



HAL
open science

Gaussian distributions on Riemannian symmetric spaces of non-positive curvature

Salem Said, Cyrus Mostajeran, Simon Heuveline

► **To cite this version:**

Salem Said, Cyrus Mostajeran, Simon Heuveline. Gaussian distributions on Riemannian symmetric spaces of non-positive curvature. Handbook of Statistics Vol 46, pp.357-400, 2022. hal-03865789

HAL Id: hal-03865789

<https://hal.science/hal-03865789>

Submitted on 22 Nov 2022

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

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

Gaussian distributions on
Riemannian symmetric spaces
of non-positive curvature

Salem Said (CNRS, Laboratoire LJK – Université Grenoble-Alpes)
Cyrus Mostajeran (Department of Engineering – University of Cambridge)
Simon Heuveline (Centre for Mathematical Sciences – University of Cambridge)

Abstract

This article aims to give a coherent presentation of the theory of Gaussian distributions on Riemannian symmetric spaces and also to report on recent original developments of this theory. The initial goal is to define a family of probability distributions, on any suitable Riemannian manifold, for which maximum-likelihood estimation, based on a finite sequence of observations, is equivalent to computation of the Riemannian barycentre of these observations. As it turns out, this goal is achievable whenever the underlying Riemannian manifold is a Riemannian symmetric space of non-positive curvature. In this case, the required Gaussian distributions are exactly the maximum-entropy distributions, for fixed barycentre and dispersion. The second step is the search for efficient means of computing the normalising factors associated with these distributions. This leads to a fascinating connection with random matrix theory, and even with theoretical physics (Chern-Simons theory), which yields a series of original results that provide exact expressions, as well as high-dimensional asymptotic expansions of the normalising factors. Another outcome of this connection with random matrix theory is the new idea of *duality*, between Gaussian distributions on Riemannian symmetric spaces of opposite curvatures. The present article also investigates Bayesian inference for Gaussian distributions on symmetric spaces. This investigation motivates original results regarding Markov-Chain Monte Carlo and convex optimisation on Riemannian manifolds. It also reveals a new open problem (roughly, this concerns the equality of *a posteriori* mode with *a posteriori* barycentre), which should be the focus of future developments.

keywords : Gaussian distribution ; Symmetric space ; Random matrix theory ; Markov-Chain Monte Carlo ; Convex optimisation ; Bayesian inference

1 Introduction

The realisation that an essentially new approach, beyond that of classical statistics, is needed in order to learn from data that live in non-Euclidean spaces, can be credited to Fréchet, who invented what we today call the Fréchet mean, back in 1948 [1].

The Fréchet mean generalises the concept of the mean (average or expectation) of a sequence of observations (x_1, \dots, x_N) , from the classical case where these observations lie in a Euclidean space, to the general setting where they belong to a non-Euclidean space.

Let us call our sample space M (the observations belong to M). If M is a Euclidean space, it has a vector space structure, and the mean of (x_1, \dots, x_N) is just the arithmetic mean $(x_1 + \dots + x_N)/N$. If M is a non-Euclidean space, it will have no vector space structure, and this definition will lose all meaning.

To salvage the concept of mean, Fréchet suggested looking at the set of global minima of the sum of squared distances (the factor $1/2$ is included for later convenience)

$$\mathcal{E}(x) = \frac{1}{2} \sum_{n=1}^N d^2(x_n, x) \quad \text{for } x \in M$$

He noted that any global minimum of \mathcal{E} deserves to be called a mean of (x_1, \dots, x_N) . In this way, the mean of a sequence of observations in a non-Euclidean space is well-defined, at the cost of eventually failing to be unique.

Fast-forward to the present, learning from data that live in Riemannian manifolds (a particular class of non-Euclidean spaces) has become central to many applications, ranging from radar signal processing to neuroscience [2, 3], and the Fréchet mean (more descriptively called the Riemannian barycentre) a very popular tool in this respect [4, 5].

Naturally, it also became important to provide a statistical (specifically, inferential) foundation for the use of this tool. One way or another, this led to the quest for a suitable definition of a Gaussian distribution on a Riemannian manifold. This appeared inevitable, already because of the intimate connection between arithmetic means and Gaussian distributions, in the classical Euclidean case (see Paragraph 2.1, for discussion).

Gaussian distributions, defined as maximum entropy distributions on a Riemannian manifold, for a given Fréchet mean and dispersion, were first introduced by Pennec [6]. For a while, it remained difficult to study these distributions, as there was no practical means of computing the associated normalising factors. However, a first breakthrough came when these factors were expressed as multiple integrals, in the case of Gaussian distributions on the space of real positive-definite matrices [7].

In [8][9], the approach of [7] was generalised to Gaussian distributions on Riemannian symmetric spaces of non-positive curvature, which include hyperbolic spaces, as well as spaces of real, complex and quaternion positive-definite matrices, and spaces of structured (Toeplitz or block-Toeplitz) positive-definite matrices. This opened the way to rigorous learning algorithms for data that live in these spaces (this is partially discussed in [9]).

The introduction of Riemannian symmetric spaces reduced normalising factors of Gaussian distributions to multiple integrals, which could be computed using Monte Carlo techniques [10]. Only very recently, it was realised that the techniques of random matrix theory made it possible to write down both analytic expressions and high-dimensional asymptotic expansions of these multiple integrals. This was studied by the theoretical physics community [11] (see also our paper, currently under review [12]).

The aim of the present article is to give a coherent presentation of the theory of Gaussian distributions on Riemannian symmetric spaces of non-positive curvature, and report on recent original developments of this theory, including (but not limited to) the ones just mentioned. Its main body (Sections 2 and 3) relies on a variety of new results, mostly contained in the habilitation thesis [13] — one advantage of this situation is that the flow of results is not interrupted by their sometimes lengthy proofs, given in [13].

In the following, Section 2 introduced Gaussian distributions and their connection with random matrix theory. Section 3 investigates Bayesian inference of these distributions. Each one of these sections opens with a description of the original results which it contains.

Another, more modest, contribution of this article is Appendix B (not based on [13]). This appendix provides new results on the convergence rates for Riemannian gradient descent, applied to strictly convex and strongly convex functions, defined on a convex subset of a Riemannian manifold. The main results are Propositions B.8 and B.9.

2 Gaussian distributions and RMT

The starting point of this section is a historical discussion of the concept of a Gaussian distribution. This leads up to the definition of Gaussian distributions, adopted in Paragraph 2.2, as a family of distributions $P(\bar{x}, \sigma)$ on a Riemannian manifold M , with parameters $\bar{x} \in M$ and $\sigma > 0$, for which maximum-likelihood estimation of \bar{x} is equivalent to computation of the Riemannian barycentre. It turns out that this definition can be pursued whenever M is a Riemannian symmetric space of non-positive curvature.

Paragraph 2.3 then gives a general expression of the normalising factor $Z(\sigma)$ of the Gaussian distribution $P(\bar{x}, \sigma)$, in the form of a multiple integral (8). When M is a space of positive-definite matrices, or when M is the so-called Siegel domain, (8) is further reduced to a kind of integral familiar in random matrix theory ((10) and (15), respectively).

Paragraph 2.4 states the existence and uniqueness of maximum-likelihood estimates of the parameters \bar{x} and σ . It also states the maximum-entropy property of the Gaussian distribution $P(\bar{x}, \sigma)$, in Proposition 2.5. Paragraph 2.5 provides expressions of the barycentre (shown to be equal to \bar{x}) and the covariance tensor of $P(\bar{x}, \sigma)$.

Paragraph 2.6 begins the series of results based on random matrix theory (RMT). These concern Gaussian distributions on the space $H(N)$ of complex positive-definite matrices. First, the analytic expression of $Z(\sigma)$ is given in Proposition 2.8. Then, an asymptotic expansion of this expression, in the limit where N goes to infinity while $t = N\sigma^2$ remains constant, is given in Proposition 2.9.

Paragraph 2.7 describes the asymptotic distribution of eigenvalues of a random positive-definite matrix in $H(N)$, drawn from the Gaussian distribution $P(I_N, \sigma)$ (I_N denotes the $N \times N$ identity matrix). This asymptotic distribution has a probability density function, whose explicit expression is provided in Proposition 2.10.

Paragraph 2.8 introduces Θ distributions. These are classical normal distributions, wrapped around the unitary group $U(N)$, which is the dual symmetric space of $H(N)$. Proposition 2.11 uncovers an unexpected relationship between Θ distributions on $U(N)$ and Gaussian distributions on $H(N)$: the normalising factors of these distributions are connected by a simple identity (38).

Proofs of the above-mentioned results may be found in Chapter 3 of [13].

2.1 From Gauss to Shannon

The story of Gaussian distributions is a story of discovery and re-discovery. Different scientists, at different times, were repeatedly lead to these distributions, through different routes. It seems the story began in 1801, on New Year's day, when Giuseppe Piazzi sighted a heavenly body (in fact, the asteroid Ceres), which he thought to be a new planet. Less than six weeks later, this “new planet” disappeared behind the sun. Using a method of least squares, Gauss predicted the area in the sky, where it re-appeared one year later. His justification of this method of least squares (cast in modern language) is that measurement errors follow a family of distributions, which satisfies

Property 1: maximum-likelihood estimation is equivalent to the least-squares problem.

In his *Theoria motus corporum coelestium (1809)*, he used this property to show that the distribution of measurement errors is (again, in modern language) a Gaussian distribution.

In 1810, Laplace studied the distribution of a quantity, which is the aggregate of a great number of elementary observations. He was lead in this (completely different) way, to the same distribution discovered by Gauss. Laplace was among the first scientists to show

Property 2: the distribution of the sum of a large number of elementary observations is (asymptotically) a Gaussian distribution.

Around 1860, Maxwell rediscovered Gaussian distributions, through his investigation of the velocity distribution of particles in an ideal gas (which he viewed as freely-colliding perfect elastic spheres). Essentially, he showed that

Property 3: the distribution of a rotationally-invariant random vector, which has independent components, is a Gaussian distribution.

Kinetic theory lead to another fascinating development, related to Gaussian distributions. Around 1905, Einstein (and, independently, Smoluchowsky) showed that

Property 4: the distribution of the position of a particle, which is undergoing a Brownian motion, is a Gaussian distribution.

In addition to kinetic theory, alternative routes to Gaussian distributions have been found in quantum mechanics, information theory, and other fields. In quantum mechanics, a Gaussian distribution is a position distribution with minimum uncertainty. That is, it achieves equality in Heisenberg's inequality. In information theory, one may attribute to Shannon the following maximum-entropy characterisation

Property 5: a probability distribution with maximum entropy, among all distributions with a given mean and variance, is a Gaussian distribution.

The above list of re-discoveries of Gaussian distributions may be extended much longer. However, the main point is the following. In a Euclidean space, identified with \mathbb{R}^d , any one of the above five properties leads to the same famous expression of a Gaussian distribution,

$$P(dx|\bar{x}, \sigma) = (2\pi\sigma^2)^{-\frac{d}{2}} \exp\left[-\frac{(x - \bar{x})^2}{2\sigma^2}\right] dx$$

as a probability distribution on \mathbb{R}^d , with mean vector $\bar{x} \in \mathbb{R}^d$ and variance parameter $\sigma > 0$ (here dx denotes the Lebesgue measure on \mathbb{R}^d).

In non-Euclidean space, each one of these properties may lead to a different distribution, which may then be called a Gaussian distribution, but only from a restricted point

of view. People interested in Brownian motion may call the heat kernel of a Riemannian manifold a Gaussian distribution on that manifold. However, statisticians will not like this definition, since it will (in general) fail to have a straightforward connection to maximum-likelihood estimation.

2.2 The “right” Gaussian

As of now, the following definition of Gaussian distributions is chosen. Gaussian distributions, on a Riemannian manifold M , are a family of distributions $P(\bar{x}, \sigma)$, parameterised by $\bar{x} \in M$ and $\sigma > 0$, such that: a maximum-likelihood estimate \hat{x}_N of \bar{x} , based on samples $(x_n; n = 1, \dots, N)$ from $P(\bar{x}, \sigma)$, is a solution of the least-squares problem

$$\text{minimise over } x \in M \quad \mathcal{E}_N(x) = \frac{1}{2} \sum_{n=1}^N d^2(x_n, x) \quad (1)$$

This means that \hat{x}_N is an empirical barycentre of the samples (x_n) . In order to construct probability distributions $P(\bar{x}, \sigma)$, which satisfy this definition, consider the density profile

$$f(x|\bar{x}, \sigma) = \exp \left[-\frac{d^2(x, \bar{x})}{2\sigma^2} \right] \quad (2)$$

and the normalising factor,

$$Z(\bar{x}, \sigma) = \int_M f(x|\bar{x}, \sigma) \text{vol}(dx) \quad (3)$$

where vol denotes Riemannian volume. If this is finite, then

$$P(dx|\bar{x}, \sigma) = (Z(\bar{x}, \sigma))^{-1} f(x|\bar{x}, \sigma) \text{vol}(dx) \quad (4)$$

is a well-defined probability distribution on M . In 2.4, below, it will be seen that $P(\bar{x}, \sigma)$, as defined by (4), is indeed a Gaussian distribution, if M is a Hadamard manifold and also a homogeneous space. The following propositions will then be helpful.

Proposition 2.1. *Let M be a Hadamard manifold, whose sectional curvatures lie in $[\kappa, 0]$, where $\kappa = -c^2$. Then, for any $\bar{x} \in M$ and $\sigma > 0$, if $Z(\bar{x}, \sigma)$ is given by (3),*

$$Z_0(\sigma) \leq Z(\bar{x}, \sigma) \leq Z_c(\sigma) \quad (5)$$

where $Z_0(\sigma) = (2\pi\sigma^2)^{\frac{d}{2}}$ and $Z_c(\sigma)$ is positive, given by (d denotes the dimension of M)

$$Z_c(\sigma) = \omega_{d-1} \frac{\sigma}{(2c)^{d-1}} \sum_{k=0}^{d-1} (-1)^k \binom{d-1}{k} \frac{\Phi((d-1-2k)\sigma c)}{\Phi'((d-1-2k)\sigma c)} \quad (6)$$

where ω_{d-1} denotes the area of the unit sphere in \mathbb{R}^d , and Φ denotes the standard normal distribution function.

Proposition 2.2. *If M is a Riemannian homogeneous space, and $Z(\bar{x}, \sigma)$ is given by (3), then $Z(\bar{x}, \sigma)$ does not depend on \bar{x} . In other words, $Z(\bar{x}, \sigma) = Z(\sigma)$.*

If M is a Hadamard manifold and also a homogeneous space, then both Propositions 2.1 and 2.2 apply to M . Indeed, if M is a Riemannian homogeneous space, then its sectional curvatures lie within a bounded subset of the real line. Therefore, Proposition 2.1 implies $Z(\bar{x}, \sigma)$ is finite for all $\bar{x} \in M$ and $\sigma > 0$. On the other hand, Proposition 2.2 implies that $Z(\bar{x}, \sigma) = Z(\sigma)$. Then, (4) reduces to

$$P(dx|\bar{x}, \sigma) = (Z(\sigma))^{-1} \exp\left[-\frac{d^2(x, \bar{x})}{2\sigma^2}\right] \text{vol}(dx) \quad (7)$$

and yields a well-defined probability distribution $P(\bar{x}, \sigma)$ on M . This will be the main focus, throughout the following.

Remark : here, readers may wish to recall the concept of a Hadamard manifold, or of a homogeneous space, from [14] (or any other good Riemannian geometry textbook). The point of appealing to these concepts is the following. The assumption that M is a Hadamard manifold implies that geodesic spherical coordinates, which cover all of M , can be introduced at any point $\bar{x} \in M$. Proposition 2.1 is obtained by writing the integral (3) in terms of these spherical coordinates, and then applying Riemannian volume comparison theorems, that state, very roughly speaking, that manifolds with more positive curvature have less volume. On the other hand, to say that M is a homogeneous space means that all points $\bar{x} \in M$ are equivalent, so changing the point of origin \bar{x} does not change the integral (3). This is the key to Proposition 2.2.

2.3 The normalising factor $Z(\sigma)$

Assume now $M = G/K$ is a Riemannian symmetric space which belongs to the non-compact case, described in Appendix A.1. In particular, M is a Hadamard manifold and also a homogeneous space. Thus, for each $\bar{x} \in M$ and $\sigma > 0$, there is a well-defined probability distribution $P(\bar{x}, \sigma)$ on M , given by (7). Here, the normalising factor $Z(\sigma)$ can be expressed as a multiple integral, using the integral formula (74) of Proposition A.1, from Appendix A.1. Applying this proposition (with $o = \bar{x}$), it is enough to note

$$f(\varphi(s, a)|\bar{x}, \sigma) = \exp\left[-\frac{\|a\|_B^2}{2\sigma^2}\right]$$

where $\|a\|_B^2 = B(a, a)$, in terms of the $\text{Ad}(G)$ -invariant symmetric bilinear form B (see Appendix A.1). Since this expression only depends on a , (74) yields the following formula

$$Z(\sigma) = \frac{\omega(S)}{|W|} \int_{\mathfrak{a}} \exp\left[-\frac{\|a\|_B^2}{2\sigma^2}\right] \prod_{\lambda \in \Delta_+} |\sinh \lambda(a)|^{m_\lambda} da \quad (8)$$

This formula expresses $Z(\sigma)$ as a multiple integral on the vector space \mathfrak{a} . Recall that the dimension of \mathfrak{a} is known as the rank of M [15].

