

# Information geometry for neural networks

Daniel Wagenaar

6th April 1998

**Information geometry is the result of applying non-Euclidean geometry to probability theory. The present work introduces some of the basics of information geometry with an eye on applications in neural network research. The Fisher metric and Amari's  $\alpha$ -connections are introduced and a proof of the uniqueness of the former is sketched. Dual connections and dual coordinate systems are discussed as is the associated divergence. It is shown how information geometry promises to improve upon the learning times for gradient descent learning. Due to the inclusion of an appendix about Riemannian geometry, this text should be mostly self-contained.**

*Information geometry for neural networks* by Daniel Wagenaar,  
Centre for Neural Networks, King's College London, April 1998

Project option for the MSc course in Information Processing and Neural Networks

Supervisor: Dr A. C. C. Coolen

Second examiner: Prof R. F. Streater

---

# CONTENTS

<b>Introduction</b>	<b>5</b>
<b>Notes on notation</b>	<b>6</b>
<b>1 Information geometry</b>	<b>7</b>
1.1 Probability distributions . . . . .	7
1.2 Families of distributions as manifolds . . . . .	7
1.3 Distances between distributions: a metric . . . . .	8
1.4 Affine connection on a statistical manifold . . . . .	10
<b>2 Duality in differential geometry</b>	<b>13</b>
2.1 Dual connections . . . . .	13
2.2 Dual flatness . . . . .	15
2.3 Dual coordinate systems . . . . .	15
2.3.1 Relation of the metrics . . . . .	16
2.3.2 Dual coordinate systems in dually flat manifolds . . . . .	19
2.4 Divergence . . . . .	19
2.5 Duality of $\alpha$ - and $-\alpha$ -connections . . . . .	21
<b>3 Optimization</b>	<b>22</b>
3.1 Minimizing a scalar function: Gradient descent . . . . .	22
3.2 Minimizing the divergence . . . . .	23
3.2.1 Between two constrained points . . . . .	23
3.2.2 Application: Recurrent neural networks . . . . .	24
<b>4 Uniqueness of the Fisher metric</b>	<b>27</b>
4.1 The finite case . . . . .	27
4.1.1 Markov embeddings . . . . .	27
4.1.2 Embeddings and vectors . . . . .	28
4.1.3 Invariant metrics . . . . .	28
4.1.4 Rationale . . . . .	29
4.2 Generalization to continuous sets . . . . .	30
4.3 The parametric case . . . . .	30
4.A Appendix: The more formal language of probability theory . . . . .	31
4.B A proof of the invariance properties of the Fisher information . . . . .	33
4.B.1 Invariance under transformations of the random variable . . . . .	33

---

4.B.2 Covariance under reparametrization . . . . .	34
<b>A Riemannian geometry</b>	<b>35</b>
A.1 Manifolds . . . . .	35
A.2 Vectors . . . . .	36
A.2.1 Tangent vectors . . . . .	36
A.2.2 Vector fields . . . . .	37
A.2.3 Transformation behaviour . . . . .	37
A.3 Tensor fields . . . . .	38
A.4 Metrics . . . . .	38
A.5 Affine and metric connection . . . . .	39
A.5.1 Affine connection and parallel transport . . . . .	39
A.5.2 Metric connection . . . . .	42
A.6 Curvature . . . . .	43
A.6.1 Intuitive introduction . . . . .	43
A.6.2 Formal definition . . . . .	44
A.6.3 Affine flatness . . . . .	44
A.7 Submanifolds . . . . .	45
<b>B Some special families of probability distributions</b>	<b>46</b>
B.1 Exponential family . . . . .	46
B.2 Mixture family . . . . .	47
<b>Bibliography</b>	<b>49</b>

---

## INTRODUCTION

Information geometry is the result of applying the ideas of non-Euclidean geometry to probability theory. Although interest in this subject can be traced back to the late 1960's, it reached maturity only through the work of Amari in the 1980's. His book [1] is still the canonical reference for anyone wishing to learn about it.

One of the fundamental questions information geometry helps to answer is: 'Given two probability distributions, is it possible to define a notion of "distance" between them?' An important application in neural networks is the gradient descent learning rule, which is used to minimize an error function by repeatedly taking small steps in parameter space. Traditionally, this space was tacitly assumed to have a trivial (flat) geometry. Information geometry shows that this assumption is false, and provides a theoretical recipe to find a gradient descent-style rule for any given network which can be shown to be optimal. In this way it promises a potentially dramatic improvement of learning time in comparison to the traditional rule.

In the present work, I describe some of the basics of information geometry, with the applicability to neural networks as a guide. Since my background in theoretical physics has given me more intuition about geometry than about probability theory, this text takes geometry rather than probability as its starting point. Having said that, I have avoided using the abstract language of modern differential geometry, and opted for the slightly less austere framework of Riemannian geometry. Appendix A introduces all the geometric concepts used in the main text, with the aim of making it accessible for readers with little previous experience with non-Euclidean geometry. That said, I have to admit that the appendix is rather compact, and that textbooks may offer a more gentle introduction into the subject. For example, [2] is a definite recommendation.

The rest of this text is laid out as follows: chapter 1 sets up the basic framework of information geometry, introducing a natural metric and a class of connections for families of probability distributions. Chapter 2 sets out to define some notions of duality in geometry which have a major impact on information geometry. In chapter 3 the results of the previous chapters are used to devise a more natural algorithm for parametric gradient descent that takes the geometry of the parameter space into account. Finally, in chapter 4 the properties of the metric introduced in chapter 1 are investigated in more detail, and a proof is sketched of the uniqueness of this metric.

I'd like to thank Ton Coolen for introducing me to this fascinating subject and for many useful discussions, Dr Streater for his thorough reading of the drafts, Dr Corcuera for sending me a preprint of his article on the classification of divergences and Dr Amari for bringing this article to my attention.

---

## NOTES ON NOTATION

- Throughout this text we shall employ the Einstein summation convention, that is, summation is implied over indices that occur once upstairs and once downstairs in an expression, unless explicitly stated. For example, we write

$$x^j y_i \equiv \sum_i x^j y_i.$$

In some cases we shall be less than scrupulous about re-using indices that are bound by implicit summation. We might for example write something like

$$\frac{d}{dt} e^{tx^j y_i} = x^j y_i e^{tx^j y_i},$$

where the three pairs of  $i$ 's are supposed to be completely independent.

- We shall sometimes use the notation

$$f|_P$$

to denote ‘the function  $f$ , evaluated at the point  $P$ ’. This is intended to leave equations more transparent than the notation  $f(P)$ , and has no other significance.

- We shall use boldface to denote vectors in  $\mathbb{R}^n$ , e.g.  $\mathbf{x} \equiv (x^i)_{i=1}^n$ . Vectors on manifolds (see appendix A.1) will be denoted by italic uppercase letters, e.g.  $X = X^\mu \hat{e}_\mu$ <sup>1</sup>).
- The distinction between downstairs and upstairs labels shall be important throughout this text, as usual when dealing with Riemannian geometry<sup>2</sup>).
- We shall use greek indices  $\mu, \nu, \dots$  to enumerate coordinates on manifolds.

---

<sup>1</sup>For the benefit of those readers not well acquainted with differential geometry, we shall avoid writing vectors as differential operators,  $X = X^\mu \partial_\mu$ . This has the unfortunate side effect of making a small number of statements less immediately obvious. Where appropriate the differential operator form of equations shall be given in footnotes.

<sup>2</sup>In appendix A we review the fundamentals of Riemannian geometry as needed for the main text.

## INFORMATION GEOMETRY

As mentioned in the introduction, information geometry is Riemannian geometry applied to probability theory. This chapter, which introduces some of the basic concepts of information geometry, does not presuppose any knowledge of the theory of probability and distributions. Unfortunately however, it does require some knowledge of Riemannian geometry. The reader is referred to appendix A for a brief introduction intended to provide the necessary background for this chapter.

### 1.1 Probability distributions

We shall begin by defining what we mean by a probability distribution. For our purposes, a probability distribution over some field (or set)  $X$  is a distribution  $p : X \rightarrow \mathbb{R}$ , such that

- $\int_X d\mathbf{x} p(\mathbf{x}) = 1$ ;
- For any finite subset  $S \subset X$ ,  $\int_S d\mathbf{x} p(\mathbf{x}) > 0$ .

In the following we shall consider families of such distributions. In most cases these families will be parametrized by a set of continuous parameters  $\boldsymbol{\theta} = (\theta^\mu)_{\mu=1}^N$ , that take values in some open interval  $M \subseteq \mathbb{R}^N$  and we write  $p_{\boldsymbol{\theta}}$  to denote members of the family. For any fixed  $\boldsymbol{\theta}$ ,  $p_{\boldsymbol{\theta}} : \mathbf{x} \mapsto p_{\boldsymbol{\theta}}(\mathbf{x})$  is a mapping from  $X$  to  $\mathbb{R}$ .

As an example, consider the Gaussian distributions in one dimension:

$$p_{(\mu, \sigma)}(x) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{1}{2}(x-\mu)^2/\sigma^2}.$$

In this case we may take  $\boldsymbol{\theta} = (\theta^1, \theta^2) = (\mu, \sigma)$  as a parametrization of the family. Note that this is not the only possible parametrization: for example  $(\theta^1, \theta^2) = (\frac{\mu}{\sigma^2}, \frac{1}{\sigma^2})$  is also commonly used.

### 1.2 Families of distributions as manifolds

In information geometry, one extends a family of distributions,  $F = \{p_{\boldsymbol{\theta}} | \boldsymbol{\theta} \in M\}$ , to a manifold  $\mathcal{M}$  such that the points  $p \in \mathcal{M}$  are in a one to one relation with the distributions  $p \in F$ . The parameters  $(\theta^\mu)$  of  $F$  can thus also be used as coordinates on  $\mathcal{M}$ . In doing so, one hopes to gain some insight into the structure of such a family. For example, one might hope to discover a reasonable measure of ‘nearness’ of two distributions in the family.

Having made the link between families of distributions and manifolds, one can try to identify which objects in the language of distributions naturally correspond to objects in the language of manifolds and vice versa. Arguably the most important objects in the language of manifolds are tangent vectors. The tangent space  $T_\theta$  at the point in  $\mathcal{M}$  with coordinates  $(\theta^\mu)$  is seen to be isomorphic to the vector space spanned by the random variables<sup>1)</sup>  $\frac{\partial \log p_\theta(\cdot)}{\partial \theta^\mu}$ ,  $\mu = 1 \dots N$ . This space is called  $T_\theta^{(1)}$ . A vector field  $A(\theta) \in T(\mathcal{M})$ :

$$A : \theta \mapsto A(\theta) = A^\mu(\theta) \hat{e}_\mu \quad (1.1)$$

thus is equivalent to a random variable  $A_\theta(\cdot) \in T^{(1)}(\mathcal{M})$ :

$$A_\theta(\mathbf{x}) = A^\mu(\theta) \frac{\partial \log p_\theta(\mathbf{x})}{\partial \theta^\mu}. \quad (1.2)$$

(Just as  $T(\mathcal{M})$  is the space of continuously differentiable mappings that assigns some vector  $A(\theta) \in T_\theta$  to each point  $\theta \in \mathcal{M}$ ,  $T^{(1)}(\mathcal{M})$  assigns a random variable  $A_\theta \in T_\theta^{(1)}$ .)

In view of the above equivalence we shall not find it necessary to distinguish between the vector field  $A$  and the corresponding random variable  $A(\cdot)$ .

(1.2) is called the *1-representation* of the vector field  $A$ . It is clearly possible to use some other basis of functionals of  $p_\theta$  instead of  $\frac{\partial \log p_\theta}{\partial \theta^\mu}$ . Our present choice has the advantage that the 1-representation of a vector has zero expectation value:

$$E\left[\frac{\partial \log p_\theta}{\partial \theta^\mu}\right] \equiv \int_X d\mathbf{x} p_\theta(\mathbf{x}) \frac{\partial \log p_\theta(\mathbf{x})}{\partial \theta^\mu} = \int_X d\mathbf{x} \frac{\partial p_\theta(\mathbf{x})}{\partial \theta^\mu} = \frac{\partial}{\partial \theta^\mu} \int_X d\mathbf{x} p_\theta(\mathbf{x}) = \frac{\partial}{\partial \theta^\mu} 1 = 0.$$

Using other functionals can be useful, and in fact the 1-representation turns out to be just one member of the family of  $\alpha$ -representations [1].

In order to simplify notation, we shall use

$$\ell(\theta) \equiv \log p_\theta$$

in the following. (The argument to  $\ell$  shall be omitted when obvious from the context.) Note that  $\ell(\theta)$  is a random variable, i.e. a function from  $X$  to  $\mathbb{R}$ :  $\ell(\theta) : \mathbf{x} \mapsto \log p_\theta(\mathbf{x})$ .

We shall also use the shorthand

$$\partial_\mu \equiv \frac{\partial}{\partial \theta^\mu}.$$

### 1.3 Distances between distributions: a metric

In several applications we are interested in distances between distributions. For example, given a distribution  $p \in \mathcal{M}$  and a submanifold  $S \subseteq \mathcal{M}$  we may wish to find the distribution  $p' \in S$  that is ‘nearest’ to  $p$  in some sense. To give a specific example: suppose  $\mathcal{M}$  is a large family of distributions containing both the gaussian and the binomial distributions of one variable as

<sup>1)</sup>A random variable in this context is just a function from  $X$  to  $\mathbb{R}$ .



subsets. One may wish to approximate a given binomial distribution by the ‘nearest’ exponential distribution. For this particular case, of course, an approximation formula has been found long ago. However, a general framework for constructing such approximations is seen to be useful. For all such tasks we need a notion of distance on manifolds of distributions. In other words we need a metric.

It turns out that the following is a suitable metric for manifolds of distributions:

**Definition:** The *Fisher metric* on a manifold of probability distributions is defined as

$$g_{\mu\nu}(\boldsymbol{\theta}) = E[\partial_\mu \ell(\boldsymbol{\theta}) \partial_\nu \ell(\boldsymbol{\theta})]. \quad (1.3)$$

Obviously, this also gives us an inner product: for vector fields  $A$  and  $B$  we have:

$$\langle A, B \rangle|_{\boldsymbol{\theta}} = g_{\mu\nu} A_\theta^\mu B_\theta^\nu = E[A_\theta^\mu \partial_\mu \ell(\boldsymbol{\theta}) B_\theta^\nu \partial_\nu \ell(\boldsymbol{\theta})] = E[A_\theta B_\theta].$$

This last form is known to statisticians as the Fisher information of the two random variables  $A_\theta$  and  $B_\theta$ . It is related to the maximum amount of information that can be inferred about  $A_\theta$  and  $B_\theta$  by a single measurement through the Cramér-Rao theorem, which we shall not discuss.

At first sight, the definition (1.3) may seem rather ad hoc. However, it has recently been proven by Corcuera and Giummolè [3] to be unique in having the following very appealing properties:

- $g_{\mu\nu}$  is invariant under reparametrizations of the sample space  $X$ ;
- $g_{\mu\nu}$  is covariant under reparametrizations of the manifold (the parameter space).

This uniqueness will be used later to prove the optimality of the gradient descent rule based on the Fisher metric. The full proof as given in [3] is rather involved. In chapter 4 we shall present an outline of the proof, plus a glossary of the terminology needed to understand [3].

Before going on, let us note that the metric may also be written as

$$g_{\mu\nu} = -E[\partial_\mu \partial_\nu \ell(\boldsymbol{\theta})],$$

since

$$\begin{aligned} E[\partial_\mu \partial_\nu \ell] &= \int_X \mathbf{d}\mathbf{x} p \partial_\mu \left( \frac{1}{p} \partial_\nu p \right) = \int_X \mathbf{d}\mathbf{x} \left\{ \partial_\mu \partial_\nu p - \frac{1}{p} \partial_\mu p \partial_\nu p \right\} \\ &= 0 - \int_X \mathbf{d}\mathbf{x} p \left( \frac{1}{p} \partial_\mu p \right) \left( \frac{1}{p} \partial_\nu p \right) = -E[\partial_\mu \ell \partial_\nu \ell]. \end{aligned} \quad (1.4)$$

**Example: Metric for a single neuron**

Consider a single  $N$ -input binary neuron with output defined as follows:

$$y(t) = \text{sgn}[\tanh \beta h(\mathbf{x}) + \eta(t)], \quad (1.5)$$

with

$$h(\mathbf{x}) = \sum_{i=1}^N J^i x_i + J^0.$$

In these formulas,  $J^i$  are connection weights,  $J^0$  is the external field or bias,  $x_i$  are the (real valued) inputs, and  $\eta(t)$  is a source of uniform random noise in  $[-1, 1]$ .

