is geodesically complete.

## 1.2 Exponential map and cut locus

From the second order differential equations theory, we know that there exists one and only one geodesic  $\gamma_{(x, \partial_v)}$  going through the point  $x \in \mathcal{M}$  at  $t = 0$  with tangent vector  $\partial_v \in T_x \mathcal{M}$ . This geodesic is theoretically defined in a sufficiently small interval around zero but since we the manifold is geodesically complete, its definition domain can be extended to  $\mathbb{R}$ .

Thus, the point  $\gamma_{(x, \partial_v)}(t)$  is defined for all vector  $\partial_v \in T_x \mathcal{M}$  and all parameter  $t$ . The *exponential map* maps each vector  $\partial_v$  to the point of the manifold reached in a unit time:

$$\text{exp}_x : \begin{array}{l} T_x \mathcal{M} \longrightarrow \mathcal{M} \\ \partial_v \longmapsto \text{exp}_x(\partial_v) = \gamma_{(x, \partial_v)}(1) \end{array} \tag{5}$$

This function realizes a local diffeomorphism from a sufficiently small neighborhood  $\mathcal{V}$  of  $0 \in T_x \mathcal{M}$  into a neighborhood  $\mathcal{U}_x$  of the point  $x \in \mathcal{M}$ . We denote by  $\log_x = \text{exp}_x^{(-1)}$  the inverse map or simply  $\overline{x\vec{y}} = \log_x(y)$ . In this chart, the geodesics going through  $x$  are the represented by the lines going through the origin:  $\log_x \gamma_{(x, \overline{x\vec{y}})}(t) = t \cdot \overline{x\vec{y}}$ . Moreover, the distance with respect to the development point  $x$  is preserved:

$$\text{dist}(x, y) = \|\overline{x\vec{y}}\| = (\langle \overline{x\vec{y}} | \overline{x\vec{y}} \rangle)^{1/2}$$

Thus, the *exponential chart at x* can be seen as the development of the manifold in the tangent space at a given point along the geodesics. This is also called a *normal coordinate system* if it has an orthonormal basis. At the origin of such a chart, the metric reduces to the identity matrix and the Christoffel symbols vanish.

Now, it is natural to search for the maximal domain where the exponential map is a diffeomorphism. If we follow a geodesic  $\gamma_{(x, \partial_v)}(t) = \text{exp}_x(t \cdot \partial_v)$  from  $t = 0$  to infinity, it is either always minimizing all along or it is minimizing up to a time  $t_0 < \infty$  and not any more after (thanks to

<sup>3</sup>In facts, geodesics are defined as the critical points of the energy functional  $\mathcal{E}(\gamma) = \frac{1}{2} \int_a^b \|\dot{\gamma}\|^2 \cdot dt$ . It turns out that they also optimize the length functional but they are moreover parameterized proportionally to arc-length.

<sup>4</sup>Einstein summation convention is the implicit sum upon each index that appear up and down in the formula.

the geodesic completeness). In this last case, the point  $z = \gamma_{(x, \partial_v)}(t_0)$  is called a *cut point* and the corresponding tangent vector  $t_0 \cdot \partial_v$  a *tangential cut point*. The set of all cut points of all geodesics starting from  $x$  is the *cut locus*  $C(x) \in \mathcal{M}$  and the set of corresponding vectors the *tangential cut locus*  $\mathcal{C}(x) \in T_x \mathcal{M}$ . Thus, we have  $C(x) = \exp_x(\mathcal{C}(x))$ , and the maximal definition domain for the exponential chart is the domain  $\mathcal{D}(x)$  containing 0 and delimited by the tangential cut locus.

It is easy to see that this domain is connected and star-shaped with respect to the origin of  $T_x \mathcal{M}$ . Its image by the exponential map covers all the manifold except the cut locus and the segment  $[0, \vec{xy}]$  is transformed into the unique minimizing geodesic from  $x$  to  $y$ . Moreover, the boundary of the domain  $\mathcal{C}(x) = \partial \mathcal{D}(x)$  is continuous and contains a dense set of points where several minimizing geodesics meet.

Hence, the exponential chart at  $x$  is a chart centered at  $x$  with a connected and star-shaped definition domain that covers all the manifold except the cut locus  $C(x)$ :

$$\begin{aligned} \mathcal{D}(x) \in \mathbb{R}^n &\longleftrightarrow \mathcal{M} - C(x) \\ \vec{xy} = \log_x(y) &\longleftrightarrow y = \exp_x(\vec{xy}) \end{aligned}$$

From a computational point of view, it is often interesting to extend this representation to include the tangential cut locus but we have to take care of the multiplicity.

**Example:** On the sphere  $\mathcal{S}_n$  (center 0 and radius 1) with the canonical Riemannian metric (induced by the ambient Euclidean space  $\mathbb{R}^{n+1}$ ), the geodesics are the great circles and the cut locus of a point  $x$  is its antipodal point  $\underline{x} = -x$ . The exponential chart is obtained by rolling the sphere onto its tangent space so that the great circles going through  $x$  become lines. The maximal definition domain is thus the open ball  $\mathcal{D} = \mathcal{B}_n(\pi)$ . On its boundary  $\partial \mathcal{D} = \mathcal{C} = \mathcal{S}_{n-1}(\pi)$ , all the points represent  $\underline{x}$ .

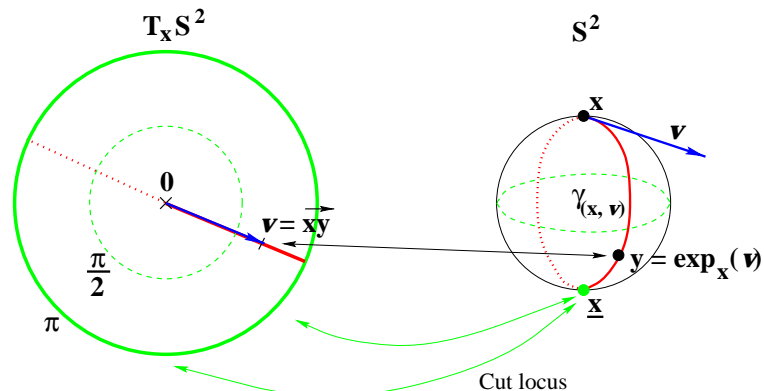


Figure 1: Exponential chart and cut locus for the sphere  $\mathcal{S}_2$  and the projective space  $\mathcal{P}_2$

For the real projective space  $\mathcal{P}_n$  (obtained by identification of antipodal points of the sphere  $\mathcal{S}_n$ ), the geodesics are still the great circles, but the cut locus of the point  $\{x, -x\}$  is now the equator of the two points where antipodal points are still identified (thus the cut locus is  $\mathcal{P}_{n-1}$ ). The definition domain of the exponential chart is the open ball  $\mathcal{D} = \mathcal{B}_n(\frac{\pi}{2})$ , and the tangential cut locus is the sphere  $\partial \mathcal{D} = \mathcal{S}_{n-1}(\frac{\pi}{2})$  where antipodal points are identified.

### 1.3 Riemannian measure or volume form

In a vector space with coordinate frame  $\mathcal{A} = (a_1, \dots, a_n)$ , the representation of the canonical metric is given by  $G = A^T \cdot A$  where  $A = [a_1, \dots, a_n]$  is the matrix of coordinates change from  $\mathcal{A}$  to the canonical orthonormal basis. Similarly, the measure (or the infinitesimal volume element) is given by the volume of the parallelepipedon spanned by the basis vectors:  $d\mathcal{V} = |\det(A)| \cdot dx = \sqrt{|\det(G)|} \cdot dx$ .

Assuming now a Riemannian manifold  $\mathcal{M}$ , we can see that the Riemannian metric  $G(x)$  induces an infinitesimal volume element on each tangent space, and thus a measure on the manifold:

$$d\mathcal{M}(x) = \sqrt{|G(x)|} dx \tag{6}$$

One can show that the cut locus has a null measure. This means that we can integrate indifferently in  $\mathcal{M}$  or in any exponential chart. If  $f$  is an integrable function of the manifold and  $f_x(\vec{x}_y) = f(\exp_x(\vec{x}_y))$  is its image in the exponential chart at  $x$ , we have:

$$\int_{\mathcal{M}} f(x) \cdot d\mathcal{M} = \int_{\mathcal{D}(x)} f_x(\vec{z}) \cdot \sqrt{G_{\vec{x}}(\vec{z})} \cdot d\vec{z}$$

### 1.4 Gradient, Hessian and Taylor expansion of a real function

**Gradient** Let  $f$  be a smooth function from  $\mathcal{M}$  to  $\mathbb{R}$  (an observable). Its Gradient  $\text{grad } f$  at point  $x$  is the linear form on  $T_x\mathcal{M}$  corresponding to the directional derivatives  $\partial_v$ :

$$\forall v \in T_x\mathcal{M} \quad \text{grad } f(v) = \partial_v f$$

Thanks to the dot product, we can identify the linear form  $d\omega$  in the tangent space  $T_x\mathcal{M}$  with the vector  $\omega$  if  $d\omega(v) = \langle \omega | v \rangle_x$  for all vector  $v \in T_x\mathcal{M}$ . All these notions can be extended to the whole manifold using vector fields.

In a local chart  $x$ , we have  $\partial_v f = \frac{\partial f(x)}{\partial x} \cdot v$  and  $\langle \omega | v \rangle_x = \omega^T \cdot G(x) \cdot v$ . Thus, the expression of the gradient in a chart is:

$$\text{grad } f = G^{(-1)}(x) \cdot \frac{\partial f}{\partial x} = g^{ij} \partial_j f$$

This definition corresponds to the classical gradient in  $\mathbb{R}^n$  even in the case of a non orthonormal basis.

**Hessian** The second covariant derivative (the *Hessian*) is a little bit more complex and makes use of the connection  $\nabla$ . We just give here its expression in a local coordinate system:

$$\text{Hess } f = \nabla df = (\partial_{ij} f - \Gamma_{ij}^k \partial_k f) dx^i dy^j$$

**Taylor expansion** Let  $f_x$  be the expression of  $f$  in a normal coordinate system at  $x$ . Its Taylor expansion around the origin is:

$$f_x(v) = f_x(0) + J_{f_x} \cdot v + \frac{1}{2} v^T \cdot H_{f_x} \cdot v + O(\|v\|^3)$$

where  $J_{f_x} = [\partial_i f]$  and  $H_{f_x} = [\partial_{ij} f]$ . Since we are in a normal coordinate system, we have  $f_x(v) = f(\exp_x(v))$ . Moreover, the metric at the origin reduces to the identity:  $J_{f_x} = \text{grad } f^T$ , and the Christoffel symbols vanished so that the matrix of second derivatives  $H_{f_x}$  corresponds to the Hessian  $\text{Hess } f$ . Thus, The Taylor expansion can be written in any coordinate system:

$$f(\exp_x(v)) = f(x) + \text{grad } f(v) + \frac{1}{2} \text{Hess } f(v, v) + O(\|v\|^3) \tag{7}$$

## 1.5 Curvature and regular geodesic balls

**Riemannian curvature** Gauss showed that the curvature of a surface can be expressed from its metric. Riemann used this property to generalize the notion of curvature to manifolds in the following way. Let  $\mathcal{V} \subset T_x\mathcal{M}$  be a 2-dimensional vector subset of the tangent space at  $x$ . The space spanned by the geodesics of  $\mathcal{M}$  starting at  $x$  with a tangent vector in  $\mathcal{V}$  is a sub-manifold with the induced metric. Thus, it is a surface and we can determine its Gaussian curvature: we call *sectional curvature* or *Riemannian curvature* this intrinsic quantity  $\kappa(x, \mathcal{V})$

It is clear that the curvature is always null in a vector space. In simple cases, the “surface” of the manifold (viewed as embedded in  $\mathbb{R}^p$ ) is always on the same side of the tangent space at any point. The curvature is in this case positive. This is in particular the case of the sphere or of  $\mathcal{S}O_n$ . For rigid transformations of  $\mathbb{R}^n$ , we add the translations (thus a vector space of null curvature), and the manifold has only a *non negative curvature*. There are also examples of manifolds with negative curvature, such as the hyperbolic paraboloid  $z = x^2 - y^2$  (the saddle surface).

**Ricci curvature** The Ricci curvature tensor is a linear operator of the tangent space (i.e. a matrix). It is obtained a partial trace of the full Riemannian curvature tensor, and can be thought of as a Laplacian of the metric: this is the property we will be interested in to make a Taylor expansion of the metric.

Let  $u \in T_x\mathcal{M}$  be a unit vector of the tangent plane, then  $\langle Ric(u) | u \rangle$  is the sum of the sectional curvatures of all the planes spanned by the vector  $u$  and one of the  $n - 1$  vector from an orthonormal frame containing  $u$ . In dimensions 2 and 3 the Ricci curvature matrix specifies all sectional curvatures of the manifold, but in higher dimensions Ricci curvature contains less information than the full curvature tensor. Using standard tensor notations and Einstein summation convention, the Riemannian curvature tensor and its contraction into the Ricci tensor are given by:

$$R_{ijk}^\gamma = \partial_k \Gamma_{ij}^\gamma + \Gamma_{ij}^\sigma \Gamma_{\sigma k}^\gamma - \partial_j \Gamma_{ik}^\gamma - \Gamma_{ik}^\sigma \Gamma_{\sigma j}^\gamma \quad \text{and} \quad Ric_{ij} = R_{ij\sigma}^\sigma$$

**Regular geodesic balls** The ball  $\mathcal{B}(x, r)$  is the set of points  $y \in \mathcal{M}$  such that the distance to the center  $x$  is less than the radius  $r$ . This ball is said to be *geodesic* if it does not meet the cut locus of its center. This means that there exists a unique minimizing geodesic from the center to any point of a geodesic ball. Let  $\kappa$  be the maximum of the Riemannian curvature in this geodesic ball. The ball is said to be *regular* if its radius verifies  $2r \cdot \sqrt{\kappa} < \pi$ .

For instance, on the sphere  $\mathcal{S}_2$  with radius one, the curvature is constant and equal to 1. A geodesic ball is regular if  $r < \pi/2$ . Such a ball can almost cover an hemisphere, but not the equator. In a Riemannian manifold with non positive curvature, a regular geodesic ball can cover the whole manifold (according to the Hadamard theorem, such a manifold is diffeomorphic to  $\mathbb{R}^n$  if it is connected).

## 2 Probabilities and statistics on a Riemannian Manifold

In this section, we do not consider measurements that are real variables or vectors depending on the outcome of a random experiment, but rather measurements of elements of a manifold. We call such a measurement a *random primitive*, reserving the names of *random transformation* and *random feature* for the particular case of transformation groups and homogeneous manifolds.

### Definition 1 (Random primitive)

Let  $(\Omega, \mathcal{B}, \Pr)$  be a probabilized space with the Borelian tribe. A random primitive in the Riemannian manifold  $\mathcal{M}$  is a Borelian function  $\mathbf{x} = \mathbf{x}(\omega)$  from  $\Omega$  to  $\mathcal{M}$ .

As in the real or vectorial case, we can now make abstraction of the original space  $\Omega$  and directly work with the induced probability measure on  $\mathcal{M}$ .

### 2.1 Probability density function

**Definition 2** Let  $\mathcal{A}$  be the Borelian tribe of  $\mathcal{M}$  (the tribe generated by the class of open sets). The random primitive  $\mathbf{x}$  has a probability density function  $p_{\mathbf{x}}$  (real, positive and integrable function) if:

$$\forall \mathcal{X} \in \mathcal{A}, \quad \Pr(\mathbf{x} \in \mathcal{X}) = \int_{\mathcal{X}} p(y).d\mathcal{M}(y) \quad \text{and} \quad \Pr(\mathcal{M}) = \int_{\mathcal{M}} p(y).d\mathcal{M}(y) = 1 \quad (8)$$

A simple example of a pdf is the *uniform pdf* in a bounded set  $\mathcal{X}$ :

$$p_{\mathcal{X}}(y) = \frac{1}{\int_{\mathcal{X}} d\mathcal{M}} \mathbf{1}_{\mathcal{X}}(y) = \frac{\mathbf{1}_{\mathcal{X}}(y)}{\mathcal{V}(\mathcal{X})}$$

where  $\mathcal{V}(\mathcal{X})$  is the “volume” of the set  $\mathcal{X}$ .

One must be careful that this pdf is uniform with respect to the measure  $d\mathcal{M}$  and is not uniform for another measure on the manifold. This problem is the basis of the Bertrand paradox for geometrical probabilities [Poincaré, 1912, Kendall and Moran, 1963, Pennec and Ayache, 1998] and raise the problem of the measure to choose on the manifold. In our case, the measure is induced by the Riemannian metric but the problem is only lifted: which Riemannian metric do we have to choose ? We will address this question in the next sections for transformation groups and homogeneous manifolds and show that an invariant metric is a good geometric choice.

**Expression of the density in a chart** Let  $\mathbf{x}$  be a random primitive of pdf  $p_{\mathbf{x}}$ . If  $x = \pi(\mathbf{x})$  is a chart of the manifold defined almost everywhere, we obtain a random vector  $\mathbf{x} = \pi(\mathbf{x})$  which pdf  $\rho_{\mathbf{x}}$  is defined with respect to the Lebesgue measure  $dx$  in  $\mathbb{R}^n$  instead of  $d\mathcal{M}$  in  $\mathcal{M}$ . Using the expression of the Riemannian measure of equation (6), the two pdfs are related by

$$\rho_{\mathbf{x}}(y) = p_{\mathbf{x}}(y) \cdot \sqrt{|G(y)|} \quad (9)$$

We shall note that the density  $\rho_{\mathbf{x}}$  depends on the chart used whereas the pdf  $p_{\mathbf{x}}$  is intrinsic to the manifold.

**Expectation of an observable** Let  $\varphi(\mathbf{x})$  be a Borelian real valued function defined on  $\mathcal{M}$  and  $\mathbf{x}$  a random primitive of pdf  $p_{\mathbf{x}}$ . Then,  $\varphi(\mathbf{x})$  is a real random variable and we can compute its expectation:

$$\mathbf{E}[\varphi(\mathbf{x})] = \mathbf{E}_{\mathbf{x}}[\varphi] = \int_{\mathcal{M}} \varphi(y) \cdot p_{\mathbf{x}}(y) \cdot d\mathcal{M}(y) \quad (10)$$

This notion of expectation corresponds to the one we defined on real random variables and vectors. However, we cannot directly extend it to define the mean value of the distribution since we have no way to generalize this integral in  $\mathbb{R}$  into an integral with value in the manifold.

## 2.2 Expectation and Mean value

We focus in this section to the notion of central value of a distribution. We will preferably use the denomination *mean value or mean primitive* than *expected primitive* to stress the difference between this notion and the expectation of a real function.

**Fréchet expectation or mean value** Let  $\mathbf{x}$  be a random vector of  $\mathbb{R}^n$ . [Fréchet, 1944, Fréchet, 1948] observed that the variance  $\sigma_{\mathbf{x}}^2(y) = \mathbf{E}[\text{dist}(\mathbf{x}, y)^2]$  is minimized for the mean vector  $\bar{\mathbf{x}} = \mathbf{E}[\mathbf{x}]$ . The major point for the generalization is that the expectation of a real valued function is well defined for our connected and geodesically complete Riemannian manifold  $\mathcal{M}$ .

### Definition 3 Variance of a random primitive

Let  $\mathbf{x}$  be a random primitive of pdf  $p_{\mathbf{x}}$ . The variance  $\sigma_{\mathbf{x}}^2(y)$  is the expectation of the squared distance between the random primitive and the fixed primitive  $y$ :

$$\sigma_{\mathbf{x}}^2(y) = \mathbf{E}[\text{dist}(y, \mathbf{x})^2] = \int_{\mathcal{M}} \text{dist}(y, z)^2 \cdot p_{\mathbf{x}}(z) \cdot d\mathcal{M}(z) \quad (11)$$

### Definition 4 Fréchet expectation of a random primitive

Let  $\mathbf{x}$  be a random primitive. If the variance  $\sigma_{\mathbf{x}}^2(y)$  is finite for all primitive  $y \in \mathcal{M}$  (which is in particular true for a density with a compact support), every primitive  $\bar{\mathbf{x}}$  minimizing this variance is called an *expected or mean primitive*. Thus, the set of the mean primitives is:

$$\mathbb{E}[\mathbf{x}] = \arg \min_{y \in \mathcal{M}} (\mathbf{E}[\text{dist}(y, \mathbf{x})^2]) \quad (12)$$

If there exists a least one mean primitive  $\bar{\mathbf{x}}$ , we call variance the minimal value  $\sigma_{\mathbf{x}}^2 = \sigma_{\mathbf{x}}^2(\bar{\mathbf{x}})$  and standard deviation its square-root.

In the same way, one define the *empirical or discrete mean primitive* of a set of measures  $x_1, \dots, x_n$  with the discrete version:

$$\mathbb{E}[\{x_i\}] = \arg \min_{y \in \mathcal{M}} (\mathbf{E}[\{\text{dist}(y, x_i)^2\}]) = \arg \min_{y \in \mathcal{M}} \left( \frac{1}{n} \sum_i \text{dist}(y, x_i)^2 \right) \quad (13)$$

If there exists a least a mean primitive  $\bar{\mathbf{x}}$ , one call *empirical variance* the minimal value  $s^2 = \frac{1}{n} \sum_i \text{dist}(\bar{\mathbf{x}}, x_i)^2$  and *empirical standard deviation* or RMS (for Root Mean Square) its square-root.

Following the same principle, one can define other types of central values. The *mean deviation at order  $\alpha$*  is

$$\sigma_{\mathbf{x}, \alpha}(y) = (\mathbf{E}[\text{dist}(y, \mathbf{x})^\alpha])^{1/\alpha} = \left( \int_{\mathcal{M}} \text{dist}(y, z)^\alpha \cdot p_{\mathbf{x}}(z) \cdot d\mathcal{M}(z) \right)^{1/\alpha} \quad (14)$$

If this function is bounded on  $\mathcal{M}$ , one call *central primitive at order  $\alpha$*  every primitive  $\bar{x}_\alpha$  minimizing it. For instance, the *modes* are obtained for  $\alpha = 0$ . Exactly like in a vector space, they are the primitives where the density is maximal on the manifold (which is generally not a maximum for the density on the charts). The *median primitive* is obtained for  $\alpha = 1$ . For  $\alpha \rightarrow \infty$ , we obtain the “barycenter” of the distribution support (which has to be compact).