Example 1 : the easiest instance of (8) arises when M is a hyperbolic space of dimension d , and constant sectional curvature equal to -1 . Then, M has rank equal to 1, so that $\mathfrak{a} = \mathbb{R}\hat{a}$ for some unit vector $\hat{a} \in \mathfrak{a}$. Since the sectional curvature is equal to -1 , there is only one positive root λ , say $\lambda(\hat{a}) = 1$, with multiplicity $m_\lambda = d - 1$. In addition,

$|W| = 2$ because there are two Weyl chambers, $C_+ = \{r\hat{a}; r > 0\}$ and $C_- = \{r\hat{a}; r < 0\}$. Accordingly, (8) reads

$$Z(\sigma) = \frac{\omega_{d-1}}{2} \int_{-\infty}^{+\infty} \exp\left[-\frac{r^2}{2\sigma^2}\right] |\sinh(r)|^{d-1} dr = \omega_{d-1} \int_0^{+\infty} \exp\left[-\frac{r^2}{2\sigma^2}\right] \sinh^{d-1}(r) dr$$

In general, if all distances are divided by $c > 0$, the sectional curvature -1 is replaced by $-c^2$. Thus, when M is a hyperbolic space of dimension d , and sectional curvature $-c^2$,

$$Z(\sigma) = \omega_{d-1} \int_0^{+\infty} \exp\left[-\frac{r^2}{2\sigma^2}\right] (c^{-1} \sinh(cr))^{d-1} dr$$

This is exactly $Z_c(\sigma)$, expressed analytically in (6).

Example 2: another example, also susceptible of analytic expression, is when M is a space of positive-definite matrices with real, complex, or quaternion coefficients. Then, $M = G/K$ with $G = \text{GL}(N, \mathbb{K})$, where $\mathbb{K} = \mathbb{R}, \mathbb{C}$ or \mathbb{H} (real numbers, complex numbers, or quaternions), and $K \subset G$ a maximal compact subgroup, $K = O(N), U(N)$ or $Sp(N)$. In each of these three cases, \mathfrak{a} is the space of $N \times N$ real diagonal matrices, and the positive roots are the linear maps $\lambda(a) = a_{ii} - a_{jj}$ where $i < j$, each one having its multiplicity $m_\lambda = \beta$, ($\beta = 1, 2$ or 4 , for $\mathbb{K} = \mathbb{R}, \mathbb{C}$ or \mathbb{H}). In addition, $\|a\|_B^2 = 4\text{tr}(a^2)$. The Weyl group W is the groupe of permutation matrices in K , so $|W| = N!$, while $S = K/T_N$ where T_N is the subgroup of all diagonal matrices in K . Replacing all of this into (8), it follows that

$$Z(\sigma) = \frac{\omega_\beta(N)}{N!} \int_{\mathfrak{a}} \prod_{i=1}^N \exp\left[-\frac{2a_{ii}^2}{\sigma^2}\right] \prod_{i<j} |\sinh(a_{ii} - a_{jj})|^\beta da \quad (9)$$

where $\omega_\beta(N)$ stands for $\omega(S)$, and $da = da_{11} \dots da_{NN}$. Introducing $x_i = \exp(2a_{ii})$,

$$Z(\sigma) = \frac{\omega_\beta(N)}{2^{NN_\beta} N!} \int_{\mathbb{R}_+^N} |V(x)|^\beta \prod_{i=1}^N \rho(x_i, 2\sigma^2) x_i^{-N_\beta} dx_i$$

where $N_\beta = (\beta/2)(N-1) + 1$, $\rho(x, k) = \exp(-\log^2(x)/k)$ and $V(x) = \prod_{i<j} (x_j - x_i)$ is the Vandermonde determinant. Finally, using the elementary identity

$$\rho(x, k) x^\alpha = \exp\left[\frac{k}{4} \alpha^2\right] \rho\left(e^{-\frac{k}{2} \alpha} x, k\right)$$

it is immediately found that

$$Z(\sigma) = \frac{\omega_\beta(N)}{2^{NN_\beta} N!} \times \exp[-NN_\beta^2(\sigma^2/2)] \times \int_{\mathbb{R}_+^N} |V(u)|^\beta \prod_{i=1}^N \rho(u_i, 2\sigma^2) du_i \quad (10)$$

For the case $\beta = 2$, the integral in (10) will be expressed analytically in 2.6, below.

Remark: curious readers will want to compute $\omega_\beta(N)$. For example, $\omega_2(N)$ can be found using the Weyl integral formula on $U(N)$ [16]. This yields $\omega_2(N) = \text{vol}(U(N))/(2\pi)^N$. The volume of the unitary group can be found by looking at the normalising factor of a

Gaussian unitary ensemble [17]. Specifically, $\text{vol}(U(N)) = (2\pi)^{(N^2+N)/2}/G(N)$, in terms of $G(N) = \Gamma(1) \times \Gamma(2) \times \dots \times \Gamma(N)$ (Γ denotes the Euler Gamma function). ■

Example 3: for this last example, let $M = D_N$ be the Siegel domain [18]. This is the set of $N \times N$ symmetric complex matrices z , such that $I_N - z^\dagger z$ is positive-definite. Here, $M = G/K$, where $G \simeq \text{Sp}(N, \mathbb{R})$ (real symplectic group) and $K \simeq U(N)$ (unitary group). Precisely, G is the group of $2N \times 2N$ complex matrices g , with $g^\dagger \Omega g = \Omega$ and $g^\dagger \Gamma g = \Gamma$, where † denotes the transpose, and where Ω and Γ are the matrices

$$\Omega = \begin{pmatrix} & I_N \\ -I_N & \end{pmatrix} ; \quad \Gamma = \begin{pmatrix} I_N & \\ & -I_N \end{pmatrix}$$

In addition, K is the group of block-diagonal matrices $k = \text{diag}(U, U^*)$ where $U \in U(N)$, and * denotes the conjugate. The action of G on M is given by Möbius transformations,

$$g \cdot z = (Az + B)(Cz + D)^{-1} \quad g = \begin{pmatrix} A & B \\ C & D \end{pmatrix} \quad (11)$$

This action preserves the Siegel metric, which is defined by

$$\langle v, v \rangle_z = \|(I_N - z z^\dagger)^{-1} v\|_B^2 \quad \|v\|_B^2 = \frac{1}{2} \text{tr}(v v^\dagger) \quad (12)$$

where each tangent vector v is identified with a symmetric complex matrix. Now [19],

$$\mathfrak{a} = \left\{ \begin{pmatrix} & a \\ a & \end{pmatrix} ; a = \text{diag}(a_{11}, \dots, a_{NN}) \right\} \quad (13)$$

The positive roots are $\lambda(a) = a_{ii} - a_{jj}$ for $i < j$, and $\lambda(a) = a_{ii} + a_{jj}$ for $i \leq j$, all with $m_\lambda = 1$. The order of the Weyl group is $|W| = 2^N N!$, and $\omega(S) = \text{vol}(U(N))/2^N$. Replacing into (8), it follows that

$$Z(\sigma) = \frac{\text{vol}(U(N))}{2^{2N} N!} \int_{\mathfrak{a}} \prod_{i=1}^N \exp \left[-\frac{a_{ii}^2}{2\sigma^2} \right] \prod_{i < j} \sinh |a_{ii} - a_{jj}| \prod_{i \leq j} \sinh |a_{ii} + a_{jj}| da \quad (14)$$

or, after introducing $u_i = \cosh(2a_{ii})$,

$$Z(\sigma) = 2^{-2N} \text{vol}(U(N)) \times \int_{C^N} V(u) \prod_{n=1}^N w(u_n, 8\sigma^2) du_n \quad (15)$$

where $C^N = \{u \in \mathbb{R}_+^N : u_1 \leq u_2 \leq u_1 \leq \dots \leq u_N\}$ and $w(u, k) = \exp(-a \cosh^2(u)/k)$, while $V(u)$ is the Vandermonde determinant, as in (10).

2.4 MLE and maximum entropy

Let M be a Hadamard manifold, which is also a homogeneous space. Consider the family of distributions $P(\bar{x}, \sigma)$ on M , given by (7) for $\bar{x} \in M$ and $\sigma > 0$. This family of distributions fits the definition of Gaussian distributions, stated at the beginning of 2.2.

Proposition 2.3. *Let $P(\bar{x}, \sigma)$ be given by (7), for $\bar{x} \in M$ and $\sigma > 0$. The maximum-likelihood estimate of the parameter \bar{x} , based on samples $(x_n; n = 1, \dots, N)$ from $P(\bar{x}, \sigma)$, is unique and equal to the empirical barycentre \hat{x}_N of the samples (x_n) .*

This proposition is almost immediate. From (7), one has the log-likelihood function

$$\ell(\bar{x}, \sigma) = -N \log Z(\sigma) - \frac{1}{2\sigma^2} \sum_{n=1}^N d^2(x_n, \bar{x}) \quad (16)$$

Since the first term does not depend on \bar{x} , one may maximise $\ell(\bar{x}, \sigma)$, first over \bar{x} and then over σ . Clearly, maximising over \bar{x} is equivalent to minimising the sum of squared distances $d^2(x_n, \bar{x})$. This is just the least-squares problem (1), whose solution is the empirical barycentre \hat{x}_N . Moreover, \hat{x}_N is unique, since M is a Hadamard manifold [20][4].

Consider now maximum-likelihood estimation of σ . This is better carried out in terms of the natural parameter $\eta = (-2\sigma^2)^{-1}$, or in terms of the moment parameter $\delta = \psi'(\eta)$, where $\psi(\eta) = \log Z(\sigma)$ and the prime denotes the derivative.

Proposition 2.4. *The function $\psi(\eta)$, just defined, is a strictly convex function, which maps the half-line $(-\infty, 0)$ onto \mathbb{R} . The maximum-likelihood estimates of the parameters η and δ are*

$$\hat{\eta}_N = (\psi')^{-1}(\hat{\delta}_N) \quad \text{and} \quad \hat{\delta}_N = \frac{1}{N} \sum_{n=1}^N d^2(x_n, \hat{x}_N) \quad (17)$$

where $(\psi')^{-1}$ denotes the reciprocal function.

Remark: $\hat{\eta}_N$ in (17) is well-defined, since the range of ψ' is equal to $(0, \infty)$. Indeed, one has the following inequalities, analogous to (5),

$$\psi'_0(\eta) \leq \psi'(\eta) \leq \psi'_c(\eta) \quad (18)$$

where $\psi_0(\eta) = \log Z_0(\sigma)$, and $\psi_c(\eta) = \log Z_c(\sigma)$. Now, $\psi'_0(\eta) = n\sigma^2$, which increases to $+\infty$ when σ increases to $+\infty$. On the other hand, since $\eta = (-2\sigma^2)^{-1}$,

$$\psi'_c(\eta) = \sigma^3 \frac{d}{d\sigma} (\log Z_c(\sigma)) \quad (19)$$

which, from (6), is $= 0$ when $\sigma = 0$. Thus, it follows from (18) that ψ' maps the half-line $(-\infty, 0)$ onto the half-line $(0, +\infty)$. ■

An alternative definition of Gaussian distributions is provided by their maximum-entropy property, stated in the following proposition. Here, entropy specifically means Shannon's differential entropy. If P is a probability distribution on M , with probability density function p , this entropy is equal to

$$S(P) = \int_M (\log p(x)) p(x) \text{vol}(dx)$$

Proposition 2.5. *The Gaussian distribution $P(\bar{x}, \sigma)$ is the unique distribution on M , having maximum Shannon entropy, among all distributions P with given barycentre \bar{x} and dispersion $\delta = \mathbb{E}_{x \sim P}[d^2(x, \bar{x})]$. Its entropy is equal to $\psi^*(\delta)$ where ψ^* is the Legendre transform of ψ .*

2.5 Barycentre and covariance

Let M be a Hadamard manifold, which is also a homogeneous space. Consider the barycentre and covariance of the Gaussian distribution $P(\bar{x}, \sigma)$ on M , given by (7).

First, it should be noted $P(\bar{x}, \sigma)$ has a well-defined Riemannian barycentre, since it has finite second-order moments. To see that this is true, it is enough to note that

$$\int_M d^2(\bar{x}, x) P(dx|\bar{x}, \sigma) < \infty$$

Ineded, this integral is just $\psi'(\eta)$ in (18).

Proposition 2.6. *Let $P(\bar{x}, \sigma)$ be given by (7), for $\bar{x} \in M$ and $\sigma > 0$. The Riemannian barycentre of $P(\bar{x}, \sigma)$ is equal to \bar{x} .*

The proof of this proposition relies on the fact that the so-called variance function

$$\mathcal{E}(x) = \frac{1}{2} \int_M d^2(x, y) P(dy|\bar{x}, \sigma) \quad (20)$$

is strongly convex (see [13], Paragraph 2.2.3). Thus, if $\text{grad}\mathcal{E}(\bar{x}) = 0$, then \bar{x} is the global minimiser of \mathcal{E} , and therefore the barycentre of $P(\bar{x}, \sigma)$. However, $\text{grad}\mathcal{E}(\bar{x}) = 0$ follows by a direct application of the following ‘‘Fisher’s identity’’,

$$\int_M (\text{grad}_{\bar{x}} \log p(x|\bar{x}, \sigma)) P(dx|\bar{x}, \sigma) = 0$$

where $\text{grad}_{\bar{x}}$ denotes the gradient with respect to \bar{x} , defined according to the Riemannian metric of M , and $p(x|\bar{x}, \sigma)$ is the probability density function, appearing in (7).

The covariance form of $P(\bar{x}, \sigma)$ is the symmetric bilinear form $C_{\bar{x}}$ on $T_{\bar{x}}M$,

$$C_{\bar{x}}(u, v) = \int_M \langle u, \text{Exp}_{\bar{x}}^{-1}(x) \rangle \langle \text{Exp}_{\bar{x}}^{-1}(x), v \rangle p(x|\bar{x}, \sigma) \text{vol}(dx) \quad u, v \in T_{\bar{x}}M \quad (21)$$

where Exp denotes the Riemannian exponential map (Exp^{-1} is well-defined, since M is a Hadamard manifold).

With $\sigma > 0$ fixed, the map which assigns to $\bar{x} \in M$ the covariance form $C_{\bar{x}}$ is a (0,2)-tensor field on M , here called the covariance tensor of $P(\bar{x}, \sigma)$. In order to compute this tensor field, consider the following situation.

Assume $M = G/K$ is a Riemannian symmetric space. Here, $K = K_o$, the stabiliser in G of $o \in M$. For $k \in K$ and $u \in T_oM$, it is clear $k \cdot u \in T_oM$. This defines a representation of K in the tangent space T_oM , called the isotropy representation. One says that M is an irreducible symmetric space, if this isotropy representation is irreducible.

If M is not irreducible, then it is a product of irreducible Riemannian symmetric spaces $M = M_1 \times \dots \times M_s$ [15] (Proposition 5.5, Chapter VIII. This is the de Rham decomposition of M). Accordingly, for $x \in M$ and $u \in T_xM$, one may write $x = (x_1, \dots, x_s)$ and $u = (u_1, \dots, u_s)$, where $x_r \in M_r$ and $u_r \in T_{x_r}M_r$. Now, looking back at (7), it may be seen that

$$p(x|\bar{x}, \sigma) = \prod_{r=1}^s p(x_r|\bar{x}_r, \sigma) \quad p(x_r|\bar{x}_r, \sigma) = (Z_r(\sigma))^{-1} \exp \left[-\frac{d^2(x_r, \bar{x}_r)}{2\sigma^2} \right] \quad (22)$$

For the following proposition, let $\eta = (-2\sigma^2)^{-1}$ and $\psi_r(\eta) = \log Z_r(\sigma)$.

Proposition 2.7. *Assume that M is a product of irreducible Riemannian symmetric spaces, $M = M_1 \times \dots \times M_s$. The covariance tensor C in (21) is given by*

$$C_{\bar{x}}(u, u) = \sum_{r=1}^s \frac{\psi'_r(\eta)}{\dim M_r} \|u_r\|_{\bar{x}_r}^2 \quad (23)$$

for $u \in T_{\bar{x}}M$ where $\bar{x} = (\bar{x}_1, \dots, \bar{x}_s)$ and $u = (u_1, \dots, u_s)$, with $\bar{x}_r \in M_r$ and $u_r \in T_{\bar{x}_r}M_r$.