From (1.5) we immediately find that

$$p_J(y|\mathbf{x}) = \frac{1}{2} + \frac{1}{2}y \tanh \beta h(\mathbf{x}) \quad (1.6)$$

Since  $p_J(y, \mathbf{x}) = p_J(y|\mathbf{x})p(\mathbf{x})$ , we find

$$\ell(\mathbf{J}) \equiv \log p_J(y, \mathbf{x}) = \log p_J(y|\mathbf{x}) = \log [1 + y \tanh \beta h(\mathbf{x})].$$

Introducing  $x_0 \equiv 1$ , we may write  $h(\mathbf{x}) = \sum_{\mu=0}^N J^\mu x_\mu \equiv \mathbf{J} \cdot \mathbf{x}$ , to find

$$\partial_\mu \ell(\mathbf{J}) = \frac{y}{1 + y \tanh \beta h(\mathbf{x})} [1 - \tanh^2 \beta h(\mathbf{x})] \beta x_\mu. \quad (1.7)$$

Therefore

$$\begin{aligned} g_{\mu\nu}(\mathbf{J}) &= \sum_{\mathbf{x} \in \{-1,1\}^N} \sum_{y \in \{-1,1\}} p(\mathbf{x}) p_J(y|\mathbf{x}) \partial_\mu \ell \partial_\nu \ell \\ &= \sum_{\mathbf{x}} \sum_y \frac{1}{2} p(\mathbf{x}) \frac{y^2}{1 + y \tanh \beta h(\mathbf{x})} [1 - \tanh^2 \beta h(\mathbf{x})]^2 \beta^2 x_\mu x_\nu. \end{aligned}$$

Noting that  $y = \pm 1$  implies  $y^2 = 1$ , and using the relation

$$\frac{1}{1 + \tanh x} + \frac{1}{1 - \tanh x} = \frac{(1 - \tanh x) + (1 + \tanh x)}{1 - \tanh^2 x} = \frac{2}{1 - \tanh^2 x},$$

this can be simplified to

$$g_{\mu\nu}(\mathbf{J}) = \sum_{\mathbf{x} \in \{-1,1\}^N} p(\mathbf{x}) [1 - \tanh^2(\beta \mathbf{J} \cdot \mathbf{x})] \beta^2 x_\mu x_\nu. \quad (1.8)$$

Unfortunately, finding the contravariant form of the metric,  $g^{\mu\nu}$ , is non-trivial since the matrix inverse of (1.8) cannot easily be computed for generic  $p(\mathbf{x})$ . ■

## 1.4 Affine connection on a statistical manifold

In this section we shall introduce a family of affine connections based on the 1-representation of vectors on a statistical manifold. These connections have been named  $\alpha$ -connections by Amari in [1].

As shown in appendix A.5, an affine connection provides a means of comparing vectors at nearby points, thus providing a non-local notion of parallelism. Through the related notion of affine geodesic, it also provides a notion of straight line between two points in a manifold. We noted that an affine connection is defined by a mapping from the tangent space at  $P'$  with coordinates  $\boldsymbol{\theta} + \delta\boldsymbol{\theta}$  to the tangent space at  $P$  with coordinates  $\boldsymbol{\theta}$ .

In the 1-representation, the former space,  $T_{\mathbf{p}'}^{(1)}$ , is spanned by the basis

$$\partial_\mu \ell(\boldsymbol{\theta} + \delta\boldsymbol{\theta}) = \partial_\mu \ell(\boldsymbol{\theta}) + \delta\theta^\nu \partial_\mu \partial_\nu \ell(\boldsymbol{\theta}) + O(\delta\theta^\nu \delta\theta^\mu). \quad (1.9)$$

In trying to construct a connection, we seek a mapping from these functions to some functions in  $T_{\mathbf{p}}^{(1)}$ , which is spanned by

$$\partial_\mu \ell(\boldsymbol{\theta}). \quad (1.10)$$

It is quite clear that the functions (1.9) cannot be expressed as linear combinations of (1.10), since the expectation values of the latter vanish, while  $E[\partial_\mu \partial_\nu \ell(\boldsymbol{\theta})]$  certainly does not vanish. There are several ways to cure this problem.

One is to add  $g_{\mu\nu} \delta\theta^\nu$  to (1.9), yielding

$$\partial_\mu \ell(\boldsymbol{\theta}) + \{\partial_\mu \partial_\nu \ell(\boldsymbol{\theta}) + g_{\mu\nu}\} \delta\theta^\nu, \quad (1.11)$$

which has vanishing expectation value, but still does not yet necessarily belong to  $T_{\mathbf{p}}^{(1)}$ . This we repair by bluntly projecting it down to  $T_{\mathbf{p}}^{(1)}$ . Since the projection of any random variable  $A_\theta(\cdot)$  down to  $T_{\mathbf{p}}^{(1)}$  is given by

$$A_\theta'(\mathbf{x}) = E[A_\theta(\mathbf{x}) \partial_\nu \ell(\mathbf{x}; \boldsymbol{\theta})] g^{\mu\nu}(\boldsymbol{\theta}) \partial_\mu \ell(\mathbf{x}; \boldsymbol{\theta}),^1$$

we find that

$$\phi : \partial_\mu \ell(\boldsymbol{\theta} + \delta\boldsymbol{\theta}) \mapsto \partial_\mu \ell(\boldsymbol{\theta}) + E[\{\partial_\mu \partial_\nu \ell(\boldsymbol{\theta}) + g_{\mu\nu}(\boldsymbol{\theta})\} \delta\theta^\nu \partial_\rho \ell(\boldsymbol{\theta})] g^{\rho\lambda}(\boldsymbol{\theta}) \partial_\lambda \ell(\boldsymbol{\theta}) \quad (1.12)$$

is a suitable projection from  $T_{\mathbf{p}'}^{(1)}$  into  $T_{\mathbf{p}}^{(1)}$ . As  $g_{\mu\nu}$  is an expectation value itself, the expectation value of the second term in braces factorizes, and since the expectation value of  $\partial_\lambda \ell$  is zero, this term actually vanishes. The resulting connection therefore is

$$\Gamma_{\mu\nu}{}^\lambda = E[\partial_\mu \partial_\nu \ell g^{\rho\lambda} \partial_\rho \ell] = E[\partial_\mu \partial_\nu \ell \partial_\rho \ell] g^{\rho\lambda}. \quad (1.13)$$

However, there are other possibilities: (1.12) is not the only projection of  $T_{\mathbf{p}'}^{(1)}$  to  $T_{\mathbf{p}}^{(1)}$ . Since the expectation value of the combination

$$\partial_\mu \partial_\nu \ell + \partial_\mu \ell \partial_\nu \ell$$

vanishes (see (1.4)), we may consider

$$\partial_\mu \ell(\boldsymbol{\theta}) + \{\partial_\mu \partial_\nu \ell(\boldsymbol{\theta}) + \partial_\mu \ell \partial_\nu \ell\} \delta\theta^\nu \quad (1.14)$$

as a replacement for (1.11)<sup>2</sup>. Projecting this down to  $T_{\mathbf{p}}^{(1)}$  yields

$$\Gamma_{\mu\nu}{}^\lambda = E[\{\partial_\mu \partial_\nu \ell + \partial_\mu \ell \partial_\nu \ell\} g^{\rho\lambda} \partial_\rho \ell] = E[\partial_\mu \partial_\nu \ell \partial_\rho \ell + \partial_\mu \ell \partial_\nu \ell \partial_\rho \ell] g^{\rho\lambda}. \quad (1.15)$$

Obviously, any linear combination of (1.11) and (1.14) also has vanishing expectation value, so we have in fact found an entire family of connections, which are called the  $\alpha$ -connections:

<sup>1</sup>cf.  $X = X^\mu \hat{e}_\mu \Leftrightarrow X^\mu = \langle X, \hat{e}_\nu \rangle g^{\mu\nu}$ .

<sup>2</sup>Note the subtle difference between adding  $g_{\mu\nu}$  and adding  $\partial_\mu \ell \partial_\nu \ell$  without taking expectation value straight-away.

**Definition:  $\alpha$ -connection**

The  $\alpha$ -connection on a statistical manifold is defined as

$$\Gamma_{\mu\nu}^{(\alpha)\lambda} = E\left[\partial_\mu \partial_\nu \ell \partial_\rho \ell + \frac{1-\alpha}{2} \partial_\mu \ell \partial_\nu \ell \partial_\rho \ell\right] g^{\rho\lambda}. \quad (1.16)$$

As an aside, we show that the metric connection is the same as the 0-connection:

$$\begin{aligned} \Gamma_{\mu\nu\rho}^{(\text{metric})} &= \frac{1}{2} \left\{ \partial_\mu g_{\nu\rho} + \partial_\nu g_{\mu\rho} - \partial_\rho g_{\mu\nu} \right\} \\ &= \frac{1}{2} \int d\mathbf{x} \left\{ \partial_\mu [p \partial_\nu \ell \partial_\rho \ell] + \partial_\nu [p \partial_\mu \ell \partial_\rho \ell] - \partial_\rho [p \partial_\mu \ell \partial_\nu \ell] \right\} \\ &= \frac{1}{2} \int d\mathbf{x} p \left\{ \frac{1}{p} \partial_\mu [p \partial_\nu \ell \partial_\rho \ell] + \frac{1}{p} \partial_\nu [p \partial_\mu \ell \partial_\rho \ell] - \frac{1}{p} \partial_\rho [p \partial_\mu \ell \partial_\nu \ell] \right\} \\ &= \frac{1}{2} \int d\mathbf{x} p \left\{ \partial_\mu \ell \partial_\nu \ell \partial_\rho \ell + 2 \partial_\mu \partial_\nu \ell \partial_\rho \ell \right\} \\ &= \frac{1}{2} E[\partial_\mu \ell \partial_\nu \ell \partial_\rho \ell] + E[\partial_\mu \partial_\nu \ell \partial_\rho \ell] = \Gamma_{\mu\nu\rho}^{(0)}. \end{aligned}$$

We may therefore write

$$\Gamma_{\mu\nu\rho}^{(\alpha)} = \Gamma_{\mu\nu\rho}^{(\text{metric})} + \alpha T_{\mu\nu\rho}, \quad (1.17)$$

where

$$T_{\mu\nu\rho} \equiv -\frac{1}{2} E[\partial_\mu \ell \partial_\nu \ell \partial_\rho \ell].$$

We round this of chapter with another example.

**Example:** Continuing our previous example, we may compute the  $\alpha$ -connection for a single binary neuron: differentiating (1.7) again, we find:

$$\partial_\mu \partial_\nu \ell = -[1 - \tanh^2 \beta h(\mathbf{x})] x_\mu x_\nu$$

after some straightforward calculations. Similar calculations yield

$$\partial_\mu \partial_\nu \ell \partial_\rho \ell = -\frac{y}{1 + y \tanh \beta h(\mathbf{x})} [1 - \tanh^2 \beta h(\mathbf{x})]^2 \beta^3 x_\mu x_\nu x_\rho,$$

and

$$\partial_\mu \ell \partial_\nu \ell \partial_\rho \ell = \frac{y}{[1 + y \tanh \beta h(\mathbf{x})]^3} [1 - \tanh^2 \beta h(\mathbf{x})]^3 \beta^3 x_\mu x_\nu x_\rho.$$

Taking expectation values gives us

$$E[\partial_\mu \partial_\nu \ell \partial_\rho \ell] = 0,$$

and

$$E[\partial_\mu \ell \partial_\nu \ell \partial_\rho \ell] = \sum_{\mathbf{x}} p(\mathbf{x}) \tanh \beta h(\mathbf{x}) [1 - \tanh^2 \beta h(\mathbf{x})] \beta^3 x_\mu x_\nu x_\rho.$$

Therefore,

$$\Gamma_{\mu\nu\rho}^{(\alpha)} = \frac{1-\alpha}{2} \sum_{\mathbf{x}} p(\mathbf{x}) \tanh \beta h(\mathbf{x}) [1 - \tanh^2 \beta h(\mathbf{x})] \beta^3 x_\mu x_\nu x_\rho. \quad (1.18)$$

Since we haven't been able to compute  $g^{\mu\nu}$ , we cannot give  $\Gamma_{\mu\nu}^{(\alpha)\lambda}$  either. ■

## DUALITY IN DIFFERENTIAL GEOMETRY

In this chapter we shall investigate the notions of dual connections and dual coordinate systems. The link with chapter 1 will be made at the very end, when we discover the duality properties of the  $\alpha$ -connections.

The key result of this chapter will be the introduction of ‘divergence’, which we shall find to be a measure of the difference between two distributions. For computational purposes the divergence as the important advantage over the Riemannian distance, that it may be calculated without integration along a geodesic.

Most of the ideas presented in this chapter may also be found in [1].

### 2.1 Dual connections

Consider vector fields  $X$  and  $Y$  that are defined on a curve  $\gamma \subset \mathcal{M}$  as the parallel transports of the vectors  $X(0) \equiv X|_{\gamma(0)}$  and  $Y(0) \equiv Y|_{\gamma(0)}$  relative to some affine connection  $\Gamma_{\mu\nu\rho}$ . Parallel transport does not necessarily preserve inner product, that is

$$\langle X(t), Y(t) \rangle \neq \langle X(0), Y(0) \rangle$$

in general. A connection for which the inner product between any pair of vectors is preserved across parallel transport is called a *metric* connection.

For non-metric connections, it may be possible to find another connection, say  $\Gamma_{\mu\nu\rho}^*$ , such that

$$\langle X(t), Y^*(t) \rangle = \langle X(0), Y(0) \rangle, \quad (2.1)$$

where  $Y^*(t)$  is the parallel transport of  $Y|_{\gamma(0)}$  relative to  $\Gamma_{\mu\nu\rho}^*$ . If (2.1) holds for any two vectors  $X$  and  $Y$ , then the connections  $\Gamma_{\mu\nu\rho}$  and  $\Gamma_{\mu\nu\rho}^*$  are said to be *dual* to each other<sup>1)</sup>. Metric connections may then be called *self-dual*.

Duality of connections may be locally defined as follows:

**Definition:** Two covariant derivatives  $\nabla$  and  $\nabla^*$  (and the corresponding connections  $\Gamma_{\mu\nu\rho}$  and  $\Gamma_{\mu\nu\rho}^*$ ) are said to be *dual* to each other when for any three vector fields  $X, Y$  and  $Z$ :

$$X^\mu \partial_\mu \langle Y, Z \rangle = \langle \nabla_X Y, Z \rangle + \langle Y, \nabla_X^* Z \rangle. \quad (2.2)$$

---

<sup>1)</sup>The notion of a dual connection was introduced in a slightly different manner originally: even in the absence of a metric, one might try to find a way to transport a 1-form  $\omega$  such that  $\omega(X)$  remains constant along any curve if  $X$  is transported parallel to that curve. It was found that this is indeed possible, and the connection that transports  $\omega$  as required was named the dual to the connection used for the parallel transport of  $X$ .

**Theorem:** There exists a dual to *any* affine connection. This dual is unique, and the dual of the dual is the connection itself.

**Proof:** Substituting  $X = \hat{e}_\mu$ ,  $Y = \hat{e}_\nu$  and  $Z = \hat{e}_\rho$  in the definition, we find

$$\partial_\mu g_{\nu\rho} = \Gamma_{\mu\nu\rho} + \Gamma_{\mu\rho\nu}^* \quad (2.3)$$

from which all three claims follow directly.  $\square$

Note that the essential step in this proof is: if (2.3) holds for two connections  $\Gamma_{\mu\nu\rho}$  and  $\Gamma_{\mu\rho\nu}^*$ , then they are indeed dual in the sense of (2.2). This equivalence between (2.3) and (2.2) will be important in the proof of the next lemma as well.

**Lemma:** Any pair of connections  $\nabla$  and  $\nabla^*$  satisfying (2.1) for any vector fields  $X$  and  $Y^*$  and any curve  $\gamma$  also satisfies the differential definition of duality, and vice versa.

**Proof:**

$\Rightarrow$  Start out with (2.1) for infinitesimal  $t$ , so that we may expand:

$$\begin{aligned} X^\mu(t) &= X^\mu|_{t=0} - t [\Gamma_{\lambda\nu}^\mu \dot{\theta}^\lambda X^\nu]_{t=0} + O(t^2), \\ Y^{*\mu}(t) &= Y^{*\mu}|_{t=0} - t [\Gamma_{\lambda\nu}^{*\mu} \dot{\theta}^\lambda Y^\nu]_{t=0} + O(t^2), \\ g_{\mu\nu}(t) &= g_{\mu\nu}|_{t=0} + t [\partial_\lambda g_{\mu\nu} \dot{\theta}^\lambda]_{t=0} + O(t^2), \end{aligned}$$

where  $\dot{\phantom{x}}$  denotes differentiation with respect to  $t$ .

Inserting these Taylor expansions into (2.1), we find

$$\begin{aligned} 0 &= \frac{d}{dt} [X^\mu(t) Y^{*\nu}(t) g_{\mu\nu}(t)] \\ &= [-\Gamma_{\lambda\rho}^\mu \dot{\theta}^\lambda X^\rho Y^\nu g_{\mu\nu} - X^\nu \Gamma_{\lambda\rho}^{*\nu} \dot{\theta}^\lambda Y^\rho g_{\mu\nu} + X^\mu Y^\nu \partial_\lambda g_{\mu\nu} \dot{\theta}^\lambda]_{t=0} + O(t) \\ &= [(-\Gamma_{\lambda\mu\nu} - \Gamma_{\lambda\nu\mu}^* + \partial_\lambda g_{\mu\nu}) X^\mu Y^\nu \dot{\theta}^\lambda]_{t=0}. \end{aligned}$$

Since this is supposed to hold for any vectors  $X$  and  $Y$  and any curve  $\gamma$ , we may conclude that the first factor on the right hand side must be zero, proving (2.3), and thereby duality of  $\nabla$  and  $\nabla^*$ .