The definition of these central values can be extended to the discrete case easily, except perhaps for the modes and for  $\alpha \rightarrow \infty$ . We note that the Fréchet expectation is defined for all metric space and not only for Riemannian manifolds.

**Existence and uniqueness: Karcher expectation** As our mean primitive is the result of a minimization, its existence is not ensured (the global minimum could be unreachable) and anyway the result is a set and no longer a single element. This is to be compared with some central values in vector spaces, for instance the modes. However, the Fréchet expectation does not define all the modes even in vector spaces: one only keep the modes of maximal intensity.

To get rid of this constraint, [Karcher, 1977] proposed to consider the local minima of the variance  $\sigma_{\mathbf{x}}^2(y)$  (equation 11) instead of the global ones. As global minima are local minima, the Fréchet expected primitives are a subset of the Karcher expected primitives. However, the use of local minima allows to characterize the Karcher means using only local derivatives of order two.

Using this extended definition, [Karcher, 1977] and [Kendall, 1990] established conditions on the manifold and the distribution to ensure the existence and uniqueness of the mean. We just recall here the results without the proofs:

**Theorem 1 (Existence and uniqueness of the Karcher mean)**

Let  $\mathbf{x}$  be a random primitive of pdf  $p_{\mathbf{x}}$ .

- [Kendall, 1990] *If the support of  $p_{\mathbf{x}}$  is included in a regular geodesic ball  $\mathcal{B}(y, r)$ , then there exists one and only one Karcher mean  $\mathbf{x}$  on this ball.*
- [Karcher, 1977] *If the support of  $p_{\mathbf{x}}$  is included in a geodesic ball  $\mathcal{B}(y, r)$  and if the ball of double radius  $\mathcal{B}(y, 2.r)$  is still geodesic and regular, then the variance  $\sigma_{\mathbf{x}}^2(z)$  is a convex function of  $z$  and has only one critical point on  $\mathcal{B}(y, r)$ , necessarily the Karcher mean.*

These conditions are relatively restrictive but ensure a correct behavior of our mean for localized distributions.

**Other possible definitions of the mean primitives** The Karcher mean is perfectly adapted for our purpose, thanks to the good properties it has for optimization (see below). However, there are other works proposing different ways to generalize the notion of mean value or barycenter of a distribution in a manifold. We review them for the sake of completeness and for their mathematical interest, but they do not seem to be practically applicable.

[Doss, 1949] uses another property of the expectation as the starting point for the generalization: if  $\mathbf{x}$  is a real random variable, the only real number  $\bar{x}$  verifying:

$$\forall y \in \mathbb{R} \quad |y - \bar{x}| \leq \mathbf{E} [ |\mathbf{x} - \bar{x}| ]$$

is the mean value  $\mathbf{E} [ \mathbf{x} ]$ . Thus, in a metric space, the *mean according to Doss* is defined as the set of primitives  $\bar{x} \in \mathcal{M}$  verifying:

$$\forall y \in \mathcal{M} \quad \text{dist}(y, \bar{x}) \leq \mathbf{E} [ \text{dist}(\mathbf{x}, \bar{x}) ]$$

Herer shows in [Herer, 1986, Herer, 1988] that this definition includes the classical expectation in a Banach space (with possibly other points) and develop on this basis a conditional expectation.

A similar definition that uses convex functions on the manifold instead of metric properties proposed in [Emery and Mokobodzki, 1991] and [Arnaudon, 1994, Arnaudon, 1995]. A function from  $\mathcal{M}$  to  $\mathbb{R}$  is convex if its restriction to all geodesic is convex (considered as a function from  $\mathbb{R}$  to  $\mathbb{R}$ ). The *convex barycenter* of a random primitive  $\mathbf{x}$  with density  $p_{\mathbf{x}}$  is the set  $\mathbb{B}(\mathbf{x})$  of primitives  $y \in \mathcal{M}$  such that  $\alpha(y) \leq \mathbf{E}[\alpha(\mathbf{x})]$  holds for every real bounded and convex function  $\alpha$  on a neighborhood of the support of  $p_{\mathbf{x}}$ .

This definition seems to be of little interest in our case since for compact manifolds, such as the sphere or  $SO_3$ , the manifold of 3D rotations, the geodesics are closed and the only convex functions on the manifold are the constant ones. Thus, every random primitive for which the support distribution is the whole manifold has the whole manifold as convex barycenter.

However, in the case where the support of the distribution is included in a *strongly convex open set*<sup>5</sup>  $\mathcal{U}$ , Emery shows that the *exponential barycenters*, defined as the critical points of the variance  $\sigma_{\mathbf{x}}^2(y)$  are subset of the convex barycenter  $\mathbb{B}(\mathbf{x})$ . Local and global minima being particular critical points, the exponential barycenters include the Karcher means that include themselves the Fréchet means.

[Picard, 1994] realizes a good synthesis of most of these notions of mean value and show that the definition of a “barycenter” (i.e. a mean value) is linked to a connector, which determines itself a connection, and thus possibly a metric. An interesting property brought by this formulation is that the distance between two barycenters (with different definitions) is of the order of  $O(\sigma_{\mathbf{x}})$ . Thus, for sufficiently centered random primitives, all these values are close.

### 2.3 Characterizing a Karcher mean

To characterize a local minimum of a twice differentiable function, we just have to require a null gradient and a negative definite Hessian matrix. The problem with the variance function  $\sigma^2(y)$  is that the integration domain (namely  $\mathcal{M} \setminus C(y)$ ) depends on the derivation point  $y$ . Thus we cannot just use the Lebesgue theorem to differentiate under the sum. Fortunately, we were able to generalize in appendix A a differentiability proof of Pr Maillot [Maillot, 1997] for the uniform distribution on compact manifolds. The theorem we obtain is the following:

#### Theorem 2 Gradient of the variance function

Let  $P$  be a probability on the Riemannian manifold  $\mathcal{M}$ . The variance function

$$\sigma^2(\mathbf{x}) = \int_{\mathcal{M}} \text{dist}(\mathbf{x}, z)^2 . dP(z)$$

is differentiable at any point  $y \in \mathcal{M}$  where the variance is finite and the cut locus  $C(y)$  has a null probability measure:

$$P(C(y)) = \int_{C(y)} dP(z) = 0 \quad \text{and} \quad \sigma^2(y) = \int_{\mathcal{M}} \text{dist}(y, z)^2 . dP(z) < \infty$$

At such a point, it has the following gradient:

$$(\text{grad } \sigma^2)(y) = -2. \int_{\mathcal{M}/C(y)} \vec{y}\vec{z} . dP(z) \tag{15}$$

<sup>5</sup>Here, strongly convex means that for every two points of  $\mathcal{U}$  there is a unique minimizing geodesic joining them that depend in a  $C^\infty$  of the two points.

Using the random primitive formalism, this result can be rewritten more simply

$$\text{grad} (\sigma_{\mathbf{x}}^2(y)) = -2 \cdot \mathbf{E} [\overrightarrow{y\mathbf{x}}] \tag{16}$$

Now, we know that the variance is continuous but not differentiable at the points where the cut locus has a non-zero probability measure. At these points, the variance can have an extremum (think for instance to  $\|x\|$  in vector spaces). Thus, the extrema of  $\sigma^2$  are characterized by  $(\text{grad } \sigma^2)(y) = 0$  if this is defined or  $P(C(y)) > 0$ .

**Corollary 1 (Characterization of Karcher means for manifolds with a cut locus)**

Assume that the random primitive  $\mathbf{x}$  has a finite variance everywhere and let  $\mathcal{A}$  be the set of points where the cut locus has a non-zero probability measure. A necessary condition  $\bar{x}$  to be a Karcher mean is:

$$\bar{x} \in \mathbb{E}[\mathbf{x}] \implies \begin{cases} \mathbf{E} [\overrightarrow{\bar{x}\mathbf{x}}] = 0 & \text{if } \bar{x} \notin \mathcal{A} \\ \text{or} \\ \bar{x} \in \mathcal{A} \end{cases} \tag{17}$$

For discrete or empirical means, the characterization is the same but we can write explicitly the set  $\mathcal{A} = \cup_i C(x_i)$  and the expectation  $\mathbf{E} [\overrightarrow{\{\bar{x}x_i\}}] = \sum_i \overrightarrow{\bar{x}x_i}$ .

If the manifold does not have a cut locus, we have no differentiation problem. One can even differentiate one step further to obtain the Hessian matrix:

$$\text{Hess} (\sigma_{\mathbf{x}}^2(y)) = -2 \cdot \text{Id}$$

In this case, there is one and only one minimum of the variance, necessarily the Fréchet mean.

**Corollary 2 (Characterization of the Fréchet mean for manifolds without a cut locus)**

Assume that the random primitive  $\mathbf{x}$  has a finite variance. Then, there is one and only one Fréchet (or Karcher) mean characterized by

$$\mathbb{E}[\mathbf{x}] = \{\bar{x}\} \iff \mathbf{E} [\overrightarrow{\bar{x}\mathbf{x}}] = 0 \tag{18}$$

For discrete or empirical means, the characterization is the similar

$$\mathbb{E}[\{x_i\}] = \{\bar{x}\} \iff \mathbf{E} [\overrightarrow{\{\bar{x}x_i\}}] = \sum_i \overrightarrow{\bar{x}x_i} = 0 \tag{19}$$

Basically, the characterization is the same as in Euclidean spaces, except that the non-zero probability cut loci induce some discontinuities in the first derivative of the variance. This corresponds to something like a Dirac on the second order derivative, which is the main reason why it is much more difficult to compute the Hessian matrix of the variance on a manifold with a cut locus.

**Example on the circle** The easiest example is probably a symmetric distribution on the circle. Let  $p = \cos(\theta)^2/\pi$  be the probability density function of our random primitive  $\theta$  on the circle. For a circle with the canonical metric, the exponential chart centered at  $\alpha$  is  $\overrightarrow{\alpha\theta} = \theta - \alpha$  for  $\theta \in ]\alpha - \pi; \alpha + \pi[$ , the distance being obviously  $\text{dist}(\alpha, \theta) = |\alpha - \theta|$  within this domain.

Let us first compute the mean points by computing explicitly the variance and its derivatives. The variance is:

$$\sigma^2(\alpha) = \int_{\alpha-\pi}^{\alpha+\pi} \text{dist}(\alpha, \theta)^2 \cdot p(\theta) \cdot d\theta = \int_{-\pi}^{\pi} \gamma^2 \frac{\cos(\gamma + \alpha)^2}{\pi} d\gamma = \frac{\pi^2}{3} - \frac{1}{2} + \cos(\alpha)^2.$$

Its derivative is rather easy to compute:  $\text{grad } \sigma^2(\alpha) = -2 \cos(\alpha) \sin(\alpha)$ , and the second order derivative is  $H(\alpha) = 4 \sin(\alpha)^2 - 2$ . Solving for  $\text{grad } \sigma^2(\alpha) = 0$ , we get four critical points:

- $\alpha = 0$  and  $\alpha = \pm\pi$  with  $H(0) = H(\pm\pi) = -2$ ,
- $\alpha = \pm\pi/2$  with  $H(\pm\pi) = +2$ .

Thus, there are two relative (and here absolute) minima:

$$\mathbb{E}[\boldsymbol{\theta}] = \{0, \pm\pi\}$$

Let us use now the general framework developed on Riemannian manifolds. According to theorem 2, the gradient of the variance is

$$\text{grad } \sigma^2(\alpha) = -2\mathbb{E}\left[\overrightarrow{\alpha\boldsymbol{\theta}}\right] = -2 \int_{\alpha-\pi}^{\alpha+\pi} \overrightarrow{\alpha\boldsymbol{\theta}}.d\theta = -2 \int_{\alpha-\pi}^{\alpha+\pi} (\theta - \alpha) \cdot \frac{\cos(\theta)^2}{\pi} .d\theta = -2 \cos(\alpha) \sin(\alpha),$$

which is in accordance with our previous computations. Now, differentiating once again under the sum, we get:

$$\int_{\alpha-\pi}^{\alpha+\pi} \frac{\partial^2 \text{dist}(\alpha, \theta)}{\partial \alpha^2} .p(\theta) .d\theta = -2 \int_{\alpha-\pi}^{\alpha+\pi} \frac{\partial \overrightarrow{\alpha\boldsymbol{\theta}}}{\partial \alpha} .p(\theta) .d\theta = 2 \int_{\alpha-\pi}^{\alpha+\pi} p(\theta) .d\theta = 2,$$

which is clearly different from our direct calculus. One way to see the problem is the following: the vector field  $\overrightarrow{\alpha\boldsymbol{\theta}}$  is continuous and differentiable on the circle except at the cut locus of  $\alpha$  (i.e. at  $\theta = \alpha \pm \pi$ ) where it has a jump of  $2\pi$ . Thus, the second order derivative of the squared distance should be  $-2(-1 + 2\pi\delta_{(\alpha \pm \pi)}(\theta))$ , where  $\delta$  is the Dirac distribution, and the integral becomes:

$$H(\alpha) = -2 \int_{\alpha-\pi}^{\alpha+\pi} (-1 + 2\pi\delta_{(\alpha \pm \pi)}(\theta)) .p(\theta) .d\theta = 2 - 4\pi p(\alpha \pm \pi) = 2 - 4 \cos(\theta)^2$$

which is this time in accordance with the direct calculus.

## 2.4 A gradient descent algorithm to obtain the mean

Gradient descent is a usual technique to compute a minimum. Moreover, as we have a canonical way to go from the tangent space to the manifold thanks to the exponential map, this iterative algorithm seems to be perfectly adapted. In this section, we assume that the conditions of theorem (2) are fulfilled.

Let  $y$  be an estimation of the mean of the random primitive  $\mathbf{x}$  and  $f(y) = \sigma_{\mathbf{x}}^2(y)$  the variance. A practical gradient descent algorithm is to minimize the second order approximation of the cost function at the current point. According to the Taylor expansion of equation (7), the second order approximation of  $f$  and  $y$  is:

$$f(\exp_y(v)) = f(y) + \text{grad } f(v) + \frac{1}{2}\text{Hess } f(v, v)$$

This is a function of the vector  $v \in T_y\mathcal{M}$ . Assuming that  $\text{Hess } f$  is positive definite, this function is concave and has thus a minimum characterized by a null gradient. Let  $H_f(v)$  denote the linear form verifying  $\langle H_f(v) | w \rangle = \text{Hess } f(v, w)$  for all  $w$  and  $H_f^{(-1)}$  denote the inverse map. The minimum is characterized by

$$\text{grad}_v f_y = 0 = \text{grad } f + H_f(v) \quad \Leftrightarrow \quad v = -H_f^{(-1)}(\text{grad } f)$$

We saw in the previous section that  $\text{grad } f = -2\mathbf{E}[\overrightarrow{y\mathbf{x}}]$ . Neglecting the “cut locus term” in the Hessian matrix gives us a perfect positive definite matrix  $\text{Hess } f \simeq 2 \cdot \text{Id}$ . Thus, the gradient descent algorithm is

$$y_{t+1} = \exp_{y_t} \left( \mathbf{E}[\overrightarrow{y_t\mathbf{x}}] \right) \tag{20}$$

In the case of the discrete or empirical mean, which is much more interesting from a statistical point of view, we have exactly the same algorithm, but with the empirical expectation:

$$y_{t+1} = \exp_{y_t} \left( \sum_i \overrightarrow{y_t x_i} \right) \tag{21}$$

We note that in the case of a vector space, these two formula simplify to  $y_{t+1} = \mathbf{E}[\mathbf{x}]$  and  $y_{t+1} = \frac{1}{n} \sum_i x_i$ , which are the definition of the mean value and the barycenter. Moreover, the algorithm converges in a single step.

An important point for this algorithm is to determine a good starting point. In the case on a set of measurements  $\{x_i\}$ , one can choose at random one of the measurements as the starting point. Another solution is to map to each point  $x_i$  its mean distance with respect to other points (or the median distance to be robust) and choose as the starting point the minimizing point. From a computer science point of view, the complexity is  $k^2$  (where  $k$  is the number of measurements) but the method can be randomized efficiently [Huber, 1981, Rousseeuw and Leroy, 1987].

To verify the uniqueness of the solution, we can repeat the algorithm from several starting points (for instance all the measurements  $x_i$ ). If we know the the Riemannian curvature of the manifold (for instance if it is constant or if there is an upper bound  $\kappa$ ), we can use the theorem (1). We just have to verify that the maximum distance between the measurements and the mean value we have found is sufficiently small so that all measurements fits into a regular godesic ball of radius:

$$r = \max_i \text{dist}(\bar{x}, x_i) < \frac{\pi}{2\sqrt{\kappa}}$$

## 2.5 Covariance matrix

With the mean value, we have a dispersion value: the variance. To go one step further, we observe that the covariance matrix of a random vector  $\mathbf{x}$  with respect to a point  $y$  is the *directional* dispersion of the “difference” vector  $\overrightarrow{y\mathbf{x}} = \mathbf{x} - y$ :

$$\text{Cov}_y(\mathbf{x}) = \mathbf{E}[\overrightarrow{y\mathbf{x}} \cdot \overrightarrow{y\mathbf{x}}^T] = \int_{\mathbb{R}^n} (\overrightarrow{y\mathbf{x}}) \cdot (\overrightarrow{y\mathbf{x}})^T \cdot p_{\mathbf{x}}(x) \cdot dx$$

This definition is easily extendible to a complete Riemannian manifold using the random vector  $\overrightarrow{y\mathbf{x}}$  in  $T_y\mathcal{M}$  and the Riemannian measure. In fact, we are usually interested in the covariance relative to the mean value:

### Definition 5 (Covariance)

Let  $\mathbf{x}$  be a random primitive and  $\bar{x} \in \mathbb{E}[\mathbf{x}]$  a mean value that we assume to be unique to simplify the notations (otherwise we have to keep a reference to the mean value). We note  $\Sigma_{\mathbf{xx}}$  and we call covariance the expression:

$$\Sigma_{\mathbf{xx}} = \text{Cov}_{\bar{x}}(\mathbf{x}) = \mathbf{E}[\overrightarrow{\bar{x}\mathbf{x}} \cdot \overrightarrow{\bar{x}\mathbf{x}}^T] = \int_{\mathcal{D}(\bar{x})} (\overrightarrow{\bar{x}\mathbf{x}}) \cdot (\overrightarrow{\bar{x}\mathbf{x}})^T \cdot p_{\mathbf{x}}(x) \cdot d\mathcal{M}(x) \tag{22}$$

The empirical covariance is defined in the same way using the discrete version of the expectation operator.

We observe that the covariance depends on the basis used for the exponential chart if we see it as a matrix, but it does not depend on it if we consider it as a bilinear form over the tangent plane.

The covariance is related to the variance just as in the vector case:

$$\text{Tr}(\Sigma_{\mathbf{x}\mathbf{x}}) = \mathbf{E} \left[ \text{Tr}(\overrightarrow{\bar{\mathbf{x}}}\overrightarrow{\bar{\mathbf{x}}}^T) \right] = \mathbf{E} [ \text{dist}(\bar{\mathbf{x}}, \mathbf{x}) ] = \sigma_{\bar{\mathbf{x}}}^2$$

This formula is still valid relatively to any fixed point:  $\text{Tr}(\text{Cov}_y(\mathbf{x})) = \sigma_{\bar{\mathbf{x}}}^2(y)$ .

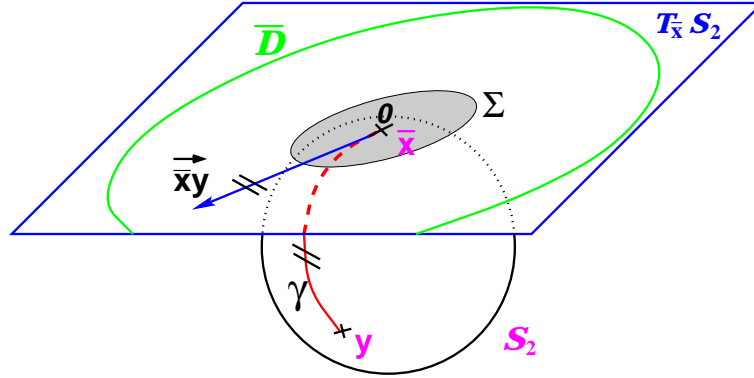


Figure 2: The covariance is defined in the tangent plane at  $S_2$  at the mean point as the classical covariance matrix of the random vector “deviation from the mean”  $\Sigma_{\mathbf{x}\mathbf{x}} = \mathbf{E} \left[ \overrightarrow{\bar{\mathbf{x}}}\overrightarrow{\bar{\mathbf{x}}}^T \right]$ .

In fact, as soon as we have found a (or the) mean value, everything appears to be similar to the case of a centered random vector by developing the manifold onto the tangent space at the mean value. Indeed,  $\overrightarrow{\bar{\mathbf{x}}}$  is a random vector of pdf  $\rho_{\bar{\mathbf{x}}}(y) = p_{\mathbf{x}}(y) \cdot \sqrt{|G(y)|}$  with respect to the Lebesgue measure in the connected and star-shaped domain  $\mathcal{D}(\bar{\mathbf{x}}) \subset T_{\bar{\mathbf{x}}}\mathcal{M}$ . We know that its expectation is  $\mathbf{E} \left[ \overrightarrow{\bar{\mathbf{x}}} \right] = 0$  and its covariance matrix is defined as usual. Thus, we could define higher order moments of the distribution by tensors on this tangent space, just as we have done for the covariance.

## 2.6 Several random primitives

Now assume that we have several simultaneous measures of the same random experiment. We quickly develop in this section the case of two random primitives  $\mathbf{x}$  and  $\mathbf{y}$  in two manifolds  $\mathcal{M}$  and  $\mathcal{N}$ . The generalization to any (finite) number of random primitives is straightforward.

Let  $\mathbf{z} = (\mathbf{x}, \mathbf{y})$  be the corresponding random primitive in the product manifold  $\mathcal{M} \times \mathcal{N}$ . The measure and metric and this product manifold is simply the product of the measures and metrics. Thus the pdf is  $p_{\mathbf{z}}(\mathbf{z}) = p_{(\mathbf{x}, \mathbf{y})}(\mathbf{x}, \mathbf{y})$  with respect to the measure  $d\mathcal{M}.d\mathcal{N}$ .