Example: let $M = \mathbf{H}(N)$, the space of $N \times N$ Hermitian positive-definite matrices, so $M = \mathrm{GL}(N, \mathbb{C})/U(N)$, with $U(N)$ the stabiliser of $o = I_N$ ($N \times N$ identity matrix). The de Rham decomposition of M is $M = M_1 \times M_2$, where $M_1 = \mathbb{R}$ and M_2 is the submanifold whose elements are those $x \in M$ such that $\det(x) = 1$. Accordingly, each $\bar{x} \in M$ is identified with the couple (\bar{x}_1, \bar{x}_2) ,

$$\bar{x}_1 = \frac{1}{N} \log \det(\bar{x}) \quad \bar{x}_2 = (\det(\bar{x}))^{-1/N} \bar{x}$$

and each $u \in T_{\bar{x}}M$ is written $u = u_1 \bar{x} + u_2$

$$u_1 = \frac{1}{N} \mathrm{tr}(\bar{x}^{-1}u) \quad u_2 = u - \frac{1}{N} \mathrm{tr}(\bar{x}^{-1}u) \bar{x}$$

These may be replaced into expression (23),

$$C_{\bar{x}}(u, u) = \psi'_1(\eta) u_1^2 + \frac{\psi'_2(\eta)}{N^2 - 1} \|u_2\|_{\bar{x}_2}^2 \quad (24)$$

where $\psi_1(\eta) = \log(2\pi\sigma^2)^{\frac{1}{2}}$, and $\psi_2(\eta) = \log Z(\sigma) - \psi_1(\eta)$ ($Z(\sigma)$ is given by (26) in 2.6, below). After a direct calculation, this can be brought under the form

$$C_{\bar{x}}(u, u) = g_1(\sigma) \mathrm{tr}^2(\bar{x}^{-1}u) + g_2(\sigma) \mathrm{tr}(\bar{x}^{-1}u)^2 \quad (25)$$

where $g_1(\sigma)$ and $g_2(\sigma)$ are certain functions of σ .

Remark: as a corollary of Proposition 2.7, the covariance tensor C is a G -invariant Riemannian metric on M . This is clear, for example, in the special case of (25), which coincides with the general expression of a $\mathrm{GL}(N, \mathbb{C})$ -invariant metric. ■

2.6 $Z(\sigma)$ from RMT

Random matrix theory is very helpful in the calculation of integrals such as (10) and (15), leading both to exact expressions and to asymptotic expansions of these integrals. Here, this is illustrated for the integral (10), with $\beta = 2$. This corresponds to $M = \mathbf{H}(N)$, the space of $N \times N$ Hermitian positive-definite matrices. In this case, it is possible to provide an analytic formula for the normalising factor $Z(\sigma)$.

Proposition 2.8. *When $M = \mathbf{H}(N)$, the normalising factor $Z(\sigma)$, given by (10) with $\beta = 2$, admits of the following analytic expression*

$$Z(\sigma) = \frac{\omega_2(N)}{2^{N^2}} (2\pi\sigma^2)^{\frac{N}{2}} \exp \left[\left(\frac{N^3 - N}{6} \right) \sigma^2 \right] \prod_{n=1}^{N-1} \left(1 - e^{-n\sigma^2} \right)^{N-n} \quad (26)$$

The proof of this proposition is a direct application of a well-known formula from random matrix theory [17]. The integral in (10) reads

$$I_N(\sigma) = \int_{\mathbb{R}_+^N} |V(u)|^\beta \prod_{i=1}^N \rho(u_i, 2\sigma^2) du_i$$

According to [17] (Chapter 5, Page 79), if $(p_n; n = 0, 1, \dots)$ are orthonormal polynomials with respect to the weight function $\rho(u, 2\sigma^2)$ on \mathbb{R}_+ , then $I_N(\sigma)$ is given by

$$I_N(\sigma) = N! \prod_{n=0}^{N-1} p_{nn}^{-2} \quad (27)$$

where p_{nn} is the leading coefficient in p_n . The required orthonormal polynomials p_n are given by $p_n = (2\pi\sigma^2)^{-\frac{1}{4}} s_n$, where s_n are the Stieltjes-Wigert polynomials [21] (Page 33). Looking up the expression of these polynomials, it is easy to find

$$p_{nn}^{-2} = (2\pi\sigma^2)^{\frac{1}{2}} \exp\left[\frac{(2n+1)^2}{2}\sigma^2\right] \prod_{m=1}^n (1 - e^{-m\sigma^2})$$

Then, working out the product (27) and replacing into (10), one is lead to (26).

Moving on, it is possible to derive an asymptotic expression of $Z(\sigma)$, valid in the limit where N goes to infinity while the product $t = N\sigma^2$ remains constant.

Proposition 2.9. *Let $Z(\sigma)$ be given by (26). If $N \rightarrow \infty$, while $t = N\sigma^2$ remains constant, then the following equivalence holds,*

$$\frac{1}{N^2} \log Z(\sigma) \sim -\frac{1}{2} \log\left(\frac{2N}{\pi}\right) + \frac{3}{4} + \frac{t}{6} - \frac{\text{Li}_3(e^{-t}) - \zeta(3)}{t^2} \quad (28)$$

where $\text{Li}_3(x) = \sum_{k=1}^{\infty} x^k/k^3$ for $|x| < 1$ (the trilogarithm), and ζ is the Riemann Zeta function.

The main idea behind (28) is that, taking the logarithm in (26), the product on the right-hand side turns into a Riemann sum for the improper integral

$$\int_0^1 (1-x) \log(1 - e^{-tx}) dx = -(\text{Li}_3(e^{-t}) - \zeta(3))/t^2$$

where the equality follows by integrating term-by-term the power series of the logarithm.

2.7 The asymptotic distribution

From the point of view of random matrix theory, a Gaussian distribution $P(I_N, \sigma)$ on $M = \text{H}(N)$ defines a unitary matrix ensemble. If x is a random matrix, drawn from this ensemble, and $(x_i; i = 1, \dots, N)$ are its eigenvalues, which all belong to $(0, \infty)$, then the empirical distribution ν_N , which is given by (as usual, δ_{x_i} is the Dirac distribution at x_i)

$$\nu_N(B) = \mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \delta_{x_i}(B) \right] \quad (29)$$

for measurable $B \subset (0, \infty)$, converges to an absolutely continuous distribution ν_t , when N goes to infinity, while the product $t = N\sigma^2$ remains constant.

Proposition 2.10. *Let $c = e^{-t}$ and $a(t) = c(1 + \sqrt{1 - c})^{-2}$ while $b(t) = c(1 - \sqrt{1 - c})^{-2}$. When N goes to infinity, while the product $t = N\sigma^2$ remains constant, the empirical distribution ν_N converges weakly to the distribution ν_t with probability density function*

$$\frac{d\nu_t}{dx}(x) = \frac{1}{\pi tx} \arctan\left(\frac{4e^t x - (x+1)^2}{x+1}\right) \mathbf{1}_{[a(t), b(t)]}(x) \quad (30)$$

where $\mathbf{1}_{[a(t), b(t)]}$ denotes the indicator function of the interval $[a(t), b(t)]$.

Remark: as one should expect, when $t = 0$ (so $\sigma^2 = 0$), $a(t) = b(t) = 1$. ■

The proof of Proposition 2.10 is a relatively direct application of a result in [22] (Page 191). In fact, the integration variables in (10) are $u_i = e^t x_i$. Let $\tilde{\nu}_N$ be the empirical distribution of the u_i (this is the same as (29), but with u_i instead of x_i). By applying [17] (Chapter 5, Page 81),

$$\tilde{\nu}_N(B) = \frac{1}{N} \int_B R_N^{(1)}(u) (du) \quad (31)$$

for measurable $B \subset (0, \infty)$, where the one-point correlation function $R_N^{(1)}(u)$ is given by

$$R_N^{(1)}(u) = \rho(u, 2\sigma^2) \sum_{n=0}^{N-1} p_n^2(u) \quad (32)$$

in the notation of 2.6 (p_n are orthonormal polynomials, with respect to the weight $\rho(u, 2\sigma^2)$). According to [23] (Page 133), $\tilde{\nu}_N$ given by (31) converges weakly to the so-called equilibrium distribution $\tilde{\nu}_t$, which minimises the electrostatic energy functional

$$E(\nu) = \frac{1}{t} \int_0^\infty \frac{1}{2} \log^2(u) \nu(du) - \int_0^\infty \int_0^\infty \log|u - v| \nu(du) \nu(dv) \quad (33)$$

over probability distributions ν on $(0, \infty)$. Also according to [23] (Page 133), this equilibrium distribution is the asymptotic distribution of the zeros of the polynomial p_N (in the limit $N \rightarrow \infty$ while $N\sigma^2 = t$). Fortunately, p_N is just a constant multiple of the Stieltjes-Wigert polynomial s_N [21] (Page 33). Therefore, the required asymptotic distribution of zeros can be read from [22] (Page 191). Finally, (30) follows by introducing the change of variables $x = e^{-t}u$.

Remark: in [24], the equilibrium distribution $\tilde{\nu}_t$ is derived directly, by searching for stationary distributions of the energy functional (33). This leads to a singular integral equation, whose solution reduces to a Riemann-Hilbert problem. Astoundingly, the Gaussian distributions on $H(N)$, as introduced in the present chapter, provide a matrix model for Chern-Simons quantum field theory (a detailed account is given in [24]). ■

2.8 Duality: the Θ distributions

Recall the Riemannian symmetric space $M = H(N)$ of 2.6. Its dual space is the unitary group $M^* = U(N)$ (the definition of *duality* may be found in Appendix A.2).

Consider now a family of distributions on M^* , which will be called Θ distributions, and which display an interesting connection with Gaussian distributions on M , studied

in 2.6. Recall Jacobi's ϑ function¹,

$$\vartheta(e^{i\phi}|\sigma^2) = \sum_{m=-\infty}^{+\infty} \exp(-m^2\sigma^2 + 2mi\phi)$$

As a function of ϕ , up to some minor modifications, this is a wrapped normal distribution (in other words, the heat kernel of the unit circle),

$$\frac{1}{2\pi} \vartheta(e^{i\phi}|\frac{\sigma^2}{2}) = \sum_{m=-\infty}^{\infty} \exp\left[-\frac{(2\phi - 2m\pi)^2}{2\sigma^2}\right]$$

Each $x \in M^*$ can be written $x = k \cdot e^{i\theta}$ where $k \in U(N)$ and $e^{i\theta} = \text{diag}(e^{i\theta_j}; j = 1, \dots, N)$. Here, $k \cdot y = kyk^\dagger$ for $y \in M^*$. With this notation, define the following matrix ϑ function,

$$\Theta(x|\sigma^2) = k \cdot \vartheta(e^{i\theta}|\frac{\sigma^2}{2}) \quad (34)$$

which is obtained from x by applying Jacobi's ϑ function to each eigenvalue of x . Further, consider the positive function,

$$f_*(x|\bar{x}, \sigma) = \det\left[(2\pi\sigma^2)^{\frac{1}{2}} \Theta(x\bar{x}^\dagger|\sigma^2)\right] \quad (35)$$

where $\bar{x} \in M^*$. This is also equal to

$$\det\left[(2\pi\sigma^2)^{\frac{1}{2}} \Theta(\bar{x}^\dagger x|\sigma^2)\right]$$

since the matrices $x\bar{x}^\dagger$ and $\bar{x}^\dagger x$ are similar. Then, let $Z_{M^*}(\sigma)$ denote the normalising constant

$$Z_{M^*}(\sigma) = \int_{M^*} f_*(x|\bar{x}, \sigma) \text{vol}(dx) \quad (36)$$

which does not depend on \bar{x} , as can be seen, by introducing the new variable of integration $z = x\bar{x}^\dagger$, and using the invariance of $\text{vol}(dx)$.

Now, define a Θ distribution $\Theta(\bar{x}, \sigma)$ as the probability distribution on M^* , whose probability density function, with respect to $\text{vol}(dx)$, is given by

$$p_*(x|\bar{x}, \sigma) = (Z_{M^*}(\sigma))^{-1} f_*(x|\bar{x}, \sigma) \quad (37)$$

Proposition 2.11. *Let $Z_M(\sigma) = Z(\sigma)$ be given by (26), and $Z_{M^*}(\sigma)$ be given by (36). Then, the following equality holds*

$$\frac{Z_M(\sigma)}{Z_{M^*}(\sigma)} = \exp\left[\left(\frac{N^3 - N}{6}\right)\sigma^2\right] \quad (38)$$

Remark: the Gaussian density (7) on M , and the Θ distribution density (37) on M^* are apparently unrelated. Therefore, it is interesting to note their normalising constants $Z_M(\sigma)$ and $Z_{M^*}(\sigma)$ scale together according to the simple relation (38). The connection between the two distributions is due to the duality between M and M^* . ■

¹To follow the original notation of Jacobi [25], this should be written $\vartheta(e^{i\phi}|q)$ where $q = e^{-\sigma^2}$. In other popular notations, this function is called ϑ_{00} or ϑ_3 .

3 Gaussian distributions and Bayesian inference

This section aims to investigate Bayesian inference for Gaussian distributions. Precisely, it aims to study Bayesian estimation of the parameter x of a Gaussian distribution $P(x, \sigma)$, when this parameter is assigned a prior density which is also Gaussian, say $P(z, \tau)$. Both the prior density $P(z, \tau)$ and the likelihood density $P(x, \sigma)$ are defined on a Riemannian symmetric space of non-positive curvature M ².

Paragraph 3.1 begins by expressing the posterior density $\pi(x)$, based on the general definition

$$\pi(x) \propto \text{prior density} \times \text{likelihood density}$$

Here, $\pi(x)$ will remain partially unknown, as the missing normalising factor cannot be determined. Then, two Bayesian estimators are studied. The maximum *a posteriori* \hat{x}_{MAP} is the mode of $\pi(x)$,

$$\hat{x}_{\text{MAP}} = \operatorname{argmax}_{x \in M} \pi(x)$$

while the minimum mean square error estimator \hat{x}_{MMS} , classically understood as the mean (expectation) of the posterior density, is here the Riemannian barycentre of $\pi(x)$.

It is seen that \hat{x}_{MAP} can be computed directly, being a geodesic convex combination of the prior barycentre z and a new observation y , with respective weights $1 - \rho$ and ρ , where $\rho = \tau^2 / (\sigma^2 + \tau^2)$. On the other hand, \hat{x}_{MMS} seems much harder to compute.

However, Proposition 3.1 states that $\hat{x}_{\text{MMS}} = \hat{x}_{\text{MAP}}$ if $\rho = 1/2$, and Proposition 3.2 states that, in the special case where M is a hyperbolic space, \hat{x}_{MMS} is a geodesic convex combination of z and y , just like \hat{x}_{MAP} , but with different weights, say $(1 - t^*)$ and t^* .

Paragraph 3.2 reports on numerical experiments which show, again in the special case where M is a hyperbolic space, that \hat{x}_{MMS} and \hat{x}_{MAP} lie very close to each other, and that they even appear to be equal (this would mean $t^* = \rho$). At present, I am unaware of any mathematical explanation of this phenomenon.

Paragraph 3.3 describes the computational tools employed in calculating \hat{x}_{MMS} . First, Proposition 3.3 provides easy-to-verify sufficient conditions, for the geometric ergodicity of an isotropic Metropolis-Hastings Markov chain, in a Riemannian symmetric space M . These conditions are shown to apply in the case of the posterior density $\pi(x)$, making it possible to generate geometrically ergodic samples $(x_n; n \geq 1)$ from this density.

The next Proposition 3.4 states that the empirical barycentre \bar{x}_N of the samples (x_1, \dots, x_N) converges almost-surely to \hat{x}_{MMS} , so \bar{x}_N may be used to approximate \hat{x}_{MMS} to any required accuracy.

Concretely, computing the empirical barycentre \bar{x}_N requires solving a strongly convex optimisation problem on the Riemannian manifold M (here, convexity is with respect to the Riemannian connection of M [26]). This has come to be called “geodesic convexity”). Appendix B is devoted to a brief but systematic study of convex optimisation on Riemannian manifolds. Specifically, it establishes the rate of convergence of Riemannian gradient descent schemes, applied to strictly convex or strongly convex cost functions.

The gradient descent schemes under consideration are retraction schemes (not limited to the Riemannian exponential) with a constant step-size. The problem is then to find the largest possible step-size which guarantees a certain rate of convergence. Proposition B.8 addresses this problem for strictly convex functions, and Proposition B.9 for strongly

²Proofs of the results stated in this section can be found in Chapter 4 of [13].

convex functions. For any strictly convex cost function, and suitable retraction, Proposition B.8 gives the largest possible step-size which guarantees a rate of convergence at least as fast as $O(1/t)$ (t is the number of iterations). Proposition B.9 does the same for strongly convex functions, but with an exponential rate of convergence. In fact, with regard to the original motivation of computing \bar{x}_N , this ensures that a gradient descent scheme, using the Riemannian exponential, converges after only a few iterations.