$\Leftarrow$  We shall show that the left hand side of (2.1) cannot differ from the right hand side if  $\Gamma_{\mu\nu\rho}$  and  $\Gamma_{\mu\rho\nu}^*$  are dual connections. Let the curve  $\gamma$  be parametrized by  $\theta = \theta(t)$ , and denote the tangent vector field to  $\gamma$  by  $\dot{\theta} \in T(\mathcal{M})$ . Since  $X$  and  $Y^*$  are defined as parallel transports along  $\gamma$ , we have

$$\nabla_{\dot{\theta}} X(t) = \nabla_{\dot{\theta}}^* Y^*(t) = 0.$$

Therefore

$$\begin{aligned} \frac{d}{dt} \langle X(t), Y^*(t) \rangle &= \frac{d\theta^\mu}{dt} \frac{\partial}{\partial\theta^\mu} \langle X(\theta(t)), Y^*(\theta(t)) \rangle \\ &\equiv \dot{\theta}^\mu \partial_\mu \langle X(\theta(t)), Y^*(\theta(t)) \rangle \\ &= \langle \nabla_{\dot{\theta}} X(t), Y^*(t) \rangle + \langle X(t), \nabla_{\dot{\theta}}^* Y^*(t) \rangle \\ &= 0 + 0 = 0. \end{aligned}$$

Taking the integral with respect to  $t$  yields (2.1).  $\square$

## 2.2 Dual flatness

A manifold is *flat* with respect to an affine connection when there is a set of coordinates such that  $\Gamma_{\mu\nu}{}^\rho = 0$ . The following theorem links flatness with respect to a connection and flatness with respect to its dual:

**Theorem: Dual flatness**

When a manifold is flat with respect to an affine connection, it is also flat with respect to its dual.

This theorem follows from the following lemma:

**Lemma:** When  $\Gamma_{\mu\nu}{}^\rho$  and  $\Gamma_{\mu\nu}^*{}^\rho$  are dual connections, their curvatures obey the following relationship:

$$R_{\mu\nu\rho\lambda} = -R_{\mu\nu\lambda\rho}^*$$

**Proof:** From its definition (A.11), we see that it is possible to write the curvature as

$$R_{\mu\nu\rho\lambda} = \langle \nabla_\mu \nabla_\nu \hat{e}_\rho, \hat{e}_\lambda \rangle - (\mu \leftrightarrow \nu),$$

since

$$\nabla_\mu (\Gamma_{\nu\rho}{}^\lambda \hat{e}_\lambda) = (\partial_\mu \Gamma_{\nu\rho}{}^\lambda) \hat{e}_\lambda + \Gamma_{\nu\rho}{}^\lambda \nabla_\mu \hat{e}_\lambda.$$

Using definition (2.2) twice, we find

$$\begin{aligned} R_{\mu\nu\rho\lambda} &= \left( \partial_\mu \langle \nabla_\nu \hat{e}_\rho, \hat{e}_\lambda \rangle - \langle \nabla_\nu \hat{e}_\rho, \nabla_\mu^* \hat{e}_\lambda \rangle \right) - (\mu \leftrightarrow \nu) \\ &= \left( \partial_\mu \partial_\nu \langle \hat{e}_\rho, \hat{e}_\lambda \rangle - \partial_\mu \langle \hat{e}_\rho, \nabla_\nu^* \hat{e}_\lambda \rangle - \partial_\nu \langle \hat{e}_\rho, \nabla_\mu^* \hat{e}_\lambda \rangle + \langle \hat{e}_\rho, \nabla_\nu^* \nabla_\mu^* \hat{e}_\lambda \rangle \right) - (\mu \leftrightarrow \nu) \\ &= \left( 0 - 0 - 0 + \langle \hat{e}_\rho, \nabla_\nu^* \nabla_\mu^* \hat{e}_\lambda \rangle \right) - (\mu \leftrightarrow \nu) \\ &= \langle \nabla_\nu^* \nabla_\mu^* \hat{e}_\lambda, \hat{e}_\rho \rangle - (\mu \leftrightarrow \nu) = R_{\nu\mu\lambda\rho}^* \\ &= -R_{\mu\nu\lambda\rho}^*. \end{aligned} \quad \square$$

## 2.3 Dual coordinate systems

On a dually flat manifold, there exist two special coordinate systems: the affine flat coordinates for each of the connections. These coordinate systems are related to one another by a duality relation of their own: they are *dual coordinate systems*:

**Definition: Dual coordinate systems**

Two coordinate systems  $(\theta^\mu)$  and  $(\tilde{\theta}^{\tilde{\nu}})$  are said to be dual to one another when their coordinate basis vectors satisfy:

$$\langle \hat{e}_\mu, \tilde{e}^{\tilde{\nu}} \rangle = \delta_\mu^{\tilde{\nu}},$$

where  $\hat{e}_\mu$  and  $\tilde{e}^{\tilde{\nu}}$  the coordinate basis vectors for the  $\theta$  and  $\tilde{\theta}$  systems respectively.

Note that we use lower indices to denote the components of  $\tilde{\theta}$ . At the present stage, this has no other significance than notational convenience: it will, for example, allow us to stick to the usual summation convention.

For  $(\theta^\mu)$  and  $(\tilde{\theta}_{\tilde{\nu}})$  to be dual to one another, they need not necessarily be affine coordinate systems. However there is no guarantee for general manifolds that a pair of dual coordinate systems exists.

We shall investigate some of the properties of dual coordinate systems.

### 2.3.1 Relation of the metrics

We may express  $\theta$  and  $\tilde{\theta}$  in terms of one another:

$$\theta = \theta(\tilde{\theta}), \quad \tilde{\theta} = \tilde{\theta}(\theta).$$

The coordinate basis vectors are therefore related by

$$\hat{e}_\mu = \frac{\partial \tilde{\theta}_{\tilde{\nu}}}{\partial \theta^\mu} \tilde{e}^{\tilde{\nu}}, \quad \tilde{e}^{\tilde{\mu}} = \frac{\partial \theta^\nu}{\partial \tilde{\theta}_{\tilde{\mu}}} \hat{e}_\nu.$$

Using these relations, we may express the metrics  $g_{\mu\nu} = \langle \hat{e}_\mu, \hat{e}_\nu \rangle$  and  $g^{\tilde{\mu}\tilde{\nu}} = \langle \tilde{e}^{\tilde{\mu}}, \tilde{e}^{\tilde{\nu}} \rangle$  induced by the coordinate systems in terms of the Jacobians:

$$g_{\mu\nu} \equiv \langle \hat{e}_\mu, \hat{e}_\nu \rangle = \frac{\partial \tilde{\theta}_{\tilde{\nu}}}{\partial \theta^\mu} \langle \tilde{e}^{\tilde{\nu}}, \hat{e}_\nu \rangle = \frac{\partial \tilde{\theta}_{\tilde{\nu}}}{\partial \theta^\mu} \delta_{\tilde{\nu}}^\nu, \quad (2.4a)$$

and

$$\tilde{g}^{\tilde{\mu}\tilde{\nu}} \equiv \langle \tilde{e}^{\tilde{\mu}}, \tilde{e}^{\tilde{\nu}} \rangle = \frac{\partial \theta^\mu}{\partial \tilde{\theta}_{\tilde{\mu}}} \langle \hat{e}_\mu, \tilde{e}^{\tilde{\nu}} \rangle = \frac{\partial \theta^\mu}{\partial \tilde{\theta}_{\tilde{\mu}}} \delta_\mu^{\tilde{\nu}}. \quad (2.4b)$$

Noting that the Jacobians  $[J_{\mu\tilde{\nu}}] \equiv \left[ \frac{\partial \tilde{\theta}_{\tilde{\nu}}}{\partial \theta^\mu} \right]$  and  $[J^{\tilde{\mu}\nu}] \equiv \left[ \frac{\partial \theta^\nu}{\partial \tilde{\theta}_{\tilde{\mu}}} \right]$  are each other's matrix inverse, we find that  $[\tilde{g}^{\tilde{\mu}\tilde{\nu}}]$  and  $[g_{\mu\nu}]$  are also each other's matrix inverse. Since the matrix inverse of  $[g_{\mu\nu}]$  is known to be the contravariant form of the metric,  $[g^{\mu\nu}]$ , we find that  $\tilde{g}^{\tilde{\mu}\tilde{\nu}} = \delta_{\tilde{\mu}}^\mu \delta_{\tilde{\nu}}^\nu g^{\mu\nu}$ . In fact, this means that any tensor  $T$  expressed in  $\theta$ -coordinates as  $T^{\mu_1 \dots \mu_m}_{\nu_1 \dots \nu_n}$ , may be re-expressed in  $\tilde{\theta}$ -coordinates as

$$T^{\tilde{\mu}_1 \dots \tilde{\mu}_m}_{\tilde{\nu}_1 \dots \tilde{\nu}_n} = \delta_{\tilde{\mu}_1}^{\mu_1} \dots \delta_{\tilde{\mu}_m}^{\mu_m} \delta_{\tilde{\nu}_1}^{\nu_1} \dots \delta_{\tilde{\nu}_n}^{\nu_n} T^{\mu_1 \dots \mu_m}_{\nu_1 \dots \nu_n}.$$

At this stage it is obvious that we may as well clean up our notation by dropping the distinction between labels with and without tildes<sup>2)</sup>.

The following theorem allows us to find the functional form of  $\theta = \theta(\tilde{\theta})$  and  $\tilde{\theta} = \tilde{\theta}(\theta)$ :

<sup>1)</sup>Again these relations are obvious in the differential operator formalism: they follow directly from  $\hat{e}_\mu \equiv \frac{\partial}{\partial \theta^\mu}$  and  $\tilde{e}^{\tilde{\nu}} \equiv \frac{\partial}{\partial \tilde{\theta}_{\tilde{\nu}}}$

<sup>2)</sup>This explains why we chose to put the contravariant labels in the  $\tilde{\theta}$  coordinate system downstairs. The transformation between  $\theta$  and  $\tilde{\theta}$  coordinates is equal to the identity only when covariant labels are replaced by contravariant ones and vice versa at the same time. The present convention means that we do not have to change upstairs labels into downstairs labels, thus avoiding a lot of confusion.



**Theorem:** When  $(\theta^\mu)$  and  $(\tilde{\theta}_\mu)$  are dual coordinate systems, there exist *potential functions*  $\Theta(\theta)$  and  $\tilde{\Theta}(\tilde{\theta})$  such that

$$\theta^\mu = \tilde{\partial}^\mu \tilde{\Theta}(\tilde{\theta}) \quad \text{and} \quad \tilde{\theta}_\mu = \partial_\mu \Theta(\theta). \quad (2.5)$$

It follows that

$$g_{\mu\nu} = \partial_\mu \partial_\nu \Theta(\theta) \quad \text{and} \quad \tilde{g}^{\mu\nu} = \tilde{\partial}^\mu \tilde{\partial}^\nu \tilde{\Theta}(\tilde{\theta}). \quad (2.6)$$

Furthermore,

$$\Theta(\theta) + \tilde{\Theta}(\tilde{\theta}) = \theta^\mu \tilde{\theta}_\mu. \quad (2.7)$$

Conversely, when a potential function  $\Theta(\theta)$  exists such that  $g_{\mu\nu} = \partial_\mu \partial_\nu \Theta(\theta)$ , (2.5) yields a coordinate system  $(\tilde{\theta}_\mu)$  which will be dual to  $(\theta^\mu)$ , and (2.7) may be used to derive the other potential function  $\tilde{\Theta}(\tilde{\theta})$ .

**Proof:** Symmetry of the metric  $g_{\mu\nu} = \partial_\mu \tilde{\theta}_\nu(\theta)$  shows that  $\partial_\mu \tilde{\theta}_\nu - \partial_\nu \tilde{\theta}_\mu = 0$ , from which we may conclude that, at least locally,  $\tilde{\theta}(\theta)$  is the derivative of some function, i.e. there exists a function  $\Theta(\theta)$  such that  $\tilde{\theta}_\mu = \partial_\mu \Theta$ . (2.6) follows directly from inserting (2.5) into (2.4a) and (2.4b). Finally, (2.7) is a general fact about Legendre transforms.

The other direction is easy: when  $g_{\mu\nu} = \partial_\mu \partial_\nu \Theta(\theta)$ , we see that  $(\tilde{\theta}_\mu) \equiv (\partial_\mu \Theta(\theta))$  is dual to  $(\theta)$  from the fact that

$$\tilde{e}^\nu = \left( \left[ \frac{\partial \tilde{\theta}}{\partial \theta} \right]^{-1} \right)^{\nu\mu} \hat{e}_\mu = ([g_{\mu\nu}]^{-1})^{\mu\nu} \hat{e}_\mu,$$

whence  $\langle \tilde{e}^\nu, \hat{e}_\mu \rangle = \delta_\mu^\nu$ , proving duality.  $\square$

### Example: Binary neuron

Recall that for a binary neuron we obtained

$$g_{\mu\nu}(\mathbf{J}) = \sum_{\mathbf{x} \in \{-1,1\}^N} p(\mathbf{x}) [1 - \tanh^2(\beta \mathbf{J} \cdot \mathbf{x})] \beta^2 x_\mu x_\nu.$$

It is not difficult to integrate this expression and find that

$$\Theta(\mathbf{J}) = \sum_{\mathbf{x} \in \{1,1\}^N} p(\mathbf{x}) \ln \cosh(\beta \mathbf{J} \cdot \mathbf{x}),$$

whence we compute

$$\tilde{J}_\mu = \sum_{\mathbf{x} \in \{1,1\}^N} p(\mathbf{x}) \tanh(\beta \mathbf{J} \cdot \mathbf{x}) \beta x_\mu.$$

In order to find  $\tilde{\Theta}(\tilde{\mathbf{J}})$ , we would have to solve this equation for  $\mathbf{J}$ . This is not easy in general, but we may note that in the case of a single input and no threshold the equations trivialize: we get

$$g_{11} = \sum_{x \in \{-1,1\}} p(x) [1 - \tanh^2(\beta J x)] \beta^2 x^2.$$

From the antisymmetry of  $\tanh$  and noting that  $x^2 \equiv 1$  this simplifies to

$$g_{11} = \beta^2 [1 - \tanh^2(\beta J)].$$

We then find  $\Theta(J) = \text{Incosh}(\beta J)$  and

$$\tilde{J} = \beta \tanh(\beta J).$$

Inverting the metric becomes easy too:

$$g^{11} = \frac{1}{\beta^2 [1 - \tanh^2(\beta J)]} = \frac{1}{\beta^2 - \tilde{J}^2}.$$

More interestingly, some headway can be made in the limit  $N \rightarrow \infty$ , if we assume the inputs to be uniformly distributed,  $p(x) = 2^{-N}$ . We may write:

$$\tilde{J}_\mu = \beta \langle x_\mu \tanh(\beta J^v x_v) \rangle = \beta \langle \tanh(\beta J^v x_v x_\mu) \rangle = \beta \left\langle \tanh \beta \left( J^\mu + \sum_{v \neq \mu} J^v x_v \right) \right\rangle,$$

where we performed a gauge transformation  $x_v \rightarrow x_v x_\mu$  in the last step.

In the limit  $N \rightarrow \infty$  the second term,  $\sum_{v \neq \mu} J^v x_v$ , approaches a Gaussian distribution with zero mean and variance:

$$\sigma^2 = \left\langle \left( \sum_{v \neq \mu} J^v x_v \right)^2 \right\rangle = 2^{-N} \sum_{x \in \{-1, 1\}} \sum_{v, \rho \neq \mu} J^v J^\rho x_v x_\rho = \sum_{v \neq \mu} J^v J^v \equiv |\mathbf{J}|^2 - (J^\mu)^2,$$

where no summation over  $\mu$  is implied and  $|\mathbf{J}|$  is the Euclidean (!) length of  $\mathbf{J}$ . We may therefore compute  $\tilde{J}_\mu$  by

$$\tilde{J}_\mu = \beta \int_{-\infty}^{\infty} \frac{dz}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} \tanh \beta \left( J^\mu + z \sqrt{|\mathbf{J}|^2 - (J^\mu)^2} \right). \quad (2.8)$$

Assuming that  $|\mathbf{J}|$  is of order 1, and that for each of the components  $|J^\mu|$  is of order  $1/\sqrt{N}$ , we may expand:

$$\begin{aligned} \tanh \beta \left( J^\mu + z \sqrt{|\mathbf{J}|^2 - (J^\mu)^2} \right) &= \tanh \beta \left( J^\mu + z |\mathbf{J}| + O(N^{-1}) \right) \\ &= \tanh(\beta z |\mathbf{J}|) + \beta J^\mu [1 - \tanh^2(\beta z |\mathbf{J}|)] + O(N^{-1}). \end{aligned}$$

Inserting this into (2.8), the first term vanishes and we are left with:

$$\tilde{J}_\mu = \beta^2 \int_{-\infty}^{\infty} \frac{dz}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} J^\mu [1 - \tanh^2(\beta z |\mathbf{J}|)].$$

While this integral cannot be computed analytically, numerical methods may be employed, or one could expand in one of the limits  $\beta \rightarrow 0$  or — with some more care —  $\beta \rightarrow \infty$ . We shall not pursue these possibilities here. ■

### 2.3.2 Dual coordinate systems in dually flat manifolds

The following theorem states that a pair of dual coordinate systems exists in a dually flat manifold:

**Theorem:** When a manifold  $\mathcal{M}$  is flat with respect to a dual pair of torsion-free connections  $\nabla$  and  $\nabla^*$ , there is a pair of dual coordinate systems  $(\theta^\mu)$  and  $(\tilde{\theta}_\mu)$ , such that  $(\theta^\mu)$  is  $\nabla$ -affine, and  $(\tilde{\theta}_\mu)$  is  $\nabla^*$ -affine.