As in the vector case, the *marginal* densities are obtained by partial integration:

$$p_{\mathbf{x}}(\mathbf{x}) = \int_{\mathcal{N}} p_{(\mathbf{x}, \mathbf{y})}(\mathbf{x}, \mathbf{y}).d\mathcal{N}(\mathbf{y}) \quad \text{et} \quad p_{\mathbf{y}}(\mathbf{y}) = \int_{\mathcal{M}} p_{(\mathbf{x}, \mathbf{y})}(\mathbf{x}, \mathbf{y}).d\mathcal{M}(\mathbf{x})$$

The mean value of the joint measure is

$$\bar{\mathbf{z}} \in \mathbb{E}[\mathbf{z}] = \mathbb{E}[\mathbf{x}] \times \mathbb{E}[\mathbf{y}] \quad \iff \quad \bar{\mathbf{x}} \in \mathbb{E}[\mathbf{x}] \quad \text{and} \quad \bar{\mathbf{y}} \in \mathbb{E}[\mathbf{y}]$$

Assuming that there is only one mean, the covariance is given by:

$$\Sigma_{\mathbf{z}\mathbf{z}} = \begin{bmatrix} \Sigma_{\mathbf{x}\mathbf{x}} & \Sigma_{\mathbf{x}\mathbf{y}} \\ \Sigma_{\mathbf{y}\mathbf{x}} & \Sigma_{\mathbf{y}\mathbf{y}} \end{bmatrix}$$

where  $\Sigma_{\mathbf{x}\mathbf{y}} = \Sigma_{\mathbf{y}\mathbf{x}}^T$  is the following cross covariance

$$\Sigma_{\mathbf{x}\mathbf{y}} = \mathbf{E} \left[ \overrightarrow{\bar{\mathbf{x}}\mathbf{x}} \cdot \overrightarrow{\bar{\mathbf{y}}\mathbf{y}}^T \right]$$

**Independent primitives** The pdf of the joint measure  $p_{\mathbf{z}}(z)$  can be factored in

$$p_{(\mathbf{x},\mathbf{y})}(\mathbf{x}, \mathbf{y}) = p_{\mathbf{x}}(\mathbf{x}) \cdot p_{\mathbf{y}}(\mathbf{y})$$

if and only if the random primitives  $\mathbf{x}$  and  $\mathbf{y}$  are independent.

The independence assumption will often be made in statistics. In this case the cross covariance can be factored in  $\Sigma_{\mathbf{x}\mathbf{y}} = \mathbf{E} \left[ \overrightarrow{\bar{\mathbf{x}}\mathbf{x}} \right] \cdot \mathbf{E} \left[ \overrightarrow{\bar{\mathbf{y}}\mathbf{y}}^T \right]$  but these two integrals are null by the mean value characterization. Thus, we have :  $\Sigma_{\mathbf{x}\mathbf{y}} = 0$ .

## 2.7 A generalization of the Normal distribution

We now present an approach based on the the information minimization to generalize the Normal distribution to a manifold. This may not be the best generalization but it allows us to approximate such a distribution by the usual Gaussian in the tangent space of the mean value in the case of small covariances. In this section the symbols  $\log$  and  $\exp$  denote the standard logarithmic and exponential functions in  $\mathbb{R}$ .

**Information and uniform law** As we can integrate a real valued function, the extension of the *entropy*  $\mathbf{H}[\mathbf{x}]$  (or its opposite, the *information*  $\mathbf{I}[\mathbf{x}]$ ) of a random primitive is straightforward:

$$\mathbf{I}[\mathbf{x}] = -\mathbf{H}[\mathbf{x}] = \mathbf{E}[\log(p_{\mathbf{x}}(\mathbf{x}))] = \int_{\mathcal{M}} \log(p_{\mathbf{x}}(\mathbf{x})) \cdot p_{\mathbf{x}}(\mathbf{x}) \cdot d\mathcal{M}(\mathbf{x}) \quad (23)$$

This definition is coherent since the pdf  $p_{\mathcal{U}}$  that minimize the information when we only know that the measure is in a compact set  $\mathcal{U}$  is the uniform density in this set:

$$p_{\mathcal{U}}(\mathbf{x}) = \mathbf{1}_{\mathcal{U}}(\mathbf{x}) \Big/ \int_{\mathcal{U}} d\mathcal{M}(y)$$

**Constrained information minimization** Now assume that we only know the mean (that we suppose to be unique) and the covariance of a random primitive: we denote it by  $\mathbf{x} \sim (\bar{\mathbf{x}}, \Sigma)$ . If we need to assume a pdf for that random primitive, it seems reasonable to choose the one which is the least informative, while fitting the mean and the covariance. The pdf is minimizing in this case the conditional information:

$$\mathbf{I}[\mathbf{x} | \bar{\mathbf{x}} \in \mathbb{E}[\mathbf{x}], \Sigma_{\mathbf{x}\mathbf{x}} = \Sigma]$$

We can express all the constraints directly in the exponential chart at the mean value. Let  $\rho(y) = p(\exp_{\bar{\mathbf{x}}}(y))$  be the density in the chart with respect to the induced Riemannian measure  $d\mathcal{M}_{\bar{\mathbf{x}}}(y) = \sqrt{|G_{\bar{\mathbf{x}}}(y)|} \cdot dy$  (we use here the Riemannian measure instead of the Lebesgue one to simplify equations below). The constraints are:

- the normalization:  $\mathbf{E} [ \mathbf{1}_{\mathcal{M}} ] = \int_{\mathcal{D}(\bar{x})} \rho(y).d\mathcal{M}_{\bar{x}}(y) = 1$
- a nul mean value:  $\mathbf{E} [ \vec{\bar{x}\bar{x}} ] = \int_{\mathcal{D}(\bar{x})} y.\rho(y).d\mathcal{M}_{\bar{x}}(y) = 0,$
- and a fixed covariance  $\Sigma$  :  $\mathbf{E} [ \vec{\bar{x}\bar{x}}.\vec{\bar{x}\bar{x}}^T ] = \int_{\mathcal{D}(\bar{x})} y.y^T.\rho(y).d\mathcal{M}_{\bar{x}}(y) = \Sigma$

To simplify the optimization, we won't consider any continuity or differentiability constraint on the cut locus  $C(\bar{x})$  (which would mean constraints on the border of the domain). This means that we can do the optimization in the exponential chart at the mean point like if we were in the open domain  $\mathcal{D}(\bar{x}) \in \mathbb{R}^n$ .

Using a scalar Lagrange multiplier  $\alpha$  for the normalization constraint, a vector multiplier  $\beta$  for the mean value and a symmetric matrix  $\Gamma$  for the covariance, the Lagrangian is:

$$\Lambda(\rho) = \int_{\mathcal{D}(\bar{x})} \left( \rho.\log(\rho) + \alpha.\rho + \beta^T.y.\rho + \frac{y^T.\Gamma.y}{2}.\rho \right).d\mathcal{M}_{\bar{x}}(y)$$

It is stationary at a critical point:

$$\frac{\partial \Lambda(\rho)}{\partial \rho} = 0 = \log(\rho) + 1 + \alpha + \beta^T.y + \frac{y^T.\Gamma.y}{2}$$

Thus we find  $\rho(y) = k.\exp\left(-\beta^T.y - \frac{y^T.\Gamma.y}{2}\right)$  with  $k = \exp(-1 - \alpha)$ . Assuming that the definition domain  $\mathcal{D}(\bar{x})$  is symmetric with respect to the origin, we find that  $\beta = 0$  ensures a null mean. Reporting in the constraints gives the following equations.

### Theorem 3 (Normal law)

We call Normal law on the manifold  $\mathcal{M}$  the pdf minimizing the information with a fixed mean value and covariance.

Assuming no continuity nor differentiability constraint on the cut locus  $C(\bar{x})$  and a symmetric domain  $\mathcal{D}(\bar{x})$ , the Normal law of mean  $\bar{x}$  (the mean value) and concentration matrix  $\Gamma$  is given by:

$$N_{(\bar{x},\Gamma)}(y) = k.\exp\left(-\frac{\vec{\bar{x}y}^T.\Gamma.\vec{\bar{x}y}}{2}\right) \quad (24)$$

where the normalization constant is

$$k^{(-1)} = \int_{\mathcal{M}} \exp\left(-\frac{\vec{\bar{x}y}^T.\Gamma.\vec{\bar{x}y}}{2}\right).d\mathcal{M}(y) \quad (25)$$

The covariance and concentration are related by:

$$\Sigma = k.\int_{\mathcal{M}} \vec{\bar{x}y}.\vec{\bar{x}y}^T.\exp\left(-\frac{\vec{\bar{x}y}^T.\Gamma.\vec{\bar{x}y}}{2}\right).d\mathcal{M}(y) \quad (26)$$

From the concentration matrix, we can compute the covariance of the random primitive, at least numerically, but the reverse is more difficult. As  $y^T.\Gamma.y = \text{Tr}(\Gamma.y.y^T)$ , the information of this pdf simplifies to:

$$\mathbf{I} [ N_{(\bar{x},\Gamma)} ] = \log(k) - \frac{1}{2}.\int_{\mathcal{M}} y^T.\Gamma.y.N_{(\bar{x},\Gamma)}.d\mathcal{M}_{\bar{x}}(y) = \log(k) - \frac{1}{2}.\text{Tr}(\Gamma.\Sigma)$$

**The vectorial case** The integrals can be entirely computed, and we find  $k^{(-1)} = \frac{(2\pi)^{\frac{n}{2}}}{\sqrt{|\Gamma|}}$  and  $\Gamma = \Sigma^{(-1)}$ . The Normal density is thus the usual Gaussian density:

$$N_{(\bar{x},\Gamma)}(x) = k \cdot \exp\left(-\frac{(x - \bar{x})^T \cdot \Gamma \cdot (x - \bar{x})}{2}\right) = \frac{1}{(2\pi)^{n/2} \cdot \sqrt{|\Sigma|}} \cdot \exp\left(-\frac{(x - \bar{x})^T \cdot \Sigma^{(-1)} \cdot (x - \bar{x})}{2}\right)$$

**Example on a simple manifold: the circle** The exponential chart for the circle of radius 1 with the canonical metric is the angle  $\theta \in \mathcal{D} = ]-\pi; \pi[$  with respect to the development point, and the measure is simply  $d\theta$ . For a circle of radius  $r$ , the exponential chart becomes  $x = r \cdot \theta$ . The domain is  $\mathcal{D} = ]-a; a[$  (with  $a = \pi \cdot r$ ) and the measure is  $dx = r \cdot d\theta$ . Thus, the normalization factor of the Normal density is:

$$k^{(-1)} = \int_{-a}^a \exp\left(-\frac{\gamma \cdot x^2}{2}\right) \cdot dx = \sqrt{\frac{2\pi}{\gamma}} \cdot \operatorname{erf}\left(\sqrt{\frac{\gamma}{2}} \cdot a\right)$$

where  $\operatorname{erf}$  is the error function  $\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \cdot \int_0^x \exp(-t^2) \cdot dt$ . The density is the truncated Gaussian  $N_{(0,\gamma)}(x) = k \cdot \exp\left(-\frac{\gamma \cdot x^2}{2}\right)$ . It is continuous but not differentiable on the cut locus  $\pi \equiv -\pi$ . The truncation introduces a bias in the relation between the variance and the concentration parameter:

$$\sigma^2 = \int_{-a}^a x^2 \cdot k \cdot \exp\left(-\frac{\gamma \cdot x^2}{2}\right) \cdot dx = \frac{1}{\gamma} \left(1 - 2 \cdot a \cdot k \cdot \exp\left(-\frac{\gamma \cdot a^2}{2}\right)\right)$$

It is interesting to have a look on limit properties: if the circle radius goes to infinity, the circle becomes the real line and we obtain the usual Gaussian with the relation  $\sigma^2 = 1/\gamma$ , as expected.

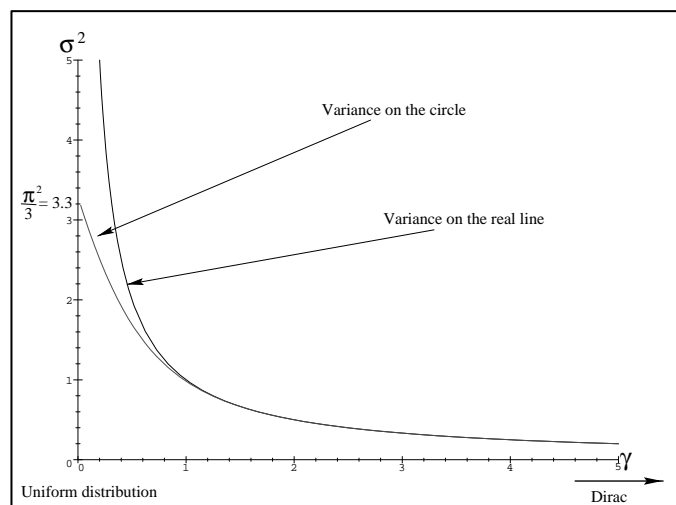


Figure 3: Variance  $\sigma^2$  with respect to the concentration parameter  $\gamma$  on the circle of radius 1 and the real line. This variance tends toward  $\sigma_0^2 = \pi^2/3$  for the uniform distribution on the circle ( $\gamma = 0$ ) whereas it tends to infinity for the uniform measure on  $\mathbb{R}$ . For a strong concentration ( $\gamma > 1$ ), the variance on the circle can be accurately approximated by  $\sigma^2 \simeq 1/\gamma$ , as in the real case.

As the circle is compact, the variance cannot become infinite as in the real case. If the concentration parameter  $\gamma$  goes to zero (which corresponds to  $\sigma^2 \rightarrow +\infty$  in the real case), a Taylor

expansion gives  $\sigma^2 = a^2/3 + O(\gamma)$ . Thus, the maximal variance on the circle is

$$\sigma_0^2 = \lim_{\gamma \rightarrow 0} \sigma^2 = \frac{a^2}{3} \quad \text{with the density} \quad N_{(0,0)}(x) = \frac{1}{2.a}$$

As expected, the Normal density of concentration 0 is the uniform density. On the other hand, if  $\gamma$  goes to infinity, the variance goes to zero and the density tends to a Dirac (see figure 3).

**Approximation for a small  $\Sigma$**  If the pdf is sufficiently concentrated (a high concentration matrix  $\Gamma$  or a small covariance matrix  $\Sigma$ ), then we can use a Taylor expansion of the metric in a normal coordinate system around the mean value to approximate the previous integrals and obtain a Taylor expansions of the normalization factor and the concentration matrix with respect to the covariance matrix.

The Taylor expansion of the metric is given by [Chavel, 1993, p84]. We easily deduce the Taylor expansion of the measure around the origin (Ric is the Ricci (or scalar) curvature matrix in the considered normal coordinate system):

$$d\mathcal{M}(y) = \sqrt{\det(G(y))} = 1 - \frac{y^T \cdot \text{Ric} \cdot y}{6} + O(\|y\|^3)$$

Reporting this Taylor expansion in the integrals and manipulating the formulas (see appendix B) leads to the following theorem.

**Theorem 4 (Approximate normal density)**

*In a complete Riemannian manifold, the normal density is  $N(y) = k \cdot \exp(-\vec{xy}^T \cdot \Gamma \cdot \vec{xy}/2)$ . Let  $r$  be the injection radius at the mean point. The normalization constant and the concentration matrices are approximated by the following expressions for a covariance matrix  $\Sigma$  of small variance  $\sigma^2 = \text{Tr}(\Sigma)$ :*

$$k = \frac{1 + O(\sigma^3) + \varepsilon\left(\frac{\sigma}{r}\right)}{\sqrt{(2\pi)^n \cdot \det(\Sigma)}} \quad \text{and} \quad \Gamma = \Sigma^{(-1)} - \frac{1}{3}\text{Ric} + O(\sigma) + \varepsilon\left(\frac{\sigma}{r}\right)$$

Here,  $\varepsilon(x)$  is a function that is a  $O(x^k)$  for any positive  $k$ , with the convention that  $\varepsilon\left(\frac{\sigma}{+\infty}\right) = \varepsilon(0) = 0$ . More precisely, this is a function such that  $\forall k \in \mathbb{R}^+$ ,  $\lim_{0+} x^{-k} \cdot \varepsilon(x) = 0$

**Discussion** The information minimization approach to generalize the normal distribution to Riemannian manifolds is interesting since we obtain a whole family of densities going from the Dirac to the uniform distribution (or the uniform measure if the manifold is only locally compact). Unfortunately, this distribution is generally not differentiable at the cut locus, and often even not continuous.

However, if the relation between the parameters and the moments of the distribution are not as simple as in the vector case (but can we expect something simpler in the general case of Riemannian manifolds?), the approximation for small covariances turns out to be rather simple. Thus, this approximate distribution can be handled quite easily for statistical purposes.

It would be interesting to see if there exists similar families of densities that are continuous and differentiable at the cut locus. A starting point could be to extend some of the distributions used in directional statistics (statistics on the sphere) as in [Bingham, 1974, Jupp and Mardia, 1989, Kent, 1992, Mardia, 1995].

## 2.8 Mahalanobis distance and $\chi^2$ law

The problem we are now tackling is to determine if a measure  $\hat{y}$  was drawn from a random primitive  $\mathbf{x}$ . From a statistical point of view, the pdf of the measurement process is often too rich an information to be estimated or handled and, in practice, one often characterizes this pdf by its moments, and more particularly by the mean and the covariance<sup>6</sup>. We denote it by

$$\mathbf{x} \sim (\bar{\mathbf{x}}, \Sigma_{\mathbf{xx}})$$

Based on these characteristics only, we want to decide if the measure  $\hat{y}$  is compatible with this measurement process.

In the vectorial case and assuming a Gaussian distribution, the  $\chi^2$  test is well adapted to do that. This test measures the probability of the Mahalanobis distance  $\chi^2 = (\hat{x} - \bar{x})^T \cdot \Sigma_{\mathbf{xx}}^{(-1)} \cdot (\hat{x} - \bar{x})$  assuming that  $\hat{x}$  is drawn from  $\mathbf{x}$ . If the probability is too small (i.e.  $\chi^2$  is too large), the hypothesis is rejected.

This definition of the Mahalanobis distance can be easily generalized to complete Riemannian manifolds with our tools. We note that it is well defined for any distribution of the random primitive and not only the normal one.

### Definition 6 (Mahalanobis distance)

We call Mahalanobis distance between a random primitive  $\mathbf{x} \sim (\bar{\mathbf{x}}, \Sigma_{\mathbf{xx}})$  and a (deterministic) point  $y$  on the manifold the value

$$\mu_{\mathbf{x}}^2(y) = \bar{\mathbf{x}}y^T \cdot \Sigma_{\mathbf{xx}}^{(-1)} \cdot \bar{\mathbf{x}}y \quad (27)$$

In fact, the Mahalanobis distance measures the distance between  $y$  and the mean value  $\bar{\mathbf{x}}$  according to the “metric”  $\Sigma_{\mathbf{xx}}^{(-1)}$ .

**Property** Since  $\mu_{\mathbf{x}}^2$  is a function from  $\mathcal{M}$  to  $\mathbb{R}$ ,  $\mu_{\mathbf{x}}^2(\mathbf{y})$  is a real random variable. The expectation of this random variable is well defined and turns out to be quite simple:

$$\begin{aligned} \mathbf{E} [\mu_{\mathbf{x}}^2(\mathbf{y})] &= \int_{\mathcal{M}} \mu_{\mathbf{x}}^2(z) \cdot p_{\mathbf{y}}(z) \cdot d\mathcal{M}(z) = \int_{\mathcal{M}} \bar{\mathbf{x}}z^T \cdot \Sigma_{\mathbf{xx}}^{(-1)} \cdot \bar{\mathbf{x}}z \cdot p_{\mathbf{y}}(z) \cdot d\mathcal{M}(z) \\ &= \text{Tr} \left( \Sigma_{\mathbf{xx}}^{(-1)} \cdot \int_{\mathcal{M}} \bar{\mathbf{x}}z \cdot \bar{\mathbf{x}}z^T \cdot p_{\mathbf{y}}(z) \cdot d\mathcal{M}(z) \right) = \text{Tr} (\Sigma_{\mathbf{xx}}^{(-1)} \cdot \text{Cov}_{\bar{\mathbf{x}}}(\mathbf{y})) \end{aligned}$$

The expectation of the Mahalanobis distance of a random primitive with itself is even simpler:

$$\mathbf{E} [\mu_{\mathbf{x}}^2(\mathbf{x})] = \text{Tr}(\Sigma_{\mathbf{xx}}^{(-1)} \cdot \Sigma_{\mathbf{xx}}) = \text{Tr}(\text{Id}_n) = n$$

**Theorem 5** The expected Mahalanobis distance of a random primitive with itself is independent of the distribution and does only depend on the dimension of the manifold:

$$\mathbf{E} [\mu_{\mathbf{x}}^2(\mathbf{x})] = n \quad (28)$$

This identity can be used to verify with a posteriori measurements that the covariance matrix has been correctly estimated. It can be compared with the expectation of the “normalized” squared distance, which is by definition:  $\mathbf{E} [\text{dist}(\mathbf{x}, \bar{\mathbf{x}})^2 / \sigma_{\mathbf{x}}^2] = 1$ .

---

<sup>6</sup>If we have to assume a density for this random primitive, the least informative one is the normal density that can be approximated by theorem 4.

**A generalized  $\chi^2$  law** Assuming that the random primitive  $\mathbf{x} \sim (\bar{\mathbf{x}}, \Sigma_{\mathbf{xx}})$  is normal, we can go one step further and compute the probability that  $\chi^2 = \mu_{\mathbf{x}}^2 < \alpha^2$  (see appendix C). This generalization of the  $\chi^2$  law turns out to be still independent of the mean value and the covariance matrix of the random primitive (at least up to the order  $O(\sigma^3)$ ):

**Theorem 6 (Approximate  $\chi^2$  law)**

*With the same hypotheses as for the approximate normal law, the  $\chi^2$  probability is*

$$\Pr\{\chi^2 \leq \alpha^2\} = (2\pi)^{-\frac{n}{2}} \int_{\|x\| \leq \alpha} \exp\left(-\frac{\|x\|^2}{2}\right) .dx + O(\sigma^3) + \varepsilon\left(\frac{\sigma}{r}\right) \quad (29)$$

while the density is

$$p_{\chi^2}(u) = \frac{1 + O(\sigma^3) + \varepsilon\left(\frac{\sigma}{r}\right)}{2 \cdot \Gamma\left(\frac{n}{2}\right)} \left(\frac{u}{2}\right)^{\frac{n}{2}-1} \exp\left(-\frac{u}{2}\right) \quad (30)$$

The  $\chi^2$  probability can be computed using the incomplete gamma function  $\Pr\{\chi^2 \leq \alpha^2\} = P\left(\frac{n}{2}, \frac{\alpha^2}{2}\right)$  (see for instance [Press et al., 1991]).