3.1 MAP versus MMS

Assume $M = G/K$ is a Riemannian symmetric space which belongs to the non-compact case, described in Appendix A.1. Recall the Gaussian distribution $P(x, \sigma)$ on M given by its probability density function (7)

$$p(y|x, \sigma) = (Z(\sigma))^{-1} \exp \left[-\frac{d^2(y, x)}{2\sigma^2} \right] \quad (39)$$

In 2.4, it was seen that maximum-likelihood estimation of the parameter x , based on samples $(y_n; n = 1, \dots, N)$, amounts to computing the empirical barycentre of these samples. The one-sample maximum-likelihood estimate, given a single observation y , is therefore $\hat{x}_{ML} = y$.

Instead of maximum-likelihood estimation, consider Bayesian estimation of x , based on the observation y . To do so, assign to x a prior density, which is also Gaussian,

$$p(x|z, \tau) = (Z(\tau))^{-1} \exp \left[-\frac{d^2(x, z)}{2\tau^2} \right] \quad (40)$$

Upon observation of y , Bayesian inference concerning x is carried out using the posterior density

$$\pi(x) \propto \exp \left[-\frac{d^2(y, x)}{2\sigma^2} - \frac{d^2(x, z)}{2\tau^2} \right] \quad (41)$$

where \propto indicates a missing (unknown) normalising factor.

In particular, the maximum *a posteriori* estimator \hat{x}_{MAP} of x is equal to the mode of the posterior density $\pi(x)$. In other words, \hat{x}_{MAP} minimises the weighted sum of squared distances $d^2(y, x)/\sigma^2 + d^2(x, z)/\tau^2$. This is expressed in the following notation³,

$$\hat{x}_{\text{MAP}} = z \#_{\rho} y \quad \text{where } \rho = \frac{\tau^2}{\sigma^2 + \tau^2} \quad (42)$$

Thus, \hat{x}_{MAP} is a geodesic convex combination of the prior barycentre z and the observation y , with respective weights $\sigma^2/(\sigma^2 + \tau^2)$ and $\tau^2/(\sigma^2 + \tau^2)$.

On the other hand, the minimum mean square error estimator \hat{x}_{MMS} is the barycentre of the posterior density $\pi(x)$. That is, \hat{x}_{MMS} is the unique global minimiser of

$$\mathcal{E}_{\pi}(w) = \frac{1}{2} \int_M d^2(w, x) \pi(x) \text{vol}(dx) \quad (43)$$

³If $p, q \in M$ and $c : [0, 1] \rightarrow M$ is a geodesic curve with $c(0) = p$ and $c(1) = q$, then $p \#_t q = c(t)$, for $t \in [0, 1]$. In other words, $p \#_t q$ is a geodesic convex combination of p and q , with respective weights $(1 - t)$ and t .

While it is easy to compute \hat{x}_{MAP} from (42), it is much harder to find \hat{x}_{MMS} , as this requires minimising the integral (43), where the density $\pi(x)$ is known only up to normalisation.

Still, there is one special case where these two estimators are equal.

Proposition 3.1. *In the above notation, if $\sigma^2 = \tau^2$ (that is $\rho = 1/2$), then $\hat{x}_{\text{MMS}} = \hat{x}_{\text{MAP}}$.*

When M is a Euclidean space, it is well-known that $\hat{x}_{\text{MMS}} = \hat{x}_{\text{MAP}}$ for any value of ρ . When M is a space of constant negative curvature, the following proposition indicates \hat{x}_{MMS} and \hat{x}_{MAP} cannot be too far away from one another.

Proposition 3.2. *In the above notation, if M is a space of constant negative curvature (hyperbolic space), then $\hat{x}_{\text{MMS}} = z \#_{t^*} y$ for some $t^* \in (0, 1)$.*

3.2 Bounding the distance

In general, one expects \hat{x}_{MMS} and \hat{x}_{MAP} to be different from one another, when $\rho \neq 1/2$. However, when M is a space of constant negative curvature, Proposition 3.2 shows the distance between these two estimators is always less than the distance between z and y .

Surprisingly (again when M is a space of constant negative curvature), numerical experiments show that \hat{x}_{MMS} and \hat{x}_{MAP} lie very close to each other, and that they even appear to be equal. I am unaware of any mathematical explanation of this phenomenon.

It is possible to bound the distance between \hat{x}_{MMS} and \hat{x}_{MAP} , using the fundamental contraction property [20] (this is an immediate application of Jensen's inequality, as explained in the proof of Theorem 6.3 in [20]).

$$d(\hat{x}_{\text{MMS}}, \hat{x}_{\text{MAP}}) \leq W(\pi, \delta_{\hat{x}_{\text{MAP}}}) \quad (44)$$

where W denotes the Kantorovich (L^1 -Wasserstein) distance, and $\delta_{\hat{x}_{\text{MAP}}}$ denotes the Dirac probability distribution concentrated at \hat{x}_{MAP} . Now, the right-hand side of (44) is equal to the first-order moment

$$m_1(\hat{x}_{\text{MAP}}) = \int_M d(\hat{x}_{\text{MAP}}, x) \pi(x) \text{vol}(dx) \quad (45)$$

Of course, the upper bound in (44) is not tight, since it is strictly positive, even when $\rho = 1/2$, as one may see from (45).

It will be shown below that a Metropolis-Hastings algorithm, with Gaussian proposals, can be used to generate geometrically ergodic samples $(x_n; n \geq 1)$ from the posterior density π . It is therefore possible to approximate (45) by an empirical average

$$\bar{m}_1(\hat{x}_{\text{MAP}}) = \frac{1}{N} \sum_{n=1}^N d(\hat{x}_{\text{MAP}}, x_n) \quad (46)$$

In addition, the samples (x_n) can be used to compute a convergent approximation of \hat{x}_{MMS} . Precisely, the empirical barycentre \bar{x}_{MMS} of the samples (x_1, \dots, x_N) converges almost-surely to \hat{x}_{MMS} (this is a result of Proposition 3.4).

Numerical experiments were conducted in the case where M is a hyperbolic space of curvature equal to -1 and of dimension d . The following table was obtained for the values $\sigma^2 = \tau^2 = 0.1$, using samples (x_1, \dots, x_N) where $N = 2 \times 10^5$.

dimension d	2	3	4	5	6	7	8	9	10
$\bar{m}_1(\hat{x}_{\text{MAP}})$	0.28	0.35	0.41	0.47	0.50	0.57	0.60	0.66	0.70
$d(\bar{x}_{\text{MMS}}, \hat{x}_{\text{MAP}})$	0.00	0.00	0.00	0.01	0.01	0.02	0.02	0.02	0.03

and the following table for $\sigma^2 = 1$ and $\tau^2 = 0.5$, again using $N = 2 \times 10^5$.

dimension d	2	3	4	5	6	7	8	9	10
$\bar{m}_1(\hat{x}_{\text{MAP}})$	0.75	1.00	1.12	1.44	1.73	1.97	2.15	2.54	2.91
$d(\bar{x}_{\text{MMS}}, \hat{x}_{\text{MAP}})$	0.00	0.00	0.03	0.02	0.02	0.03	0.04	0.03	0.12

The first table confirms Proposition 3.1. The second table, more surprisingly, shows that \hat{x}_{MMS} and \hat{x}_{MAP} can be quite close to each other, even when $\rho \neq 1/2$.

In both of these tables, $d(\bar{x}_{\text{MMS}}, \hat{x}_{\text{MAP}})$ is an approximation of $d(\hat{x}_{\text{MMS}}, \hat{x}_{\text{MAP}})$, based on using the empirical barycentre \bar{x}_{MMS} instead of \hat{x}_{MMS} . The main source of error affecting this approximation is the fact that the samples (x_1, \dots, x_N) follow from a Metropolis-Hastings algorithm, and not directly from the posterior density π .

Other values of σ^2 and τ^2 lead to similar orders of magnitude for $\bar{m}_1(\hat{x}_{\text{MAP}})$ and $d(\bar{x}_{\text{MMS}}, \hat{x}_{\text{MAP}})$. While $\bar{m}_1(\hat{x}_{\text{MAP}})$ increases with the dimension d , $d(\bar{x}_{\text{MMS}}, \hat{x}_{\text{MAP}})$ does not appear sensitive to increasing dimension.

Based on these experimental results, one is tempted to conjecture that $\hat{x}_{\text{MMS}} = \hat{x}_{\text{MAP}}$, even when $\rho \neq 1/2$. Of course, numerical experiments do not equate to a mathematical proof.

3.3 Computing the MMS

3.3.1 Metropolis-Hastings algorithm

A crucial step in Bayesian inference is sampling from the posterior density. Here, this is $\pi(x)$ given by (41). Since $\pi(x)$ is known only up to normalisation, a suitable sampling method is afforded by the Metropolis-Hastings algorithm. This algorithm generates a Markov chain $(x_n; n \geq 1)$, with transition kernel [27]

$$Pf(x) = \int_M \alpha(x, y)q(x, y)f(y)\text{vol}(dy) + \rho(x)f(x) \quad (47)$$

for any bounded measurable function $f : M \rightarrow \mathbb{R}$, where $\alpha(x, y)$ is the probability of accepting a transition from x to dy , and $\rho(x)$ is the probability of staying at x , and where $q(x, y)$ is the proposed transition density

$$q(x, y) \geq 0 \text{ and } \int_M q(x, y)\text{vol}(dy) = 1 \quad \text{for } x \in M \quad (48)$$

In the following, (x_n) will always be an isotropic Metropolis-Hastings chain, in the sense that $q(x, y) = q(d(x, y))$, so $q(x, y)$ only depends on the distance $d(x, y)$. In this case, the acceptance probability $\alpha(x, y)$ is given by $\alpha(x, y) = \min\{1, \pi(y)/\pi(x)\}$.

The aim of the Metropolis-Hastings algorithm is to produce a Markov chain (x_n) which is geometrically ergodic. Geometric ergodicity means the distribution π_n of x_n converges to π , with a geometric rate, in the sense that there exist $\beta \in (0, 1)$ and $R(x_1) \in (0, \infty)$, as well as a function $V : M \rightarrow \mathbb{R}$, such that (in the following, $\pi(dx) = \pi(x)\text{vol}(dx)$)

$$V(x) \geq \max\{1, d^2(x, x^*)\} \text{ for some } x^* \in M \quad (49)$$

$$\left| \int_M f(x)(\pi_n(dx) - \pi(dx)) \right| \leq R(x_1)\beta^n \quad (50)$$

for any function $f : M \rightarrow \mathbb{R}$ with $|f| \leq V$. If the chain (x_n) is geometrically ergodic, then it satisfies the strong law of large numbers [28]

$$\frac{1}{N} \sum_{n=1}^N f(x_n) \longrightarrow \int_M f(x)\pi(dx) \quad (\text{almost-surely}) \quad (51)$$

as well as a corresponding central limit theorem (see Theorem 17.0.1, in [28]). Then, in practice, the Metropolis-Hastings algorithm can be used to generate samples (x_n) from the posterior density $\pi(x)$.

The following general statement can be proved, concerning the geometric ergodicity of isotropic Metropolis-Hastings chains. The proof (see [13], Section 4.6) is a generalisation of the one carried out in the special case where M is a Euclidean space [29].

Proposition 3.3. *Let M be a Riemannian symmetric space, which belongs to the non-compact case. Assume $(x_n; n \geq 1)$ is a Markov chain in M , with transition kernel given by (47), with proposed transition density $q(x, y) = q(d(x, y))$, and with strictly positive invariant density π .*

The chain (x_n) satisfies (49) and (50), if the following assumptions hold,

(a1) there exists $x^ \in M$, such that $r(x) = d(x^*, x)$ and $\ell(x) = \log \pi(x)$ satisfy*

$$\limsup_{r(x) \rightarrow \infty} \frac{\langle \text{grad } r, \text{grad } \ell \rangle_x}{r(x)} < 0$$

(a2) if $n(x) = \text{grad } \ell(x) / \|\text{grad } \ell(x)\|$, then $n(x)$ satisfies

$$\limsup_{r(x) \rightarrow \infty} \langle \text{grad } r, n \rangle_x < 0$$

(a3) there exist $\delta_q > 0$ and $\varepsilon_q > 0$ such that $d(x, y) < \delta_q$ implies $q(x, y) > \varepsilon_q$

Remark : the posterior density π in (41) verifies Assumptions (a1) and (a2). To see this, let $x^* = z$, and write

$$\text{grad } \ell(x) = -\frac{1}{\tau^2} r(x) \text{grad } r(x) - \frac{1}{\sigma^2} \text{grad } f_y(x)$$

where $f_y(x) = d^2(y, x)/2$. Then, taking the scalar product with $\text{grad } r$,

$$\langle \text{grad } r, \text{grad } \ell \rangle_x = -\frac{1}{\tau^2} r(x) - \frac{1}{\sigma^2} \langle \text{grad } r, \text{grad } f_y \rangle_x \quad (52)$$

since $\text{grad } r(x)$ is a unit vector, for all $x \in M$. Now, $\text{grad } f_y(x) = -\text{Exp}_x^{-1}(y)$ (see [30]). But, since $r(x)$ is a convex function of x , it follows, by (85) in Appendix B.1, that

$$\langle \text{grad } r, \text{Exp}_x^{-1}(y) \rangle \leq r(y) - r(x)$$

for any $y \in M$. Thus, the right-hand side of (52) is strictly negative, as soon as $r(x) > r(y)$, and Assumption (a1) is indeed verified. That Assumption (a2) is also verified can be proved by a similar reasoning. ■

Remark: on the other hand, Assumption (a3) holds, if the proposed transition density $q(x, y)$ is a Gaussian density, $q(x, y) = p(y|x, \tau_a)$. With this choice of $q(x, y)$, all the assumptions of Proposition 3.3 are verified, for the posterior density π in (41). Therefore, Proposition 3.3 implies that the Metropolis-Hastings algorithm generates geometrically ergodic samples $(x_n; n \geq 1)$, from this posterior density. ■

3.3.2 The empirical barycentre

Let $(x_n; n \geq 1)$ be a Metropolis-Hastings Markov chain in M , with its transition kernel (47), and invariant density π . Assume the chain (x_n) is geometrically ergodic, so it satisfies the strong law of large numbers (51).

Then, let \bar{x}_N denote the empirical barycentre of the first N samples (x_1, \dots, x_N) . This is the unique global minimum of the variance function

$$\mathcal{E}_N(w) = \frac{1}{2N} \sum_{n=1}^N d^2(w, x_n) \quad (53)$$

Let \hat{x} denote the Riemannian barycentre of the invariant density π . It turns out that \bar{x}_N converges almost-surely to \hat{x} .

Proposition 3.4. *Let (x_n) be any Markov chain in a Hadamard manifold M , with invariant distribution π . Denote \bar{x}_N the empirical barycentre of (x_1, \dots, x_N) , and \hat{x} the Riemannian barycentre of π . If (x_n) satisfies the strong law of large numbers (51), then \bar{x}_N converges to \hat{x} , almost-surely.*

The proof of Proposition 3.4 is nearly a word-for-word repetition of the proof in [31] (that of Theorem 2.3).

According to the remarks after Proposition 3.3, the Metropolis-Hastings Markov chain (x_n) , whose invariant density is the posterior density $\pi(x)$, given by (41), is geometrically ergodic. Therefore, by Proposition 3.4, the empirical barycentre \bar{x}_{MMS} , of the samples (x_1, \dots, x_N) , converges almost-surely to the minimum mean square error estimator \hat{x}_{MMS} (since this is just the barycentre of the posterior density π). This provides a practical strategy for approximating \hat{x}_{MMS} . Indeed, \bar{x}_{MMS} can be computed using the Riemannian gradient descent method (this method is discussed in Appendix B.4).

Remark: this strategy for approximating \hat{x}_{MMS} provided the numerical results discussed

in Paragraph 3.2. For an additional, visual illustration, consider (as in Paragraph 3.2) the case where M is a space of constant negative curvature -1 , and of dimension $d = 2$. Figure 1 represents M in the shape of the Poincaré disc. The prior barycentre z is designated by a square \square and the observation y by a circle \circ . Grey crosses \times mark the last 1000 out of $N = 100000$ samples x_n generated using the Metropolis-Hastings kernel (47), and the empirical barycentre \bar{x}_{MMS} is designated by a black circle \bullet . In both of the Subfigures 1a and 1b, \bar{x}_{MMS} is seen to lie on the geodesic connecting z and y , here the dashed circle arc. Note that 1a corresponds to Proposition 3.1 and 1b to Proposition 3.2.