**Proof:**  $\nabla$ -flatness allows us to introduce a coordinate system  $(\theta^\mu)$  in which  $\Gamma_{\mu\nu\rho} = 0$ . According to (2.3) this means that  $\Gamma_{\mu\nu\rho}^* = -\partial_\mu g_{\nu\rho}$  (still in  $\theta$  coordinates). Since we assumed that  $\nabla^*$  is torsion free, we have  $\Gamma_{\mu\nu\rho}^* = \Gamma_{\nu\mu\rho}^*$ , and therefore  $\partial_\mu g_{\nu\rho} = \partial_\nu g_{\mu\rho}$ . Combining this with the fact that  $g_{\mu\nu} = g_{\nu\mu}$ , we may conclude that (again, at least locally) a potential function  $\Theta$  exists such that  $g_{\mu\nu} = \partial_\mu \partial_\nu \Theta$ . This allows us to introduce a coordinate system  $(\tilde{\theta}_\mu)$  dual to  $(\theta^\mu)$  defined by  $\tilde{\theta}_\nu = \partial_\nu \Theta(\theta)$ .

In order to show that  $(\tilde{\theta}_\mu)$  is a  $\nabla^*$ -affine coordinate system as claimed, we note that for any  $\mu$ :

$$\partial_\mu \langle \hat{e}_\nu, \tilde{e}^\rho \rangle = \partial_\mu \delta_\nu^\rho = 0,$$

since  $(\tilde{\theta}_\mu)$  is dual to  $(\theta^\mu)$ . On the other hand, (2.2) shows that

$$\begin{aligned} \partial_\mu \langle \hat{e}_\nu, \tilde{e}^\rho \rangle &= \langle \nabla_{\hat{e}_\mu} \hat{e}_\nu, \tilde{e}^\rho \rangle + \langle \hat{e}_\nu, \nabla_{\hat{e}_\mu}^* \tilde{e}^\rho \rangle \\ &= g^{\rho\lambda} \langle \nabla_{\hat{e}_\mu} \hat{e}_\nu, \hat{e}_\lambda \rangle + g_{\nu\lambda} g_{\mu\sigma} \langle \tilde{e}^\lambda, \nabla_{\tilde{e}^\sigma}^* \tilde{e}^\rho \rangle \\ &\equiv g^{\rho\lambda} \Gamma_{\mu\nu\lambda} + g_{\nu\lambda} g_{\mu\sigma} (\Gamma^*)^{\sigma\rho\lambda}, \end{aligned}$$

where  $(\Gamma^*)^{\sigma\rho\lambda}$  is the connection corresponding to  $\nabla^{*1}$ .

Since both the left-hand side and the first term on the right are zero, we conclude that  $(\Gamma^*)^{\sigma\rho\lambda} = 0$ , proving that  $(\tilde{\theta}_\mu)$  is a  $\nabla^*$ -affine coordinate system.  $\square$

## 2.4 Divergence

On a manifold with dual connections we define the divergence between two points as

$$D(P, Q) = \Theta(\theta_P) + \tilde{\Theta}(\tilde{\theta}_Q) - \theta_P^\mu \tilde{\theta}_{Q,\mu}. \quad (2.9)$$

At first sight this definition may seem meaningless, but in fact the divergence has rather nice properties: it behaves very much like the square of a distance, and it is obviously very easy to compute: one does not have to evaluate integrals as in the calculation of the Riemannian distance on a curved manifold.

More specifically, the properties of the divergence can be stated as follows:

1.  $D(P, Q) \geq 0$ , with equality iff  $P = Q$ ;<sup>2)</sup>

<sup>1)</sup>Note that for  $(\Gamma^*)^{\mu\nu\rho}$  we need not make the distinction between the covariant form in  $\tilde{\theta}$  coordinates and the contravariant form in  $\theta$  coordinates, just as for tensors.

<sup>2)</sup>Assuming that the metric is positive definite, which the Fisher information certainly is.

2.  $\frac{\partial}{\partial \theta_P^\mu} D(P, Q) \Big|_{P=Q} = \frac{\partial}{\partial \theta_Q^\mu} D(P, Q) \Big|_{P=Q} = 0;$
3.  $\frac{\partial}{\partial \theta_P^\mu} \frac{\partial}{\partial \theta_P^\nu} D(P, Q) \Big|_{P=Q} = g_{\mu\nu}(P);$  (in fact  $\frac{\partial}{\partial \theta_P^\mu} \frac{\partial}{\partial \theta_P^\nu} D(P, Q) = g_{\mu\nu}(P)$  for any  $P, Q$ );
4. Given three points  $P, Q$  and  $R$ , then

$$D(P, R) \geq D(P, Q) + D(Q, R)$$

if the angle between the tangent vectors at  $Q$  of the  $\nabla$ -geodesic joining  $P$  and  $Q$ , and the  $\nabla^*$ -geodesic joining  $Q$  and  $R$  is greater than, equal to, or less than  $90^\circ$ . (This angle is labelled  $\varphi$  in the picture below.)

Properties 2 and 3 follow directly from differentiating the definition. Property 1 then follows from these by noting that  $D(P, Q)$  is strictly convex in  $\theta_Q - \theta_P$ , since the metric is strictly positive definite.

Property 4 – which can be viewed as a generalized Pythagoras law – can be proved as follows: Let  $\gamma_{PQ}$  be the  $\nabla$ -geodesic joining  $P$  and  $Q$ , and  $\gamma_{QR}$  the  $\nabla^*$ -geodesic joining  $Q$  and  $R$ . Being geodesics, these curves can be written in terms of the affine coordinates  $\theta$  and  $\tilde{\theta}$  as follows:

$$\gamma_{PQ} : t \mapsto \theta_P + (\theta_Q - \theta_P)t$$

and

$$\gamma_{QR} : t \mapsto \tilde{\theta}_Q + (\tilde{\theta}_R - \tilde{\theta}_Q)t.$$

By definition, we have

$$D(P, R) = \Theta(\theta_P) + \tilde{\Theta}(\tilde{\theta}_R) - \theta_P^\mu \tilde{\theta}_{R,\mu}.$$

On the other hand,

$$D(P, Q) + D(Q, R) = \Theta(\theta_P) + \tilde{\Theta}(\tilde{\theta}_Q) - \theta_P^\mu \tilde{\theta}_{Q,\mu} + \Theta(\theta_Q) + \tilde{\Theta}(\tilde{\theta}_R) - \theta_Q^\mu \tilde{\theta}_{R,\mu}.$$

Inserting (2.7) and collecting terms this can be rewritten as:

$$\begin{aligned} D(P, Q) + D(Q, R) &= D(P, R) + \theta_Q^\mu \tilde{\theta}_{Q,\mu} - \theta_P^\mu \tilde{\theta}_{Q,\mu} - \theta_Q^\mu \tilde{\theta}_{R,\mu} + \theta_P^\mu \tilde{\theta}_{R,\mu} \\ &= D(P, R) - (\theta_Q^\mu - \theta_P^\mu)(\tilde{\theta}_{R,\mu} - \tilde{\theta}_{Q,\mu}) \\ &= D(P, R) - \langle \dot{\gamma}_{PQ}, \dot{\gamma}_{QR} \rangle_Q \\ &= D(P, R) - \|\dot{\gamma}_{PQ}\| \|\dot{\gamma}_{QR}\| \cos(\pi - \varphi), \end{aligned}$$

proving the statement.

There is an intimate relation between the divergence just defined and the Kullback-Leibler distance  $D(p||q) = \int_X dx p(x) \log \frac{p(x)}{q(x)}$ : locally they are equivalent, and for some special families of distributions one may show that they are globally equal.

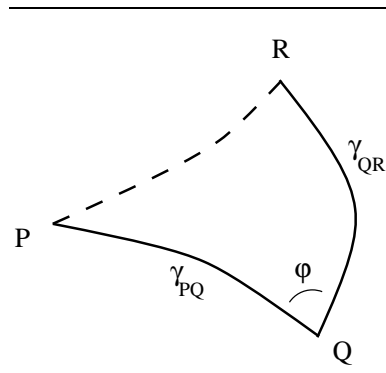


Figure 2.1: *The extension of Pythagoras' law.*

## 2.5 Duality of $\alpha$ - and $-\alpha$ -connections

We end this chapter by establishing a link between the dualities found above in differential geometry and the information geometry introduced in the previous chapter. This link consists of the following:

**Theorem:** The  $\alpha$ - and  $-\alpha$ -connections are dual to one another.

**Proof:** From (1.17) we find that

$$\Gamma_{\mu\nu\rho}^{(\alpha)} + \Gamma_{\mu\rho\nu}^{(-\alpha)} = \Gamma_{\mu\nu\rho}^{(\text{metric})} + \alpha T_{\mu\nu\rho} + \Gamma_{\mu\rho\nu}^{(\text{metric})} - \alpha T_{\mu\rho\nu}.$$

Since  $T_{\mu\nu\rho}$  is fully symmetric in its indices, the two terms that involve  $\alpha$  cancel. The metric connection being self dual by construction, the remaining terms add up to  $\partial_\mu g_{\nu\rho}$  (using (2.3)). We have thus established that

$$\Gamma_{\mu\nu\rho}^{(\alpha)} + \Gamma_{\mu\rho\nu}^{(-\alpha)} = \partial_\mu g_{\nu\rho}.$$

According to (2.3), this means that  $\Gamma_{\mu\nu\rho}^{(\alpha)}$  and  $\Gamma_{\mu\rho\nu}^{(-\alpha)}$  are dual to one another.  $\square$

**Corollary:**  $\alpha$ -flatness implies  $-\alpha$ -flatness.

## OPTIMIZATION

## 3.1 Minimizing a scalar function: Gradient descent

Minimizing a scalar function is a very common task. For example in neural networks one often wishes to find the point in weight space where the *generalization error* is minimal: In a feed forward network with  $n$  inputs  $\mathbf{x}$  and  $m$  outputs  $\mathbf{f} = \mathbf{f}(\mathbf{x}; \mathbf{J})$ , ( $\mathbf{J}$  contains all the connection weights defining the network), the generalization error is given by

$$\varepsilon_g(\mathbf{J}) = \left\langle \overline{\|\mathbf{T}(\mathbf{x}) - \mathbf{f}(\mathbf{x}; \mathbf{J})\|^2} \right\rangle,$$

where  $\overline{\dots}$  indicates averaging over the noise in  $\mathbf{f}$ ,  $\langle \dots \rangle$  indicates averaging over the inputs, and  $\mathbf{T}(\mathbf{x})$  is the function that the network is supposed to learn.

In many applications, people use the following *gradient descent learning rule* (in which  $\eta$  is called the *learning rate*):

$$J^\mu \rightarrow J^\mu - \eta \frac{\partial}{\partial J^\mu} \varepsilon_g(\mathbf{J}).$$

From our study of differential geometry we see immediately that something funny is going on here: in the first term on the right,  $\mu$  appears as an upstairs index, while in the second term it appears downstairs. Thus the subtraction doesn't yield a proper vector.

The way to fix this problem is obvious: use the inverse metric to raise the index on the second term. This leads us to the following corrected learning rule

$$J^\mu \rightarrow J^\mu - \eta g^{\mu\nu} \frac{\partial}{\partial J^\nu} \varepsilon_g(\mathbf{J}).$$

Alternatively, this result can be obtained from first principles, by taking a step back, and considering what we are really trying to achieve by gradient descent: we wish to find a  $\delta\mathbf{J}$  that maximizes

$$-\frac{\varepsilon_g(\mathbf{J} + \delta\mathbf{J}) - \varepsilon_g(\mathbf{J})}{d(\mathbf{J} + \delta\mathbf{J}, \mathbf{J})}. \quad (3.1)$$

For infinitesimal  $\delta\mathbf{J}$  appendix A.4 tells us that

$$d(\mathbf{J} + \delta\mathbf{J}, \mathbf{J}) = \sqrt{g_{\mu\nu} \delta J^\mu \delta J^\nu}.$$

Inserting this into (3.1) and extremizing the result with respect to  $\delta\mathbf{J}$  shows that  $\delta J^\mu$  should be taken to be a scalar multiple of  $g^{\mu\nu} \frac{\partial}{\partial J^\nu} \varepsilon_g(\mathbf{J})$  as claimed.

A number of studies [4, 5, 6] indicate that this modified learning rule converges much more rapidly than the original, and in particular, that many of the plateaus phases encountered in flat gradient descent are avoided or substantially reduced when one acknowledges that the weight space is not in fact flat.

### 3.2 Minimizing the divergence

Another important optimization problem is to find a distribution that approximates another, given distribution as closely as possible under some constraints. An example would be where we wish to find a weight vector for a stochastic neural network, such that it most closely approximates a given distribution, which need not be exactly representable by the network.

Mathematically, this translates into the following: given a manifold  $\mathcal{M}$  and a submanifold  $\mathcal{S}$ , and a point  $P$  in  $\mathcal{M}$ , find the point  $P'$  in  $\mathcal{S}$  that is nearest to  $P$ . Instead of minimizing the Riemann distance — which would involve integrating the metric along the metric geodesic between  $P$  and  $P'$  at every iteration step — we aim to minimize the divergence introduced in chapter 2.

When  $\mathcal{M}$  is flat under the dual connections  $\Gamma_{\mu\nu\rho}$  and  $\Gamma_{\mu\nu\rho}^*$ , and  $\mathcal{S}$  is convex<sup>1)</sup> with respect to  $\Gamma_{\mu\nu\rho}^*$ , this problem may be solved as follows:

- Introduce affine coordinates  $\theta$  and  $\tilde{\theta}$  on  $\mathcal{M}$ .
- Introduce any convenient coordinate system  $\vartheta$  on  $\mathcal{S}$ .
- The task then reduces to minimizing

$$D(P, P') \equiv D(\theta_P, \tilde{\theta}(\vartheta_{P'})) \quad (3.2)$$

with respect to  $\vartheta$ . Since (3.2) defines a scalar function on  $\mathcal{S}$ , this is easy: simply iterate

$$\vartheta^\alpha \rightarrow \vartheta^\alpha - \eta g_S^{\alpha\beta} \frac{\partial}{\partial \vartheta^\beta} D(\theta_P, \tilde{\theta}(\vartheta)),$$

where  $g_S^{\alpha\beta}$  is the metric on  $\mathcal{S}$  induced by the metric on  $\mathcal{M}$ .

When  $\mathcal{S}$  is not convex, the procedure may still be used, but the minimum reached does not necessarily correspond to the global minimum.

#### 3.2.1 Between two constrained points

Sometimes we may not wish to approximate a specific distribution in  $\mathcal{M}$  by a distribution in  $\mathcal{S}$ , but rather to find a distribution in  $\mathcal{S}$  that is closest to any distribution in  $\mathcal{M}$  that satisfies some property. Mathematically this translates to the possibility that there may be a partitioning of  $\mathcal{M}$  such that all distributions in one set are in some sense equivalent. This situation typically occurs in recurrent neural networks with hidden layers: only the states of non-hidden neurons are relevant from an external point of view.

Suppose then that we wish to find a distribution in  $\mathcal{S}$  that is closest to any distribution in another submanifold  $\mathcal{N}$  of  $\mathcal{M}$ . In other words, we wish to find the points  $P \in \mathcal{N}$  and  $P' \in \mathcal{S}$  that minimize  $D(P, P') = D(\theta(\tau_P), \tilde{\theta}(\vartheta_{P'}))$ .

Introducing coordinates  $\tau = (\tau^a)$  on  $\mathcal{N}$ , this problem may again be solved by gradient descent<sup>2)</sup>: just iterate

$$(\vartheta, \tau) \rightarrow (\vartheta', \tau'),$$

<sup>1)</sup>A submanifold  $\mathcal{S}$  is said to be convex with respect to  $\nabla$  when the  $\nabla$ -geodesic between any two points in  $\mathcal{S}$  lies entirely within  $\mathcal{S}$ .

<sup>2)</sup>The solution presented here is based on [7]

where

$$\vartheta'^{\alpha} = \vartheta^{\alpha} - \eta g_{\mathcal{S}}^{\alpha\beta} \frac{\partial}{\partial \vartheta^{\beta}} D(\boldsymbol{\theta}(\boldsymbol{\tau}), \tilde{\boldsymbol{\theta}}(\boldsymbol{\vartheta})),$$

and

$$\tau'^a = \tau^a - \eta g_{\mathcal{N}}^{ab} \frac{\partial}{\partial \tau^b} D(\boldsymbol{\theta}(\boldsymbol{\tau}), \tilde{\boldsymbol{\theta}}(\boldsymbol{\vartheta})),$$

in which  $g_{\mathcal{N}}^{ab}$  is the metric on  $\mathcal{N}$  induced by the metric on  $\mathcal{M}$ .

Again this procedure converges to a unique global minimum if  $\mathcal{S}$  is convex with respect to  $\Gamma_{\mu\nu\rho}^*$  and  $\mathcal{N}$  is convex with respect to  $\Gamma_{\mu\nu\rho}$ , but the procedure may be used even in the general case if one accepts that it may not find the global minimum.