In practice, one often use this law to test if a measurement  $\hat{\mathbf{x}}$  has been drawn from a random primitive  $\mathbf{x}$ : if the hypothesis is true, the value  $\mu_{\mathbf{x}}^2(\hat{\mathbf{x}})$  will be less than  $\alpha^2$  with a probability  $\gamma = \Pr\{\chi^2 \leq \alpha^2\}$ . Thus, one choose a confidence level  $\gamma$  (for instance 95% or 99%), then we find the value  $\alpha(\gamma)$  such that  $\gamma = \Pr\{\chi^2 \leq \alpha^2\}$  and accept the hypothesis if  $\mu_{\mathbf{x}}^2(\hat{\mathbf{x}}) \leq \alpha^2$ .

### 3 The case of connected Lie groups

A Lie group being a differentiable manifold, we are just left with the choice of a Riemannian metric to apply our statistical framework. Firstly, we detail the two canonical Riemannian metrics (left- and right-invariant metrics). Then, choosing the left invariant metric, we show that all the exponential could be identified to the exponential chart at the identity (called the principal chart). From a computer science point of view, this means that we only have only one chart to deal with in order to cover the whole group (with special cases of the points of the cut locus). In Section 3.3, we investigate the propagation of the pdfs of random transformations for the group operations: left and right translation by a deterministic transformation, inversion, composition of two random transformations. Interestingly, this last pdf is based on a generalisation of the convolution product used for the addition of two random variables in  $\mathbb{R}^n$ . Section 3.4 details the simplifications of our gradient descent algorithm to obtain the mean using the group operations and the properties of the left invariant metric, while Section 3.5 investigates the propagation of the Karcher mean values under the group operations: if the left translation of a random transformation (by a deterministic transformation) is stable, we need to require the left and right invariance of the metric in order to ensure the stability of the mean values with respect to the right translation and the inversion. An example with 2D rigid transformations, where the metric is only left invariant, shows that the empirical mean value of the inverse is generally not the inverse of the mean value. Last but not least, we show in section 3.6 that the Normal distribution is stable by left translation.

#### 3.1 Left and Right Invariant Metrics

We have two canonical way to compare the tangent spaces at different points of the manifold: the left translation  $L_g(f) = g \circ f$  and the right translation  $R_g(f) = f \circ g$ . These applications can obviously

be differentiated to provide linear operators from  $T_f\mathcal{G}$  to  $T_{(g \circ f)}\mathcal{G}$  and  $T_{(f \circ g)}\mathcal{G}$ . Let us denote by  $J_L$  and  $J_R$  the expression of the differentials (in a given chart) at the identity:

$$J_L(f) = \left. \frac{\partial(f \circ e)}{\partial e} \right|_{e=\text{Id}} \quad \text{and} \quad J_R(f) = \left. \frac{\partial(e \circ f)}{\partial e} \right|_{e=\text{Id}}$$

Let  $v$  be a tangent vector at the identity. It corresponds to a tangent vector of  $T_f\mathcal{M}$   $v_l = J_L(f).v$  by left translation and  $v_r = J_R(f).v$  by right translation.

Let now  $\langle \cdot | \cdot \rangle$  be any dot product on the tangent space at the identity, characterised by a symmetric matrix  $Q$  in our chart. We can translate this dot product in any tangent space by left or right translation. If  $v$  and  $w$  are tangent vectors at  $f$ , we define the left-invariant metric by:

$$\langle v | w \rangle_f = \langle J_L(f)^{(-1)}.v | J_L(f)^{(-1)}.w \rangle = v^T . J_L(f)^{(-T)} . Q . J_L(f)^{(-1)}.w,$$

which means that the expression of the metric is in our chart:  $G_L(f) = J_L(f)^{(-T)} . Q . J_L(f)^{(-1)}$ . The right-invariant metric is defined similarly with  $G_R(f) = J_R(f)^{(-T)} . Q . J_R(f)^{(-1)}$ .

We note that the volume element associated to these invariant metrics are the Haar measures of the group ( $|J| = |\det(J)|$ ) [Pennec and Ayache, 1998]:

$$d_L\mathcal{G}(f) = \sqrt{|G_L(f)|}.df = \lambda. \frac{df}{|J_L(f)|} \quad \text{and} \quad d_R\mathcal{G}(f) = \sqrt{|G_R(f)|}.df = \lambda. \frac{df}{|J_R(f)|}$$

Since the metric is left (resp. right) invariant, the geodesics are globally conserved by left (resp. right) translation and one can determine only the geodesics starting from the identity. One should be careful that left and right invariant metrics are generally different. However, if the group is compact (for instance rotations), there exists a left and right invariant metric [Spivak, 1979, Carmo, 1992] and the geodesics for this metric (starting from the identity) are the one parameter sub-groups, i.e. the curves verifying  $\forall(s, t) \in \mathbb{R}^2, \gamma(s + t) = \gamma(s) \circ \gamma(t) = \gamma(t) \circ \gamma(s)$ . This equation is often easier to solve than the standard second order differential equations system, but it is not valid for only locally compact groups (such as rigid transformations for instance).

### 3.2 Principal chart

From now on, we consider that the group is geodesically complete and that geodesics are determined for the left-invariant metric. The same developments could of course be done (with different results) with the right invariant one.

We call *principal chart* the exponential chart at the identity and we denote by  $\vec{f}$  the representation of  $f$  in this chart. To simplify the notations, we will write simply  $\exp$  and  $\log$  the exponential map and its inverse function at the identity. As the metric matrix  $Q$  is symmetric and positive, it can be written  $Q = A^T.A$ . Thus, changing the coordinate system into  $\vec{f}' = A. \vec{f}$  allows to obtain a principal chart where  $Q = \text{Id}$ . We will consider  $Q = \text{Id}$  in the sequel.

Let  $\gamma_{(\text{Id}, \vec{f})}$  be a geodesic starting at the identity with tangent vector  $\vec{f}$ . We have  $\gamma(0) = \text{Id}$  and  $\gamma(1) = f$ . Translating this geodesic by  $g$ , we obtain a new geodesic  $\delta = g \circ \gamma_{(\text{Id}, \vec{f})}$  starting at  $g$  with tangent vector  $J_L(\vec{g}). \vec{f}$  and arriving at time one at  $g \circ f$ . Thus, by definition of the exponential map, we get  $\exp_g(J_L(\vec{g}). \vec{f}) = g \circ f$ . Taking  $v = J_L(\vec{g}). \vec{f}$ , we obtain the following formula for the exponential map at any point:

$$\exp_g(v) = g \circ \exp(J_L(\vec{g})^{(-1)}.v) \tag{31}$$

Taking now  $v = J_L(\vec{g}).(\vec{g}^{(-1)} \circ \vec{f})$ , we obtain the expression  $\vec{g}\vec{f}$  of the  $f$  in the exponential chart at  $g$  (expressed in the basis induced by the principal chart):

$$\vec{g}\vec{f} = \log_{\vec{g}}(f) = J_L(\vec{g}).(\vec{g}^{(-1)} \circ \vec{f}) \quad (32)$$

The main interest of all these developments is that, from a computer science point of view, we can use only one chart (the principal chart) and obtain all the others by the left translation.

The distance between two transformations  $f$  and  $g$  is the length of the minimising geodesic joining the two points. By definition, such a geodesic starts at  $g$  with tangent vector  $\vec{g}\vec{f}$ :

$$\text{dist}(f, g)^2 = \|\vec{g}\vec{f}\|_{\vec{g}}^2 = \vec{g}\vec{f}^T \cdot Q(\vec{g}).\vec{g}\vec{f} = (\vec{g}^{(-1)} \circ \vec{f})^T \cdot (\vec{g}^{(-1)} \circ \vec{f}) = \|\vec{g}^{(-1)} \circ \vec{f}\|_{\text{Id}}^2 = \text{dist}(g^{(-1)} \circ f, \text{Id})^2 \quad (33)$$

### 3.3 Propagation of the pdfs for some simple group operations

Since we chose the left invariant metric on our Lie group, the probability density functions are defined using the left Haar measure:

$$\Pr(\mathbf{f} \in \mathcal{Y}) = \int_{\mathcal{Y}} p(y).d_L\mathcal{G}(y)$$

As already mentioned, the left and right Haar measures are generally different. Their difference is quantified by the *module*  $\Delta\mathcal{G}(g)$ , defined by the relation:  $d_L\mathcal{G}(g) = \Delta\mathcal{G}(g).d_R\mathcal{G}(g)$ . With our notations, we have  $\Delta\mathcal{G}(g) = |J_R(g)|/|J_L(g)|$ . The group is said unimodular if  $\Delta\mathcal{G}(g)$  is constant at each point of the group (i.e. equal to 1 up to a normalisation factor). A compact group is always unimodular but only locally compact group may have different left and right Haar measures (they may have in this case different left and right invariant metrics).

#### Theorem 7 (Translation of a random transformation)

Let  $\mathbf{f}$  be a random transformation of pdf  $p_{\mathbf{f}}$ , and  $f_o$  a deterministic transformation. The pdf of  $\mathbf{g}_l = f_o \circ \mathbf{f}$  (left translation of  $\mathbf{f}$ ) is:

$$p_{(f_o \circ \mathbf{f})}(g) = p_{\mathbf{f}}(f_o^{(-1)} \circ g) \quad (34)$$

whereas the pdf of the right translation is:

$$p_{(\mathbf{f} \circ f_o)}(g) = \frac{\Delta\mathcal{G}(g \circ f_o^{(-1)})}{\Delta\mathcal{G}(g)}.p_{\mathbf{f}}(g \circ f_o^{(-1)}) \quad (35)$$

**Proof:** The probability of  $\mathbf{g}_l = f_o \circ \mathbf{f}$  to be in a set  $\mathcal{Y} \subset \mathcal{G}$  is:

$$P(f_o \circ \mathbf{f} \in \mathcal{Y}) = P(\mathbf{f} \in f_o^{(-1)} \circ \mathcal{Y}) = \int_{f_o^{(-1)} \circ \mathcal{Y}} p_{\mathbf{f}}(h).d_L\mathcal{G}(h) \int_{\mathcal{Y}} p_{\mathbf{f}}(f_o^{(-1)} \circ g).d_L\mathcal{G}(f_o^{(-1)} \circ g)$$

where the last equality is obtained by the change of variable  $h = f_o^{(-1)} \circ g$ . The first result is obtained using the left invariance of the measure.

Now, for the right translation, the probability of  $\mathbf{g}_r = \mathbf{f} \circ f_o$  to be in a set  $\mathcal{Y} \subset \mathcal{G}$  is:

$$\Pr((\mathbf{f} \circ f_o) \in \mathcal{Y}) = \Pr(\mathbf{f} \in (\mathcal{Y} \circ f_o^{(-1)})) = \int_{(\mathcal{Y} \circ f_o^{(-1)})} p_{\mathbf{f}}(h).d_L\mathcal{G}(h)$$

With the change of variable  $h = g \circ f_o^{(-1)}$ , we have

$$\Pr((\mathbf{f} \circ f_o) \in \mathcal{Y}) = \int_{\mathcal{Y}} p_{\mathbf{f}}(g \circ f_o^{(-1)}) \cdot d_L\mathcal{G}(g \circ f_o^{(-1)}) = \int_{\mathcal{Y}} \Delta\mathcal{G}(g \circ f_o^{(-1)}) \cdot p_{\mathbf{f}}(g \circ f_o^{(-1)}) \cdot d_R\mathcal{G}(g \circ f_o^{(-1)})$$

Since  $d_R\mathcal{G}$  is the right invariant measure, we can simplify and come back to the left invariant measure:

$$\Pr((\mathbf{f} \circ \mathbf{f}_o) \in \mathcal{Y}) = \int_{\mathcal{Y}} \Delta\mathcal{G}(\mathbf{g} \circ \mathbf{f}_o^{(-1)}) \cdot p_{\mathbf{f}}(\mathbf{g} \circ \mathbf{f}_o^{(-1)}) \cdot d_R\mathcal{G}(\mathbf{g}) = \int_{\mathcal{Y}} \frac{\Delta\mathcal{G}(\mathbf{g} \circ \mathbf{f}_o^{(-1)})}{\Delta\mathcal{G}(\mathbf{g})} \cdot p_{\mathbf{f}}(\mathbf{g} \circ \mathbf{f}_o^{(-1)}) \cdot d_L\mathcal{G}(\mathbf{g})$$

This gives the second result. ■

### Theorem 8 Composition of random transformations

Let  $\mathbf{f}_1$  and  $\mathbf{f}_2$  be two random transformations of pdf  $p_{\mathbf{f}_1}$  and  $p_{\mathbf{f}_2}$ . The pdf of their composition is:

$$p_{(\mathbf{f}_1 \circ \mathbf{f}_2)}(\mathbf{f}) = \int_{\mathcal{G}} p_{\mathbf{f}_1}(\mathbf{g}) \cdot p_{\mathbf{f}_2}(\mathbf{g}^{(-1)} \circ \mathbf{f}) \cdot d_L\mathcal{G}(\mathbf{g}) \quad (36)$$

**Proof:** Let  $\mathcal{F} \subset \mathcal{G}$  be a set of transformation. The set of couples  $(\mathbf{g}_1, \mathbf{g}_2)$  such that  $\mathbf{g}_1 \circ \mathbf{g}_2 \in \mathcal{F}$  is  $\mathcal{A} = \left\{ (\mathbf{g}_1, \mathbf{g}_2) / \mathbf{g}_1 \in \mathcal{G}, \mathbf{g}_2 \in \mathbf{g}_1^{(-1)} \circ \mathcal{F} \right\}$ . Thus, the probability of  $\mathbf{f} = \mathbf{f}_1 \circ \mathbf{f}_2$  to be in  $\mathcal{F}$  is:

$$\begin{aligned} \Pr(\mathbf{f}_1 \circ \mathbf{f}_2 \in \mathcal{F}) &= \int_{\mathcal{G}} \left( \int_{\mathbf{g}_1^{(-1)} \circ \mathcal{F}} p_{\mathbf{f}_2}(\mathbf{g}_2) \cdot d_L\mathcal{G}(\mathbf{g}_2) \right) p_{\mathbf{f}_1}(\mathbf{g}_1) \cdot d_L\mathcal{G}(\mathbf{g}_1) \\ &= \int_{\mathcal{G}} \left( \int_{\mathcal{F}} p_{\mathbf{f}_2}(\mathbf{g}_1^{(-1)} \circ \mathbf{g}_2) \cdot d_L\mathcal{G}(\mathbf{g}_2) \right) p_{\mathbf{f}_1}(\mathbf{g}_1) \cdot d_L\mathcal{G}(\mathbf{g}_1) \\ &= \int_{\mathcal{F}} \left( \int_{\mathcal{G}} p_{\mathbf{f}_2}(\mathbf{g}_1^{(-1)} \circ \mathbf{g}_2) \cdot p_{\mathbf{f}_1}(\mathbf{g}_1) \cdot d_L\mathcal{G}(\mathbf{g}_1) \right) d_L\mathcal{G}(\mathbf{g}_2) \end{aligned}$$

■

**Analogy with the convolution product** Equation (36) is remarkably similar to the classical convolution product obtain for the addition of two random vectors in  $\mathbb{R}^n$ :

$$p_{(\mathbf{x}+\mathbf{y})}(z) = p_{\mathbf{x}} \otimes p_{\mathbf{y}}(z) = \int_{\mathbb{R}^n} p_{\mathbf{x}}(t) \cdot p_{\mathbf{y}}(z - t) \cdot dt$$

This formula can be found as a particular case of equation (36) for the group of translations of  $\mathbb{R}^n$  for which we have  $x^{(-1)} = -x$ ,  $x \circ y = x + y$  and  $d_L\mathcal{G}(x) = dx$ .

### Theorem 9 (Inversion of a random transformation)

Let  $\mathbf{f}$  be a random transformation of density  $p_{\mathbf{f}}$ . The density of the inverse transformation is:

$$p_{\mathbf{f}^{(-1)}}(\mathbf{g}) = \Delta\mathcal{G}(\mathbf{g}^{(-1)}) \cdot p_{\mathbf{f}}(\mathbf{g}^{(-1)}) \quad (37)$$

**Proof:** The first step is to determine the relation between  $d_L\mathcal{G}(\mathbf{g}^{(-1)})$  and  $d_L\mathcal{G}(\mathbf{g})$ . We have:

$$\begin{aligned} d_L\mathcal{G}(\mathbf{g}^{(-1)}) &= |J_L(\mathbf{g}^{(-1)})|^{(-1)} \cdot d(\mathbf{g}^{(-1)}) = |J_L(\mathbf{g}^{(-1)})|^{(-1)} \cdot \left| \frac{\partial \mathbf{g}^{(-1)}}{\partial \mathbf{g}} \right| \cdot d\mathbf{g} = \left| \frac{\partial \mathbf{g}}{\partial \mathbf{g}^{(-1)}} \cdot \frac{\partial \mathbf{g}^{(-1)} \circ \mathbf{e}}{\partial \mathbf{e}} \right|_{\mathbf{e}=\text{Id}}^{(-1)} \cdot d\mathbf{g} \\ &= \left| \frac{\partial \mathbf{e}^{(-1)} \circ \mathbf{g}}{\partial \mathbf{e}} \right|_{\mathbf{e}=\text{Id}}^{(-1)} \cdot d\mathbf{g} = |J_R(\mathbf{g})|^{(-1)} \cdot \left| \frac{\partial \mathbf{e}^{(-1)}}{\partial \mathbf{e}} \right|_{\mathbf{e}=\text{Id}}^{(-1)} \cdot d\mathbf{g} = \left| \frac{J_L(\mathbf{g})}{J_R(\mathbf{g})} \right| \cdot d_L\mathcal{G}(\mathbf{g}) \end{aligned}$$

the last equality being due to  $\partial \mathbf{e}^{(-1)} / \partial \mathbf{e} = -\text{Id}$  at  $\mathbf{e} = \text{Id}$ . Thus, we obtain:

$d_L\mathcal{G}(\mathbf{g}^{(-1)}) = \Delta\mathcal{G}(\mathbf{g}^{(-1)}) \cdot d_L\mathcal{G}(\mathbf{g})$ . Now, if  $\mathcal{F}$  is a set of transformations, we have:

$$\begin{aligned} \Pr(\mathbf{f}^{(-1)} \in \mathcal{F}) &= \Pr(\mathbf{f} \in \mathcal{F}^{(-1)}) = \int_{\mathcal{F}^{(-1)}} p_{\mathbf{f}}(\mathbf{g}) \cdot d_L\mathcal{G}(\mathbf{g}) \\ &= \int_{\mathcal{F}} p_{\mathbf{f}}(\mathbf{g}^{(-1)}) \cdot d_L\mathcal{G}(\mathbf{g}^{(-1)}) = \int_{\mathcal{F}} \Delta\mathcal{G}(\mathbf{g}^{(-1)}) p_{\mathbf{f}}(\mathbf{g}^{(-1)}) \cdot d_L\mathcal{G}(\mathbf{g}) \end{aligned}$$

which gives the proposed density. ■

### 3.4 Obtaining the Karcher mean values

The Karcher means are minimising the variance function defined with the left invariant distance as defined in the previous section. There are few simplifications due to the left invariant structure: thanks to formula (32), we can write:

$$\mathbf{E} \left[ \vec{g\mathbf{f}} \right] = J_L(\vec{g}) \cdot \int_{G/C(g)} (\vec{g}^{(-1)} \circ \vec{f}) \cdot dP(\mathbf{f})$$

Thus, assuming that the conditions of theorem (2) (i.e. that the variance is finite and the cut locus has a null measure), our gradient descent algorithm to obtain the mean can be entirely written in the principal chart (we focus here on the practical case: the empirical mean value):

$$\vec{g}_{t+1} = \vec{g}_t \circ \left( \frac{1}{n} \sum_i \vec{g}_t^{(-1)} \circ \vec{f}_i \right) \quad (38)$$

which simply corresponds to the following algorithm: (left) translate all measured values to the origin using the inverse of the current transformation estimation; compute the barycenter of the transformations in the principal chart; use that barycenter as the (right) increment for the current value.

**Example with 3D rigid rotations** We have shown in [Pennec, 1996] that the principal chart is the rotation vector (the unit axis of rotation multiplied by the angle of rotation) with an norm (angle) less than  $\pi$ . Thus, the mean rotation vector satisfies<sup>7</sup>:

$$\mathbf{E} \left[ \vec{r}^{(-1)} \circ \mathbf{r} \right] = \mathbf{0}$$

If the mean rotation  $\vec{r}$  is sufficiently centered in the principal chart (i.e. if the angle of rotation is small), and if the distribution is sufficiently peaked, then a Taylor expansion shows that the classical expectation is an approximation of the Karcher one:

$$\vec{r}^{(-1)} \circ \mathbf{r} = \mathbf{r} - \vec{r} + O(\|\vec{r}\|^2) + O(\|\vec{r} - \mathbf{r}\|^2)$$

However, as soon as the Karcher mean value goes farther from the origin, the difference with the classical expectation is more and more visible (see figure 4).

### 3.5 Properties of the Karcher expectation for the group operations

Let us now turn to the properties of the mean transformations for the group operations.