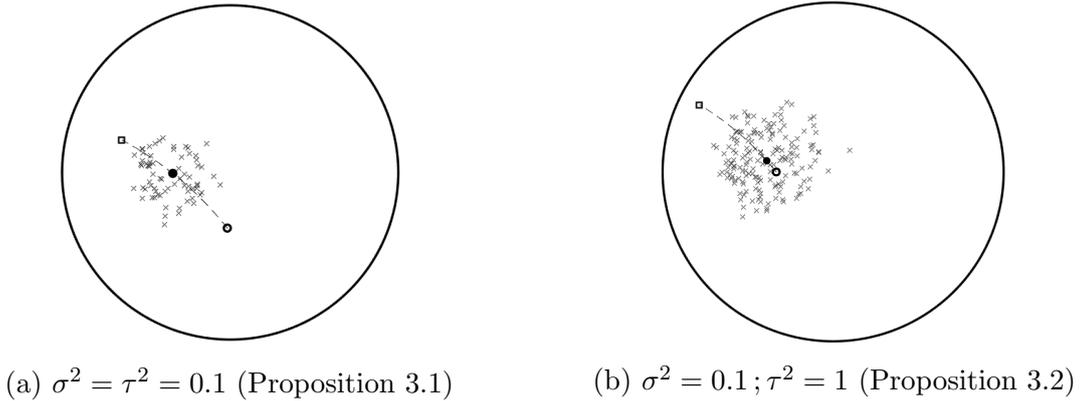


Figure 1: Poincaré disc with $z = \square; y = \circ; \bar{x}_{\text{MMS}} = \bullet$

3.4 Proof of Proposition 3.2

Neither this proposition, nor its proof, appeared in [13]. The proof is here given in a series of lemmas. Recall that M is now a hyperbolic space (simply-connected space of constant negative curvature -1).

Lemma 3.1. *Let $\gamma : \mathbb{R} \rightarrow M$ denote the geodesic curve with $\gamma(0) = z$ and $\gamma(1) = y$. Then, \hat{x}_{MMS} lies on this geodesic curve γ .*

Proof: recall from [32] (Section 3) that there exists an isometry $\sigma : M \rightarrow M$, such that $\sigma \circ \sigma$ is the identity (σ is an involution), and the set of fixed points of σ is exactly the geodesic curve γ . The key point in the following, which can be seen from (41), is that

$$\pi(\sigma(x)) = \pi(x) \quad \text{for } x \in M \quad (54)$$

In other words, σ leaves invariant the posterior density π . Let \mathcal{E}_π be the function in (43). Then, note that

$$(\mathcal{E}_\pi \circ \sigma)(w) = \frac{1}{2} \int_M d^2(w, \sigma(x)) \pi(x) \text{vol}(dx) = \frac{1}{2} \int_M d^2(w, x) \pi(\sigma(x)) \sigma^*(\text{vol})(dx)$$

where the first equality follows from (41), because σ is an isometry and an involution, and the second equality by a change of variables ($\sigma^*(\text{vol})$ denotes the pullback of the volume form vol by σ). Using (54) and the fact that σ preserves the volume, it now follows that

$$(\mathcal{E}_\pi \circ \sigma)(w) = \mathcal{E}_\pi(w) \quad \text{for } w \in M \quad (55)$$

Finally, taking $w = \hat{x}_{\text{MMS}}$ and recalling that \hat{x}_{MMS} is the unique global minimiser of \mathcal{E}_π , it follows that $\sigma(\hat{x}_{\text{MMS}}) = \hat{x}_{\text{MMS}}$, so that \hat{x}_{MMS} indeed lies on γ . ■

Lemma 3.2. *There exists a continuous function $c : \mathbb{R} \rightarrow \mathbb{R}$ such that*

$$\text{grad } \mathcal{E}_\pi(\gamma(t)) = c(t)\dot{\gamma}(t) \quad \text{for } t \in \mathbb{R} \quad (56)$$

Remark : taking covariant derivatives in (56),

$$\text{Hess } \mathcal{E}_\pi(\gamma(t)) \cdot \dot{\gamma}(t) = c'(t)\dot{\gamma}(t) \quad (57)$$

where $c'(t) = dc(t)/dt$. Because \mathcal{E}_π is 1-strongly convex (see Item (iii) of Proposition B.2, Appendix B.1), it follows that $c'(t) \geq 1$. In particular, $c(t)$ is strictly increasing.

Proof : let $w = \gamma(t)$ and take the gradient of (55). This yields

$$d\sigma \cdot \text{grad } \mathcal{E}_\pi(\gamma(t)) = \text{grad } \mathcal{E}_\pi(\gamma(t)) \quad (58)$$

However, the derivative $d\sigma : T_{\gamma(t)}M \rightarrow T_{\gamma(t)}M$ is equal to 1 on vectors parallel to $\dot{\gamma}(t)$ and to -1 on vectors orthogonal to $\dot{\gamma}(t)$. Thus, (58) implies that $\text{grad } \mathcal{E}_\pi(\gamma(t))$ should be parallel to $\dot{\gamma}(t)$. This is equivalent to (56). ■

The next step in the proof will be to compute $c(0)$ and $c(1)$. This will show that $c(0)$ is negative and $c(1)$ is positive. Computing $c(0)$ and $c(1)$ requires taking a closer look at the posterior density π .

Lemma 3.3. *Let $Z(z, y, \tau, \sigma)$ denote the missing normalising factor in (41). Then, $Z(z, y, \tau, \sigma) = Z(\delta, \tau, \sigma)$, where $\delta = d^2(z, y)$.*

Proof : a hyperbolic space is a two-point homogeneous space [15] (Page 355). This means if $z', y' \in M$ have $d(z', y') = d(z, y)$, then there exists an isometry $g : M \rightarrow M$ such $g(z) = z'$ and $g(y) = y'$. Now, since g is an isometry,

$$\begin{aligned} Z(z, y, \tau, \sigma) &= \int_M \exp \left[-\frac{d^2(y, x)}{2\sigma^2} - \frac{d^2(x, z)}{2\tau^2} \right] \text{vol}(dx) \\ &= \int_M \exp \left[-\frac{d^2(g(y), g(x))}{2\sigma^2} - \frac{d^2(g(x), g(z))}{2\tau^2} \right] \text{vol}(dx) \end{aligned}$$

Thus, introducing the change of variables $w = g(x)$,

$$Z(z, y, \tau, \sigma) = \int_M \exp \left[-\frac{d^2(y', w)}{2\sigma^2} - \frac{d^2(w, z')}{2\tau^2} \right] \text{vol}(dw) = Z(z', y', \tau, \sigma)$$

In other words, $Z(z, y, \tau, \sigma)$ only depends on the distance between z and y . ■

It is now possible to compute $c(0)$ and $c(1)$.

Lemma 3.4. *There exists a positive constant $\psi(\delta, \tau, \sigma)$, such that*

$$c(0) = -\psi(\delta, \tau, \sigma)\tau^2 \quad \text{and} \quad c(1) = \psi(\delta, \tau, \sigma)\sigma^2 \quad (59)$$

In the notation of (56).

Proof: for any value of the parameters (z, y, τ, σ) ,

$$\int_M \exp \left[-\frac{d^2(y, x)}{2\sigma^2} - \frac{d^2(x, z)}{2\tau^2} - \log Z(\delta, \tau, \sigma) \right] \text{vol}(dx) = 1$$

where $\delta = d^2(z, y)$. Taking the gradient of this identity with respect to z , and using

$$\text{grad}_z d^2(x, z) = -2\text{Exp}_z^{-1}(x) \quad ; \quad \text{grad}_z d^2(z, y) = -2\text{Exp}_z^{-1}(y)$$

where grad_z denotes the gradient with respect to z , it follows that

$$\frac{1}{\tau^2} \int_M \text{Exp}_z^{-1}(x) \pi(x) \text{vol}(dx) - \psi(\delta, \tau, \sigma) \text{Exp}_z^{-1}(y) = 0 \quad (60)$$

where $\psi(\delta, \tau, \sigma) = -2 \times \partial \log Z(\delta, \tau, \sigma) / \partial \delta$. However, here one has,

$$\int_M \text{Exp}_z^{-1}(x) \pi(x) \text{vol}(dx) = -\text{grad } \mathcal{E}_\pi(z) \quad ; \quad \text{Exp}_z^{-1}(y) = \dot{\gamma}(0) \quad (61)$$

Thus, replacing (61) into (60), it follows that

$$c(0) = -\psi(\delta, \tau, \sigma) \tau^2$$

which is the first part of (59). The second part can be proved in the same way, taking the gradient with respect to y rather than z . The fact that $\psi(\delta, \tau, \sigma) > 0$ follows because $c(t)$ is strictly increasing (see the remark after Lemma 3.2), and has at most one zero (because \mathcal{E}_π has exactly one stationary point). ■

It is now possible to complete the proof of Proposition 3.2. Lemma 3.4, shows that $c(0)$ is negative and $c(1)$ is positive. Therefore, $c(t^*) = 0$ for some $t^* \in (0, 1)$. By (56) of Lemma 3.2, $\text{grad } \mathcal{E}_\pi(\gamma(t^*)) = 0$. Since \mathcal{E}_π is strongly convex, it follows immediately that $\gamma(t^*)$ is the unique global minimiser of \mathcal{E}_π . In other words, $\hat{x}_{\text{MMS}} = \gamma(t^*)$, as required. ■

A Riemannian symmetric spaces

A Riemannian symmetric space is a Riemannian manifold M , such that, for each $x \in M$, there exists an isometry $s_x : M \rightarrow M$, with $s_x(x) = x$ and $ds_x(x) = -\text{Id}_x$. This isometry s_x is called the geodesic symmetry at x [15].

Let G denote the identity component of the isometry group of M , and $K = K_o$ be the stabiliser in G of some point $o \in M$. Then, $M = G/K$ is a Riemannian homogeneous space. The mapping $\theta : G \rightarrow G$, $\theta(g) = s_o \circ g \circ s_o$ is an involutive isomorphism of G .

Let \mathfrak{g} denote the Lie algebra of G , and consider the Cartan decomposition, $\mathfrak{g} = \mathfrak{k} + \mathfrak{p}$, where \mathfrak{k} is the $+1$ eigenspace of $d\theta$ and \mathfrak{p} is the -1 eigenspace of $d\theta$. One clearly has the commutation relations,

$$[\mathfrak{k}, \mathfrak{k}] \subset \mathfrak{k} ; [\mathfrak{k}, \mathfrak{p}] \subset \mathfrak{p} ; [\mathfrak{p}, \mathfrak{p}] \subset \mathfrak{k} \quad (62)$$

In addition, it turns out that \mathfrak{k} is the Lie algebra of K , and that \mathfrak{p} may be identified with T_oM , in a natural way.

The Riemannian metric of M may always be expressed in terms of an $\text{Ad}(K)$ -invariant scalar product Q on \mathfrak{g} . If $x \in M$ is given by $x = g \cdot o$ for some $g \in G$ (where $g \cdot o = g(o)$), then

$$\langle u, v \rangle_x = Q(g^{-1} \cdot u, g^{-1} \cdot v) \quad (63)$$

where the vectors $g^{-1} \cdot u$ and $g^{-1} \cdot v$, which belong to T_oM , are identified with elements of \mathfrak{p} . Here, by an abuse of notation, $dg^{-1} \cdot u$ is denoted $g^{-1} \cdot u$.

Let $\exp : \mathfrak{g} \rightarrow G$ denote the Lie group exponential. If $v \in T_oM$, then the Riemannian exponential $\text{Exp}_o(v)$ is given by

$$\text{Exp}_o(v) = \exp(v) \cdot o \quad (64)$$

Moreover, if Π_0^t denotes parallel transport along the geodesic $c(t) = \text{Exp}_o(tv)$, then

$$\Pi_0^t(u) = \exp(tv) \cdot u \quad (65)$$

for any $u \in T_oM$ (note that the identification $T_oM \simeq \mathfrak{p}$ is always made, implicitly). Using (65), one can derive the following expression for the Riemann curvature tensor at o ,

$$R_o(v, u)w = -[[v, u], w] \quad v, u, w \in T_oM \quad (66)$$

A fundamental property of symmetric spaces is that the curvature tensor is parallel: $\nabla R = 0$. This is often used to solve the Jacobi equation [15][14], and then express the derivative of the Riemannian exponential,

$$d\text{Exp}_x(v)(u) = \exp(v) \cdot \text{sh}(R_v)(u) \quad (67)$$

where $\text{sh}(R_v) = \sum_{n=0}^{\infty} (R_v)^n / (2n+1)!$ for the self-adjoint curvature operator $R_v(u) = [v, [v, u]]$. As a result of (67), since $\exp(v)$ is an isometry, the following expression of the Riemannian volume is immediate

$$\text{Exp}_o^*(\text{vol}) = |\det(\text{sh}(R_v))| dv \quad (68)$$

where dv denotes the volume form on T_oM , associated with the restriction of the scalar product Q to \mathfrak{p} .

Expression (68) yields applicable integral formulae, when \mathfrak{g} is a reductive Lie algebra ($\mathfrak{g} = \mathfrak{z} + \mathfrak{g}_{ss}$: \mathfrak{z} the centre of \mathfrak{g} and \mathfrak{g}_{ss} semisimple). If \mathfrak{a} is a maximal Abelian subspace of \mathfrak{p} , any $v \in \mathfrak{p}$ is of the form $v = \text{Ad}(k)a$ for some $k \in K$ and $a \in \mathfrak{a}$ (see [15], Lemma 6.3, Chapter V). Moreover, using the fact that $\text{Ad}(k)$ is an isomorphism of \mathfrak{g} ,

$$\text{Ad}(k^{-1}) \circ R_v \circ \text{Ad}(k) = R_a = \sum_{\lambda \in \Delta_+} (\lambda(a))^2 \Pi_\lambda \quad (69)$$

where each $\lambda \in \Delta_+$ is a linear form $\lambda : \mathfrak{a} \rightarrow \mathbb{R}$, and Π_λ is the orthogonal projector onto the corresponding eigenspace of R_a . Here, Δ_+ is the set of positive roots of \mathfrak{g} with respect to \mathfrak{a} [15] (see Lemma 2.9, Chapter VII).

It is possible to use the diagonalisation (69), in order to evaluate the determinant (68). To obtain a regular parameterisation, let $S = K/K_{\mathfrak{a}}$, where $K_{\mathfrak{a}}$ is the centraliser of \mathfrak{a} in K . Then, let $\varphi : S \times \mathfrak{a} \rightarrow M$ be given by $\varphi(s, a) = \text{Exp}_o(\beta(s, a))$ where $\beta(s, a) = \text{Ad}(s)a$. Now, by (68) and (69),

$$\varphi^*(\text{vol}) = \prod_{\lambda \in \Delta_+} \left| \frac{\sinh \lambda(a)}{\lambda(a)} \right|^{m_\lambda} \beta^*(dv)$$

where m_λ is the multiplicity of λ (the rank of Π_λ). On the other hand, one may show

$$\beta^*(dv) = \prod_{\lambda \in \Delta_+} |\lambda(a)|^{m_\lambda} da \omega(ds) \quad (70)$$

where da is the volume form on \mathfrak{a} , and ω is the invariant volume induced onto S from K .

Finally, the Riemannian volume, in terms of the parameterisation φ , can be expressed in the following way

$$\varphi^*(\text{vol}) = \prod_{\lambda \in \Delta_+} |\sinh \lambda(a)|^{m_\lambda} da \omega(ds) \quad (71)$$

Using (71), it will be possible to write down integral formulae for Riemannian symmetric spaces, either non-compact or compact.

A.1 The non-compact case

This is the case where \mathfrak{g} admits an $\text{Ad}(G)$ -invariant, non-degenerate, symmetric bilinear form B , such that $Q(u, z) = -B(u, d\theta(z))$ is an $\text{Ad}(K)$ -invariant scalar product on \mathfrak{g} .

In this case, B is negative-definite on \mathfrak{k} and positive-definite on \mathfrak{p} . Moreover, the linear map $\text{ad}(z) = [z, \cdot]$ is skew-symmetric or symmetric (with respect to Q), according to whether $z \in \mathfrak{k}$ or $z \in \mathfrak{p}$.

If $u_1, u_2 \in \mathfrak{p}$ are orthonormal, the sectional curvature of $\text{Span}(u_1, u_2)$ is found from (66), $\kappa(u_1, u_2) = -\|[u_1, u_2]\|_o^2 \leq 0$. Therefore, M has non-positive sectional curvatures.