### 3.2.2 Application: Recurrent neural networks

As a specific example, consider a stochastic recurrent neural network with symmetric interaction matrix and without self-interactions<sup>1</sup>. If the network consists of  $N$  cells, we may represent the states of these cells as an  $N$ -dimensional vector  $\boldsymbol{\sigma} \in \{-1, 1\}^N$ . The sequential evolution rule picks one of these cells at random at every timestep and updates its state according to

$$\sigma_i(t+1) = \text{sgn} \left( \tanh(\beta h_i(\boldsymbol{\sigma}(t))) + z_i(t) \right),$$

where  $z_i(t)$  is a uniform random variable taking values in  $[-1, 1]$  and

$$h_i(\boldsymbol{\sigma}) = \sum_{j=1}^N w_{ij} \sigma_j + \theta_i.$$

Here  $w_{ij}$  is the interaction matrix and  $\theta_i$  are the biases.

Iterating this step yields a Markov process, which can be written as

$$P[\boldsymbol{\sigma}(t+1) = \boldsymbol{\sigma}] = \sum_{\boldsymbol{\sigma}'} W[\boldsymbol{\sigma}; \boldsymbol{\sigma}'] P[\boldsymbol{\sigma}(t) = \boldsymbol{\sigma}'],$$

where

$$W[\boldsymbol{\sigma}; \boldsymbol{\sigma}'] = \delta_{\boldsymbol{\sigma}, \boldsymbol{\sigma}'} + \frac{1}{N} \sum_{i=1}^N (\mathcal{W}_i(F_i \boldsymbol{\sigma}) \delta_{\boldsymbol{\sigma}, F_i \boldsymbol{\sigma}'} - \mathcal{W}_i(\boldsymbol{\sigma}) \delta_{\boldsymbol{\sigma}, \boldsymbol{\sigma}'}).$$

Here  $F_i$  denotes the  $i$ -th spin-flip operator:  $(F_i \boldsymbol{\sigma})_j = \sigma_j - 2\sigma_j \delta_{ij}$ , and

$$\mathcal{W}_i(\boldsymbol{\sigma}) = \frac{1}{2} - \frac{1}{2} \sigma_i \tanh(\beta h_i(\boldsymbol{\sigma})).$$

If the weight matrix is symmetric ( $w_{ij} = w_{ji}$ ) and self-interactions are absent ( $w_{ii} = 0$ ), this process obeys detailed balance and the equilibrium distribution is given by

$$P[\boldsymbol{\sigma}] = e^{-\beta H(\boldsymbol{\sigma}) + c_0(\boldsymbol{w}, \boldsymbol{\theta})}, \quad (3.3)$$

where

$$H(\boldsymbol{\sigma}) = -\frac{1}{2} \sum_{i \neq j} \sigma_i w_{ij} \sigma_j - \sum_i \theta_i \sigma_i, \quad (3.4)$$

---

<sup>1</sup>This example is also treated in [7].



and  $c_0$  is a normalisation constant.

If we assume that the network contains  $H$  hidden cells and  $V = N - H$  visible cells, the energy function may be rewritten as

$$H(\boldsymbol{\sigma}^V; \boldsymbol{\sigma}^H) = - \sum_{i < j} \sigma_i^H w_{ij}^H \sigma_j^H - \sum_{i < j} \sigma_i^V w_{ij}^V \sigma_j^V - \sum_{i,j} \sigma_i^H w_{ij}^{HV} \sigma_j^V,$$

in which we have implicitly introduced  $\sigma_0^H = \sigma_0^V = 1$  and replaced  $\theta_i^H$  and  $\theta_i^V$  by  $w_{0i}^H$  and  $w_{0i}^V$  respectively. Note that the number of degrees of freedom in the interaction matrix is  $\frac{1}{2}H(H+1)$  in  $\{w_{ij}^H\}$ ,  $\frac{1}{2}V(V+1)$  in  $\{w_{ij}^V\}$  and  $HV$  in  $\{w_{ij}^{HV}\}$ . (Recall that  $w_{ij} = w_{ji}$  and  $w_{ij} = 0$ .)

The full space of probability distributions on  $X = \{(\boldsymbol{\sigma}^V, \boldsymbol{\sigma}^H)\}$  is a manifold  $\mathcal{M}$  of dimension  $2^{H+V} - 1$ . Only a small part of these distributions can be represented exactly by the neural network: the interaction matrix parametrizes a submanifold  $\mathcal{S}$  of  $\mathcal{M}$  of dimension  $\frac{1}{2}H(H+1) + \frac{1}{2}V(V+1) + VH$ . On the other hand, from an external point of view only the distribution on the visible cells is important. This distribution is defined by  $2^V - 1$  parameters. This partitions  $\mathcal{M}$  into submanifolds  $\mathcal{N}_{Q^V}$  of dimension  $2^V(2^H - 1)$  each, containing distributions that cannot be distinguished by looking at the visible cells only:  $Q^V$  labels the probability distribution  $Q^V : \{-1, 1\}^V \rightarrow [0, 1]$  on the visible cells to which all distributions in  $\mathcal{N}_{Q^V}$  are equivalent from an external point of view. Our task is to find a neural network that most closely approximates a given  $Q^V$ . In other words: we seek points  $P \in \mathcal{N}_{Q^V}$  and  $P' \in \mathcal{S}$  that are as ‘close’ to each other as possible. We shall show how this may be achieved by minimizing the divergence.

Minimization is quite straightforward if we notice that  $\mathcal{S}$  is an exponential family (see appendix B) and that  $\mathcal{N}_{Q^V}$  are mixture families. The first fact follows from (3.3) with (3.4), since it is clearly possible to introduce alternative random variables  $\mathbf{s} = \mathbf{s}(\boldsymbol{\sigma})$  such that (3.3) may be rewritten as:

$$P[\mathbf{s}] = e^{\sum w_{ij} s_{ij} + c_0(\mathbf{w})}.$$

Having thus established that  $\mathcal{S}$  is 1-flat<sup>1)</sup>, we can immediately conclude that  $\mathcal{S}$  is 1-convex, since the parameter space  $\{w_{ij}\}$  stretches out infinitely in every direction.

To see that  $\mathcal{N}_{Q^V}$  is a mixture family, consider the following:  $\mathcal{M}$  clearly is a mixture family, since it is possible to write any probability distribution  $p = p(\boldsymbol{\sigma}) \in \mathcal{M}$  as

$$p(\boldsymbol{\sigma}) = \sum_{\boldsymbol{\sigma}'} \theta^{\boldsymbol{\sigma}'} \delta_{\boldsymbol{\sigma}, \boldsymbol{\sigma}'},$$

(where each of the  $\theta^{\boldsymbol{\sigma}'}$ s must be in  $[0, 1]$  and  $\sum_{\boldsymbol{\sigma}} \theta^{\boldsymbol{\sigma}} = 1$  for  $p$  to be a properly normalized probability distribution. Since  $\mathcal{N}_{Q^V}$  is defined from  $\mathcal{M}$  by *linear* constraints:

$$\mathcal{N}_{Q^V} = \left\{ p \in \mathcal{M} \left| \sum_{\boldsymbol{\sigma}^H} p(\boldsymbol{\sigma}^V; \boldsymbol{\sigma}^H) = Q^V(\boldsymbol{\sigma}^V) \quad (\forall \boldsymbol{\sigma}^V) \right. \right\},$$

it follows that  $\mathcal{N}_{Q^V}$  is also a mixture family, and is therefore  $-1$ -flat. Thus  $\mathcal{N}_{Q^V}$  is  $-1$ -convex.

Our optimization problem can therefore be solved by introducing  $-1$ -affine coordinates ( $\theta^\mu$ ) and  $+1$ -affine coordinates ( $\tilde{\theta}_\mu$ ) on  $\mathcal{M}$  and minimizing the divergence

$$D(P, P') = D(\boldsymbol{\theta}(\boldsymbol{\tau}), \tilde{\boldsymbol{\theta}}(\boldsymbol{\vartheta}))$$

<sup>1)</sup>In appendix B it is shown that the exponential family is 1-flat.

with respect to  $\vartheta$  and  $\tau$  simultaneously. (As before,  $(\vartheta^\alpha)$  are coordinates on  $\mathcal{S}$  and  $(\tau^a)$  are coordinates on  $\mathcal{X}$ .)

The beauty of this technique is that it is not necessary to do the evolution and find the moments of  $\sigma^H$  at each timestep, since the divergence is given by  $D(\tau, \vartheta) = \Theta(\theta(\tau) + \tilde{\Theta}(\tilde{\theta}(\vartheta))) - \theta^\mu \tilde{\theta}_\mu$ , which is specified in terms of the weights only. This makes the optimization process much less time consuming.

If minimizing the divergens seems an arbitrary thing to do, it may help to know that in the case considered in this example, the divergence equals the Kullback-Leibler distance

$$D(P\|Q) = \sum_{\sigma} P(\sigma) \log \frac{P(\sigma)}{Q(\sigma)}. \quad (3.5)$$

We shall prove that  $D(p_{\mathbf{w}}, p_{\mathbf{w}'}) = D(p_{\mathbf{w}'}\|p_{\mathbf{w}})$ . To reduce the notational clutter, we shall relabel the weights by a single vector  $\mathbf{w} = (w^\mu) \equiv (w_{ij}^H, w_{ij}^V, w_{ij}^{HV})$ .

From the definition of the divergence, we have:

$$D(p_{\mathbf{w}}, p_{\mathbf{w}}) = \left. \frac{\partial}{\partial w^\mu} D(p_{\mathbf{w}}, p_{\mathbf{w}'}) \right|_{\mathbf{w}'=\mathbf{w}} = 0, \quad \text{and} \quad \frac{\partial}{\partial w^\mu} \frac{\partial}{\partial w^\nu} D(p_{\mathbf{w}}, p_{\mathbf{w}'}) = g_{\mu\nu}(\mathbf{w}).$$

In this equation,  $(w^\mu)$  are supposed to be  $-1$ -affine coordinates.

It is not difficult to see that

$$D(p_{\mathbf{w}}\|p_{\mathbf{w}}) = \left. \frac{\partial}{\partial w^\mu} D(p_{\mathbf{w}'}\|p_{\mathbf{w}}) \right|_{\mathbf{w}'=\mathbf{w}} = 0.$$

All that remains to be shown then is that  $\frac{\partial}{\partial w^\mu} \frac{\partial}{\partial w^\nu} D(p_{\mathbf{w}'}\|p_{\mathbf{w}}) = g_{\mu\nu}(\mathbf{w})$  for all  $\mathbf{w}$  and  $\mathbf{w}'$ .

Differentiating (3.5) twice, we find:

$$\begin{aligned} \frac{\partial}{\partial w^\mu} \frac{\partial}{\partial w^\nu} D(p_{\mathbf{w}'}\|p_{\mathbf{w}}) &= \sum_{\sigma} p_{\mathbf{w}'}(\sigma) \frac{\partial}{\partial w^\mu} \frac{\partial}{\partial w^\nu} p_{\mathbf{w}}(\sigma) \\ &= \sum_{\sigma} p_{\mathbf{w}}(\sigma) \frac{\partial}{\partial w^\mu} \frac{\partial}{\partial w^\nu} p_{\mathbf{w}}(\sigma) - \sum_{\sigma} [p_{\mathbf{w}'}(\sigma) - p_{\mathbf{w}}(\sigma)] \frac{\partial}{\partial w^\mu} \frac{\partial}{\partial w^\nu} p_{\mathbf{w}}(\sigma). \end{aligned} \quad (3.6)$$

The first term on the right will be recognized as  $g_{\mu\nu}(\mathbf{w})$ . The second term vanishes, as is shown by the following argument: we may expand

$$p_{\mathbf{w}'}(\sigma) - p_{\mathbf{w}}(\sigma) = (w'^\mu - w^\mu) \partial_\mu p_{\mathbf{w}} + \frac{1}{2} (w'^\mu - w^\mu) (w'^\nu - w^\nu) \partial_\mu \partial_\nu p_{\mathbf{w}} + \dots$$

Since  $(w^\mu)$  as a  $-1$ -affine coordinate system on a mixture family, the probability distributions can be written in the form

$$p_{\mathbf{w}} = \sum_{\mu} w^\mu f_{\mu}(\sigma).$$

Therefore, all higher order terms vanish identically. However, since  $\mathcal{M}$  is a mixture family, it is  $-1$ -flat and thus also  $+1$ -flat. Hence the  $1$ -connection vanishes:

$$\Gamma_{\mu\nu\rho}^{(1)} \equiv E\left[\frac{1}{p} \partial_\rho p \partial_\mu \partial_\nu p\right] = 0.$$

Therefore the entire second term of (3.6) vanishes, and the equality of the divergence and the Kullback-Leibler distance is proven.

## UNIQUENESS OF THE FISHER METRIC

The uniqueness of the Fisher information as a metric on statistical manifolds has recently been proven by Corcuera and Giummolè [3]. In this chapter we present an outline of their proof in a form that is intended to be accessible for those who have less intimate knowledge of probability theory than Corcuera and Giummolè presuppose.

The proof consists of three distinct steps:

1. Classification of metrics on manifolds of probability distributions over a finite number of ‘atoms’;
2. Generalization to infinite (continuous) sets;
3. Application to parametrized distributions.

The first of these was first performed by Čencov in [8]. The other two steps have been taken by Corcuera and Giummolè in [3].

### 4.1 The finite case

We are ultimately looking for metrics on manifolds of parametrized probability distributions that are invariant under transformation of the random variables, and covariant under reparametrization. Starting out with finite sets only helps if we know what invariances we should require in the finite case to end up with the invariances we are looking for in the parametrized case. It turns out that invariance under *Markov embeddings* is the finite case equivalent of reparametrization covariance. We shall therefore start by looking at those.

#### 4.1.1 Markov embeddings

Consider a set  $X$  (for example an interval in  $\mathbb{R}^N$ ) and a partition  $\mathbf{A} = \{A_1, \dots, A_m\}$  of  $X$ . (A collection of subsets is a partition if (1)  $\forall i, j : i \neq j \Rightarrow A_i \cap A_j = \emptyset$ , and (2)  $\bigcup_i A_i = X$ .) Furthermore, let  $\mathbf{B} = \{B_1, \dots, B_n\}$  be a subpartition of  $\mathbf{A}$ , that is, a partition of  $X$  such that each  $B_i$  is contained entirely within one of the  $A_j$ ’s. In other words, there exists a partition  $\mathbf{I} = \{I_1, \dots, I_m\}$  of  $\{1, \dots, n\}$  such that

$$A_i = \bigcup_{j \in I_i} B_j.$$

One picture says more than a hundred mathematical symbols: figure 4.1 should make the setup clear.

Let  $\mathcal{A} \cong \mathbb{R}_+^m$  be the collection of non-negative distributions on  $A$ <sup>1)</sup> and similarly let  $\mathcal{B} \cong \mathbb{R}_+^n$  be the collection of distribution on  $B$ . A *Markov embedding* is a subpartition  $\mathbf{B}$  of  $A$  together with a mapping  $f$  from  $\mathcal{A}$  to  $\mathcal{B}$ , with  $f(\boldsymbol{\xi})$  given by

$$f_j(\boldsymbol{\xi}) = \sum_{i=1}^m q_{ij} \xi_i,$$

where  $q_{ij} = 0$  unless  $j \in I_i$ . (Thus only one term of the sum is non-zero.) Note that consistency requires that none of the  $q_{ij}$ 's be negative, and that  $\sum_j q_{ij} = 1$  for all  $i$ .

Associated with  $f$  is a mapping  $\bar{f}$  from  $\mathcal{B}$  to  $\mathcal{A}$ :

$$\bar{f}_i(\boldsymbol{\eta}) = \sum_{j \in I_i} \eta_j.$$

Note that  $\bar{f} \circ f = \mathbb{1} : \mathcal{A} \rightarrow \mathcal{A}$ , but  $f \circ \bar{f}$  is not necessarily equal to the identity mapping on  $\mathcal{B}$ .

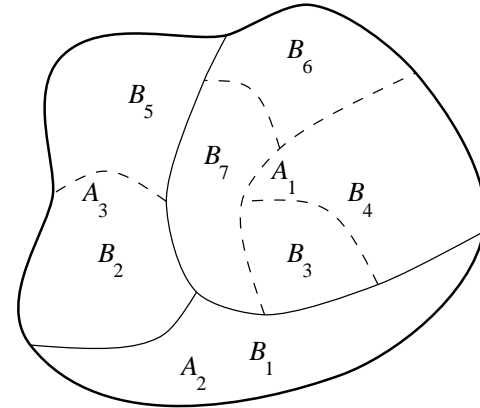


Figure 4.1: A set  $X$  with a partition  $A$  and a subpartition  $B$ , with dashed lines marking the boundaries of the elements of the subpartition.

#### 4.1.2 Embeddings and vectors

We may consider  $\mathcal{A}$  as an  $m$ -dimensional manifold on which we may introduce vectors. In particular we may introduce the coordinate basis  $\{\hat{i} \mid i = 1, \dots, m\}$ , where  $\hat{i}$  is the  $i$ -th coordinate basis vector in  $\mathbb{R}_+^m$ . A mapping  $f : \mathcal{A} \rightarrow \mathcal{B}$  induces a mapping  $f_* : T(\mathcal{A}) \rightarrow T(\mathcal{B})$  defined by

$$f_* : \hat{i} \mapsto f_*(\hat{i}) = \sum_{j=1}^n q_{ij} \tilde{j},$$

where  $\tilde{j}$  is the  $j$ -th coordinate basis vector in  $\mathbb{R}_+^n$ . (These vector relations are the ones that are slightly more involved on the set of probability distributions  $\{\xi \in \mathbb{R}_+^m \mid \sum_i \xi_i = 1\}$ , since the linear constraint reduces the dimension of the tangent space too.)