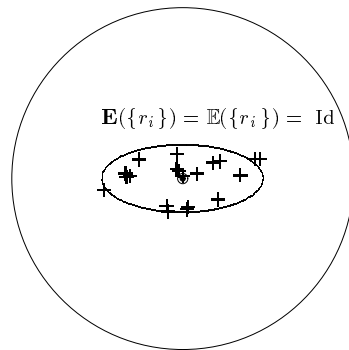
**Theorem 10 (Left translation of a random transformation)**

$$\mathbf{E} [ \mathbf{g} \circ \mathbf{f} ] = \mathbf{g} \circ \mathbf{E} [ \mathbf{f} ] \quad \text{and} \quad \mathbf{E} [ \{ \mathbf{g} \circ \mathbf{f}_i \} ] = \mathbf{g} \circ \mathbf{E} [ \{ \mathbf{f}_i \} ] \quad (39)$$

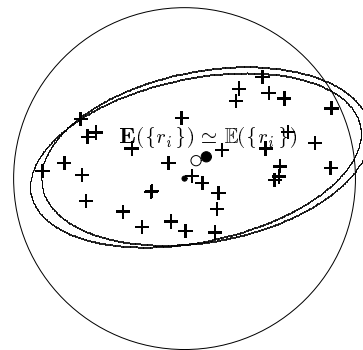
*This result also holds for the sets of central primitives of any order.*

---

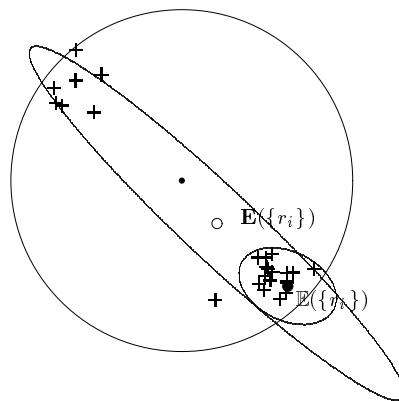
<sup>7</sup>the formulas to compute the composition, inversion and all the necessary derivatives are provided in [Pennec, 1996].



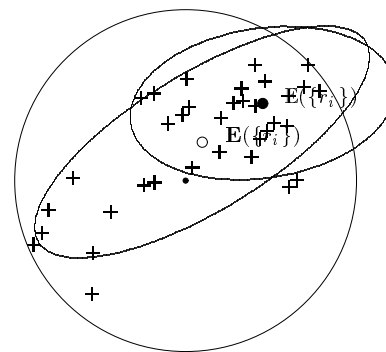
(a) Standard and Fréchet expectations are identical at the origin.



(b) When the Fréchet expectation is near the origin, the standard expectation can be considered as a first order approximation.



(c) Near the Domain boundary, standard and Fréchet expectations greatly differ. We note that with this relatively small noise, it would be possible to detect the boundary effect.



(d) With a greater amount of noise, it is no longer possible to guess the correct clustering to avoid the boundary effect.

Figure 4: Behaviour of expected rotations and their covariance matrices: projection of the measured rotation vectors, their standard and Fréchet expectation and the corresponding uncertainty ellipsoid at  $\chi^2 = 15$  onto the  $(r_x, r_y)$  plane. The circle represent the domain boundary for the rotation vector ( $\theta = \|r\| = \pi$ ). Remember that when we cross this boundary on one side at point  $r = \theta.n$ , we reenter the domain at the symmetric point  $r' = -\theta.n$ . In fact, the write way to visualise the covariance matrices would be to left-translate all measurements with respect to the “mean” point, as in figure (4(a)) so that the representation corresponds to the exponential chart at this point.

**Proof:** Thanks to the (left) invariant distance, we have:

$$\sigma_{g \circ f}^2(h) = \mathbf{E} [ \text{dist}(g \circ f, h)^2 ] = \mathbf{E} [ \text{dist}(f, g^{(-1)} \circ h)^2 ] = \sigma_f^2(g^{(-1)} \circ h)$$

Thus,  $\bar{f}$  minimises  $\sigma_f^2(f)$  if and only if  $\bar{h} = g \circ \bar{f}$  minimises  $\sigma_{g \circ f}^2(h)$ . The same argument holds for other types of central features (since only the invariance of the distance was used) and for the empirical mean. ■

**Theorem 11 (Right translation and inversion of a random transformation)**

*If the distance is left and right invariant, we have:*

$$\mathbf{E} [ \mathbf{f} \circ g ] = \mathbf{E} [ \mathbf{f} ] \circ g \quad \text{and} \quad \mathbf{E} [ \{f_i\} \circ g ] = \mathbf{E} [ \{f_i\} ] \circ g \quad (40)$$

$$\mathbf{E} [ \mathbf{f}^{(-1)} ] = \mathbf{E} [ \mathbf{f} ]^{(-1)} \quad \text{and} \quad \mathbf{E} [ \{f_i\}^{(-1)} ] = \mathbf{E} [ \{f_i\} ]^{(-1)} \quad (41)$$

*These results also holds for the sets of central primitives of any order.*

**Proof:** The first result is obtained as in theorem (10) but using the right invariance. For the inversion, let us first notice that

$$\text{dist}(g^{(-1)}, f) = \text{dist}(\text{Id}, g \circ f) = \text{dist}(f^{(-1)}, g)$$

where the first equality is obtained using the left invariance of the distance and the second one using the right invariance. Thus, we have

$$\sigma_{f^{(-1)}}^2(g) = \mathbf{E} [ \text{dist}(f^{(-1)}, g)^2 ] = \mathbf{E} [ \text{dist}(f, g^{(-1)})^2 ] = \sigma_f^2(g^{(-1)})$$

The result is then obtained as before. The same argument holds for other types of central features (since only the invariance of the distance was used) and for the empirical mean. ■

Since a compact group has a left and right invariant distance, the stability of the set of Karcher means by the group operations is ensured. It may seem weird that this is generally not the case for a non compact group. We present below an example on 2D rigid transformations where the stability is observed for the rotation part (a compact group) but not for the translation part which is responsible for the non compactity.

**Example with 2D rigid transformations** A 2D rigid transformation  $f$  is characterised by an angle of rotation  $\theta \in [-\pi; \pi[$  and a translation vector  $t \in \mathbb{R}^2$ . The composition of two transformations is given by:

$$f_1 \circ f_2 = \begin{cases} (\theta_1 + \theta_2) [2\pi] \\ R(\theta_1).t_2 + t_1 \end{cases} \quad \text{with} \quad R(\theta) = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix}$$

The modulus on the angle mean that we keep the value between  $-\pi$  and  $+\pi$ .

Excluding the rotations of  $\pm\pi$ , this representation is the principal chart for the left invariant distance. Indeed, we have:

$$J_L(f) = \left. \frac{\partial(f \circ e)}{\partial e} \right|_{e=\text{Id}} = \begin{bmatrix} 1 & 0 \\ 0 & R(\theta) \end{bmatrix}$$

Choosing an identity metric at the identity, the representation of the left invariant metric is  $G_L(f) = J_L(f)^{(-T)}.J_L(f)^{(-1)} = \text{Id}$ . Thus, the Christoffel symbols vanish and the geodesics starting from 0 (the identity) are straight lines.

Now, let  $f_1 = (\pi/4; -\sqrt{2}/2; \sqrt{2}/2)$ ,  $f_2 = (0; \sqrt{2}; 0)$  and  $f_3 = (-\pi/4; -\sqrt{2}/2; -\sqrt{2}/2)$  be three transformations in our principal chart. One easily verify that their barycenter is the origin. Thus,  $\bar{f} = (0; 0; 0) = \text{Id}$  is a Karcher mean. In this case, this is the only one:

$$\mathbb{E}[f_1, f_2, f_3] = \{(0; 0; 0)\}$$

The inverse of these transformations is easy to compute:  $f_1^{(-1)} = (-\pi/4; 0; -1)$ ,  $f_2^{(-1)} = (0; -\sqrt{2}; 0)$  and  $f_3^{(-1)} = (+\pi/4; 0; 1)$ . This time, the barycenter  $\bar{f}^{(-1)} = (0; -\sqrt{2}/3; 0)$  is null for the rotation part, but not for the translation: the mean value of the inverse transformations is not the identity

$$\mathbb{E}[f_1^{(-1)}, f_2^{(-1)}, f_3^{(-1)}] = \{(0; -\sqrt{2}/3; 0)\}$$

### 3.6 A word on the propagation of the Normal law parameters

#### Theorem 12 (Left translation of a Normal random transformation)

If  $\mathbf{f}$  is a random transformation following the Normal law  $N_{(\bar{\mathbf{f}}, \Gamma_{\mathbf{f}})}$  (in the conditions of theorem 3), then  $\mathbf{g} \circ \mathbf{f}$  is a random transformation following the Normal law  $N_{(\mathbf{g} \circ \bar{\mathbf{f}}, \Gamma_{\mathbf{g} \circ \mathbf{f}})}$  with

$$\Gamma_{\mathbf{g} \circ \mathbf{f}} = J^{(-T)} \cdot \Gamma_{\mathbf{f}} \cdot J^{(-1)} \quad \text{where} \quad J = \left. \frac{\partial(\bar{\mathbf{g}} \circ \bar{\mathbf{f}})}{\partial \bar{\mathbf{f}}} \right|_{\bar{\mathbf{f}} = \bar{\mathbf{f}}}$$

Moreover, the covariances matrices are related by  $\Sigma_{\mathbf{g} \circ \mathbf{f}} = J \cdot \Sigma_{\mathbf{f}} \cdot J^T$ .

**Proof:** According to eq. (34), the random transformation  $\mathbf{g} \circ \mathbf{f}$  has the pdf:

$$p(\mathbf{h}) = N_{(\bar{\mathbf{f}}, \Gamma_{\mathbf{f}})}(\mathbf{g}^{(-1)} \circ \mathbf{h}) = k(\bar{\mathbf{f}}, \Gamma_{\mathbf{f}}) \cdot \exp\left(-\frac{1}{2} \cdot \left(\bar{\mathbf{f}}^{(-1)} \circ (\bar{\mathbf{g}}^{(-1)} \circ \mathbf{h})\right)^T \cdot J_L(\bar{\mathbf{f}})^T \cdot \Gamma_{\mathbf{f}} \cdot J_L(\bar{\mathbf{f}}) \cdot \left(\bar{\mathbf{f}}^{(-1)} \circ (\bar{\mathbf{g}}^{(-1)} \circ \mathbf{h})\right)\right)$$

But we have

$$J_L(\mathbf{g} \circ \bar{\mathbf{f}}) = \left. \frac{\partial(\bar{\mathbf{g}} \circ \bar{\mathbf{f}}) \circ \bar{\mathbf{e}}}{\partial \bar{\mathbf{e}}} \right|_{\bar{\mathbf{e}} = \text{Id}} = \left. \frac{\partial(\bar{\mathbf{g}} \circ \bar{\mathbf{f}})}{\partial \bar{\mathbf{f}}} \right|_{\bar{\mathbf{f}} = \bar{\mathbf{f}}} \cdot \left. \frac{\bar{\mathbf{f}} \circ \bar{\mathbf{e}}}{\partial \bar{\mathbf{e}}} \right|_{\bar{\mathbf{e}} = \text{Id}} = J \cdot J_L(\bar{\mathbf{f}})$$

with the definition of  $J$  given in the theorem. Thus, we have  $J_L(\bar{\mathbf{f}}) = J^{(-1)} \cdot J_L(\mathbf{g} \circ \bar{\mathbf{f}})$  and

$$p(\mathbf{h}) = k(\bar{\mathbf{f}}, \Gamma_{\mathbf{f}}) \cdot \exp\left(-\frac{1}{2} \cdot \left((\bar{\mathbf{g}} \circ \bar{\mathbf{f}})^{(-1)} \circ \mathbf{h}\right)^T \cdot J_L(\mathbf{g} \circ \bar{\mathbf{f}})^T \cdot \Gamma_{\mathbf{g} \circ \mathbf{f}} \cdot J_L(\mathbf{g} \circ \bar{\mathbf{f}}) \cdot \left((\bar{\mathbf{g}} \circ \bar{\mathbf{f}})^{(-1)} \circ \mathbf{h}\right)\right)$$

Now, since the measure is left invariant, the change of variable  $\mathbf{h} \rightarrow \mathbf{g}^{(-1)} \circ \mathbf{h}$  in the definition of the normalisation coefficient gives (eqn. 25):

$$k(\bar{\mathbf{f}}, \Gamma_{\mathbf{f}})^{(-1)} = \int_G \exp\left(-\frac{1}{2} \cdot \left(\bar{\mathbf{f}}^{(-1)} \circ \mathbf{h}\right)^T \cdot J_L(\bar{\mathbf{f}})^T \cdot \Gamma_{\mathbf{f}} \cdot J_L(\bar{\mathbf{f}}) \cdot \left(\bar{\mathbf{f}}^{(-1)} \circ \mathbf{h}\right)\right) \cdot d_L \mathcal{G}(\mathbf{h}) = k(\mathbf{g} \circ \bar{\mathbf{f}}, \Gamma_{\mathbf{g} \circ \mathbf{f}})$$

Thus, the pdf  $p$  is the Normal law of mean  $(\mathbf{g} \circ \bar{\mathbf{f}})$  and concentration matrix  $\Gamma_{\mathbf{g} \circ \mathbf{f}}$ . The relation between the covariance matrices relies on the fact that

$$\overrightarrow{(\mathbf{g} \circ \bar{\mathbf{f}}) \mathbf{h}} = J_L(\bar{\mathbf{g}} \circ \bar{\mathbf{f}}) \cdot \left( \left(\bar{\mathbf{g}} \circ \bar{\mathbf{f}}\right)^{(-1)} \circ \bar{\mathbf{h}} \right) = J \cdot J_L(\bar{\mathbf{f}}) \cdot \left( \bar{\mathbf{f}}^{(-1)} \circ \left(\bar{\mathbf{g}}^{(-1)} \circ \bar{\mathbf{h}}\right) \right) = J \cdot \overrightarrow{(\bar{\mathbf{f}}^{(-1)} \circ (\bar{\mathbf{g}}^{(-1)} \circ \mathbf{h}))}$$

Then, the change of variable  $\mathbf{h} \rightarrow \mathbf{g}^{(-1)} \circ \mathbf{h}$  in the definition of the covariance matrix of the random transformation  $\mathbf{g} \circ \mathbf{f}$  (eqn. 26) gives the result. ■

One could obtain similar results for the right composition and the inversion in the case of a left and right invariant distance. The proof for the right translation would be obtained by the right equivalent of the construction we have done with the left invariant distance. The proof for the inversion is more tricky, but the properties of the left and right invariant metric should be sufficient.

## 4 The case of connected homogeneous manifolds

A manifold on which acts a transformation group is said to be *transitive* or *homogeneous* with respect to this group if, for any two points of the manifold, there exists a transformation of the group that transform one point to the other. Let  $o$  be a particular point of the manifold that we call *the origin* and let us denote by  $g \star x$  the action of a transformation  $g$  on a point  $x$  of the manifold. The manifold is homogeneous if  $\mathcal{G} \star o = \{g \star o / g \in \mathcal{G}\} = \mathcal{M}$ .

### 4.1 Invariant Riemannian metric

In the case of an homogeneous manifold, we can identify the elements  $x \in \mathcal{M}$  with subsets of of the group  $\mathcal{G}$ . Let  $\mathcal{H}$  be the *isotropy subgroup* of the origin, i.e. the set of transformations that leave it unchanged:

$$\mathcal{H} = \{h \in \mathcal{G} / h \star o = o\}$$

Let now  $f_x$  be a transformation that bring the origin into the point  $x = f_x \star o$ . The set of all transformation that bring the origin into the point  $x$  is the left translation of  $\mathcal{H}$  by  $f_x$ :

$$\mathcal{F}_x = \{g \in \mathcal{G} / g \star o = x\} = f_x \circ \mathcal{H}$$

These sets are in fact the elements of the quotient space  $\mathcal{G}/\mathcal{H}$  are thus called the *coset*  $\mathcal{F}_x$  of  $x$ .

To compare the tangent spaces  $T_o\mathcal{M}$  and  $T_x\mathcal{M}$  at the origine and at point  $x$ , we now have a whole set of transformations from  $\mathcal{F}_x$ . Let us chose a *placement function*  $f_x$ , i.e. a representant of each coset. In a local coordinate system, the vector  $v \in T\mathcal{M}$  is transformed by the placement function into

$$v_x = J(f_x).v \quad \text{with} \quad J(f_x) = \left. \frac{\partial(f_x \star e)}{\partial e} \right|_{e=o}$$

This transportation can be used to define the dot product from the tangent space at the origin on each tangent space  $T_x\mathcal{M}$ :

$$\langle y_1 | y_2 \rangle_x = \langle J(f_x)^{(-1)}.y_1 | J(f_x)^{(-1)}.y_2 \rangle \quad \text{avec} \quad J(f_x) = \left. \frac{\partial(f_x \star e)}{\partial e} \right|_{e=o}$$

Thus, the expression of the metric in a local chart is

$$G(x) = J(f_x)^{(-T)}.Q.J(f_x)^{(-1)} \quad (42)$$

A necessary condition to obtain an invariant metric is obviously the stability of the metric with respect to the choice of the placement function. From the definition of the cosets, this reduces to the invariance of the metric  $Q$  at the origin by the action of the isotropy group  $\mathcal{H}$ :

$$\forall h \in \mathcal{H} \quad J(h)^T.Q.J(h) = Q \quad (43)$$

The dot product  $\langle \cdot | \cdot \rangle_x$  is in this case a continuous function of  $x$ : we have obtained an invariant metric (the condition above is sufficient). Note that there does not always exists an invariant metric: think for instance to points under the affine group.

Assuming that there exists an invariant metric, the associated infinitesimal volume element is given by:

$$d\mathcal{M}(x) = \sqrt{|G(x)|}.dx = \lambda. \frac{dx}{|J(f_x)|} \quad \text{with} \quad \lambda = \sqrt{|Q|} \quad \text{and} \quad f_x \in \mathcal{F}_x$$

This measure was the invariant measure determined in [Pennec and Ayache, 1998]. However, we have a stronger constraint on the existance: an invariant distance implies an invariant measure, but the converse is false.

## 4.2 Principal chart

In the sequel, we assume that there exists an invariant distance and that the manifold is geodesically complete for this metric. We call principal chart the exponential chart at the origin with  $Q = \text{Id}$ . The action of a transformation is defined in this chart by  $f \star \vec{x} \equiv \log(f \star x) = \log(f \star \exp(\vec{x}))$ .

Since geodesics are globally invariant by the action of transformation, we can easily express the exponential chart at any other point. If  $\vec{xy}$  is a tangent vector of  $T_x\mathcal{M}$  expressed in the basis induced by the principal chart, we have:

$$\exp_{\vec{x}}(\vec{xy}) = f_{\vec{x}} \star \exp(J(f_{\vec{x}})^{(-1)} \cdot \vec{y})$$

The inverse function gives (in the star-shaped domain delimited by the cut-locus):

$$\vec{xy} = \log_{\vec{x}}(y) = J(f_{\vec{x}}) \cdot (f_{\vec{x}}^{(-1)} \star \vec{y}) \quad (44)$$

The functions are independent of the chosen placement function. Indeed, let  $f'_x$  be another placement function. Then, by definition of the cosets, for each point  $x$  there exists a transformation  $h_x \in \mathcal{H}$  such that  $f'_x = f_x \circ h_x$ . Now, the geodesics starting from the origin are straight lines in the exponential chart and are globally invariant: they are transformed into straight lines in the exponential chart by the isotropy group. This means that the action of the isotropy group is linear in the principal chart: we can write  $h \star \vec{x} = J(h) \cdot \vec{x}$ . Thus,  $f_x^{(-1)} \star \vec{y} = h_x^{(-1)} \star (f_x \star \vec{y}) = J(h_x)^{(-1)} \cdot (f_x \star \vec{y})$ . Since  $J(f'_x) = J(f_x) \cdot J(h_x)$ , we get the sought equality  $J(f'_x) \cdot (f_x^{(-1)} \star \vec{y}) = J(f'_x) \cdot (f_x^{(-1)} \star \vec{y})$ . The invariance of the exponential is similar.

The distance between two points  $x$  and  $y$  is by definition the length of the minimizing geodesic joining them, and by definition of the exponential chart the length of  $\vec{xy}$ :

$$\text{dist}(x, y)^2 = \|\log_{\vec{x}}(y)\|_{\vec{x}}^2 = \vec{xy}^T \cdot Q(\vec{x}) \cdot \vec{xy} = (f_x^{(-1)} \star \vec{y})^T \cdot (f_x^{(-1)} \star \vec{y}) \quad (45)$$

## 4.3 Propagation of pdfs

### Theorem 13 (Action of a deterministic transformation on a random feature)

Let  $\mathbf{x}$  be a random feature of pdf  $p_{\mathbf{x}}$  and  $f \in \mathcal{G}$  a deterministic transformation. Then, the pdf of the random feature  $f \star \mathbf{x}$  is

$$p_{(f \star \mathbf{x})}(y) = p_{\mathbf{x}}(f^{(-1)} \star y) \quad (46)$$

**Proof:** The probability of  $\mathbf{y} = f \star \mathbf{x}$  to be in a set  $\mathcal{Y}$  is

$$P(f \star \mathbf{x} \in \mathcal{Y}) = P(\mathbf{x} \in f^{(-1)} \star \mathcal{Y}) = \int_{(f^{(-1)} \star \mathcal{Y})} p_{\mathbf{x}}(z) \cdot d\mathcal{M}(z) = \int_{\mathcal{Y}} p_{\mathbf{x}}(f^{(-1)} \star y) \cdot d\mathcal{M}(y)$$

The last equality is obtained using the change of variable  $z = f^{(-1)} \star y$  and the invariance of the measure. ■

### Theorem 14 (Action of a random transformation on the origin)

Let  $\mathbf{f}$  be a random transformation of pdf  $p_{\mathbf{f}}$ . Its action on the origin  $\circ$  determines a random feature  $\mathbf{y} = \mathbf{f} \star \circ$  that has the pdf

$$p_{(\mathbf{f} \star \circ)}(y) = \int_{\mathcal{H}} p_{\mathbf{f}}(f_y \circ h) \cdot d\mathcal{H}(h) \quad (47)$$

**Proof:** The probability of  $\mathbf{x} = \mathbf{f} \star \circ$  to be in a set  $\mathcal{X} \in \mathcal{M}$  is:

$$\Pr(\mathbf{f} \star \circ \in \mathcal{X}) = \int_{\mathcal{F}_{\mathcal{X}}} p_{\mathbf{f}}(\mathbf{g}).d\mathcal{G}(\mathbf{g}) \quad \text{avec} \quad \mathcal{F}_{\mathcal{X}} = \{\mathbf{g} \in \mathcal{G} / \mathbf{g} \star \circ \in \mathcal{X}\}$$

Using a placement function, the set of transformations  $\mathcal{F}_{\mathcal{X}}$  is  $\mathcal{F}_{\mathcal{X}} = \{\mathbf{f}_x \circ \mathbf{h} / \mathbf{x} \in \mathcal{X}, \mathbf{h} \in \mathcal{H}\}$ . Thus:

$$\Pr(\mathbf{x} \in \mathcal{X}) = \int_{\mathcal{X}, \mathcal{H}} p_{\mathbf{f}}(\mathbf{f}_x \circ \mathbf{h}).d\mathcal{M}(\mathbf{x}).d\mathcal{H}(\mathbf{h}) = \int_{\mathcal{X}} \left( \int_{\mathcal{H}} p_{\mathbf{f}}(\mathbf{f}_x \circ \mathbf{h}).d\mathcal{H}(\mathbf{h}) \right).d\mathcal{M}(\mathbf{x})$$

■

### Theorem 15 (Action of a random transformation of a deterministic feature)

Let  $\mathbf{f}$  be a random transformation of pdf  $p_{\mathbf{f}}$ . Its action on the deterministic feature  $\mathbf{x}$  determines a random feature  $\mathbf{y} = \mathbf{f} \star \mathbf{x}$  that has the pdf:

$$p_{(\mathbf{f} \star \mathbf{x})}(\mathbf{y}) = \int_{\mathcal{H}} p_{(\mathbf{f} \circ \mathbf{f}_x)}(\mathbf{f}_y \circ \mathbf{h}).d\mathcal{H}(\mathbf{h}) = \int_{\mathcal{H}} \frac{\Delta\mathcal{G}(\mathbf{f}_y \circ \mathbf{h} \circ \mathbf{f}_x^{(-1)})}{\Delta\mathcal{G}(\mathbf{f}_y \circ \mathbf{h})}.p_{\mathbf{f}}(\mathbf{f}_y \circ \mathbf{h} \circ \mathbf{f}_x^{(-1)}).d\mathcal{H}(\mathbf{h}) \quad (48)$$

**Proof:** Choosing  $\mathbf{f}_x \in \mathcal{F}_x$  and writing  $\mathbf{y} = \mathbf{f} \star \mathbf{x} = (\mathbf{f} \circ \mathbf{f}_x) \star \circ$ , we apply equation 35 to obtain the pdf of  $\mathbf{f} \circ \mathbf{f}_x$  and then theorem 47. ■

### Theorem 16 (Action of a random transformation on a random feature)

Let  $\mathbf{f}$  be a random transformation of pdf  $p_{\mathbf{f}}$  and  $\mathbf{x}$  be a random feature of pdf  $p_{\mathbf{x}}$ . The pdf of the random feature  $\mathbf{y} = \mathbf{f} \star \mathbf{x}$  is:

$$p_{\mathbf{f} \star \mathbf{x}}(\mathbf{y}) = \int_{\mathcal{G}} p_{\mathbf{f}}(\mathbf{g}).p_{\mathbf{x}}(\mathbf{g}^{(-1)} \star \mathbf{y}).d_L\mathcal{G}(\mathbf{g}) \quad (49)$$

Once again, this is very similar to a convolution product.