In fact, M is a Hadamard manifold. It is geodesically complete by (64). It is moreover simply connected, because $\text{Exp}_o : \mathfrak{p} \rightarrow M$ is a diffeomorphism [15] (Theorem 1.1, Chapter VI). Thus, (68) yields a first integral formula,

$$\int_M f(x) \text{vol}(dx) = \int_{\mathfrak{p}} f(\text{Exp}_o(v)) |\det(\text{sh}(R_v))| dv \quad (72)$$

To obtain an integral formula from (71), one should first note that $\beta : S \times \mathfrak{a} \rightarrow \mathfrak{p}$ is not regular, nor one-to-one. Recall the following:

- the hyperplanes $\lambda(a) = 0$, where $\lambda \in \Delta_+$, divide \mathfrak{a} into finitely many connected components, which are open and convex sets, known as Weyl chambers. From (70), β is regular on each Weyl chamber.
- let $K'_\mathfrak{a}$ denote the normaliser of \mathfrak{a} in K . Then, $W = K'_\mathfrak{a}/K_\mathfrak{a}$ is a finite group of automorphisms of \mathfrak{a} , called the Weyl group, which acts freely transitively on the set of Weyl chambers [15] (Theorem 2.12, Chapter VII).

Then, for each Weyl chamber C , β is regular and one-to-one, from $S \times C$ onto its image in \mathfrak{p} . Moreover, if \mathfrak{a}_r is the union of the Weyl chambers ($a \in \mathfrak{a}_r$ if and only if $\lambda(a) \neq 0$ for any $\lambda \in \Delta_+$), then β is regular and $|W|$ -to-one from $S \times \mathfrak{a}_r$ onto its image in \mathfrak{p} . To obtain the desired integral formula, it only remains to note that φ is a diffeomorphism from $S \times C$ onto its image in M . However, this image is the set M_r of regular values of φ . By Sard's lemma, its complement is negligible [33].

Proposition A.1. *Let $M = G/K$ be a Riemannian symmetric space, which belongs to the “non-compact case”, just described. Then, for any bounded continuous function $f : M \rightarrow \mathbb{R}$,*

$$\int_M f(x) \text{vol}(dx) = \int_{C_+} \int_S f(\varphi(s, a)) \prod_{\lambda \in \Delta_+} (\sinh \lambda(a))^{m_\lambda} da \omega(ds) \quad (73)$$

$$= \frac{1}{|W|} \int_{\mathfrak{a}} \int_S f(\varphi(s, a)) \prod_{\lambda \in \Delta_+} |\sinh \lambda(a)|^{m_\lambda} da \omega(ds) \quad (74)$$

Here, C_+ is the Weyl chamber $C_+ = \{a \in \mathfrak{a} : \lambda \in \Delta_+ \Rightarrow \lambda(a) > 0\}$.

A.2 The compact case

In this case, \mathfrak{g} admits an $\text{Ad}(G)$ -invariant scalar product Q . Therefore, $\text{ad}(z)$ is skew-symmetric, with respect to Q , for each $z \in \mathfrak{g}$. Using (66), it follows that M is compact, with non-negative sectional curvature.

In fact, the compact case may be obtained from the previous non-compact case by *duality*. Denote $\mathfrak{g}_\mathbb{C}$ the complexification of \mathfrak{g} , and let $\mathfrak{g}^* = \mathfrak{k} + \mathfrak{p}_*$ where $\mathfrak{p}_* = i\mathfrak{p}$. Then, \mathfrak{g}^* is a compact real form of $\mathfrak{g}_\mathbb{C}$ (that is, \mathfrak{g}^* is a compact Lie algebra, and its complexification is equal to $\mathfrak{g}_\mathbb{C}$). Denote G^* the connected Lie group with Lie algebra \mathfrak{g}^* .

If $M = G/K$ is a Riemannian symmetric space which belongs to the non-compact case, then $M^* = G^*/K$ is a Riemannian symmetric space which belongs to the compact case. Formally, to pass from the non-compact case to the compact case, all one has to do is replace a by ia . Applying this recipe to (71), one obtains

$$\varphi^*(\text{vol}) = \prod_{\lambda \in \Delta_+} |\sin \lambda(a)|^{m_\lambda} da \omega(ds) \quad (75)$$

where da is the volume form on $\mathfrak{a}_* = i\mathfrak{a}$, and ω is the invariant volume induced onto S from K . Note that the image under Exp_o of \mathfrak{a}_* is the torus $T_* = \mathfrak{a}_*/\mathfrak{a}_K$, where \mathfrak{a}_K is the lattice given by $\mathfrak{a}_K = \{a \in \mathfrak{a}_* : \text{Exp}_o(a) = o\}$. Recall the following:

- $\varphi(s, a)$ only depends on $t = \text{Exp}_o(a)$. Thus, φ may be considered as a map from $S \times T_*$ to M .
- if $a \in \mathfrak{a}_K$ then $\exp(2a) = e$ (the identity element in G^*). Thus, $\lambda(a) \in i\pi\mathbb{Z}$ for all $\lambda \in \Delta_+$ [15] (Page 383). Therefore, there exists a function $D : T \rightarrow \mathbb{R}$, such that

$$D(t) = \prod_{\lambda \in \Delta_+} |\sin \lambda(a)|^{m_\lambda} \quad \text{whenever } t = \text{Exp}_o(a)$$

Now, T_* is a totally flat submanifold of M . Therefore, $\text{Exp}^*(dt) = da$, where dt denotes the invariant volume induced onto T_* from M . With a slight abuse of notation, (75) now reads,

$$\varphi^*(\text{vol}) = D(t) dt \omega(ds) \quad (76)$$

Denote $(T_*)_r$ the set of $t \in T_*$ such that $D(t) \neq 0$. By the same arguments as in the non-compact case, φ is a regular $|W|$ -to-one map from $S \times (T_*)_r$ onto M_r , the set of regular values of φ .

Proposition A.2. *Let $M = G^*/K$ be a Riemannian symmetric space, which belongs to the “compact case”, just described. For any bounded continuous function $f : M \rightarrow \mathbb{R}$,*

$$\int_M f(x) \text{vol}(dx) = \frac{1}{|W|} \int_{T_*} \int_S f(\varphi(t, a)) D(t) dt \omega(ds) \quad (77)$$

A.3 Example of Propositions A.1 and A.2

consider $M = \text{H}(N)$ the space of $N \times N$ Hermitian positive-definite matrices. Here, $G = \text{GL}(N, \mathbb{C})$ and $K = U(N)$. Moreover, $B(u, z) = \text{Re}(\text{tr}(uz))$ and $d\theta(z) = -z^\dagger$. Thus, \mathfrak{p} is the space of $N \times N$ Hermitian matrices, and one may choose \mathfrak{a} the space of $N \times N$ real diagonal matrices. The positive roots are the linear maps $\lambda(a) = a_{ii} - a_{jj}$ where $i < j$, and each one has its multiplicity equal to 2. The Weyl group W is the group of permutation matrices in $U(N)$ (so $|W| = N!$). Finally, $S = U(N)/T_N \equiv S_N$, where T_N is the torus of diagonal unitary matrices. By (74),

$$\int_{\text{H}(N)} f(x) \text{vol}(dx) = \frac{1}{N!} \int_{\mathfrak{a}} \int_{S_N} f(s \exp(2a) s^\dagger) \prod_{i < j} \sinh^2(a_{ii} - a_{jj}) da \omega(ds) \quad (78)$$

where $da = da_{11} \dots da_{NN}$. Now, assume f is a class function: $f(k \cdot x) = f(x)$ for $k \in K$ and $x \in \text{H}(N)$. That is, $f(x)$ depends only on the eigenvalues $x_i = e^{r_i}$ of x . By (78),

$$\int_{\text{H}(N)} f(x) \text{vol}(dx) = \frac{\omega(S_N)}{2^N N!} \int_{\mathbb{R}^N} f(\exp(r)) \prod_{i < j} \sinh^2((r_i - r_j)/2) dr \quad (79)$$

or, by introducing the eigenvalues x_i as integration variables,

$$\int_{\text{H}(N)} f(x) \text{vol}(dx) = \frac{\omega(S_N)}{2^{N^2} N!} \int_{\mathbb{R}_+^N} f(x_1, \dots, x_N) |V(x)|^2 \prod_{i=1}^N x_i^{-N} dx_i \quad (80)$$

where $V(x) = \prod_{i < j} (x_j - x_i)$ is the Vandermonde determinant.

The dual of $H(N)$ is the unitary group $U(N)$. Here, $G^* = U(N) \times U(N)$ and $K \simeq U(N)$, is the diagonal group $K = \{(x, x); x \in U(N)\}$. The Riemannian metric is given by the trace scalar product $Q(u, z) = -\text{tr}(uz)$. Moreover, $T_* = T_N$ and $S = S_N$ (this is $U(N)/T_N$). The positive roots are $\lambda(ia) = a_{ii} - a_{jj}$ where $i < j$ and where a is $N \times N$, real and diagonal⁴. By writing the integral over T_N as a multiple integral, (77) reads,

$$\int_{U(N)} f(x) \text{vol}(dx) = \frac{1}{N!} \int_{[0, 2\pi]^N} \int_{S_N} f(s \exp(2ia) s^\dagger) \prod_{i < j} \sin^2(a_{ii} - a_{jj}) \omega(ds) da \quad (81)$$

where $da = da_{11} \dots da_{NN}$.

Now, assume f is a class function. That is, $f(x)$ depends only on eigenvalues $e^{i\theta_i}$ of x . Integrating out s , from (81), it follows,

$$\int_{U(N)} f(x) \text{vol}(dx) = \frac{\omega(S_N)}{2^N N!} \int_{[0, 2\pi]^N} f(\exp(i\theta)) \prod_{i < j} \sin^2((\theta_i - \theta_j)/2) d\theta \quad (82)$$

or, after an elementary manipulation,

$$\int_{U(N)} f(x) \text{vol}(dx) = \frac{\omega(S_N)}{2^{N^2} N!} \int_{[0, 2\pi]^N} f(\theta_1, \dots, \theta_N) |V(e^{i\theta})|^2 d\theta_1 \dots d\theta_N \quad (83)$$

where $V(e^{i\theta}) = \prod_{i < j} (e^{i\theta_j} - e^{i\theta_i})$ is the Vandermonde determinant.

Integrals such as (80) and (83) are familiar in random matrix theory [17][34]. The resemblance between these integrals (for example, in the rôle played by the Vandermonde determinant) is at the origin of the sort of “duality” described in Paragraph 2.8.

⁴Please do not confuse the imaginary number i with the subscript i .

B Convex optimisation

B.1 Convex sets and functions

In Euclidean geometry, a convex set A is any set which satisfies the definition: if points x and y belong to A , then the straight line segment between x and y lies entirely in A . One hopes to extend this definition to Riemannian geometry, by letting geodesics play the rôle of straight lines. However, this does not lead to one, but to multiple definitions of a convex set. The present article will focus on the following [30].

Definition B.1. *A subset A of a complete Riemannian manifold M is called strongly convex if, whenever points x and y belong to A , there exists a unique length-minimising geodesic $\gamma_{x,y}$ connecting x and y , and $\gamma_{x,y}$ lies entirely in A .*

Remark: as an example of a different way of defining a convex set, consider the following. A subset A of M is called weakly convex if, whenever points x and y belong to A , there exists a unique geodesic γ in M , such that γ connects x and y , and γ lies entirely in A (this coincides with the definition in [30], because γ is then the unique length-minimising curve, among all curves that connect x and y and lie entirely in A). ■

In Euclidean geometry, a ball of any radius is convex. In a Riemannian manifold, a ball may fail to be strongly (or even weakly) convex, if its radius is too large. On the other hand, a ball with sufficiently small radius is always strongly convex [30][14].

Proposition B.1. *Assume the sectional curvatures of M are bounded above by $\kappa_{\max}^2 \geq 0$. Then, denoting $\text{inj}(M)$ the injectivity radius of M , let*

$$R_c(M) = \min \left\{ \frac{1}{2} \text{inj}(M), \frac{\pi}{2\kappa_{\max}} \right\} \quad (84)$$

For any $x \in M$ and $R < R_c(M)$, the open ball $B(x, R)$ is strongly convex (if $\kappa_{\max} = 0$, it should be understood that division by zero yields infinity).

Remark: if $\frac{1}{2} \text{inj}(M)$ is replaced by $\text{inj}(M)$ in (84), then $R \leq R_c(M)$ implies $B(x, R)$ is weakly convex [30][35]. ■

There is a certain class of Riemannian manifolds, where balls of any radius are strongly convex. Namely, these are Hadamard manifolds. Recall that a Hadamard manifold is a simply connected, complete Riemannian manifold with non-positive sectional curvatures. In particular [36], this implies $\text{inj}(M) = \infty$ and $\kappa_{\max} = 0$, so that $R_c(M) = \infty$.

The following definition of a convex function on a Riemannian manifold directly extends the usual, well-known definition of a convex function on a Euclidean space.

Definition B.2. *Let A be a strongly convex subset of a complete Riemannian manifold M , and $f : A \rightarrow \mathbb{R}$. Then, f is called convex (respectively, strictly convex) if $f(\gamma_{x,y}(t))$ is a convex (respectively, strictly convex) function of the time parameter t , for all $x, y \in A$. Further, if there exists $\alpha > 0$ such that $f(\gamma_{x,y}(t))$ is an α -strongly convex function of t , for all $x, y \in A$, then f is called α -strongly convex.*

For differentiable functions, it is possible to write down first-order and second-order characterisations of convexity [26]. Recall that $\dot{\gamma}_{x,y}(0) = \text{Exp}_x^{-1}(y)$, for $x, y \in A$, where

the dot denotes the time derivative and Exp the Riemannian exponential map [36]. In addition, let $\text{grad}f$ and $\text{Hess}f$ denote the gradient and Hessian of a function f on M (defined with respect to the Riemannian metric and Levi-Civita connection of M).

Proposition B.2. *Let A be a strongly convex subset of a complete Riemannian manifold M , and $f : A \rightarrow \mathbb{R}$.*

(i) *assume f is differentiable. Then, f is convex, if and only if*

$$f(y) - f(x) \geq \langle \text{grad}f(x), \text{Exp}_x^{-1}(y) \rangle_x \quad \text{for all } x, y \in A \quad (85)$$

Moreover, f is strictly convex if and only if the above inequality is strict whenever $y \neq x$.

(ii) *assume f is differentiable. Then f is α -strongly convex, if and only if,*

$$f(y) - f(x) \geq \langle \text{grad}f(x), \text{Exp}_x^{-1}(y) \rangle_x + (\alpha/2)d^2(y, x) \quad \text{for all } x, y \in A \quad (86)$$

(iii) *assume f is twice differentiable. Then f is convex if and only if $\text{Hess}f(x) \succeq 0$, and strictly convex if and only if $\text{Hess}f(x) \succ 0$, for all $x \in A$. Moreover, f is α -strongly convex if and only if $\text{Hess}f(x) \succeq \alpha g(x)$, for all $x \in A$.*

Here, $\langle \cdot, \cdot \rangle$ and $d(\cdot, \cdot)$ denote the Riemannian scalar product and distance, associated with the Riemannian metric tensor g of M . Moreover, \succ stands for the Loewner order. A straightforward consequence of (ii) in Proposition B.2 is the so-called PL inequality (PL stands for Polyak-Lojasiewicz [37]). This inequality will be used in Paragraph B.4.2.

Proposition B.3. *Let $f : A \rightarrow \mathbb{R}$ be a twice differentiable, α -strongly convex function. If f has its minimum at $x^* \in A$, then*

$$\|\text{grad}f(x)\|_x^2 \geq 2\alpha(f(x) - f(x^*)) \quad \text{for all } x \in A \quad (87)$$

B.2 Second-order Taylor formula

Consider the second-order Taylor formula, for a twice-differentiable function $f : M \rightarrow \mathbb{R}$ (as usual, M is a complete Riemannian manifold). For $x \in M$ and $v \in T_x M$,

$$f(\text{Exp}_x(v)) = f(x) + \langle \text{grad}f(x), v \rangle_x + \frac{1}{2} \text{Hess}f_{\gamma(t^*)}(\dot{\gamma}, \dot{\gamma}) \quad (88)$$

where γ is the geodesic curve $\gamma(t) = \text{Exp}_x(tv)$ and $t^* \in (0, 1)$. Formula (88) is the first-order Taylor expansion, with Lagrange remainder, of the function $f(\gamma(t))$, at $t = 0$ [13]. This formula will be the starting point for the study of Riemannian gradient descent in Paragraph B.4. There, it will be applied with $v = -\mu \text{grad}f(x)$ and $\mu \in (0, 1]$.