#### 4.1.3 Invariant metrics

A sequence of inner products  $\langle \cdot, \cdot \rangle_{(m)}$  on  $T(\mathbb{R}_+^m)$  for  $m = 2, 3, \dots$  is said to be *embedding-invariant* if for any Markov embedding  $f : \mathbb{R}_+^m \rightarrow \mathbb{R}_+^n$  (with  $n > m$ ) the following holds:

$$\langle X, Y \rangle_{(m)}(\boldsymbol{\xi}) = \langle f_*(X), f_*(Y) \rangle_{(n)}(f(\boldsymbol{\xi})).$$

Čencov proved that the only metrics that satisfy this invariance on the submanifold of probability distributions,  $\{\xi \in \mathbb{R}_+^m \mid \sum_i \xi_i = 1\}$  are given by

$$g_{ij}^{(m)}(\boldsymbol{\xi}) = c_0 \left\{ \frac{\delta_{ij}}{\xi_i} + \frac{1}{\xi_m} \right\}, \quad \text{for } i, j = 1, \dots, m-1,$$

where  $c_0$  is a constant.

<sup>1)</sup>It turns out that the theory is much easier if we do not restrict ourselves to probability distributions, i.e. if we do not require that  $\sum_i \xi_i = 1$  for  $\boldsymbol{\xi} \in \mathcal{A}$ . (This idea is taken from [9], in which a variant of Čencov's result is proven pertaining to such non-negative distributions.)

#### 4.1.4 Rationale

Why are embedding invariant metrics relevant? The solution lies in the fact that a classification of metrics leads to a classification of *divergences*, and for divergences embedding invariance is a very natural property to require.

A divergence, in this context, is any object  $D$  that takes two probability distributions over the same partition  $A$  of  $X$ , and associates a real number to the pair. If we consider a divergence as a measure of the difference between two distributions  $P$  and  $Q$ , it is natural to require that

$$D(K(P), K(Q)) \leq D(P, Q)$$

for any mapping  $K$  from one space of probability distributions to another, since such a mapping cannot add new information, and should therefore not be expected to increase the distance between distributions. Divergences which satisfy this property are called *monotone*.

If there is a mapping  $\bar{K}$  such that  $\bar{K}(K(P)) = P$  for all  $P$ , and  $D$  is monotone, then

$$D(K(P), K(Q)) = D(P, Q),$$

because

$$D(P, Q) = D(\bar{K}(K(P)), \bar{K}(K(Q))) \leq D(K(P), K(Q)) \leq D(P, Q),$$

Since Markov embeddings are invertible in this sense, monotone divergences are invariant under Markov embeddings.

One other fact about monotone divergences plays a role in our present quest: since one may consider a mapping  $K_0$  that maps any distribution down to some fixed  $P_0$ , we have for any  $P, Q$ :

$$D(P_0, P_0) = D(K_0(P), K_0(Q)) \leq D(P, Q).$$

In particular,  $D(P, P)$  is a minimum of  $D(P, Q)$ . We may as well limit ourselves to divergences for which  $D(P, P) = 0$  and  $D(P, Q) > 0$  unless  $P = Q$ , consolidating the interpretation of divergences as a measure of the difference between two probability distributions.

For any partition  $A$  of  $X$  we can then expand

$$D^{(A)}(P, Q) = \sum_{i,j=1}^{\#A} D_{ij}^{(A)}(P) (Q(A_i) - P(A_i)) (Q(A_j) - P(A_j)) + O([Q - P]^3), \quad (4.1)$$

where  $D_{ij}^{(A)}(P)$  is a strictly positive definite symmetric rank 2 tensor field.

Now we see why embedding invariant metrics are interesting: any metric can be used to play the role of  $D_{ij}(P)$  here, but only embedding invariant metrics will yield monotone divergences. Conversely, since any strictly positive definite symmetric rank 2 tensor field can be used as a metric, Čencov's theorem gives a classification of monotone divergences upto second order. It can be shown that the only monotone divergences are those for which  $D_{ij}^{(A)}(P) = c_0 \frac{\delta_{ij}}{P(A_i)}$ . (A proof can be found in [3].)

## 4.2 Generalization to continuous sets

In the previous section we found classifications of metrics and of divergences for finite partitions  $A$ , and we found that those classifications are tightly related. The next step is to show that a similar relation can be established when  $X$  is partitioned into infinitely many sets. We shall not go into the details, but merely state the results and the conditions under which they apply.

One may consider probability distributions on  $X$  as probability distributions on an infinitely fine partition of  $X$ , which can be viewed as a limit of finer and finer finite partitions. If  $P$  is a probability distribution on  $X$ , we denote by  $P_{\mathcal{A}}$  the restriction of  $P$  to  $\mathcal{A}$ , that is  $P_{\mathcal{A}}(A_i) = \int_{A_i} dx P(x)$ . Let  $(A_{(n)})$  be a sequence of partitions of  $X$ , with  $\#A_{(n)} = n$ . A divergence  $D$  is said to be *regular* if

$$\lim_{n \rightarrow \infty} D(P_{\mathcal{A}_{(n)}}, Q_{\mathcal{A}_{(n)}}) = D(P, Q)$$

for any probability distributions  $P$  and  $Q$  on  $X$ , and irrespective of the sequence  $(A_{(n)})$ .

Corcuera and Giummolè show that the only monotone and regular divergences for which  $D(P, P) = 0$ , are given by

$$D(P, Q) = c_0 \int_X dx \frac{[Q(x) - P(x)]^2}{P(x)} + O([Q - P]^3).$$

Although the mathematics is more complicated, the end result turns out to be a straightforward generalization of the finite case.

We shall not try to write down a metric in this case, since that would force us to deal with the intricacies of infinite dimensional manifolds.

## 4.3 The parametric case

Instead, we move straight on to our final goal, showing that the Fisher information is unique. To this end, consider a subspace of parametrized probability distributions on  $X$ . If  $P$  is in this subspace, we may write  $P = p_{\theta}(x)$ , where  $\theta$  is supposed to take values in some interval of  $\mathbb{R}^N$  for some  $N \in \mathbb{N}$ . Any monotone and regular divergence between  $P = p_{\theta}(x)$  and  $Q = p_{\theta'}(x)$  can then be written as

$$D(P, Q) = D(\theta, \theta') = A \int_X dx \frac{[p_{\theta'}(x) - p_{\theta}(x)]^2}{p_{\theta}(x)} + O([p_{\theta'} - p_{\theta}]^3), \quad (4.2)$$

(where we still take  $D(\theta, \theta)$  to be zero).

In writing (4.2), we have implicitly assumed that monotonicity in the parametric case is equivalent to invariance under transformations of  $x$  and  $\theta$ : those transformations correspond to the invertible mappings of §4.1.4. We shall come back to this crucial point later.

For the moment, we shall proceed assuming that (4.2) is indeed a full classification of invariant divergences upto the stated order: the final term in (4.2),  $O([p_{\theta'} - p_{\theta}]^3)$ , is supposed to indicate the fact that we are still making expansions similar to the one in (4.1).

Taking this expansion a bit more seriously, we have

$$p_{\theta'}(x) - p_{\theta}(x) = (\theta'^{\mu} - \theta^{\mu}) \partial_{\mu} p_{\theta}(x) + O([\theta' - \theta]^2).$$

Inserting this into (4.2) gives:

$$D(\boldsymbol{\theta}, \boldsymbol{\theta}') = c_0 \int_X dx \frac{1}{p_{\boldsymbol{\theta}}(x)} \partial_{\mu} p_{\boldsymbol{\theta}}(x) \partial_{\nu} p_{\boldsymbol{\theta}}(x) (\theta'^{\mu} - \theta^{\mu})(\theta'^{\nu} - \theta^{\nu}) + O([\boldsymbol{\theta}' - \boldsymbol{\theta}]^3).$$

Even now, the fact that any metric could be used as a prefactor for the  $(\theta'^{\mu} - \theta^{\mu})(\theta'^{\nu} - \theta^{\nu})$  term remains unchanged. We conclude that the only metrics that give rise to regular and monotone divergences in the parametric case are

$$g_{\mu\nu}(\boldsymbol{\theta}) = c_0 \int_X dx \frac{1}{p_{\boldsymbol{\theta}}(x)} \partial_{\mu} p_{\boldsymbol{\theta}}(x) \partial_{\nu} p_{\boldsymbol{\theta}}(x), \quad (4.3)$$

since invariance of  $D$  is equivalent to covariance of the prefactor for the  $(\theta'^{\mu} - \theta^{\mu})(\theta'^{\nu} - \theta^{\nu})$  term. The metrics (4.3) are just the scalar multiples of the Fisher information.

Note that the uniqueness of the Fisher information as a covariant metric also proves the optimality of using the Fisher information in gradient descent learning: since the optimal learning rule should certainly be reparametrization invariant, and the only covariant gradient is  $g^{\mu\nu} \partial_{\nu}$ , this gradient cannot but yield the optimal rule, apart from a possible time-dependent scalar pre-factor.

One final word of warning seems to be in order: strictly speaking, we have only proven that the Fisher information is the only metric that can be used to build *regular* divergences. While it is clear that any parametrized divergence of the form (4.2) can be used to construct a monotone and regular divergence for the space of all probability distributions, it is not quite as obvious that *all* parameter invariant divergence must of necessity be extendable in this way. Only when this point is cleared up will the classification of monotone and regular divergences constitute a full mathematical proof of the uniqueness of the Fisher metric.

#### 4.A Appendix: The more formal language of probability theory

Mathematicians have come up with a paradigm to avoid our admittedly vague term ‘probability distributions over finite or infinite partitions of a set  $X$ ’. In this appendix we shall give a very brief introduction to this language, to serve as an introduction to the literature.

The following definitions have been taken from [10].

**Definition:** A collection  $\mathcal{B}$  of subsets of a set  $X$  is called a  $\sigma$ -*algebra* if it satisfies

1.  $\emptyset \in \mathcal{B}$ ;
2. If  $A \in \mathcal{B}$ , then also  $X \setminus A \in \mathcal{B}$ ;
3. If  $A_1, A_2, A_3, \dots \in \mathcal{B}$ , then  $\bigcup_i A_i \in \mathcal{B}$ .

A pair  $(X, \mathcal{B})$  where  $\mathcal{B}$  is a  $\sigma$ -algebra over  $X$  is called a *borel space*.

**Definition:** A map  $m : \mathcal{B} \rightarrow [0, 1]$  is called a *probability distribution* if

1. It is *countably additive*, that is

$$m\left(\bigcup_i A_i\right) = \sum_i m(A_i)$$

for any (countable) collection  $\{A_i\}$  of pairwise disjoint sets in  $\mathcal{B}$ , and

2.  $m(X) = 1$ .

**Definition:** Let  $(X_1, \mathcal{B}_1)$  and  $(X_2, \mathcal{B}_2)$  be borel spaces, and let  $f : X_1 \rightarrow X_2$  be a map. Then  $f$  is called a *measurable map* if

$$\forall A \in \mathcal{B}_2 : f^{-1}(A) \equiv \{x \in X_1 \mid f(x) \in A\} \in \mathcal{B}_1.$$

A map  $f : X_1 \rightarrow Y \subseteq \mathbb{R}$  is similarly called a measurable map if

$$\forall y \in Y : f^{-1}((-\infty, y]) \equiv \{x \in X_1 \mid f(x) \leq y\} \in \mathcal{B}_1.$$

**Definition:** Let  $(X, \mathcal{B})$  be a borel space. A map  $f : X \rightarrow Y$  is called a *simple function* if  $f$  takes only a finite number of values, and for every  $f^{-1}(\{y\}) \in \mathcal{B}$  for all  $y \in Y$ .

**Definition: Integral**

Let  $s$  be any non-negative simple function on  $(X, \mathcal{B})$ . Then there exists a partition of  $X$  into disjoint sets  $A_1, \dots, A_k$  all of which belong to  $\mathcal{B}$ , and  $k$  numbers  $a_1, \dots, a_k$  in  $[0, \infty]$  such that  $s = \sum_{i=1}^k a_i \chi_{A_i}$ , where  $\chi_A$  is the indicator function:  $\chi_A(x) = 1$  if  $x \in A$  and zero otherwise.

Let  $P$  be a probability distribution on  $X$ . We then define the *integral*  $\int_X s dP$  of  $s$  with respect to  $P$  by

$$\int_X s dP = \sum_i a_i P(A_i).$$

If  $f$  is a non-negative borel function from  $X$  to  $\mathbb{R}_+^{1)}$ , we define its integral with respect to  $P$  by

$$\int_X f dP = \sup \left\{ \int_X s dP \mid s \text{ is a simple function on } (X, \mathcal{B}), \text{ and } s(x) \leq f(x) \text{ for a } x \in X \right\}.$$

This leads to a more precise definition of mappings between spaces of probability distributions<sup>2)</sup>:

**Definition:** Given two borel spaces  $(X_1, \mathcal{B}_1)$  and  $(X_2, \mathcal{B}_2)$ ,  $K : X_1 \times \mathcal{B}_2 \rightarrow [0, 1]$  is called a *Markov kernel* if

1.  $K(\cdot, A) : X_1 \rightarrow [0, 1]$  is a measurable map for any  $A \in \mathcal{B}_2$ , and
2.  $K(x, \cdot) : \mathcal{B}_2 \rightarrow [0, 1]$  is a probability distribution on  $(X_2, \mathcal{B}_2)$ .

<sup>1)</sup>A borel function is a basically a measurable map, with the exception that it need not be defined on subsets  $A \in \mathcal{B}$  for which  $P(A) = 0$ .

<sup>2)</sup>based on [3].



If  $P$  is a probability distribution on  $(X_1, \mathcal{B}_1)$ , then  $K$  induces a probability distribution  $KP$  on  $(X_2, \mathcal{B}_2)$  defined by

$$KP(A) = \int_{X_1} K(x, A) dP,$$

where the integration is well-defined since  $K(\cdot, A)$  is a measurable map.

The only further definition that is needed to understand [3] is the following:

**Definition:** Let  $(X, \mathcal{B})$  be a borel space. A *finite sub- $\sigma$ -field* is a subset  $\mathcal{A} \subset \mathcal{B}$  such that

1.  $\mathcal{A}$  is a  $\sigma$ -algebra over  $X$  and
2.  $\mathcal{A}$  is finite.

## 4.B A proof of the invariance properties of the Fisher information

For the sake of completeness we shall show that the Fisher metric is indeed invariant under transformations of the random variable and covariant under reparametrizations. Both proofs are straightforward.

### 4.B.1 Invariance under transformations of the random variable

Suppose that our probability distributions are defined in terms of a random variable  $\mathbf{x}$  taking values in  $X \subseteq \mathbb{R}^n$ . Then

$$g_{\mu\nu}(\boldsymbol{\theta}) = \int_X d\mathbf{x} \frac{1}{p_{\boldsymbol{\theta}}(\mathbf{x})} \partial_{\mu} p_{\boldsymbol{\theta}}(\mathbf{x}) \partial_{\nu} p_{\boldsymbol{\theta}}(\mathbf{x}).$$

We can re-express this in terms of another random variable  $\mathbf{y}$  taking values in  $Y \subseteq \mathbb{R}^n$ , if we suppose that  $\mathbf{y} = \mathbf{f}(\mathbf{x})$  is an invertible mapping. We clearly have:

$$\tilde{p}_{\boldsymbol{\theta}}(\mathbf{y}) = \int_X d\mathbf{x} p_{\boldsymbol{\theta}}(\mathbf{x}) \delta(\mathbf{y} - \mathbf{f}(\mathbf{x})). \quad (4.4)$$

If  $\mathbf{f}$  is invertible, then we can use the relation

$$\delta(\mathbf{y} - \mathbf{f}(\mathbf{x})) = \frac{1}{\left| \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right|} \delta(\mathbf{f}^{-1}(\mathbf{y}) - \mathbf{x}),$$

to find that

$$\tilde{p}_{\boldsymbol{\theta}}(\mathbf{y}) = \int_X d\mathbf{x} p_{\boldsymbol{\theta}}(\mathbf{x}) \frac{1}{\left| \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right|} \delta(\mathbf{f}^{-1}(\mathbf{y}) - \mathbf{x}) = \left[ \frac{1}{\left| \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right|} p_{\boldsymbol{\theta}}(\mathbf{x}) \right]_{\mathbf{x}=\mathbf{f}^{-1}(\mathbf{y})} \quad {}^1). \quad (4.5)$$

---

<sup>1)</sup>Historically, this expression is in fact older than (4.4). However, at present the properties of the  $\delta$ -function seem to be wider known than the properties of distributions in general.

If we further note that  $\left| \frac{\partial f}{\partial x} \right|$  does not depend on  $\theta$ , we see that

$$\begin{aligned} \int_Y dy \frac{1}{\tilde{p}_\theta(y)} \partial_\mu \tilde{p}_\theta(y) \partial_\nu \tilde{p}_\theta(y) &= \int_Y dy \left[ \frac{1}{\left| \frac{\partial f}{\partial x} \right| p(x)} \left| \frac{\partial f}{\partial x} \right| \partial_\mu p_\theta(x) \left| \frac{\partial f}{\partial x} \right| \partial_\nu p_\theta(x) \right]_{x=f^{-1}(y)} \\ &= \int_X dx \frac{1}{p(x)} \partial_\mu p_\theta(x) \partial_\nu p_\theta(x), \end{aligned}$$

since  $\int_Y dy = \int_X dx \left| \frac{\partial f}{\partial x} \right|$ .