**Proof:** Let  $\mathcal{X} \subset \mathcal{M}$  be a set of features. The set  $\mathcal{A}$  of couples  $(\mathbf{f}, \mathbf{x})$  such that  $\mathbf{f} \star \mathbf{x} \in \mathcal{X}$  can be written:

$$\mathcal{A} = \{(\mathbf{f}, \mathbf{x}) / \mathbf{f} \in \mathcal{G}, \mathbf{x} \in \mathbf{f}^{(-1)} \star \mathcal{X}\}$$

Thus, the probability of  $\mathbf{y} = \mathbf{f} \star \mathbf{x}$  to be in  $\mathcal{X}$  is

$$\begin{aligned} \Pr(\mathbf{f} \star \mathbf{x} \in \mathcal{X}) &= \int_{\mathcal{G}} \left( \int_{\mathbf{g}^{(-1)} \star \mathcal{X}} p_{\mathbf{x}}(\mathbf{y}).d\mathcal{M}(\mathbf{y}) \right).p_{\mathbf{f}}(\mathbf{g}).d_L\mathcal{G}(\mathbf{g}) \\ &= \int_{\mathcal{G}} \left( \int_{\mathcal{X}} p_{\mathbf{x}}(\mathbf{g}^{(-1)} \star \mathbf{z}).d\mathcal{M}(\mathbf{z}) \right).p_{\mathbf{f}}(\mathbf{g}).d_L\mathcal{G}(\mathbf{g}) \\ &= \int_{\mathcal{X}} \left( \int_{\mathcal{G}} p_{\mathbf{x}}(\mathbf{g}^{(-1)} \star \mathbf{z}).p_{\mathbf{f}}(\mathbf{g}).d_L\mathcal{G}(\mathbf{g}) \right).d\mathcal{M}(\mathbf{z}) \end{aligned}$$

■

## 4.4 Obtaining the Karcher mean values

As for the Lie groups, we have few simplifications due to the left invariant structure in the Karcher means characterization: thanks to formula (44), we can write:

$$\mathbf{E} [\vec{\mathbf{y}\mathbf{x}}] = J(\mathbf{f}_{\vec{\mathbf{y}}}). \int_{\mathcal{M}/C(\mathbf{y})} (\mathbf{f}_{\vec{\mathbf{y}}}^{(-1)} \star \vec{\mathbf{z}}).dP(\vec{\mathbf{z}})$$

Thus, assuming that the conditions of theorem (2) (i.e. that the variance is finite and the cut locus has a null measure), our gradient descent algorithm to obtain the mean can be entirely written in the principal chart (we focus here on the practical case: the empirical mean value):

$$\bar{x}_{t+1} = f_{\bar{x}_t} \star \left( \frac{1}{n} \sum_i (f_{\bar{x}_t}^{(-1)} \star \vec{x}_i) \right) \tag{50}$$

which simply corresponds to the following algorithm: translate all measured values to the origin using the inverse of the placement function at the current estimation; compute the barycenter of the points in the principal chart; use that barycenter as the increment for the current value.

#### 4.5 Stability of the Fréchet expectation and Normal law

Let us now turn to the stability of the Fréchet expectation with respect to the action of transformations.

##### Theorem 17 (Action of a deterministic transformation on a random feature)

$$\mathbb{E}[g \star \mathbf{x}] = g \star \mathbb{E}[\mathbf{x}] \quad et \quad \mathbb{E}[\{g \star x_i\}] = g \star \mathbb{E}[\{x_i\}] \tag{51}$$

*This result also holds for the set of central features of any order.*

**Proof:** Let  $\mathbf{z} = g \star \mathbf{x}$  be the random feature obtained by the action of the deterministic transformation  $g$  on the random feature  $\mathbf{x}$ . We have:

$$\sigma_{\mathbf{z}}^2(y) = \mathbf{E}[\text{dist}(g \star \mathbf{x}, y)^2] = \int_{\mathcal{M}} \text{dist}(y, g \star \mathbf{x})^2 \cdot p_{\mathbf{x}}(\mathbf{x}) \cdot d\mathcal{M}(\mathbf{x})$$

Thanks to the invariance of the distance, we get:

$$\sigma_{\mathbf{z}}^2(y) = \int_{\mathcal{M}} \text{dist}(g^{(-1)} \star y, \mathbf{x})^2 \cdot p_{\mathbf{x}}(\mathbf{x}) \cdot d\mathcal{M}(\mathbf{x}) = \sigma_{\mathbf{x}}^2(g^{(-1)} \star y)$$

Thus, the feature  $\bar{x}$  minimizes  $\sigma_{\mathbf{x}}^2(\mathbf{x})$  if and only if  $\bar{z} = g \star \bar{x}$  minimizes  $\sigma_{\mathbf{z}}^2(\mathbf{z})$ , which may be rewritten  $\mathbb{E}[\mathbf{z}] = g \star \mathbb{E}[\mathbf{x}]$ . Moreover, variances are equal:  $\sigma_{\mathbf{z}} = \sigma_{\mathbf{x}}$ .

The same demonstration holds for the stability of the central feature of any order, as well as for the empirical mean feature. ■

However, if the transformation turns out to be probabilistic, then we have no simple expression for  $\mathbb{E}[\mathbf{f} \star \mathbf{x}]$ .

##### Theorem 18 (Action of a deterministic transformation on a normal random feature)

*Let  $\mathbf{x}$  be a random Normal feature of law  $N_{(\bar{x}, \Gamma_{\mathbf{x}})}$  (in the conditions of theorem 3), and  $\mathbf{f}$  be a deterministic transformation. Then,  $\mathbf{f} \star \mathbf{x}$  is a random Normal feature following the Normal law  $N_{(\mathbf{f} \star \bar{x}, \Gamma_{\mathbf{f} \star \mathbf{x}})}$  with*

$$\Gamma_{\mathbf{f} \star \mathbf{x}} = J^{(-T)} \cdot \Gamma_{\mathbf{x}} \cdot J^{(-1)} \quad where \quad J = \left. \frac{\partial(\mathbf{f} \star \vec{x})}{\partial \vec{x}} \right|_{\vec{x}=\bar{x}}$$

*Moreover, the covariances matrices are related by  $\Sigma_{\mathbf{f} \star \mathbf{x}} = J \cdot \Sigma_{\mathbf{x}} \cdot J^T$ .*

The demonstration of this theorem is essentially the same as the one for Lie groups (Section 3.6).

## 5 Discussion

On a (geodesically complete) Riemannian manifold, it is easy to define probability density functions associated to random features, thanks to the availability of a metric. However, as soon as the expectation is concerned, we may only define the expectation of an observable (a real or vectorial function of the random feature). Thus, the definition of a mean value for a random feature is much more complex than for the vectorial case and it requires a distance-based variational formulation: the Fréchet or Karcher expected features basically minimize globally (or locally) the variance. As the mean is now defined through a minimization procedure, its existence and uniqueness are not ensured any more (except under some specific conditions). In practice, one mean value almost always exists, and it is unique as soon as the distribution is sufficiently peaked. The properties of the mean are very similar to those of the modes (that can be defined as central Karcher values of order 0) in the vectorial case. To compute the mean value, we designed an original Gauss-Newton gradient descent algorithm that essentially alternates the computation of the barycenter in the exponential chart centered at the current estimation of the mean value, and a re-centering step of the chart at the point of the manifold that corresponds to the computed barycenter.

To define higher moments of the distribution, we used the exponential chart at the mean point (which may be seen as the development of the manifold onto its tangent space at this point along the geodesics): the random feature is thus represented as a random vector with null mean in a star-shaped and symmetric domain. With this representation, there is no more problem to define the covariance matrix and potentially higher order moments. Based on this covariance matrix, we defined a Mahalanobis distance between a random and a deterministic feature that basically weights the distance between the deterministic feature and the mean feature using the inverse of the covariance matrix. Interestingly, the expected Mahalanobis distance of a random primitive with itself is independent of the distribution and is equal to the dimension of the manifold, as in the vectorial case.

Like for the mean, we chose a variational approach to generalize the Normal law: we define it as the distribution that minimizes the information knowing the mean and the covariance. Neglecting the cut-locus constraints, we show that it amounts to consider a truncated Gaussian distribution on the exponential chart centered at the mean point. However, the relation between the concentration matrix (the “metric” used in the exponential of the pdf) and the covariance matrix is slightly more complex than the simple inversion of the vectorial case, as it has to be corrected for the curvature of the manifold.

Last but not least, using the Mahalanobis distance of a Normally distributed random feature, we can generalize the  $\chi^2$  law: we were able to show that it has the same density as in the vectorial case up to an order 3 in  $\sigma$ . This opens the way to the generalization of many other statistical tests, as we may expect similarly simple approximations for sufficiently centered distributions.

In this work, all definitions are derived from the Riemannian metric of the manifold. More generally, we could conceive other definitions of the mean value using the notion of connector introduced in [Picard, 1994]. This connector formalizes a relationship between the manifold and its tangent space at one point, exactly in the way we used the exponential map of the Riemannian metric. Thus, we could generalize easily the higher order moments and the other statistical operations we defined by replacing everywhere the exponential map with the connector. One important point is that these connectors can model extrinsic distances (like the Euclidean distance on unit vectors), and could lead to very efficient approximations of the mean value for sufficiently peaked distributions. For instance, the “barycenter / re-centering” algorithm we designed will most probably converge toward a first order approximation of the Riemannian mean if we use any chart that is consistent to the

exponential chart at the first order (e.g. Euler's angle on re-centered 3D rotations). We believe that this research track may become one of the most productive from the practical point of view.

When we deal with geometric problems, we often have transformation groups (Lie groups in our case) and manifolds that are subject to the action of such a transformation group. The question that arise is how to chose an adapted Riemannian metric on these structures? In the case of a Lie group, we have two canonical Riemannian structures given by the left and right invariant metrics. Choosing one of these metrics lead to important simplifications in the propagation of the pdf of random transformations. One may notice the remarkable similarity between the pdf of composition of two random transformations and the convolution product used to compute the pdf of the sum of two random vectors. However, the propagation formulas are not symmetric, and a Riemannian metric that is only left invariant only guarantees the stability of the mean, covariance (and Normal parameters) under the action of a deterministic transformation on the left. The stability by the action on both side and by inversion is only guaranteed if we have a left and right invariant metric. Unfortunately, such a bi-invariant metric does not exists in general for non-compact Lie groups. Is there another way to find a metric that leads to symmetric properties for non-compact Lie groups (or simply for locally compact such as rigid transformation)?

The second kind of geometrical objects we investigated is homogeneous manifolds, i.e. manifolds that have no invariant w.r.t. the action of a transformation group. The basic idea is to ensure the stability of our statistical operations w.r.t. the group action by choosing an appropriate (i.e. invariant) Riemannian metric. We show that a necessary condition is the invariance of the metric at any given point by the isotropy group of this point. Such a metric does not always exist (e.g. for points under the action of the linear or affine group). If it exists, then we can write quite simply the propagation of the pdf of a random feature subject to the action of a random or deterministic transformation, and show that the mean, covariance and Normal parameters are stable under the action of a deterministic transformation.

A third kind of important geometrical objects we have not investigated yet are the so-called "shape spaces", i.e. the quotient of a manifold by a transformation group. The difficult part here is that we have to remove the assumption of geodesical completeness and deal with manifolds with boundaries. Indeed, the shape manifolds often present some singularities [Small, 1996]: for instance, the shape of pairs of points (subject to rigid transformations) is the distance between them, which is positive but can reach the singular point zero. On a geometric aspect, the question arise whether we can find compatible metrics on the feature, transformation and shape spaces, i.e. such that the metric on the feature space be the direct product of the metrics on the shape and transformation spaces.

In conclusion, we believe that we could derive most of the interesting statistical operations on a finite-dimensional manifold from a Riemannian metric (or perhaps only a connector following [Picard, 1994]), even if there the theory still has to be generalized to manifolds with boundaries. In the case of geometric objects (Lie groups or homogeneous manifolds), requiring the invariance of the metric ensures the stability of our operations w.r.t. the structural operations (composition or action, inversion), but such a metric does not always exists. Would there be some weaker conditions than invariance that could still guaranty the stability of our operations? At that point, it seems to us that the choice of the Riemannian metric is the key point to investigate. Another interesting but difficult research direction would concern the generalization of our statistical framework to infinite-dimensional groups and transformations, such as manifolds of curves or surfaces, and groups of diffeomorphisms. For this purpose, the work of Younes at al [Younes, 1998, Miller and Younes, 2001], may provide a thorough basis.

## References

- [Arnaudon, 1994] Arnaudon, M. (1994). Espérances conditionnelles et  $c$ -martingales dans les variétés. In J. Azema, P.A. Meyer, M. Y., editor, *Séminaire de probabilités XXVIII*, volume 1583 of *Lect. Notes in Math.*, pages 300–311. Springer-Verlag.
- [Arnaudon, 1995] Arnaudon, M. (1995). Barycentres convexes et approximations des martingales continues dans les variétés. In J. Azema, P.A. Meyer, M. Y., editor, *Séminaire de probabilités XXIX*, volume 1613 of *Lect. Notes in Math.*, pages 70–85. Springer-Verlag.
- [Bingham, 1974] Bingham, C. (1974). An antipodally symmetric distribution on the sphere. *The Annals of Statistics*, 2(6):1201–1225.
- [Carmo, 1992] Carmo, M. d. (1992). *Riemannian Geometry*. Mathematics. Birkhäuser, Boston, Basel, Berlin.
- [Chavel, 1993] Chavel, I. (1993). *Riemannian geometry - A modern introduction*, volume 108 of *Cambridge tracts in mathematics*. Cambridge university press.
- [Doss, 1949] Doss, S. (1949). Sur la moyenne d'un élément aléatoire dans un espace distancié. *Bull. Sc. Math.*, 73:48–72.
- [Emery and Mokobodzki, 1991] Emery, M. and Mokobodzki, G. (1991). Sur le barycentre d'une probabilité dans une variété. In J. Azema, P.A. Meyer, M. Y., editor, *Séminaire de probabilités XXV*, volume 1485 of *Lect. Notes in Math.*, pages 220–233. Springer-Verlag.
- [Fréchet, 1944] Fréchet, M. (1944). L'intégrale abstraite d'une fonction abstraite d'une variable abstraite et son application à la moyenne d'un élément aléatoire de nature quelconque. *Revue Scientifique*, pages 483–512.
- [Fréchet, 1948] Fréchet, M. (1948). Les éléments aléatoires de nature quelconque dans un espace distancié. *Ann. Inst. H. Poincaré*, 10:215–310.
- [Herer, 1986] Herer, W. (1986). Espérance mathématique au sens de Doss d'une variable aléatoire à valeur dans un espace métrique. *C. R. Acad. Sc. Paris, Série I*, t.302(3):131–134.
- [Herer, 1988] Herer, W. (1988). Espérance mathématique d'une variable aléatoire à valeur dans un espace métrique à courbure négative. *C. R. Acad. Sc. Paris, Série I*, t.306:681–684.
- [Huber, 1981] Huber, P. (1981). *Robust Statistics*. John Wiley, New York.
- [Jupp and Mardia, 1989] Jupp, P. and Mardia, K. (1989). A unified view of the theory of directional statistics, 1975-1988. *Int. Statistical Review*, 57(3):261–294.
- [Karcher, 1977] Karcher, H. (1977). Riemannian center of mass and mollifier smoothing. *Comm. Pure Appl. Math.*, 30:509–541.
- [Kendall and Moran, 1963] Kendall, M. and Moran, P. (1963). *Geometrical probability*. Number 10 in Griffin's statistical monographs and courses. Charles Griffin & Co. Ltd.
- [Kendall, 1990] Kendall, W. (1990). Probability, convexity, and harmonic maps with small image i: uniqueness and fine existence. *Proc. London Math. Soc.*, 61(2):371–406.

- [Kent, 1992] Kent, J. (1992). *The art of Statistical Science*, chapter 10 : New Directions in Shape Analysis, pages 115–127. John Wiley & Sons. K.V. Mardia, ed.
- [Klingenberg, 1982] Klingenberg, W. (1982). *Riemannian Geometry*. Walter de Gruyter, Berlin, New York.
- [Maillot, 1997] Maillot, H. (1997). Différentielle de la variance et centrage de la plaque de coupure sur une variété riemannienne compacte. Communication personnelle.
- [Mardia, 1995] Mardia, K. (1995). Directional statistics and shape analysis. Research Report STAT95/24, University of Leeds, UK.
- [Miller and Younes, 2001] Miller, M. and Younes, L. (2001). Group actions, homeomorphisms, and matching: A general framework. *International Journal of Computer Vision*, 41(1/2):61–84.
- [Neveu, 1990] Neveu, J. (1990). *Introduction aux probabilités*. Ecole Polytechnique. Cours de l’Ecole Polytechnique.
- [Papoulis, 1991] Papoulis, A. (1991). *Probability, Random Variables, and Stochastic Processes*. McGraw-Hill, Inc.
- [Pelat, 1992] Pelat, D. (1992). *Bruits et Signaux*. Cours de l’Ecole doctorale d’astrophysique d’Ile de France.
- [Pennec, 1996] Pennec, X. (1996). *L’Incertitude dans les Problèmes de Reconnaissance et de Recalage – Applications en Imagerie Médicale et Biologie Moléculaire*. PhD thesis, Ecole Polytechnique, Palaiseau (France).
- [Pennec and Ayache, 1998] Pennec, X. and Ayache, N. (1998). Uniform distribution, distance and expectation problems for geometric features processing. *Journal of Mathematical Imaging and Vision*, 9(1):49–67.
- [Pennec and Thirion, 1997] Pennec, X. and Thirion, J.-P. (1997). A framework for uncertainty and validation of 3D registration methods based on points and frames. *Int. Journal of Computer Vision*, 25(3):203–229.
- [Picard, 1994] Picard, J. (1994). Barycentres et martingales sur une variété. *Annales de l’institut Poincaré - Probabilités et Statistiques*, 30(4):647–702.
- [Poincaré, 1912] Poincaré, H. (1912). *Calcul des probabilités*. 2nd edition, Paris.
- [Press et al., 1991] Press, W., Flannery, B., Teukolsky, S., and Vetterling, W. (1991). *Numerical Recipes in C*. Cambridge Univ. Press.
- [Rousseeuw and Leroy, 1987] Rousseeuw, P. and Leroy, A. (1987). *Robust Regression and Outliers Detection*. Wiley series in prob. and math. stat. J. Wiley and Sons.
- [Small, 1996] Small, C. (1996). *The Statistical Theory of Shapes*. Springer series in statistics. Springer.
- [Spivak, 1979] Spivak, M. (1979). *Differential Geometry*, volume 1. Publish or Perish, Inc., 2nd edition.
- [Younes, 1998] Younes, L. (1998). Computable elastic distances between shapes. *SIAM Journal on Applied Mathematics*, 58(2):565–586.

## A Gradient of the variance

This proof is a generalization of a differentiability proof for the uniform distribution on compact manifolds by Pr Maillot [Maillot, 1997]. One of the main difficulty was to remove the compactness hypothesis.