To apply (88), it is quite helpful to control the second-order term in its right-hand side. One says that f is L -smooth on $B \subset M$, if there exists $L \geq 0$ with $|\text{Hess}f_y(u, u)| \leq L\|u\|_y^2$ for all $y \in B$ and $u \in T_y M$. Then, if $\gamma(t) = \text{Exp}_x(tv)$ belongs to B for all $t \in (0, 1)$,

$$f(\text{Exp}_x(v)) \leq f(x) + \langle \text{grad}f(x), v \rangle_x + (L/2)\|v\|_x^2 \quad (89)$$

Inequality (89) yields the following Proposition B.4. To state this proposition, consider a C^2 function $f : M \rightarrow \mathbb{R}$, and assume the sublevel set $B_c = \{x : f(x) \leq c\}$ is compact (and not empty), for some real c . Of course, B_c is contained in some closed ball $B = \bar{B}(z, R)$. Let G be the maximum of $\|\text{grad}f(x)\|_x$, taken over $x \in B$, and $B' = \bar{B}(z, R + G)$. Now, by compactness of B' , f is L_c -smooth on B' , for some $L_c \geq 0$.

Proposition B.4. *Let $f : M \rightarrow \mathbb{R}$ be a C^2 function, with B_c and L_c defined as above, and let $y = \text{Exp}_x(-\mu \text{grad}f(x))$ for some $\mu \in (0, 1]$. If $\mu \leq 1/L_c$, then*

$$f(y) \leq f(x) - (\mu/2) \|\text{grad}f(x)\|_x^2 \quad \text{for all } x \in B_c \quad (90)$$

In particular, $x \in B_c$ implies $y \in B_c$.

Remark : as a consequence of (90), if $x^* \in B_c$ is such that $f(x^*)$ is the minimum of $f(x)$, taken over $x \in B_c$, then

$$2L_c(f(x) - f(x^*)) \geq \|\text{grad}f(x)\|_x^2 \quad \text{for all } x \in B_c \quad (91)$$

which is complementary to (87). ■

B.3 Taylor with retractions

It is customary, in practical applications, to approximate the Riemannian exponential map by another so-called retraction map, which is easier to compute [38]. In this context, it is helpful to derive new versions of Formula (88) and of Proposition B.4, which apply when the exponential map Exp is replaced with a retraction Ret .

Recall that a retraction is a smooth map $\text{Ret} : TM \rightarrow M$ (denoted $\text{Ret}(x, v) = \text{Ret}_x(v)$ for $x \in M$ and $v \in T_xM$), such that

$$\text{Ret}_x(0_x) = x \quad \text{and} \quad d\text{Ret}_x(0_x) = \text{Id}_x \quad (92)$$

for all $x \in M$. Here, 0_x is the zero element in T_xM and Id_x is the identity map of T_xM . Most retractions, encountered in practical applications, are regular retractions, in the following sense [13].

Definition B.3. *A retraction $\text{Ret} : TM \rightarrow M$ is regular, if there exists a smooth bundle map $\Phi : TM \rightarrow TM$, such that*

$$\text{Ret}_x(v) = \text{Exp}_x(\Phi_x(v)) \quad \text{for all } x \in M \text{ and } v \in T_xM \quad (93)$$

Here, Φ is denoted $\Phi(x, v) = \Phi_x(v)$ (“bundle map” means $\Phi_x(v) \in T_xM$ for all $v \in T_xM$).

If $\text{Ret} : TM \rightarrow M$ is a regular retraction and $f : M \rightarrow \mathbb{R}$ is a twice-differentiable function, then (88) and (93) directly imply

$$f(\text{Ret}_x(v)) = f(x) + \langle \text{grad}f, \Phi_x(v) \rangle_x + \frac{1}{2} \text{Hess}f_{\gamma(t^*)}(\dot{\gamma}, \dot{\gamma}) \quad (94)$$

where γ is the geodesic curve $\gamma(t) = \text{Exp}_x(t\Phi_x(v))$ and $t^* \in (0, 1)$. Formula (94) is the required new version of (88).

The Riemannian exponential Exp is a regular retraction, with $\Phi_x = \text{Id}_x$ for $x \in M$. For a general regular retraction Ret , each map $\Phi_x : T_xM \rightarrow T_xM$ still agrees with Id_x up to second-order terms [13].

Proposition B.5. *Let $\text{Ret} : TM \rightarrow M$ be a regular retraction, with $\Phi : TM \rightarrow TM$ given by (93). Then, for each $x \in M$, the map $\Phi_x : T_xM \rightarrow T_xM$ verifies*

- (a) $\Phi_x(0_x) = 0_x$ and $\Phi'_x(0_x) = \text{Id}_x$ (the prime denotes the Fréchet derivative).
- (b) $\Phi''_x(0_x)(v, v) = \ddot{c}(0)$, where the curve $c(t)$ is given by $c(t) = \text{Ret}_x(tv)$.

The retraction Ret is called geodesic when $\Phi_x''(0_x)(v, v) = 0$ for $x \in M$ and $v \in T_x M$. In this case, Φ_x agrees with Id_x up to third order terms. For geodesic regular retractions, the following Proposition B.6 provides a new, general version of Proposition B.4.

Proposition B.6. *Let $f : M \rightarrow \mathbb{R}$ be a C^2 function, with $B_c = \{x : f(x) \leq c\}$ compact (and not empty), and let $\text{Ret} : TM \rightarrow M$ be a geodesic regular retraction. There exist constants $\beta_c, \delta_c, H_c \geq 0$, which depend on f and Ret , such that, for all $x \in B_c$,*

$$f(y) \leq f(x) - \mu \left(1 - (\beta_c H_c / 2) \mu - (\delta_c \|\text{grad} f(x)\|_x^2) \mu^2\right) \|\text{grad} f(x)\|_x^2 \quad (95)$$

whenever $y = \text{Ret}_x(-\mu \text{grad} f(x))$ for some $\mu \in (0, 1]$. In particular,

$$\frac{1}{2} - (\beta_c H_c / 2) \mu - (\delta_c \|\text{grad} f(x)\|_x^2) \mu^2 \geq 0 \implies f(y) \leq f(x) - (\mu/2) \|\text{grad} f(x)\|_x^2 \quad (96)$$

Therefore, $x \in B_c$ implies $y \in B_c$.

Remark : the application of Proposition B.6 is somewhat simplified when the retraction Ret , in addition to being regular and geodesic, is contractive and uniformly geodesic. Here, contractive means that

$$\|\Phi_x(v)\|_x \leq \|v\|_x \quad \text{for } x \in M \text{ and } v \in T_x M \quad (97)$$

In this case, it is always possible to put $\beta_c = 1$ and $H_c = L_c$, where L_c is the same constant as in Proposition B.4. Uniformly geodesic means there exists $\delta \geq 0$, such that

$$\|\Phi_x(v) - v\|_x \leq \delta \|v\|_x^3 \quad \text{for } x \in M \text{ and } v \in T_x M \quad (98)$$

In this case, it is always possible to put $\delta_c = \delta$ (independent of c and even of f). ■

The widely-used projection retractions for spheres, unitary groups and Grassmann manifolds, are examples of contractive, uniformly geodesic regular retractions [13] (see Sections 1.5 and 1.6). Here is another example, for positive-definite matrices.

Example : let $M = \text{P}(N)$, the space of symmetric positive-definite $N \times N$ matrices, equipped with its usual affine-invariant metric [6]. For $x \in \text{P}(N)$, the tangent space $T_x \text{P}(N)$ is identified with the space $\text{S}(N)$ of symmetric $N \times N$ matrices. Then, recall the Riemannian exponential map $\text{Exp}_x(v) = x \exp(x^{-1}v)$ (where \exp denotes the matrix exponential), and consider the retraction $\text{Ret}_x(v) = x + v + (1/2)vx^{-1}v$. The point of using this retraction is that the eigenvalues of $\text{Ret}_x(v)$ will always be greater than $1/2$. In addition, it is a contractive, uniformly geodesic regular retraction. ■

B.4 Riemannian gradient descent

Let $f : M \rightarrow \mathbb{R}$ be a C^2 function, and $\text{Ret} : TM \rightarrow M$ a retraction. Together, these yield the Riemannian gradient descent scheme, where $\mu \in (0, 1]$ is called the step-size,

$$x^{t+1} = \text{Ret}_{x^t}(-\mu \text{grad} f(x^t)) \quad t = 0, 1, \dots \quad (99)$$

Assume f has a compact (non-empty) sublevel set B_c , and choose a constant $L_c \geq 0$ as in Proposition B.4. Also, assume Ret is a contractive, uniformly geodesic regular retraction, as in the remark after Proposition B.6. Specifically, Ret verifies (97) and (98), for some constant $\delta \geq 0$. This allows for $\text{Ret} = \text{Exp}$, the Riemannian exponential, in which case $\delta = 0$. Now, the following lemma is a direct consequence of (96) in Proposition B.6.

Lemma B.1. *Under the two assumptions just described, on f and Ret , if $x^0 \in B_c$, then*

$$\frac{1}{2} - (L_c/2)\mu - (\delta\|\text{grad}f\|_{B_c}^2)\mu^2 \geq 0 \implies f(x^{t+1}) \leq f(x^t) - (\mu/2)\|\text{grad}f(x^t)\|_{x^t}^2 \quad (100)$$

for all $t \geq 0$, so that $x^t \in B_c$ for all $t \geq 0$. Here, $\|\text{grad}f\|_{B_c}$ is the maximum of $\|\text{grad}f(x)\|_x$, taken over $x \in B_c$.

This lemma immediately yields the convergence of the Riemannian gradient descent scheme (99), to the stationary points of f in B_c . Here, μ_c^* is the infimum of $\mu \in (0, 1]$ such that $\frac{1}{2} - (L_c/2)\mu - (\delta\|\text{grad}f\|_{B_c}^2)\mu^2 < 0$.

Proposition B.7. *Under the assumptions on f and Ret , made in Lemma B.1, if $\mu \leq \mu_c^*$, then $x^0 \in B_c$ implies the sequence (x^t) generated by (99) converges to the set of stationary points of f , in the sublevel set B_c .*

The proof of this proposition is straightforward. From Lemma B.1, if $\mu \leq \mu_c^*$, then $x^0 \in B_c$ implies $x^t \in B_c$ and $(\mu/2)\|\text{grad}f(x^t)\|_{x^t}^2 \leq f(x^t) - f(x^{t+1})$ for all $t \geq 0$. Adding these inequalities, for $t = 0, \dots, T$,

$$(\mu/2) \sum_{t=0}^T \|\text{grad}f(x^t)\|_{x^t}^2 \leq f(x^0) - f(x^{T+1})$$

Then, since f is bounded below on the compact set B_c , the series $\sum_{t=0}^{\infty} \|\text{grad}f(x^t)\|_{x^t}^2$ must converge. Finally, compactness of B_c ensures every subsequence of (x^t) has a further subsequence that converges to a stationary point of f in B_c .

Proposition B.7 holds without any convexity assumptions, made on the function f . Consider now the case of a strictly convex, and then of a strongly convex f .

B.4.1 Strictly convex case

Now, assume the function f is strictly convex on some strongly convex subset A of M . Moreover, assume f has a compact (non-empty) sublevel set $B_c \subset A$. Then, choose a constant $L_c \geq 0$ as in Proposition B.4, and let the retraction Ret be as in Lemma B.1.

Assume that f has a unique minimum at $x^* \in B_c$. Let R be the radius of the smallest ball $B(x^*, R)$ such that $B_c \subset B(x^*, R)$, D the maximum of $\|\text{grad}f(x)\|_x$ for $x \in B(x^*, R)$, and $R_c = R + D$. The following Lemma ensures that (99) ‘‘contracts the distance to x^* ’’.

Lemma B.2. *Assume that the sectional curvatures of M lie in the interval $[-\kappa_{\min}^2, \kappa_{\max}^2]$, and that $R_c < \text{inj}(x^*)$. If $x^0 \in B_c$ and $\mu \leq \mu_c^*$ (with μ_c^* as in Proposition B.7), then*

$$1 - \kappa(R_c)L_c\mu - (2\delta R_c\|\text{grad}f\|_{B_c})L_c\mu^2 \geq 0 \implies d(x^{t+1}, x^*) \leq d(x^t, x^*) \quad (101)$$

where $\kappa(R_c) = \max\{\kappa_{\min}R_c \coth(\kappa_{\min}R_c), |\kappa_{\max}R_c \cot(\kappa_{\max}R_c)|\}$ and $\|\text{grad}f\|_{B_c}$ is the maximum of $\|\text{grad}f(x)\|_x$ for $x \in B_c$.

Remark: $\text{inj}(x^*)$ denotes the injectivity radius of M at x^* [14]. The condition that $R_c < \text{inj}(x^*)$ guarantees the x^t stay away from the cut locus $\text{Cut}(x^*)$. This condition is introduced because the distance function $x \mapsto d(x, x^*)$ is not differentiable on $\text{Cut}(x^*)$, where its Hessian may even diverge to $-\infty$. ■

Let μ_d^* denote the infimum of μ such that $1 - \kappa(R_c)L_c\mu - (2\delta R_c\|\text{grad}f\|_{B_c})L_c\mu^2 < 0$.

Proposition B.8. *Under the same assumptions of Lemma B.2, if $\mu \leq \min\{\mu_c^*, \mu_d^*\}$, then*

$$f(x^{t+1}) - f(x^*) \leq \frac{2d^2(x^0, x^*)}{\mu(t+1)} \quad (102)$$

for all $t \geq 0$. In particular, the sequence (x^t) converges to x^* .

Remark: the quality of the convergence in (102) depends above all on the step-size μ . The smaller this is, the slower the convergence. From the definitions of μ_c^* and μ_d^* , it is clear that there are two reasons why μ would be smaller: a larger constant δ , and a larger curvature (in absolute value) κ_{\min} . In theory, one can always make $\delta = 0$ by using the retraction $\text{Ret} = \text{Exp}$, but this requires the ability to compute the Riemannian exponential Exp with sufficient accuracy. ■

Remark: the rate of convergence stated in (102) is a partial generalisation of the rate found in [39], for gradient descent in a Euclidean space. In the Euclidean setting, $\delta = 0$ and $\kappa_{\min} = \kappa_{\max} = 0$. It then follows from Proposition B.8, that (102) obtains whenever $\mu \leq 1/L_c$. Essentially, this is Corollary 2.1.2 (Page 81) in [39]. However, note the restriction $\mu \in (0, 1]$, which is necessary in a curved Riemannian manifold. ■

Here is an optimisation problem, which falls under the scope of Proposition B.8.

Example: let M be a Hadamard manifold, with sectional curvatures bounded below by $-\kappa_{\min}^2 \leq 0$. Fix a cutoff parameter $q > 0$, and define

$$V_y(x) = q^2 \left[1 + (d(x, y)/q)^2 \right]^{\frac{1}{2}} - q^2 \quad \text{for } x, y \in M \quad (103)$$

In [13], it was proved that $V_y : M \rightarrow \mathbb{R}$ is strictly convex, but not strongly convex, and that it is $(1 + q\kappa_{\min})$ -smooth on M .

Now, let π be a probability distribution on M and consider the problem of minimising

$$V_\pi(x) = \int_M V_y(x) \pi(dy) \quad (104)$$

Note that V_π is strictly convex (but not strongly convex), and $(1 + q\kappa_{\min})$ -smooth on M , because the same is true of each function V_y . In fact, V_π has compact sublevel sets whenever the distribution π has finite first-order moments [13]. In this case, $V_\pi(x)$ is guaranteed to achieve its minimum at some $x^* \in M$. This x^* is called the robust Riemannian barycentre of π (the adjective ‘‘robust’’ comes from the field of robust statistics [40]).

When applying Lemma B.2 and Proposition B.8 to the present example (with $f = V_\pi$), note that $\text{inj}(x^*) = \infty$, since M is a Hadamard manifold, and $L_c = (1 + q\kappa_{\min})$ does not depend on c .

B.4.2 Strongly convex case

Here, assume the function f is α -strongly convex on some strongly convex subset $A \subset M$. Let $B_c \subset A$ be a sublevel set of f (where $c > \inf_x f(x)$). Because f is strongly convex, B_c is compact, and it is possible to choose a constant $L_c \geq 0$, as in Proposition B.4. Then, let μ_c^* be given as in Proposition B.7. As usual, f has a unique minimum at $x^* \in B_c$.

Proposition B.9. *Under the assumptions just described, if $\mu \leq \mu_c^*$ and $x^0 \in B_c$, then*

$$f(x^t) - f(x^*) \leq (1 - \mu\alpha)^t (f(x^0) - f(x^*)) \quad (105)$$

for all $t \geq 0$. In particular, the sequence (x^t) converges to x^* .

The proof of this proposition follows by replacing Inequality (87) into Lemma B.1. Indeed, if $\mu \leq \mu_c^*$ and $x^0 \in B_c$, then (100) in Lemma B.1 immediately implies

$$f(x^{t+1}) - f(x^*) \leq f(x^t) - f(x^*) - (\mu/2) \|\text{grad} f(x^t)\|_{x^t}^2 \quad \text{for } t \geq 0$$

Thus, replacing (87) into the right-hand side,

$$f(x^{t+1}) - f(x^*) \leq (1 - \mu\alpha)(f(x^t) - f(x^*))$$

and (105) can be obtained by induction.