This proves the invariance under transformation of the random variable.

#### 4.B.2 Covariance under reparametrization

Suppose that  $(\tilde{\theta}^\mu)$  is a new set of coordinates, specified in terms of the old set through the invertible relationship  $\tilde{\theta} = \tilde{\theta}(\theta)$ . Defining  $\tilde{p}_{\tilde{\theta}}(\mathbf{x}) \equiv p_{\theta(\tilde{\theta})}(\mathbf{x})$ , we are then able to compute

$$\tilde{g}_{\mu\nu}(\tilde{\theta}) \equiv \int_X dx \frac{1}{\tilde{p}_{\tilde{\theta}}(\mathbf{x})} \frac{\partial}{\partial \tilde{\theta}^\mu} \tilde{p}_{\tilde{\theta}}(\mathbf{x}) \frac{\partial}{\partial \tilde{\theta}^\nu} \tilde{p}_{\tilde{\theta}}(\mathbf{x}),$$

in terms of  $g_{\mu\nu}(\theta)$ : since

$$\frac{\partial}{\partial \tilde{\theta}^\mu} \tilde{p}_{\tilde{\theta}} = \frac{\partial \theta^\nu}{\partial \tilde{\theta}^\mu} \frac{\partial}{\partial \theta^\nu} p_{\theta(\tilde{\theta})},$$

we may directly conclude that

$$\tilde{g}_{\mu\nu}(\tilde{\theta}) = \left[ \frac{\partial \theta^\rho}{\partial \tilde{\theta}^\mu} \frac{\partial \theta^\lambda}{\partial \tilde{\theta}^\nu} g_{\rho\lambda}(\theta) \right]_{\theta=\theta(\tilde{\theta})}.$$

This is precisely the covariance we claimed.

## RIEMANNIAN GEOMETRY

This appendix introduces the basics of Riemannian geometry for the purpose of the main text. It is not intended as competition for the excellent textbooks<sup>1)</sup> that are available on Riemannian geometry, either in terms of depth or in terms of educational value.

### A.1 Manifolds

For the purpose of this text, we shall not need to be completely formal in the definition of a manifold<sup>2)</sup>. For our purposes, a manifold is a set of points that allows the notion of connecting any two points by a smooth (continuously differentiable) curve. We shall also require that in the neighbourhood of any given point it is possible to define coordinates. These coordinate functions need not be global, and indeed there are many common manifolds that cannot be covered by any single coordinate patch. In such cases we just require the coordinate functions to be compatible: if  $\kappa_1$ ,  $\kappa_2$  and  $\kappa_3$  are three coordinate functions defined in the vicinity of a point P, we insist that the *transition functions*  $\phi_{ij} \equiv \kappa_i \circ \kappa_j^{-1}$  are differentiable and satisfy  $\phi_{ij} \circ \phi_{jk} = \phi_{ik}$  for any  $i$ ,  $j$  and  $k$ .

**Example:** The unit sphere  $S^2 = \{\mathbf{x} \in \mathbb{R}^3 \mid \|\mathbf{x}\| = 1\}$  is a manifold. It cannot be covered by a single coordinate function, but we may take  $\kappa_{\pm} : \mathbf{x} \mapsto (\frac{x}{1 \pm z}, \frac{y}{1 \pm z})$ , which are valid on the entire  $S^2$  except on the south and north pole respectively. Together they clearly cover the sphere. It is to show that  $\phi_{+-} \circ \phi_{-+}$  is the identity on  $S^2 \setminus \{(0, 0, 1), (0, 0, -1)\}$ . ■

In the present work we shall not worry too much about such global issues, and we shall assume that a single coordinate patch covers the entire manifold.

We will therefore use the following – very sloppy – definition of a manifold:

**Definition: Manifold**

A *manifold*  $\mathcal{M}$  is a set of points that allows the notion of smooth curves connecting any two points, and for which we can define a continuously differentiable injection  $\kappa : \mathcal{M} \rightarrow \mathbb{R}^n$  that is invertible on its co-domain.

---

<sup>1)</sup>[2] is a good introduction to general relativity. [11] is much more elaborate and mathematically precise. I found [12] and in particular [13] very useful sources on differential topology. However, they are both in Dutch and may not be easy to acquire.

<sup>2)</sup>In fact, giving a mathematically precise definition of ‘manifold’ seems to be quite difficult: one very useful introduction to differential geometry [13] gives a definition that starts with ‘A manifold is a (para-)compact Hausdorff space satisfying the following properties:’. Even so the author warns that his definition is an ‘intuitive’ one!

## A.2 Vectors

In Euclidean geometry one usually thinks of a vector as the straight line connecting two points. This definition doesn't make sense on manifolds, since we haven't yet defined the notion of straightness: a curve that appears to be straight in  $\mathbb{R}^n$  in one coordinate mapping will generally not be straight in  $\mathbb{R}^n$  in another.

### A.2.1 Tangent vectors

To avoid these complications, we wish to define vectors locally, and independently of coordinates. Instead of thinking of a vector as the line segment connecting two points, we shall use the following:

**Definition:** The *tangent vector*  $X_P$  in a point  $P$  to a curve  $\gamma(t)$  for which  $\gamma(0) = P$  is defined as

$$X_P = \left. \frac{d\gamma}{dt} \right|_{t=0}.$$

Given a set of coordinate functions, that is a mapping

$$\kappa : P \mapsto (\theta^1(P), \theta^2(P), \dots, \theta^n(P)),$$

we may introduce a basis of coordinate vectors as follows: let  $\gamma_{(\mu)}$  be the curve through  $P$  that satisfies  $\kappa(\gamma_{(\mu)}(t)) = \kappa(P) + t\hat{\mu}$ , where  $\hat{\mu}$  is the  $\mu$ -th unit vector in  $\mathbb{R}^n$ . We then define  $\hat{e}_\mu$ , for  $\mu = 1 \dots n$ , by

$$\hat{e}_\mu = \left. \frac{d\gamma_{(\mu)}}{dt} \right|_{t=0}. \quad (\text{A.1})$$

Note that in the expression  $\hat{e}_\mu$ ,  $\mu$  does *not* label the components of a vector  $\hat{e}$ , rather, for each  $\mu$ ,  $\hat{e}_\mu$  is a different vector.

The set of all tangent vectors at a point  $P$  is denoted by  $T_P$ . Any vector  $X \in T_P$  can be decomposed with respect to the basis (A.1):

$$X = X^\mu \hat{e}_\mu,$$

where summation over  $\mu$  is left implicit, in line with the Einstein summation convention which we shall employ throughout this appendix.

If  $X$  is the tangent vector to a curve  $\gamma$ , then  $X^\mu = \left. \frac{d}{dt} \theta^\mu(\gamma(t)) \right|_{t=0}$ .

Note that at this stage we have not yet introduced an inner product, so we have no way of talking about the 'length' of a vector. In particular, we should note that it makes no sense to define  $\|X\| \equiv \sqrt{\sum_\mu (X^\mu)^2}$ , since this would depend critically on the choice of coordinates.

**Remark:** Mathematicians take it one step further. Given a manifold  $\mathcal{M}$ , and a curve  $\gamma(t) \subset \mathcal{M}$  with  $\gamma(0) = P$ , they consider the vector  $X_P$  tangent to  $\gamma$  at  $t = 0$  as the derivative operator that takes the derivative of functions  $f : \mathcal{M} \rightarrow \mathbb{R}$  in the direction of  $\gamma$ :

$$Xf \equiv \nabla_X f = \left. \frac{d}{dt} f(\gamma(t)) \right|_{t=0}.$$

Using the shorthand  $\partial_\mu f$  to denote the functions  $\partial_\mu f : \mathbf{P} \mapsto \frac{\partial}{\partial \theta^\mu} [f \circ \kappa^{-1}(\boldsymbol{\theta})] \Big|_{\boldsymbol{\theta}=\kappa(\mathbf{P})}$ , they write  $X_{\mathbf{P}} f = X^\mu \partial_\mu f$ , or

$$X_{\mathbf{P}} \equiv X_{\mathbf{P}}^\mu \partial_\mu \text{ } ^1).$$

In particular, the basis vectors  $\hat{e}_\mu$  defined above may be identified with the derivative operators  $\partial_\mu$ .

These are the first building blocks of the beautiful formalism of differential topology. In this thesis I shall avoid it as much as possible, in the hope of expanding my potential audience. This approach clearly has its disadvantages, and in places where the derivative operator formalism allows shortcuts or useful insights, these will be mentioned in footnotes. █

### A.2.2 Vector fields

A vector-valued function defined on a manifold is called a vector field:

**Definition:** A (*contravariant*) *vector field*  $X$  over a manifold  $\mathcal{M}$  is a mapping that assigns a vector  $X_{\mathbf{P}}$  to each point  $\mathbf{P} \in \mathcal{M}$ . The set of all smooth vector fields over  $\mathcal{M}$  is denoted by  $T(\mathcal{M})$ .

#### Definition: Coordinate basis

Given a manifold  $\mathcal{M}$  and a coordinate function  $\kappa : \mathcal{M} \rightarrow \mathbb{R}^n$ , the  $\mu$ -th coordinate basis vector field  $\hat{e}_\mu$  is the vector field that assigns the  $\mu$ -th basis vector  $\hat{e}_\mu \Big|_{\mathbf{P}}$  to each point  $\mathbf{P} \in \mathcal{M}$ .

### A.2.3 Transformation behaviour

Given two sets of coordinates,  $(\theta^\mu)$  and  $(\tilde{\theta}^{\tilde{\mu}})$  two bases for  $T_{\mathcal{M}}$  are induced:  $\{\hat{e}_\mu\}$  and  $\{\hat{e}_{\tilde{\mu}}\}$ . Vector fields may then be decomposed in either of two ways:

$$X = X^\mu \hat{e}_\mu \quad \text{or} \quad X = X^{\tilde{\mu}} \hat{e}_{\tilde{\mu}}.$$

The components are related by

$$X^\mu = J^\mu_{\tilde{\mu}} X^{\tilde{\mu}}, \tag{A.2}$$

where  $J^\mu_{\tilde{\mu}} \equiv \frac{\partial \theta^\mu}{\partial \tilde{\theta}^{\tilde{\mu}}}$ <sup>2)</sup> is the Jacobian of the transformation. Some texts actually take this transformation behaviour as the defining property of a vector field.

---

<sup>1)</sup>This definition of the concept *vector* may seem to be far removed from the original, yet it is clear that both definitions contain the same information, and are therefore equally valid.

<sup>2)</sup>This is most easily seen when identifying  $\hat{e}_\mu \equiv \partial_\mu$ , since it is quite obvious that

$$\partial_\mu = \frac{\partial \tilde{\theta}^{\tilde{\mu}}}{\partial \theta^\mu} \partial_{\tilde{\mu}}.$$

Requiring that  $X^\mu \hat{e}_\mu = X^{\tilde{\mu}} \hat{e}_{\tilde{\mu}}$  then implies (A.2).

### A.3 Tensor fields

Vector fields are not the only objects that have ‘sensible’ transformation behaviour in terms of the Jacobian. More generally, one may consider rank  $(n, m)$ -tensor fields. For the purpose of this text, these are simply objects  $T^{\mu_1 \dots \mu_m}_{\nu_1 \dots \nu_n}$  which transform according to

$$T^{\mu_1 \dots \mu_m}_{\nu_1 \dots \nu_n} = J^{\mu_1}_{\tilde{\mu}_1} \dots J^{\mu_m}_{\tilde{\mu}_m} J_{\nu_1}^{\tilde{\nu}_1} \dots J_{\nu_n}^{\tilde{\nu}_n} T^{\tilde{\mu}_1 \dots \tilde{\mu}_m}_{\tilde{\nu}_1 \dots \tilde{\nu}_n},$$

where  $J_{\mu}^{\tilde{\mu}}$  is the matrix inverse of  $J^{\mu}_{\tilde{\mu}}$ . Because of this transformation behaviour, tensors may be multiplied in arbitrary ways to yield new tensors as long as one sticks to the summation convention. For example, if  $A$  is a  $(2, 0)$ -tensor field  $A_{\mu\nu}$  and  $X$  and  $Y$  are vector fields, then  $A(X, Y) \equiv A_{\mu\nu} X^{\mu} Y^{\nu}$  is a scalar field in its transformation behaviour: the Jacobians cancel one another.

**Remark:** More formally, a rank  $(n, m)$ -tensor field is a linear mapping that maps  $m$  covariant<sup>1)</sup> and  $n$  contravariant vector fields into real valued functions on  $\mathcal{M}$ :

$$T(A^{(1)}, \dots, A^{(m)}, X_{(1)}, \dots, X_{(n)}) = T^{\mu_1 \dots \mu_m}_{\nu_1 \dots \nu_n} A_{\mu_1}^{(1)} \dots A_{\mu_m}^{(m)} X_{(1)}^{\nu_1} \dots X_{(n)}^{\nu_n}.$$

The space of  $(n, m)$ -tensor fields over  $\mathcal{M}$  is written as  $T_{(n, m)}(\mathcal{M})$ .

An  $(n, m)$ -tensor field can also be viewed as mapping  $(0, n)$ -tensor fields into  $(m, 0)$ -tensor fields, or as mapping  $m - 1$  1-forms and  $n$  vector fields into vector fields. To give a specific example, consider a  $(2, 1)$ -tensor field  $T$ . If  $X$  and  $Y$  are vector fields, and  $A$  is a 1-form, then  $T$  may be used as any of the following mappings:

$$\begin{aligned} T : (A, X, Y) &\mapsto T^{\mu}_{\nu\rho} A_{\mu} X^{\nu} Y^{\rho} \in \mathbb{R} \\ T : A &\mapsto [T^{\mu}_{\nu\rho} A_{\mu}] \in T_{(2,0)}(\mathcal{M}) \\ T : (X, Y) &\mapsto T^{\mu}_{\nu\rho} X^{\nu} Y^{\rho} \hat{e}_{\mu} \in T(\mathcal{M}). \end{aligned}$$

If  $T$  is symmetric in its covariant indices, one may also define  $T(A, X)$  and  $T(X)$  without ambiguity. ▮

### A.4 Metrics

A metric is an object  $g_{\mu\nu}$  that assigns length to vectors and to curves in the following way: given a vector  $A_P$  at some point  $P$  in a manifold  $\mathcal{M}$ , the length of this vector is defined as

$$\|A_P\| \equiv \sqrt{g(A, A)} \Big|_P \equiv \sqrt{g_{\mu\nu} A^{\mu} A^{\nu}} \Big|_P$$

in any coordinate system.

<sup>1)</sup>A covariant vector field or 1-form  $A$  is a linear map from contravariant vectors fields into functions:  $A(X) = A_{\mu} X^{\mu}$ . 1-forms transform according to  $A_{\mu} = J_{\mu}^{\tilde{\mu}} A_{\tilde{\mu}}$ . 1-forms may be expanded in terms of basis 1-forms  $\omega^{\mu}$  which are defined by  $\omega^{\mu}(\hat{e}_{\nu}) = \delta_{\nu}^{\mu}$  in terms of a basis for contravariant vector fields.

The length  $s$  of a curve  $\gamma : [0, 1] \rightarrow \mathcal{M}$  is given by

$$s = \int dt \sqrt{g_{\mu\nu}(\gamma(t)) \frac{d\theta^\mu(\gamma(t))}{dt} \frac{d\theta^\nu(\gamma(t))}{dt}}. \quad (\text{A.3})$$

A metric also induces an inner product on  $T(\mathcal{M})$ : for any two vectors  $A$  and  $B$  defined at the same point  $P \in \mathcal{M}$ , we may define the inner product  $\langle A, B \rangle$  as

$$\langle A, B \rangle = g(A, B) = g_{\mu\nu} A^\mu B^\nu.$$

A metric must be symmetric and transform as a proper tensor field: given two sets of coordinates  $(\theta^\mu)$  and  $(\tilde{\theta}^{\tilde{\mu}})$  we must have

$$g_{\tilde{\mu}\tilde{\nu}} = \frac{\partial\theta^\mu}{\partial\tilde{\theta}^{\tilde{\mu}}} \frac{\partial\theta^\nu}{\partial\tilde{\theta}^{\tilde{\nu}}} g_{\mu\nu},$$

since otherwise the induced inner product would not be coordinate independent or symmetric.

We shall sometimes have use for a contravariant form of the metric,  $g^{\mu\nu}$ , which is defined as the matrix inverse of  $g_{\mu\nu}$ , i.e.  $g^{\mu\nu} g_{\nu\rho} = \delta_\rho^\mu$ .

Metrics may be used to raise and lower tensor indices, ie given a  $(n, m)$ -tensor  $T^{\mu_1 \dots \mu_m}_{\nu_1 \dots \nu_n}$  we may define the  $(n+1, m-1)$ -tensor  $T_{\mu_1 \dots \mu_m}_{\nu_1 \dots \nu_n} \equiv g_{\mu_1 \mu_2} T^{\mu_1 \dots \mu_m}_{\nu_1 \dots \nu_n}$ , and the  $(n-1, m+1)$ -tensor  $T^{\mu_1 \dots \mu_m \nu}_{\nu_1 \dots \nu_n} \equiv g^{\nu \nu_1} T^{\mu_1 \dots \mu_m}_{\nu_1 \dots \nu_n}$ , etcetera.