### Hypotheses

Let  $P$  be a probability on the Riemannian manifold  $\mathcal{M}$ . We assume that the cut locus  $C(y)$  of the derivation point  $y \in \mathcal{M}$  has a null measure with respect to this probability (as it has with the Riemannian measure) and that the variance is finite at that point:

$$P(C(y)) = \int_{C(y)} dP(z) = 0 \quad \text{and} \quad \sigma^2(y) = \int_{\mathcal{M}} \text{dist}(y, z)^2 .dP(z) < \infty \quad (52)$$

### Goal and problem

Let now  $\vec{g}(y) = \int_{\mathcal{M} \setminus C(y)} \vec{y}\vec{z} .dP(z)$ . As  $\|\vec{y}\vec{z}\| = d(z, y) \leq 1 + d(z, y)^2$ , and using the null probability of the cut locus, we have:

$$\|\vec{g}(y)\| \leq \int_{\mathcal{M} \setminus C(y)} \|\vec{y}\vec{z}\| .dP(z) \leq \int_{\mathcal{M} \setminus C(y)} (1 + d(z, y)^2) .dP(z) = 1 + \sigma^2(y) < \infty$$

Thus  $\vec{g}(y)$  is well defined everywhere. Let  $h(y) = d(z, y)^2 = \|\vec{y}\vec{z}\|^2$ . For a fixed  $z \notin C(y)$  (which is equivalent to  $y \notin C(z)$ ), we have  $(\text{grad } h)(y) = -2\vec{y}\vec{z}$ . Thus the proposition:

$$(\text{grad } \sigma^2)(y) = -2 \cdot \int_{\mathcal{M} \setminus C(y)} \vec{y}\vec{z} .dP(z)$$

corresponds to a derivation under the sum, but the usual conditions of the Lebesgue theorem are not fulfilled: the zero measure set  $C(y)$  varies with  $y$ . Thus, we have to come back to the original definition of the gradient.

### Definition of the gradient

Let  $\gamma(t)$  be a curve with  $\gamma(0) = y$  and  $\dot{\gamma}(0) = w$ . By definition, the function  $\sigma^2 : \mathcal{M} \rightarrow R$  is derivable if there exists a vector  $(\text{grad } \sigma^2)(y) \in T_y \mathcal{M}$  (the gradient) such that:

$$\forall w \in T_y \mathcal{M} \quad \langle (\text{grad } \sigma^2)(y) \mid w \rangle = \partial_w \sigma^2(y) = \lim_{t \rightarrow 0} \frac{\sigma^2(\gamma(t)) - \sigma^2(y)}{t}$$

Since tangent vectors are defined as equivalent classes, we can choose the geodesic curve  $\gamma(t) = \text{exp}_y(t.w)$ . Using  $v = t.w$ , the above condition can then be rewritten:

$$\forall v \in T_y \mathcal{M} \quad \lim_{\|v\| \rightarrow 0} \frac{\sigma^2(\text{exp}_y(v)) - \sigma^2(y) - \langle (\text{grad } \sigma^2)(y) \mid v \rangle}{\|v\|} = 0$$

which can be rephrased as: for all  $\eta \in \mathbb{R}_+$ , there exists  $\varepsilon$  sufficiently small such that:

$$\forall v \in T_y \mathcal{M}, \|v\| < \varepsilon \quad \left| \sigma^2(\text{exp}_y(v)) - \sigma^2(y) - \langle (\text{grad } \sigma^2)(y) \mid v \rangle \right| \leq \eta \cdot \|v\| \quad (53)$$

### General idea

Let  $\Delta(z, v)$  be the integrated function (for  $z \notin C(y)$ ):

$$\Delta(z, v) = \text{dist}(\exp_y(v), z)^2 - \text{dist}(y, z)^2 + 2 \langle \overrightarrow{yz} \mid v \rangle$$

and  $H(v) = \int_{\mathcal{M} \setminus C(y)} \Delta(z, v).dP(z)$  be the function to bound:

$$\begin{aligned} H(v) &= \sigma^2(\exp_y(v)) - \sigma^2(y) - \langle (\text{grad } \sigma^2)(y) \mid v \rangle \\ &= \int_{\mathcal{M}} \text{dist}(\exp_y(v), z)^2.dP(z) - \int_{\mathcal{M}} \text{dist}(y, z)^2.dP(z) + 2 \int_{\mathcal{M} \setminus C(y)} \langle \overrightarrow{yz} \mid v \rangle .dP(z) \end{aligned}$$

The idea is to split this integral in two in order to bound  $\Delta$  on a small neighbourhood  $W$  around the cut locus of  $y$  and to use the standard Lebesgue theorem to bound the integral of  $\Delta$  on  $\mathcal{M} \setminus W$ .

#### Lemma 1 A basis of neighbourhoods of $C(y)$

$W_\varepsilon = \bigcup_{x \in \mathcal{B}(y, \varepsilon)} C(y)$  is a continuous series of included and decreasing open sets all containing  $C(y)$  and converging toward it.

Let us first reformulate the definition of  $W_\varepsilon$ :

$$z \in W_\varepsilon \Leftrightarrow \exists x \in \mathcal{B}(y, \varepsilon) / z \in C(x) \Leftrightarrow \exists x \in \mathcal{B}(y, \varepsilon) / x \in C(z) \Leftrightarrow C(z) \cap \mathcal{B}(y, \varepsilon) \neq \emptyset$$

Going to the limit, we have:  $z \in W_0 = \lim_{\varepsilon \rightarrow 0} W_\varepsilon \Leftrightarrow C(z) \cap \{y\} \neq \emptyset \Leftrightarrow z \in C(y)$ . Thus, we have a continuous series of included and decreasing sets all containing  $C(y)$  and converging toward it. Now, let us prove that  $W_\varepsilon$  is an open set.

Let  $U = \{u = (x, v) \in \mathcal{M} \times T_x \mathcal{M} ; \|v\| = 1\}$  be the unit tangent bundle of  $\mathcal{M}$  and  $\rho : U \rightarrow \bar{\mathbb{R}}_+ = \mathbb{R}_+ \cup \{+\infty\}$  be the cutting abscissa of the geodesic starting at  $x$  with the tangent vector  $v$ . Let now  $\underline{U} = \rho^{(-1)}(\mathbb{R}_+)$  be the subset of the unit tangent bundle where  $\rho(u) < +\infty$ . The function  $\rho$  being continuous on  $U$ , the subspace  $\underline{U} = \rho^{(-1)}(+\infty)$  is open. Let  $\pi : U \rightarrow \mathcal{M}$  be the canonical projection along the fiber  $\pi((x, v)) = x$  (it is obviously continuous) and let us denote by  $U_x = \pi^{(-1)}(\{x\})$  the unit tangent bundle at point  $x$  and  $\underline{U}_x = \pi^{(-1)}(\{x\}) \cap \underline{U}$  its subset that lead to a cutting point of  $x$ .

Consider  $u = (x, v) \in \underline{U}_x$  and the geodesic  $\exp_x(t.\rho(u).v)$  for  $t \in [0; 1]$ : it is starting from  $x$  with tangent vector  $v$  and arriving at  $z = \exp_x(\rho(u).v)$  with the same tangent vector by parallel transportation. Reverting the time course ( $t \rightarrow 1 - t$ ), we have a geodesic starting at  $z$  with tangent vector  $-v$  and arriving at  $x$  with the same tangent vector. By definition of the cutting function, we have  $\rho(u) = \rho(u')$  with  $u' = (z, -v)$ . Thus, the function  $q(u) = (\exp_x(\rho(u).v), -v)$  is a continuous bijection from  $\underline{U}$  into itself with  $q^{(-1)} = q$ .

Let  $z \in W_\varepsilon$ . By definition of  $W_\varepsilon$ ,  $z$  is in the cut locus of a point  $x \in \mathcal{B}(y, \varepsilon)$ : there exists a unit tangent vector  $v$  at that point such that  $z = \pi(q(x, v))$ . Conversely, the cut locus of  $z$  intersects  $\mathcal{B}(y, \varepsilon)$ : there exists a unit tangent vector  $v$  at  $z$  such that  $x = \pi(q(z, v)) \in \mathcal{B}(y, \varepsilon)$ . Thus, we can rewrite  $W_\varepsilon = \pi(U_\varepsilon)$  where  $U_\varepsilon = q^{(-1)}(\pi^{(-1)}(\mathcal{B}(y, \varepsilon)) \cap \underline{U})$ . The functions  $\pi$  and  $q$  being continuous, this set is open.

Now, we are left the the proof that  $\pi(U_\varepsilon) = W_\varepsilon$  is open. Assume that  $W_\varepsilon$  is not open. Then,  $\overline{W_\varepsilon} = \mathcal{M} \setminus W_\varepsilon$  is not closed: there exists a series of points  $z_i \in \overline{W_\varepsilon}$  converging toward  $z \in W_\varepsilon$ . By definition of this set, there exists  $u = (z, v) \in U_\varepsilon$  such that  $\pi(u) = z$ . Conversely, there is not unit tangent vector at  $z_i \notin W_\varepsilon$  with that property:  $U_{z_i} \cap U_\varepsilon = \emptyset$ . Consider the parallel transportation  $v_i$  of  $v$  from  $z$  to  $z_i$ . The series  $u_i = (z_i, v) \in \overline{U_\varepsilon}$  is obviously converging toward  $u = (z, v) \in U_\varepsilon$ . Since  $U_\varepsilon$  is open,  $\overline{U_\varepsilon}$  is closed and the series converges within this set, which is in contradiction with the previous statement.

**Lemma 2** *Let  $y \in \mathcal{M}$  and  $\alpha > 0$ . Then there exists an open neighborhood  $W_\varepsilon$  of  $C(y)$  such that*

- (i) *For all  $x \in \mathcal{B}(y, \varepsilon)$ ,  $C(x) \in W_\varepsilon$ ,*
- (ii)  *$P(W_\varepsilon) = \int_{W_\varepsilon} dP(z) < \alpha$*
- (iii)  *$\int_{W_\varepsilon} \text{dist}(y, z).dP(z) < \alpha$*

By hypothesis, the cut locus  $C(y)$  has a null measure for the measures  $dP(z)$ . The distance being a measurable function, its measure is null on the cut locus:  $\int_{C(y)} \text{dist}(y, z).dP(z) = 0$ . Thus, there exists an open set  $W \in \mathcal{M}$  containing  $C(y)$  and having an arbitrarily small measures for the identity and the distance function (say less than  $\alpha$ ). Since  $W_\varepsilon$  is a basis of neighborhoods of  $C(y)$ , there exists  $\varepsilon$  sufficiently small such that  $W_\varepsilon \in W$ .

### Bounding $\Delta$ on $W_\varepsilon$

Let  $W_\varepsilon$ ,  $\varepsilon$  and  $\alpha$  verifying the conditions of lemma 2 and  $x, x' \in \mathcal{B}(y, \varepsilon)$ . We have  $\text{dist}(z, x) \leq \text{dist}(z, x') + \text{dist}(x', x)$ . Thus:

$$\text{dist}(z, x)^2 - \text{dist}(z, x')^2 \leq \text{dist}(x, x')(\text{dist}(x, x') + 2 \text{dist}(z, x'))$$

Using the symmetry of  $x$  and  $x'$  and the inequalities  $\text{dist}(x, x') \leq 2\varepsilon$  and  $\text{dist}(z, x') \leq \text{dist}(z, y) + \varepsilon$ , we have:  $|\text{dist}(z, x)^2 - \text{dist}(z, x')^2| \leq 2 \text{dist}(x, x') \cdot (2\varepsilon + \text{dist}(z, y))$ . Applying this bound to  $x = \exp_y(v)$  and  $x' = y$ , we obtain:

$$|\text{dist}(z, \exp_y(v))^2 - \text{dist}(z, y)^2| \leq 2(2\varepsilon + \text{dist}(z, y)) \cdot \|v\|$$

Now, the last term of  $\Delta(z, v)$  is easily bounded by:  $\langle \vec{yz} | v \rangle \leq \text{dist}(y, z) \cdot \|v\|$ . Thus, we have:  $|\Delta(z, v)| \leq 4(\varepsilon + \text{dist}(z, y)) \cdot \|v\|$  and its integral over  $W_\varepsilon$  is bounded by:

$$\int_{W_\varepsilon} \Delta(z, v).dP(z) \leq 4\|v\| \cdot \int_{W_\varepsilon} (\varepsilon + \text{dist}(z, y)).dP(z) < 8\alpha\|v\|$$

### Bounding $\Delta$ on $\mathcal{M} \setminus W_\varepsilon$ :

Let  $x = \exp_y(v) \in \mathcal{B}(y, \varepsilon)$ . We know from lemma 2 that the cut locus  $C(x)$  of such a point belong to  $W_\varepsilon$ . Thus, the integration domain  $\mathcal{M} \setminus W_\varepsilon$  is now independent of  $y$  and we can use the usual Lebesgue theorem to differentiate under the sum:

$$\text{grad} \left( \int_{\mathcal{M} \setminus W_\varepsilon} \text{dist}(y, z)^2 .dP(z) \right) = -2 \int_{\mathcal{M} \setminus W_\varepsilon} \vec{yz} .dP(z)$$

By definition, this means that for  $\|v\|$  small enough, we have:

$$\int_{\mathcal{M} \setminus W_\varepsilon} \Delta(z, v).dP(z) < \alpha\|v\|$$

### Conclusion

Thus, for  $\|v\|$  small enough, we have  $\int_{\mathcal{M}} \Delta(z, v).dP(z) < 9\alpha\|v\|$ , which means that the variance has a derivative at the point  $y$ :

$$(\text{grad } \sigma^2)(y) = -2 \cdot \int_{\mathcal{M}/C(y)} \vec{yz} .dP(z)$$

## B Approximation of the generalized Normal density

In this section, we only work with a normal coordinate system at the mean value of the considered normal law. This allows us to simplify the notations. The density is

$$N(y) = k. \exp\left(-\frac{y^T \cdot \Gamma \cdot y}{2}\right) \quad \text{with} \quad k^{(-1)} = \int_{\mathcal{M}} \exp\left(-\frac{y^T \cdot \Gamma \cdot y}{2}\right) \cdot d\mathcal{M}(y)$$

The covariance and concentration are related by:

$$\Sigma = k. \int_{\mathcal{M}} y \cdot y^T \cdot \exp\left(-\frac{y^T \cdot \Gamma \cdot y}{2}\right) \cdot d\mathcal{M}(y)$$

Since the concentration matrix  $\Gamma$  is symmetric and positive definite, we can diagonalize it in the form  $\Gamma = A^T \cdot \Lambda^2 \cdot A$  with  $A$  being orthogonal and  $\Lambda$  being the diagonal matrix of positive square roots of eigenvalues. With the change of variable  $z = \Lambda \cdot A \cdot y$ , we can rewrite the above integrals as:

$$k^{(-1)} = \det(\Lambda)^{(-1)} \cdot \int_{\mathcal{D}'} \exp\left(-\frac{\|z\|^2}{2}\right) \cdot \sqrt{\det(G(A^T \cdot \Lambda^{(-1)} \cdot z))} \cdot dz$$

$$J = k^{(-1)} \cdot \det(\Lambda) \cdot \Lambda \cdot A \cdot \Sigma \cdot A^T \cdot \Lambda = \int_{\mathcal{D}'} z \cdot z^T \cdot \exp\left(-\frac{\|z\|^2}{2}\right) \cdot \sqrt{\det(G(A^T \cdot \Lambda^{(-1)} \cdot z))} \cdot dz$$

where the new definition domain is  $\mathcal{D}' = \Lambda \cdot A \cdot \mathcal{D}$ .

From [Chavel, 1993, section 2.8, corollary 2.3, p.84], we know the Taylor expansion of the metric around the origin in a normal coordinate system:

$$\det(g_{ij}(\exp v)) = 1 - \text{Ric}(v, v)/3 + O(\|v\|^3)$$

Thus:

$$d\mathcal{M}(y) = \sqrt{\det(G(y))} = 1 - \frac{y^T \cdot \text{Ric} \cdot y}{6} + O(\|y\|^3) \tag{54}$$

where  $\text{Ric}$  is the expression of the Ricci tensor of scalar curvatures in the exponential chart. Let  $T = \Lambda^{(-1)} \cdot A \cdot \text{Ric} \cdot A^T \cdot \Lambda^{(-1)}$ . Since  $\|A^T \cdot \Lambda^{(-1)} \cdot z\|^2 = \|\Lambda^{(-1)} \cdot z\|^2 = \sum_i z_i^2 / \lambda_i^2 \leq \|z\|^2 / \lambda_{min}^2$ , we obtain:

$$\sqrt{\det(G(A^T \cdot \Lambda^{(-1)} \cdot z))} = 1 - \frac{z^T \cdot T \cdot z}{6} + O\left(\frac{\|z\|^3}{\lambda_{min}^3}\right)$$

### B.1 Manifolds of non positive curvature at the mean point

At such a point, the cut locus is void and the definition domain of the exponential chart is  $\mathcal{D} = \mathbb{R}^n$ . This greatly simplifies the integrals. We have:

$$\det(\Lambda) \cdot k^{(-1)} = \int_{\mathbb{R}^n} \exp\left(-\frac{\|z\|^2}{2}\right) \cdot \left\{1 - \frac{1}{6} z^T \cdot T \cdot z + O\left(\frac{\|z\|^3}{\lambda_{min}^3}\right)\right\} \cdot dz$$

The first term is

$$\int_{\mathbb{R}^n} \exp\left(-\frac{\|z\|^2}{2}\right) \cdot dz = \prod_i \int_{\mathbb{R}} \exp\left(-\frac{z_i^2}{2}\right) \cdot dz_i = (2\pi)^{n/2}$$

The third term in simple: as the integral converges to a constant, we obtain:

$$\int_{\mathbb{R}^n} O\left(\frac{\|z\|^3}{\lambda_{min}^3}\right) \cdot \exp\left(-\frac{\|z\|^2}{2}\right) \cdot dz = O(\lambda_{min}^{-3})$$

The second term is

$$\begin{aligned} \int_{\mathbb{R}^n} z^T \cdot T \cdot z \cdot \exp\left(-\frac{\|z\|^2}{2}\right) \cdot dz &= \sum_{i,j} T_{ij} \int z_i \cdot z_j \exp\left(-\frac{\|z\|^2}{2}\right) \cdot dz \\ &= \sum_i T_{ii} \left( \int z_i^2 \cdot \exp\left(-\frac{z_i^2}{2}\right) \cdot dz_i \right) \cdot \left( \prod_{j \neq i} \int \exp\left(-\frac{z_j^2}{2}\right) \cdot dz_j \right) \\ &= \sum_i T_{ii} \cdot (\sqrt{2\pi}) \cdot ((2\pi)^{(n-1)/2}) \\ &= (2\pi)^{n/2} \cdot \text{Tr}(T) \\ &= (2\pi)^{n/2} \cdot \text{Tr}(\Gamma^{(-1)} \cdot \text{Ric}) \end{aligned}$$

since  $\text{Tr}(A \cdot B) = \text{Tr}(B \cdot A)$ . Thus we obtain:

$$k^{(-1)} = \frac{(2\pi)^{n/2}}{\sqrt{\det(\Gamma)}} \left( 1 - \frac{\text{Tr}(\Gamma^{(-1)} \cdot \text{Ric})}{6} + O(\lambda_{min}^{-3}) \right)$$

As  $\text{Tr}(\Gamma^{(-1)} \cdot \text{Ric})$  is a term in  $\lambda_{min}^{-2}$ , we have:

$$k = \frac{\sqrt{\det(\Gamma)}}{(2\pi)^{n/2}} \left( 1 + \frac{\text{Tr}(\Gamma^{(-1)} \cdot \text{Ric})}{6} + O(\lambda_{min}^{-3}) \right) \quad (55)$$

Now, for the covariance matrix, we have:

$$J = \int_{\mathbb{R}^n} z \cdot z^T \cdot \exp\left(-\frac{\|z\|^2}{2}\right) \cdot \left( 1 - \frac{1}{6} \text{Tr}(T \cdot z \cdot z^T) + O\left(\frac{\|z\|^3}{\lambda_{min}^3}\right) \right) \cdot dz$$

We consider first the off diagonal elements  $J_{ij}$  (with  $i \neq j$ ): the first term of the integral vanished because we integrate antisymmetric functions over  $\mathbb{R}$ :  $\int z_i \cdot z_j \cdot \exp(\|z\|^2/2) \cdot dz = 0$ . The integral converges to a constant for the third term and the result is thus  $O(\lambda_{min}^{-3})$ . Since  $\text{Tr}(T \cdot z \cdot z^T) = \sum_{k,l} T_{kl} \cdot z_k \cdot z_l$ , we are left with:

$$J_{ij} = -\frac{1}{6} \sum_{k,l} T_{kl} \cdot \int z_k \cdot z_l \cdot z_i \cdot z_j \cdot \exp\left(-\frac{\|z\|^2}{2}\right) \cdot dz$$

For  $i \neq k$  and  $j \neq l$  or  $i \neq l$  and  $j \neq k$ , we integrate an antisymmetric function and the result is zero. Since the Ricci curvature matrix is symmetric, the matrix  $T$  is also symmetric and the sum is reduced to a single term:

$$J_{ij} = -\frac{1}{6} (T_{ij} + T_{ji}) \cdot \int z_i^2 \cdot z_j^2 \cdot \exp\left(-\frac{\|z\|^2}{2}\right) \cdot dz = -\frac{T_{ij}}{3} (2\pi)^{n/2}$$

Now, we consider the diagonal elements  $J_{ij}$ . The first term is

$$\begin{aligned} \int_{\mathbb{R}^n} z_i^2 \cdot \exp(-\|z\|^2/2) \cdot dz &= \left( \int_{\mathbb{R}} z_i^2 \cdot \exp(-z_i^2/2) \cdot dz_i \right) \cdot \prod_{j \neq i} \left( \int_{\mathbb{R}} \exp(-z_j^2/2) \cdot dz_j \right) \\ &= (2\pi)^{1/2} \cdot (2\pi)^{(n-1)/2} = (2\pi)^{n/2} \end{aligned}$$

The second term is

$$\alpha_i = -\frac{1}{6} \int_{\mathbb{R}^n} z_i^2 \cdot z^T \cdot T \cdot z \cdot \exp(-\|z\|^2/2) \cdot dz = -\frac{1}{6} \sum_{k,l} T_{kl} \int_{\mathbb{R}^n} z_i^2 \cdot z_k \cdot z_l \cdot \exp(-\|z\|^2/2) \cdot dz$$