Remark: (105) shows that $f(x^t)$ converges to the minimum $f(x^*)$, at an exponential rate. In practice, this can still be quite slow, if $\mu\alpha$ is very small. Indeed, one should attempt to use μ as large as possible, in order to benefit from the exponential rate (105). From the definition of μ_c^* , one cannot have μ any larger than $1/L_c$, and $\mu = 1/L_c$ is only possible if $\delta = 0$, which corresponds to using $\text{Ret} = \text{Exp}$. ■

Example: let π be a probability distribution on a complete Riemannian manifold M , and define

$$E_\pi(x) = \frac{1}{2} \int_M d^2(x, y) \pi(dy) \quad \text{for } x \in M \quad (106)$$

If the support of π is contained in a ball $B(z, R)$, where $R < R_c(M)$ given by (84), then E_π is C^2 on $B(z, R)$, and has a unique global minimum $x^* \in M$, such that $x^* \in B(z, R)$ [4] (x^* is the Riemannian barycentre of π). In addition, if $R < R_c(M)/2$, then E_π is α -strongly convex on $B(z, R)$, with α equal to $2\kappa_{\max}R \cot(2\kappa_{\max}R)$ ($= 1$ if $\kappa_{\max} = 0$).

In this case, it is possible to apply Proposition B.9 to the present example (with $f = E_\pi$). If M has positive sectional curvatures, it is always possible to choose $L_c = 1$. On the other hand, if M has negative sectional curvatures $L_c = 1 + 4\kappa_{\min}R$ always works.

C Proofs for Section B

Proof of Proposition B.3: write inequality (86) under the equivalent form

$$f(y) - f(x) \geq \langle \text{grad}f(x), \text{Exp}_x^{-1}(y) \rangle_x + (\alpha/2) \|\text{Exp}_x^{-1}(y)\|_x^2$$

With $y = x^*$, this becomes

$$f(x) - f(x^*) \leq -\langle \text{grad}f(x), \text{Exp}_x^{-1}(x^*) \rangle_x - (\alpha/2) \|\text{Exp}_x^{-1}(x^*)\|_x^2$$

or, by completing the square on the right-hand side,

$$f(x) - f(x^*) \leq -\frac{1}{2\alpha} \|\text{grad}f(x) + \text{Exp}_x^{-1}(x^*)\|_x^2 + \frac{1}{2\alpha} \|\text{grad}f(x)\|_x^2$$

Then, (87) follows immediately, by noting the first term on the right-hand side is negative.

Proof of Proposition B.4: the proof employs the notation introduced before the proposition. Let $x \in B_c$ and $v = -\mu \text{grad}f(x)$. Then, note that $\|v\|_x \leq \|\text{grad}f(x)\|_x \leq G$. This implies $\gamma(t) = \text{Exp}_x(tv)$ belongs to B' for all $t \in (0, 1)$. From the definition of L_c , it now follows by (89) that

$$f(y) \leq f(x) + \langle \text{grad}f(x), v \rangle_x + (L_c/2) \|v\|_x^2$$

and, by recalling $v = -\mu \text{grad}f(x)$,

$$f(y) \leq f(x) - \mu(1 - (L_c/2)\mu) \|\text{grad}f(x)\|_x^2$$

Then, (90) follows because $\mu \leq 1/L_c$ implies the expression in parentheses is $\geq 1/2$.

Proof of Proposition B.5: this is given in [13], Section 1.5.

Proof of Proposition B.6: assume Ret is a regular geodesic retraction, and let Φ be the corresponding map in (93). Since B_c is compact, there exist $\beta_c, \delta_c \geq 0$ such that

$$\sup \{ \|\Phi'_x(u)\|_{\text{op}}; x \in B_c \text{ and } u \in T_x M, \|u\|_x \leq \|\text{grad}f(x)\|_x \} \leq \beta_c^{1/2} \quad (107)$$

$$\sup \{ \|\Phi'''_x(u)\|_{\text{op}}; x \in B_c \text{ and } u \in T_x M, \|u\|_x \leq \|\text{grad}f(x)\|_x \} \leq \delta_c \quad (108)$$

where $\|\cdot\|_{\text{op}}$ denotes the operator norm of the linear map $\Phi'_x(u) : T_x M \rightarrow T_x M$, or of the tri-linear map $\Phi'''_x(u) : T_x M \times T_x M \times T_x M \rightarrow \mathbb{R}$. In terms of these constants β_c and δ_c ,

$$\|\Phi_x(-\mu \text{grad}f(x))\|_x^2 \leq (\beta_c \mu^2) \|\text{grad}f(x)\|_x^2 \quad \text{for } x \in B_c \quad (109)$$

$$\|\Phi_x(-\mu \text{grad}f(x)) + \mu \text{grad}f(x)\|_x \leq (\delta_c \mu^3) \|\text{grad}f(x)\|_x^3 \quad \text{for } x \in B_c \quad (110)$$

Furthermore, let B_c be contained in a closed geodesic ball $B = \bar{B}(z, R)$. Denote G the maximum of $\|\text{grad}f(x)\|_x$ taken over $x \in B$, and $B' = \bar{B}(z, R + \beta_c^{1/2} G)$. By compactness of B' , there exists $H_c \geq 0$ such that f is H_c -smooth on B' .

Now, in order to prove (95), note that $\gamma(t) = \text{Exp}_x(t\Phi_x(-\mu \text{grad}f(x)))$ belongs to B' for all $t \in (0, 1)$. It follows from (94) that (similarly to (89)),

$$f(y) \leq f(x) + \langle \text{grad}f, \Phi_x(-\mu \text{grad}f(x)) \rangle_x + (H_c/2) \|\Phi_x(-\mu \text{grad}f(x))\|_x^2$$

Then, using (109) and (110),

$$f(y) \leq f(x) - \mu \|\text{grad}f(x)\|_x^2 + (\beta_c H_c/2) \mu^2 \|\text{grad}f(x)\|_x^2 + \delta_c \mu^3 \|\text{grad}f(x)\|_x^4$$

which is the same as (95). Finally, (96) is an immediate consequence of (95).

Proof of Lemma B.1: (100) can be obtained immediately, upon replacing $\beta_c = 1$, $\delta_c = \delta$ and $H_c = L_c$ into (96).

Proof of Proposition B.7: the proof has already been summarised, right after the proposition.

Proof of Lemma B.2: let $L(x) = d^2(x, x^*)/2$. If $R_c < \text{inj}(x^*)$, then $L(x)$ is $\kappa(R_c)$ -smooth on $B(x^*, R_c)$ [14]. Note that $x^t \in B_c$ for all $t \geq 0$, because $\mu \leq \mu_c^*$ as in Proposition B.7. Then, from the Taylor expansion (94) of L , and since Ret is contractive ((109) holds with $\beta_c = 1$),

$$L(x^{t+1}) \leq L(x^t) + \langle \text{grad}L(x^t), \Phi_{x^t}(-\mu \text{grad}f(x^t)) \rangle_{x^t} + (\kappa(R_c)/2)\mu^2 \|\text{grad}f(x^t)\|_{x^t}^2$$

However, applying (91) to the third term on the right-hand side, this implies

$$L(x^{t+1}) \leq L(x^t) + \langle \text{grad}L(x^t), \Phi_{x^t}(-\mu \text{grad}f(x^t)) \rangle_{x^t} + \kappa(R_c)L_c\mu^2(f(x^t) - f(x^*)) \quad (111)$$

Now, consider the second term on the right-hand side, since $\text{grad}L(x^t) = -\text{Exp}_{x^t}^{-1}(x^*)$, this second term is equal to

$$\mu \langle \text{Exp}_{x^t}^{-1}(x^*), \text{grad}f(x^t) \rangle_{x^t} - \langle \text{Exp}_{x^t}^{-1}(x^*), \Phi_{x^t}(-\mu \text{grad}f(x^t)) + \mu \text{grad}f(x^t) \rangle_{x^t} \quad (112)$$

Applying (85) and (110) (with $\delta_c = \delta$, since Ret is uniformly geodesic),

$$(112) \leq -\mu(f(x^t) - f(x^*)) + (\delta\mu^3) \|\text{Exp}_{x^t}^{-1}(x^*)\|_{x^t} \|\text{grad}f(x^t)\|_{x^t}^3$$

Using (91) once again, along with $\|\text{Exp}_{x^t}^{-1}(x^*)\|_{x^t} \leq R_c$ and $\|\text{grad}f(x^t)\|_{x^t} \leq \|\text{grad}f\|_{B_c}$,

$$(112) \leq -\mu(f(x^t) - f(x^*)) + (2\delta R_c \|\text{grad}f\|_{B_c})L_c\mu^3(f(x^t) - f(x^*)) \quad (113)$$

Finally, from (111) and (113),

$$L(x^{t+1}) \leq L(x^t) - \mu [1 - \kappa(R_c)L_c\mu - (2\delta R_c \|\text{grad}f\|_{B_c})L_c\mu^2] (f(x^t) - f(x^*))$$

Since $f(x^t) \geq f(x^*)$, whenever the expression in square brackets is positive, one has $L(x^{t+1}) \leq L(x^t)$. However, this directly yields (101).

Proof of Proposition B.8: note from Lemma B.1 that $\mu \leq \mu_c^*$ implies

$$f(x^{t+1}) - f(x^*) \leq f(x^t) - f(x^*) - (\mu/2)\|\text{grad}f(x^t)\|_{x^t}^2$$

On the other hand, note that

$$\|\text{grad}f(x^t)\|_{x^t} \geq \frac{f(x^t) - f(x^*)}{d(x^t, x^*)} \geq \frac{f(x^t) - f(x^*)}{d(x^0, x^*)}$$

where the first inequality follows by applying Cauchy-Schwarz to (85), and the second one from Lemma B.2, since $\mu \leq \mu_d^*$. Letting $\varepsilon(t) = f(x^t) - f(x^*)$, it is now clear that

$$\varepsilon(t+1) \leq \varepsilon(t) - (\mu/2) (\varepsilon(t)/d(x^0, x^*))^2$$

so that (102) can be proved by a straightforward induction.

Proof of Proposition B.9: the proof was summarised after the proposition.

References

- [1] M. Fréchet, “Les éléments aléatoires de nature quelconque dans un espace distancié,” *Annales de l’I.H.P.*, vol. 10, no. 4, pp. 215–210, 1948.
- [2] Y. Cabanes, “Multidimensional complex stationary centered gaussian regressive time series classification: Application for audio and dar clutter machine learning in hyperbolic and siegel spaces,” Ph.D. dissertation, University of Bordeaux, 2021.
- [3] M. Congedo, A. Barachant, and R. Bhatia, “Riemannian geometry for EEG-based brain-computer interfaces; a primer and a review,” *Brain-Computer Interfaces*, vol. 4, no. 3, pp. 155–174, 2017.
- [4] B. Afsari, “Riemannian L^p center of mass: existence, uniqueness and convexity,” *Proceedings of the American Mathematical Society*, vol. 139, no. 2, pp. 655–673, 2010.
- [5] S. Said and J. H. Manton, “Riemannian barycentres of Gibbs distributions: new results on concentration and convexity,” *Information Geometry*, vol. 4, no. 2, 2021.
- [6] X. Pennec, “Intrinsic statistics on Riemannian manifolds: Basic tools for geometric measurements,” *Journal of Mathematical Imaging and Vision*, vol. 25, no. 1, pp. 127–154, 2006.
- [7] G. Cheng and B. C. Vemuri, “A novel dynamic system in the space of SPD matrices with applications to appearance tracking,” *SIAM Journal on Imaging Science*, vol. 6, no. 1, pp. 592–615, 2013.
- [8] S., L. Bombrun, Y. Berthoumieu, and J. H. Manton, “Riemannian gaussian distributions on the space of symmetric positive definite matrices,” *IEEE Transactions on Information Theory*, vol. 63, no. 4, pp. 2153–2170, 2017.
- [9] S. Said, H. Hajri, L. Bombrun, and B. C. Vemuri, “Gaussian distributions on Riemannian symmetric spaces: statistical learning with structured covariance matrices,” *IEEE Transactions on Information Theory*, vol. 64, no. 2, pp. 752–772, 2018.
- [10] P. Zanini, M. Congedo, C. Jutten, S. Said, and Y. Berthoumieu, “Parameters estimate of Riemannian gaussian distribution in the manifold of covariance matrices,” in *Sensor Array and Multichannel Signal Processing*.
- [11] L. Santilli and M. Tierz, “Riemannian gaussian distributions, random matrix ensembles and diffusion kernels,” *Nuclear Physics B*, vol. 973, 2021.
- [12] S. Heuveline, S. Said, and C. Mostajeran, “Gaussian distributions on Riemannian symmetric spaces, random matrices, and planar feynman diagrams,” *arXiv:2106.08953*, 2021.
- [13] S. Said, “Statistical models and probabilistic methods on riemannian manifolds,” *arXiv:2101.10855*, 2021.
- [14] P. Petersen, *Riemannian geometry (2nd edition)*, Springer Science, 2006.

- [15] S. Helgason, *Differential geometry and symmetric spaces*. New York and London: Academic Press, 1962.
- [16] A. W. Knap, *Lie groups, beyond an introduction (2nd edition)*. Birkhauser, 2002.
- [17] M. L. Mehta, *Random matrices (3rd edition)*. Elsevier Ltd., 2004.
- [18] C. L. Siegel, “Symplectic geometry,” *American Journal of Mathematics*, vol. 65, no. 1, pp. 1–86, 1943.
- [19] A. Terras, *Harmonic analysis on symmetric spaces and applications, Vol. II*. Springer-Verlag, 1988.
- [20] K. T. Sturm, “Probability measures on metric spaces of nonpositive curvature,” *Contemporary mathematics*, vol. 338, pp. 1–34, 2003.
- [21] G. Szegő, *Orthogonal Polynomials (1st edition)*. American Mathematical Society, 1939.
- [22] A. B. J. Kuijlaars and W. Van Assche, “The asymptotic zero distribution of orthogonal polynomials with varying recurrence coefficients,” *Journal of Approximation theory*, vol. 99, pp. 167–197, 1999.
- [23] P. Deift, *Orthogonal polynomials and random matrices: a Riemann-Hilber approach*, American Mathematical Society, 1998.
- [24] M. Mariño, *Chern-Simons theory, matrix models, and topological strings*, Oxford University Press, 2005.
- [25] E. T. Whittaker and G. N. Watson, *A course of modern analysis (4th edition)*, Cambridge University Press, 1950.
- [26] C. Udriste, *Convex functions and optimization methods on Riemannian manifolds*. Springer Science, 1994.
- [27] R. O. Roberts and J. S. Rosenthal, “General state-space Markov chains and MCMC algorithms,” *Probability Surveys*, vol. 1, pp. 20–71, 2004.
- [28] S. Meyn and R. L. Tweedie, *Markov chains and stochastic stability*, Cambridge University Press, 2008.
- [29] S. F. Jarner and E. Hansen, “Geometric ergodicity of Metropolis algorithms,” *Stochastic Processes and Applications*, vol. 58, pp. 341–361, 1998.
- [30] I. Chavel, *Riemannian geometry, a modern introduction*. Cambridge University Press, 2006.
- [31] R. Bhattacharya and V. Patrangenaru, “Large sample theory of intrinsic and extrinsic sample means on manifolds I,” *The annals of statistics*, vol. 31, no. 1, pp. 1–29, 2003.

- [32] D. V. Alekseevskij, E. B. Vinberg, and A. S. Solodovnikov, *Geometry of spaces of constant curvature (EMS vol. 29)*. Springer-Verlag, 1993.
- [33] V. I. Bogachev, *Measure Theory, Volume I*. Springer-Verlag, 2007.
- [34] E. S. Meckes, *The random matrix theory of the classical compact groups*. Cambridge University Press, 2019.
- [35] W. S. Kendall, “Probability, convexity, and harmonic maps with small image I: uniqueness and fine existence,” *Proceedings of the London Mathematical Society*, vol. 61, no. 2, pp. 371–406, 1990.
- [36] J. M. Lee, *Introduction to smooth manifolds (2nd edition)*, Springer Science, 2012.
- [37] H. Karimi, J. Nutini, and M. Schmidt, “Linear convergence of gradient and proximal-gradient methods under the Polyak-Lojasiewicz condition,” in *Machine Learning and Knowledge Discovery in Databases*, 2016.
- [38] P. A. Absil, R. Mahony, and R. Sepulchre, *Optimization algorithms on matrix manifolds*. Princeton University Press, 2008.
- [39] Y. Nesterov, *Lectures on convex optimization*. Springer Switzerland, 2018.
- [40] P. J. Huber and E. M. Ronchetti, *Robust statistics (2nd edition)*. Wiley-Blackwell, 2009.