For  $k \neq l$ , we integrate antisymmetric functions: the result is zero. Thus, we are left with:

$$\begin{aligned} \alpha_i &= -\frac{1}{6} \sum_{k \neq i} T_{kk} \int_{\mathbb{R}^n} z_i^2 \cdot z_k^2 \cdot \exp(-\|z\|^2/2) \cdot dz - \frac{1}{6} T_{ii} \int_{\mathbb{R}^n} z_i^4 \cdot \exp(-\|z\|^2/2) \cdot dz \\ &= -\frac{1}{6} \sum_{k \neq i} T_{kk} \left( \int_{\mathbb{R}} z_i^2 \cdot \exp(-z_i^2/2) \cdot dz_i \right) \cdot \left( \int_{\mathbb{R}} z_k^2 \cdot \exp(-z_k^2/2) \cdot dz_k \right) \cdot \left( \int_{\mathbb{R}} \exp(-z_j^2/2) \cdot dz_j \right)^{(n-2)} \\ &\quad - \frac{1}{6} T_{ii} \left( \int_{\mathbb{R}} z_i^4 \cdot \exp(-z_i^2/2) \cdot dz_i \right) \cdot \left( \int_{\mathbb{R}} \exp(-z_j^2/2) \cdot dz_j \right)^{(n-1)} \\ &= -\frac{1}{6} \sum_{k \neq i} T_{kk} \cdot (2\pi)^{n/2} - \frac{1}{6} T_{ii} \cdot 3 \cdot (2\pi)^{n/2} \\ &= -\frac{1}{6} (\text{Tr}(T) + 2 \cdot T_{ii}) = -\frac{1}{6} (\text{Tr}(\Gamma^{(-1)} \cdot \text{Ric}) + 2 \cdot T_{ii}) \end{aligned}$$

For the third term, the integral converges to a constant and the result is  $O(\lambda_{min}^{-3})$ . Then, the diagonal element are:

$$J_{ii} = (2\pi)^{n/2} \cdot \left( 1 - \frac{\text{Tr}(\Gamma^{(-1)} \cdot \text{Ric})}{6} - \frac{T_{ii}}{3} + O(\lambda_{min}^{-3}) \right)$$

Combining with off diagonal terms, we get:

$$J = (2\pi)^{n/2} \cdot \left\{ \text{Id} \cdot \left( 1 - \frac{\text{Tr}(\Gamma^{(-1)} \cdot \text{Ric})}{6} \right) - \frac{T}{3} + O(\lambda_{min}^{-3}) \right\}$$

As  $J = k^{(-1)} \cdot \det(\Lambda) \cdot \Lambda \cdot A \cdot \Sigma \cdot A^T \cdot \Lambda$ , we have (using equation 55 for the expansion of  $k$ ):

$$\begin{aligned} \Lambda \cdot A \cdot \Sigma \cdot A^T \cdot \Lambda &= (k / \det(\Lambda)) \cdot J \\ &= \left\{ 1 + \frac{1}{6} \text{Tr}(\Gamma^{(-1)} \cdot \text{Ric}) + O(\lambda_{min}^3) \right\} \cdot \left\{ \text{Id} \cdot \left( 1 - \frac{1}{6} \text{Tr}(\Gamma^{(-1)} \cdot \text{Ric}) \right) - \frac{1}{3} T + O(\lambda_{min}^{-3}) \right\} \\ &= \text{Id} - \frac{1}{3} T + O(\lambda_{min}^{-3}) \end{aligned}$$

We obtain thus the following relation between the covariance and the concentration matrices:

$$\Sigma = \Gamma^{(-1)} - \frac{1}{3} \Gamma^{(-1)} \cdot \text{Ric} \cdot \Gamma^{(-1)} + O(\lambda_{min}^{-5}) \tag{56}$$

This relation can be inverted to obtain the Taylor expansion of the concentration matrix with respect to the covariance. First, we shall note that from the above equation,  $O(\lambda_{min}^{-5}) = O(\sigma_{max}^5)$  where  $\sigma_{max}$  is the square root of the largest eigenvalue of  $\Sigma$ . However, a global variable such as the variance  $\sigma^2 = \text{Tr}(\Sigma)$  is more appropriate. Since  $\sigma_{max}^2 < \sum_i \sigma_i^2 = \sigma^2$ , we have  $O(\sigma_{max}) = O(\sigma)$  and the Taylor expansion of  $\Gamma^{(-1)}$  is:

$$\Gamma^{(-1)} = \Sigma + \frac{1}{3} \Sigma \cdot \text{Ric} \cdot \Sigma + O(\sigma^5)$$

Then, we have

$$\Gamma \cdot \Sigma = \text{Id} - \frac{1}{3} \text{Ric} \cdot \Gamma^{(-1)} + O(\lambda_{min}^{-3}) = \text{Id} - \frac{1}{3} \text{Ric} \cdot \Sigma + O(\sigma^3)$$

This means that:

$$\Gamma = \Sigma^{(-1)} - \frac{1}{3} \text{Ric} + O(\sigma) \tag{57}$$

To express  $k$  with respect to  $\Sigma$  instead of  $\Gamma$ , we have to compute  $\text{Tr}(\Gamma^{(-1)}.\text{Ric})$  and  $\sqrt{\det(\Gamma)}$ .

$$\text{Tr}(\Gamma^{(-1)}.\text{Ric}) = \text{Tr}(\Sigma.\text{Ric} + \frac{1}{3}\Sigma.\text{Ric}.\Sigma.\text{Ric} + O(\sigma^5)) = \text{Tr}(\Sigma.\text{Ric}) + O(\sigma^3)$$

For the determinant, we observe that

$$\frac{\partial \det(A(\alpha))}{\partial \alpha} = \det(A(\alpha)).\text{Tr}\left(\frac{\partial A(\alpha)}{\partial \alpha}.A(\alpha)^{(-1)}\right)$$

Thus, we have the following Taylor expansion (with  $A(\alpha) = \text{Id} + \alpha.B$ ):

$$\begin{aligned} \det(\text{Id} + \alpha.B) &= \det(\text{Id}) + \alpha.\left.\frac{\partial \det(\text{Id} + \alpha.B)}{\partial \alpha}\right|_{\alpha=0} + O(\alpha^2) \\ &= 1 + \alpha.\det(\text{Id}).\text{Tr}(B.(\text{Id} + B.0)^{(-1)}) + O(\alpha^2) \\ &= 1 + \text{Tr}(\alpha.B) + O(\alpha^2) \end{aligned}$$

Since  $\Gamma.\Sigma = \text{Id} - \frac{1}{3}\text{Ric}.\Sigma + O(\sigma^3)$  and  $\Sigma$  is a term in  $O(\sigma^2)$ , we have

$$\det(\Gamma) = \det(\Sigma)^{(-1)}.\det\left(\text{Id} - \frac{1}{3}\text{Ric}.\Sigma + O(\sigma^3)\right) = \det(\Sigma)^{(-1)}\left(1 - \frac{1}{3}.\text{Tr}(\Sigma.\text{Ric}) + O(\sigma^3)\right)$$

and thus

$$\sqrt{\det(\Gamma)} = \det(\Sigma)^{-1/2}.\left(1 - \frac{1}{6}.\text{Tr}(\Sigma.\text{Ric}) + O(\sigma^3)\right)$$

Reporting this expression in equation 55, we obtain:

$$k = \frac{1}{\sqrt{(2\pi)^n.\det(\Sigma)}}\left(1 - \frac{1}{6}.\text{Tr}(\Sigma.\text{Ric}) + O(\sigma^3)\right).\left(1 + \frac{1}{6}.\text{Tr}(\Sigma.\text{Ric}) + O(\sigma^3)\right)$$

Which simplifies in:

$$k = \frac{1 + O(\sigma^3)}{\sqrt{(2\pi)^n.\det(\Sigma)}} \tag{58}$$

**Summary of the approximate normal density:** In a manifold of non positive curvature, the normal density in a normal coordinate system at the mean value is:

$$N(y) = k.\exp\left(-\frac{y^T.\Gamma.y}{2}\right)$$

The normalization constant and the concentration matrices are approximated by the following expressions for a covariance matrix  $\Sigma$  of small variance  $\sigma^2 = \text{Tr}(\Sigma)$ :

$$k = \frac{1 + O(\sigma^3)}{\sqrt{(2\pi)^n.\det(\Sigma)}} \quad \text{and} \quad \Gamma = \Sigma^{(-1)} - \frac{1}{3}\text{Ric} + O(\sigma)$$

## B.2 Manifolds with a cut locus at the mean point

Now, we have to integrate all the integrals over the definition domain  $\mathcal{D}$  of the exponential chart at the mean point. We are going to bound each term of all the integrals from above and below to show that the same results still holds with just a slight modification of the Taylor expansion bounding term. Let  $r$  be the injective radius at the mean point. The open ball  $\mathcal{B}(r)$  is the greatest geodesic ball included in the definition domain:  $\mathcal{B}(r) \subset \mathcal{D} \subset \mathbb{R}^n$ . The idea is to bound the integrals of positive terms by integrating over these three sets.

In fact, with the change in variable  $z = \Lambda.A.y$ , the ball is transformed into the set  $\mathcal{B}_z = \{z/\|\Lambda^{(-1)}.z\| < r\}$ . Since  $\|\Lambda^{(-1)}\|^2 \leq \lambda_{min}^{-2}.\|z\|^2$ , we can use the smaller ball  $\mathcal{B} = \mathcal{B}(\lambda_{min}.r)$  for the integrations in polar coordinates. For the integrals that are separable along axes, it will be simpler to use the maximal cube  $\mathcal{C} = \{z/\lambda_{min}.r/\sqrt{n} < z_i < \lambda_{min}.r/\sqrt{n}\}$  included in  $\mathcal{B}$ .

In the following, we bound the different integrals used in the last section and summarized afterward the modifications it implies in the results of the last section. The first integral is bounded above by the previous value:

$$\int_{\mathcal{D}'} \exp\left(-\frac{\|z\|^2}{2}\right).dz \leq \int_{\mathbb{R}^n} \exp\left(-\frac{\|z\|^2}{2}\right).dz = (2\pi)^{n/2}$$

and below by an integration over  $\mathcal{C}$ :

$$\begin{aligned} \int_{\mathcal{D}'} \exp\left(-\frac{\|z\|^2}{2}\right).dz &\geq \prod_i \left( \int_{-\lambda_{min}.r/\sqrt{n}}^{\lambda_{min}.r/\sqrt{n}} \exp\left(-\frac{z_i^2}{2}\right).dz_i \right) \\ &\geq (2\pi)^{n/2} \cdot \left(1 - \operatorname{erfc}\left(\frac{\lambda_{min}.r}{\sqrt{2n}}\right)\right)^n \end{aligned}$$

Here,  $\operatorname{erfc}(x) = 1 - \frac{2}{\sqrt{\pi}} \int_0^x \exp(-t^2).dt$  is the complement of the error function. An interesting property is that this function (like the exponential) tends toward zero faster than any fraction  $1/x^k$  as  $x$  goes to infinity. Putting things the other way, like for Taylor expansions, we will denote by  $\varepsilon(x)$  a function that is an  $O(x^k)$  for every positive  $k$ :

$$\lim_{0+} \frac{\varepsilon(x)}{x^k} = \lim_{0+} \frac{\operatorname{erfc}(1/x)}{x^k} = \lim_{0+} \frac{\exp(-1/x)}{x^k} = 0$$

We can summarize the value of the first integral by:

$$\int_{\mathcal{D}'} \exp\left(-\frac{\|z\|^2}{2}\right).dz = (2\pi)^{n/2} + \varepsilon(\lambda_{min}^{-1}.r^{-1})$$

Now, for the following integral, we have

$$\begin{aligned} \int_{-\lambda_{min}.r/\sqrt{n}}^{\lambda_{min}.r/\sqrt{n}} z_i^2 \cdot \exp\left(-\frac{z_i^2}{2}\right).dz_i &= \sqrt{2\pi} - \sqrt{2\pi} \cdot \operatorname{erfc}\left(\frac{\lambda_{min}.r}{\sqrt{2n}}\right) - \frac{2\sqrt{n}}{\lambda_{min}.r} \exp\left(-\frac{\lambda_{min}^2.r^2}{2}\right) \\ &= \sqrt{2\pi} + \varepsilon(\lambda_{min}^{-1}.r^{-1}) \end{aligned}$$

We obtain thus

$$\int_{\mathcal{D}'} z_i^2 \cdot \exp\left(-\frac{\|z\|^2}{2}\right).dz = (2\pi)^{n/2} + \varepsilon(\lambda_{min}^{-1}.r^{-1})$$

In fact, it is not hard to show that every integral we computed over  $\mathbb{R}^n$  in the previous paragraph has the same value over  $\mathcal{D}'$  plus a term whose absolute value is of the order of  $\varepsilon(\lambda_{min}^{-1}.r^{-1})$ . Thus, we can directly generalize the previous results by replacing  $O(\sigma^k)$  with  $O(\sigma^k) + \varepsilon\left(\frac{\sigma}{r}\right)$ .

**Summary of the approximate normal density:** In a manifold with injectivity radius  $r$  at the mean point, the normal density in a normal coordinate system at this mean value is:

$$N(y) = k. \exp\left(-\frac{y^T \cdot \Gamma \cdot y}{2}\right)$$

The normalization constant and the concentration matrices are approximated by the following expressions for a covariance matrix  $\Sigma$  of small variance  $\sigma^2 = \text{Tr}(\Sigma)$ :

$$k = \frac{1 + O(\sigma^3) + \varepsilon\left(\frac{\sigma}{r}\right)}{\sqrt{(2\pi)^n \cdot \det(\Sigma)}} \quad \text{and} \quad \Gamma = \Sigma^{(-1)} - \frac{1}{3}\text{Ric} + O(\sigma) + \varepsilon\left(\frac{\sigma}{r}\right)$$

## C Approximated generalized $\chi^2$ law

Assuming that the random primitive  $\mathbf{x} \sim (\bar{\mathbf{x}}, \Sigma)$  follows a normal law, we want to compute the probability that

$$\chi^2 = \mu_{\bar{\mathbf{x}}}^2 = \bar{\mathbf{x}}\bar{\mathbf{x}}^T \cdot \Sigma^{(-1)} \cdot \bar{\mathbf{x}}\bar{\mathbf{x}} \leq \alpha^2$$

Let  $\mathcal{B}_\Sigma(\alpha)$  be the ‘‘elliptic ball’’ of covariance  $\Sigma$  and radius  $\alpha$ :  $\mathcal{B}_\Sigma(\alpha) = \{y / y^T \cdot \Sigma^{(-1)} \cdot y \leq \alpha^2\}$ . Assuming that there is not cut locus, this probability can be written in a normal coordinate system:

$$\Pr\{\chi^2 \leq \alpha^2\} = \int_{\chi^2 \leq \alpha^2} N(y) \cdot d\mathcal{M}y = \int_{\mathcal{B}_\Sigma(\alpha)} k. \exp(-y^T \cdot \Gamma \cdot y / 2) \cdot \sqrt{G(y)} \cdot dy$$

Since  $\Gamma = \Sigma^{(-1)} - \frac{1}{3}\text{Ric} + O(\sigma)$  and  $\sqrt{\det(G(y))} = 1 - \frac{y^T \cdot \text{Ric} \cdot y}{6} + O(\|y\|^3)$  we have:

$$\Pr\{\chi^2 \leq \alpha^2\} = k \int_{\mathcal{B}_\Sigma(\alpha)} \exp\left(-\frac{y^T \left(\Sigma^{(-1)} - \frac{1}{3}\text{Ric} + O(\sigma)\right) y}{2}\right) \left(1 - \frac{y^T \text{Ric} y}{6} + O(\|y\|^3)\right) dy$$

Let  $\Theta$  be a positive square root of  $\Sigma$  such that  $\Sigma = \Theta \cdot \Theta^T$ . Which the change of variable  $y = \Theta \cdot x$ , we have  $dy = \sqrt{\det(\Sigma)} \cdot dx$  and the probability is now:

$$\Pr\{\chi^2 \leq \alpha^2\} = k \cdot \sqrt{\det(\Sigma)} \int_{\|x\| \leq \alpha} \exp\left(-\frac{x^T (\text{Id} - S + O(\sigma^3)) x}{2}\right) \left(1 - \frac{x^T S x}{2} + O(\|x\|^3 \sigma^3)\right) dx$$

where  $S = \frac{1}{3}\Theta^T \text{Ric} \Theta$ . As this is a term of the order  $O(\sigma^2)$ , we have the following Taylor expansion:

$$\exp\left(-\frac{x^T (\text{Id} - S + O(\sigma^3)) x}{2}\right) = \exp\left(-\frac{\|x\|^2}{2}\right) \left(1 + \frac{x^T \cdot S \cdot x}{2} + O(\sigma^3)\right)$$

Reporting in the probability, we find that:

$$\Pr\{\chi^2 \leq \alpha^2\} = k \cdot \sqrt{\det(\Sigma)} \int_{\|x\| \leq \alpha} \exp\left(-\frac{\|x\|^2}{2}\right) (1 + O(\sigma^3) + \|x\|^2 O(\sigma^3)) dx$$

Now, we have:

$$\int_{\|x\| \leq \alpha} \exp\left(-\frac{\|x\|^2}{2}\right) (1 + \|x\|^2) \cdot O(\sigma^3) \cdot dx \leq O(\sigma^3) \int_{\mathbb{R}^n} (1 + \|x\|^2) \exp\left(-\frac{\|x\|^2}{2}\right) dx = O(\sigma^3)$$

Thus, using equation 58 for the value of  $k$ , we find the usual  $\chi^2$  probability, which is independent of the covariance  $\Sigma$ , up to the order  $\sigma^3$ :

$$\Pr\{\chi^2 \leq \alpha^2\} = (2\pi)^{-\frac{n}{2}} \int_{\|x\| \leq \alpha} \exp\left(-\frac{\|x\|^2}{2}\right) .dx + O(\sigma^3) \quad (59)$$

This integral can be computed in polar coordinates in  $\mathbb{R}^n$ : if  $r = \|x\|$  is the radius and  $n$  is the corresponding unit vector ( $x = r.n$ ), we have  $dx = r^{n-1}.dr.dn$  and thus:

$$\begin{aligned} \int_{\|x\| \leq \alpha} \exp\left(-\frac{\|x\|^2}{2}\right) .dx &= \left( \int_{S_{n-1}} dn \right) . \left( \int_0^\alpha r^{n-1} . \exp(-r^2/2) .dr \right) \\ &= \left( \frac{2\pi^{n/2}}{\Gamma(n/2)} \right) . \left( \frac{1}{2} \int_0^{\alpha^2} t^{n/2-1} . \exp\left(-\frac{t}{2}\right) .dt \right) \end{aligned}$$

with the change of variable  $t = r^2$ . In this formula,  $\Gamma$  is the Gamma function which can be recursively computed from  $\Gamma(x + 1) = x.\Gamma(x)$  with  $\Gamma(1) = 1$  and  $\Gamma(\frac{1}{2}) = \sqrt{\pi}$ . For half integer values, the result is:

$$\Gamma\left(\frac{n}{2}\right) = (k - 1)! \quad \text{if } n = 2.k \quad \text{and} \quad \Gamma\left(\frac{n}{2}\right) = \sqrt{\pi} . \prod_{i=0}^{k-1} \left(i + \frac{1}{2}\right) \quad \text{if } n = 2.k + 1$$

Thus, the probability density function of a  $\chi^2$  is:

$$p_{\chi^2}(u) = \frac{1 + O(\sigma^3)}{2.\Gamma\left(\frac{n}{2}\right)} \left(\frac{u}{2}\right)^{\frac{n}{2}-1} \exp\left(-\frac{u}{2}\right) \quad (60)$$

If there is a cut locus, we have to replace  $O(\sigma^i)$  by  $O(\sigma^i) + \varepsilon\left(\frac{\sigma}{r}\right)$  and we should only integrate on  $\mathcal{B}_\Sigma(\alpha) \cap \mathcal{D}$ . After the change of coordinates, we should integrate on  $\mathcal{B}(\alpha) \cap \mathcal{D}'$ . Thus, we have:

$$\Pr\{\chi^2 \leq \alpha^2\} = (2\pi)^{-\frac{n}{2}} \int_{\mathcal{B}(\alpha) \cap \mathcal{D}'} \exp\left(-\frac{\|x\|^2}{2}\right) .dx + O(\sigma^3) + \varepsilon\left(\frac{\sigma}{r}\right)$$

As we have  $\mathcal{B}(\lambda_{min}.r) \subset \mathcal{D}' \subset \mathbb{R}^n$ , we can enclose the researched domain into:  $\mathcal{B}(\min(\lambda_{min}.r, \alpha)) \subset \mathcal{B}(\alpha) \cap \mathcal{D}' \subset \mathcal{B}(\alpha)$ . For  $\alpha \leq \lambda_{min}.r$ , there is not problem but for  $\alpha > \lambda_{min}.r$ , we have:

$$\Pr\{\chi^2 \leq \alpha^2\} \geq (2\pi)^{-\frac{n}{2}} \int_{\mathcal{B}(\lambda_{min}.r)} \exp\left(-\frac{\|x\|^2}{2}\right) .dx + O(\sigma^3) + \varepsilon\left(\frac{\sigma}{r}\right)$$

and we have already seen that this integral is  $1 + \varepsilon\left(\frac{\sigma}{r}\right)$ . As  $\alpha > \lambda_{min}.r$ , the same integral is itself of the same order, and we obtain in all cases the same result as with no cut locus where  $O(\sigma^3)$  is replaced by  $O(\sigma^3) + \varepsilon\left(\frac{\sigma}{r}\right)$ .



---

Unité de recherche INRIA Sophia Antipolis  
2004, route des Lucioles - BP 93 - 06902 Sophia Antipolis Cedex (France)

Unité de recherche INRIA Futurs : Parc Club Orsay Université - ZAC des Vignes  
4, rue Jacques Monod - 91893 ORSAY Cedex (France)

Unité de recherche INRIA Lorraine : LORIA, Technopôle de Nancy-Brabois - Campus scientifique  
615, rue du Jardin Botanique - BP 101 - 54602 Villers-lès-Nancy Cedex (France)

Unité de recherche INRIA Rennes : IRISA, Campus universitaire de Beaulieu - 35042 Rennes Cedex (France)

Unité de recherche INRIA Rhône-Alpes : 655, avenue de l'Europe - 38334 Montbonnot Saint-Ismier (France)

Unité de recherche INRIA Rocquencourt : Domaine de Voluceau - Rocquencourt - BP 105 - 78153 Le Chesnay Cedex (France)

---

Éditeur  
INRIA - Domaine de Voluceau - Rocquencourt, BP 105 - 78153 Le Chesnay Cedex (France)

<http://www.inria.fr>

ISSN 0249-6399