## A.5 Affine and metric connection

One is often interested in the rate of change of various fields as one moves around a manifold. In the case of a scalar field (a function), this is no problem: one just takes the values at two nearby points to compute the gradient. However, for vector fields (or tensor fields in general), this is not so easy, since the tangent spaces at two different points are unequal: whereas for a real valued function  $f$  the difference  $f(P') - f(P)$  is again a real number, for a vector field  $X$ , we cannot meaningfully compute  $X(P') - X(P)$ , since the basis vectors at  $P'$  are not in  $T_P$  and vice versa. Is there a way around this? Yes and no.

No, because there is no unique way to solve this problem: given a manifold  $\mathcal{M}$  there is no unique way to define the rate of change of vector fields. However, if one is willing to accept that, the problem may be solved by introducing a linear mapping  $\Phi$  that takes vectors at  $P'$  and maps them into  $T_P$ , allowing comparison between  $\Phi(X(P'))$  and  $X(P)$ .

### A.5.1 Affine connection and parallel transport

We shall be slightly more precise: consider again the set of curves  $\{\gamma_{(\mu)}\}$  passing through  $P$  and satisfying  $\dot{\gamma}_{(\mu)}(P) = \hat{e}_\mu(P)$ , as in §A.1. Let  $P'_{(\mu, \delta t)}$  be points on  $\gamma_{(\mu)}$  near  $P$ , in the sense that  $P'_{(\mu, \delta t)} = \gamma_{(\mu)}(\delta t)$  for infinitesimal  $\delta t$ , and let  $\Phi_{(\mu, \delta t)}$  be linear mappings from  $T_{P'_{(\mu, \delta t)}}$  to  $T_P$ , that reduce to the identity as  $\delta t \rightarrow 0$ . Linearity means that  $\Phi_{(\mu, \delta t)}$  are completely defined by their actions on the coordinate vectors as follows:

$$\Phi_{(\mu, \delta t)} : \hat{e}_\rho^{(\mu, \delta t)} \mapsto \Phi_{(\mu, \delta t)}^\nu (\hat{e}_\rho^{(\mu, \delta t)}) \hat{e}_\nu,$$

where  $\{\hat{e}_\rho^{(\mu, \delta t)}\}$  and  $\{\hat{e}_\nu\}$  are the coordinate bases at  $P'_{(\mu, \delta t)}$  and  $P$  respectively.

If  $\Phi_{(\mu, \delta t)}$  are to reduce to the identity as  $\delta t \rightarrow 0$ , we can write

$$\Phi_{(\mu, \delta t)}(\hat{e}_\nu^{(\mu, \delta t)}) - \hat{e}_\nu = \delta t \Gamma_{\mu\nu}{}^\rho \hat{e}_\rho.$$

for small  $\delta t$ . The constants  $\Gamma_{\mu\nu}{}^\rho$  are called the coefficients of the *affine connection*.

Just as one defines

$$\partial_\mu f(P) = \lim_{\delta t \rightarrow 0} \frac{f(P'_{(\mu, \delta t)}) - f(P)}{\delta t}$$

for scalar functions  $f$ , we may now define the *covariant derivatives* of  $\hat{e}_\nu$  as

$$\nabla_\mu \hat{e}_\nu = \lim_{\delta t \rightarrow 0} \frac{\Phi_{(\mu, \delta t)}(\hat{e}_\nu^{(\mu, \delta t)}) - \hat{e}_\nu}{\delta t} = \Gamma_{\mu\nu}{}^\rho \hat{e}_\rho. \quad (\text{A.4})$$

Note that for any pair  $(\mu, \nu)$ ,  $\nabla_\mu \hat{e}_\nu$  is a vector.

#### A.5.1.1 Formal derivation

In this section we shall derive the action of the covariant derivative on vector fields and general tensor fields. Those who are not interested in the mathematical details, may wish to skip it.

We define the covariant derivative of a function to be the ordinary derivative:

$$\nabla_\mu f \equiv \partial_\mu f, \quad (\text{A.5})$$

and demand that  $\nabla$  behaves like a proper derivative operator, that is:

1.  $\nabla_\mu(\alpha T) = \alpha \nabla_\mu(T)$  (for any tensor  $T$  and constant  $\alpha$ ),
2.  $\nabla_\mu(T + S) = \nabla_\mu(T) + \nabla_\mu(S)$  (for any two tensors  $T$  and  $S$  of the same rank),
3.  $\nabla_\mu(T \cdot S) = T \cdot \nabla_\mu(S) + \nabla_\mu(T) \cdot S$  (for any two tensors  $T$  and  $S$ ),

where  $\cdot$  may represent either tensor multiplication or some contraction.

These properties and (A.4) allow us to conclude that for general vector fields  $X$ :

$$\nabla_\mu X = \nabla_\mu(X^\nu \hat{e}_\nu) = X^\nu \nabla_\mu(\hat{e}_\nu) + \nabla_\mu(X^\nu) \hat{e}_\nu = X^\nu \Gamma_{\mu\nu}{}^\rho \hat{e}_\rho + \partial_\mu X^\nu \hat{e}_\nu. \quad (\text{A.6})$$

They also allow us to conclude that

$$\nabla_\mu \omega^\nu = -\Gamma_{\mu\rho}{}^\nu \omega^\rho, \quad (\text{A.7})$$

since  $\omega^\nu(\hat{e}_\rho) = \delta_\rho^\nu$  implies that

$$0 = \partial_\mu(\delta_\rho^\nu) = \omega^\nu(\nabla_\mu \hat{e}_\rho) + (\nabla_\mu \omega^\nu)(\hat{e}_\rho) = \omega^\nu(\Gamma_{\mu\rho}{}^\lambda \hat{e}_\lambda) + (\nabla_\mu \omega^\nu)(\hat{e}_\rho) = \Gamma_{\mu\rho}{}^\lambda \delta_\lambda^\nu + (\nabla_\mu \omega^\nu)(\hat{e}_\rho),$$

whence

$$(\nabla_\mu \omega^\nu)_\rho = -\Gamma_{\mu\rho}{}^\nu,$$

as claimed.



Finally, we may combine (A.4), (A.5) and (A.7) to find that for a general rank  $(n, m)$ -tensor  $T$ :

$$\begin{aligned} (\nabla_\rho T)_{\mu_1 \dots \mu_m}^{v_1 \dots v_n} &= \partial_\rho T_{\mu_1 \dots \mu_m}^{v_1 \dots v_n} + \Gamma_{\rho\lambda}{}^{v_1} T_{\mu_1 \dots \mu_m}^{\lambda v_2 \dots v_n} + \dots + \Gamma_{\rho\lambda}{}^{v_n} T_{\mu_1 \dots \mu_m}^{v_1 \dots v_{n-1} \lambda} \\ &\quad - \Gamma_{\rho\mu_1}{}^\lambda T_{\lambda \mu_2 \dots \mu_m}^{v_1 \dots v_n} - \dots - \Gamma_{\rho\mu_m}{}^\lambda T_{\mu_1 \dots \mu_{m-1} \lambda}^{v_1 \dots v_n}. \end{aligned} \quad (\text{A.8})$$

Luckily we shall scarcely ever need to take the covariant derivative of anything but functions and vector fields.

### A.5.1.2 Working definitions

For the purposes of the rest of this text the following ‘working definitions’ of the covariant derivative are sufficient:

**Definition:** The *covariant derivative of a scalar function* is defined to be the ordinary derivative:

$$\nabla_\mu f \equiv \partial_\mu f.$$

**Definition:** The *covariant derivative of a vector field* is given by

$$\nabla_\mu X \equiv [\partial_\mu X^\rho + \Gamma_{\mu\nu}{}^\rho X^\nu] \hat{e}_\rho.$$

#### Definition: Covariant derivative in the direction of a vector

The covariant derivative of a function  $f$  in the direction of the vector  $X$  is defined by:

$$\nabla_X f \equiv X^\mu \nabla_\mu f = X^\mu \partial_\mu f.$$

The covariant derivative of a vector field  $Y$  in the direction of a vector  $X$  is defined by

$$\nabla_X Y \equiv X^\mu \nabla_\mu Y.$$

As mentioned before,  $\Gamma_{\mu\nu}{}^\rho$  are called the components of the affine connection:

**Definition:** An *affine connection* is an object  $\Gamma_{\mu\nu}{}^\rho$  that induces a covariant derivative as in the definition above. Requiring that  $\nabla_Y X$  be a vector implies that the affine connection is not a  $(2, 1)$ -tensor field. In fact, if  $(\theta^\mu)$  and  $(\tilde{\theta}^{\tilde{\mu}})$  are two sets of coordinates, one may show that

$$\Gamma_{\tilde{\mu}\tilde{\nu}}{}^{\tilde{\rho}} = \frac{\partial\theta^\mu}{\partial\tilde{\theta}^{\tilde{\mu}}} \frac{\partial\theta^\nu}{\partial\tilde{\theta}^{\tilde{\nu}}} \frac{\partial\tilde{\theta}^{\tilde{\rho}}}{\partial\theta^\rho} \Gamma_{\mu\nu}{}^\rho - \frac{\partial\theta^\mu}{\partial\tilde{\theta}^{\tilde{\mu}}} \frac{\partial\theta^\nu}{\partial\tilde{\theta}^{\tilde{\nu}}} \frac{\partial^2\tilde{\theta}^{\tilde{\rho}}}{\partial\theta^\mu\partial\theta^\nu}. \quad (\text{A.9})$$

This transformation behaviour can be regarded as the defining property for affine connections<sup>1)</sup>.

The antisymmetric part of an affine connection is called the torsion tensor:

---

<sup>1)</sup>One should perhaps write ‘an affine connection is an object  $\Gamma$  with components  $\Gamma_{\mu\nu}{}^\rho$ ’, but  $\Gamma$  not being a tensor we shall always write it in component form.

**Definition:** For any affine connection, the *torsion tensor* is defined as

$$T_{\mu\nu}^{\rho} \equiv \Gamma_{\mu\nu}^{\rho} - \Gamma_{\nu\mu}^{\rho}.$$

From (A.9) we see that the torsion tensor *does* transform as a proper  $(2, 1)$ -tensor.

An affine connection may be used to define what we mean by parallel transport:

**Definition: Parallel transport**

A vector field  $X$  is said to be *parallelly transported*, or *parallelly propagated*, along a curve  $\gamma$  with tangent vector field  $Y$  if

$$\nabla_Y X = 0.$$

A curve for which the tangent vector is propagated parallel to itself is called an affine geodesic:

**Definition:** A curve  $\gamma$  with tangent vector field  $X$  is called an *affine geodesic* if there exists a parametrisation of the curve such that

$$\nabla_X X = 0$$

at all points along the curve.

In terms of coordinates, we may write  $X^\mu(\gamma(t)) = \dot{x}^\mu(t)$ , where  $x^\mu(t)$  are the coordinates of  $\gamma(t)$ , and the geodesic equation becomes

$$\ddot{x}^\rho + \Gamma_{\mu\nu}^{\rho} \dot{x}^\mu \dot{x}^\nu = 0.$$

### A.5.2 Metric connection

A metric induces a special affine connection: the *metric connection*, or *Riemannian connection*. It is derived from the notion of *metric geodesic*:

**Definition:** A *metric geodesic* connecting two points is a<sup>1)</sup> shortest curve between those points in the sense of (A.3).

It is possible to show (see eg. [2]) that metric geodesics are also affine geodesics for the following connection:

$$\Gamma_{\mu\nu}^{\rho} = \frac{1}{2} g^{\rho\lambda} [\partial_\mu g_{\nu\lambda} + \partial_\nu g_{\mu\lambda} - \partial_\lambda g_{\mu\nu}] \quad (\text{A.10})$$

and vice versa. The connection (A.10) is called the metric connection. It satisfies  $\Gamma_{\mu\nu}^{\rho} = \Gamma_{\nu\mu}^{\rho}$ , and so the torsion tensor vanishes: the metric connection is *torsion free*. (Another important property of the metric connection is that it satisfies  $\nabla_\rho g_{\mu\nu} = 0$ .)

On a space with metric connection, a global notion of distance is provided by the *Riemannian distance*: the length of the metric geodesic between two points.

<sup>1)</sup>For infinitesimally separated points this geodesic is unique, but more generally there may be more than one local minimum. As an example, consider walking through a hilly area. When the hills are not very steep, there will be one shortest path between two points that are nearby. However, in the vicinity of a fairly high hill, there may be two locally-shortest paths: one around each side of the hill.

## A.6 Curvature

To top of the *exposé* of Riemannian geometry presented in this appendix, we introduce just one more important tensor that will be used in the main text: the curvature tensor. We shall start with an intuitive introduction, followed by a more formal definition.

### A.6.1 Intuitive introduction

In general, when one transports a vector parallelly from one point to another point in a curved space, the result is not path independent. The canonical example is carrying a spear around the earth: start of from the north pole pointing a spear in the direction of the Greenwich meridian and walk all the way down to the equator, still pointing the spear horizontally along the meridian. Not changing the direction of the spear, turn around and walk along the equator until you reach the 90° meridian. Still not rotating the spear, walk back to the north pole. The spear now points in the direction of the 90° meridian, in contrast to the original situation.

How does this happen? The answer must lie in the observation that parallel transport according to the metric connection on the earth is not the same as parallel transport in the three dimensional space in which the earth is embedded. On a sphere one finds that the angle of deviation is equal to the covered spherical angle. We shall wish to find a quantitative result for more general manifolds though.

To this end, consider carrying a vector  $X$  from a point  $P$  to a nearby point  $Q$  along two paths (see figure). Let's assign coordinates  $\theta = (\theta^\mu)$  to  $P$ ,  $\theta + \delta_1$  to  $R$ ,  $\theta + \delta_2$  to  $S$  and  $\theta + \delta_1 + \delta_2$  to  $Q$ . We shall assume that all components of both  $\delta_1$  and  $\delta_2$  are infinitesimal.

First carry  $X$  from  $P$  to  $R$ :

$$X_R = X_P - \Gamma_{\mu\nu}^\rho|_P \delta_1^\mu X_P^\nu \hat{e}_\rho.$$

(This follows from requiring that the covariant derivative of  $X$  in the direction of  $\delta_1$  should vanish. Integrating

$0 = \nabla_\mu X \equiv \partial_\mu X + \Gamma_{\mu\nu}^\rho X^\nu \hat{e}_\rho$  in the direction of  $\delta_1$ , we find that the first term integrates to  $X_R - X_P$ , while the second term yields  $\Gamma_{\mu\nu}^\rho|_P \delta_1^\mu X_P^\nu \hat{e}_\rho$  to first order in  $\delta_1$ .)

Expanding

$$\Gamma_{\mu\nu}^\rho|_R = \Gamma_{\mu\nu}^\rho|_P + \delta_1^\lambda \partial_\lambda \Gamma_{\mu\nu}^\rho|_P + O([\delta_1]^2),$$

we next carry  $X$  from  $R$  to  $Q$ :

$$\begin{aligned} X_{Q(\text{via } R)} &= X_R - \Gamma_{\mu\nu}^\rho|_R \delta_2^\mu X_R^\nu \hat{e}_\rho \\ &= (X_P - \Gamma_{\mu\nu}^\rho|_P \delta_1^\mu X_P^\nu \hat{e}_\rho) - (\Gamma_{\mu\nu}^\rho|_P + \delta_1^\lambda \partial_\lambda \Gamma_{\mu\nu}^\rho|_P) \delta_2^\mu (X_P^\nu - \Gamma_{\sigma\tau}^\mu|_P \delta_1^\sigma X_P^\tau) \hat{e}_\rho \\ &= X_P - \Gamma_{\mu\nu}^\rho|_P (\delta_1^\mu + \delta_2^\mu) X_P^\nu \hat{e}_\rho + \Gamma_{\mu\nu}^\rho|_P \Gamma_{\sigma\tau}^\nu|_P \delta_2^\mu \delta_1^\sigma X_P^\tau \hat{e}_\rho + \\ &\quad - (\partial_\lambda \Gamma_{\mu\nu}^\rho|_P) \delta_1^\lambda \delta_2^\mu X_P^\nu \hat{e}_\rho + O([\delta_a]^3). \end{aligned}$$

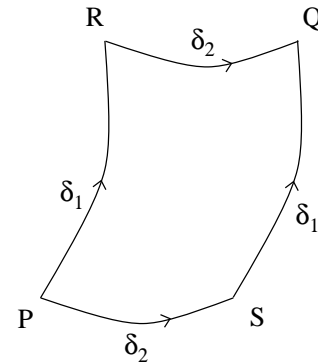


Figure A.1: *Parallel transport from a point  $P$  to a point  $Q$  along two paths. On a curved manifold, the result of parallel transport is path-dependent.*

$X_{Q(\text{via S})}$  may be found from scratch, or more easily by interchanging  $\delta_1$  and  $\delta_2$  in the expression for  $X_{Q(\text{via R})}$ . What is interesting, of course, is the difference between these two:

$$X_{Q(\text{via S})} - X_{Q(\text{via R})} = \left( \Gamma_{\mu\sigma}{}^\lambda \Big|_P \Gamma_{\nu\rho}{}^\sigma \Big|_P - \Gamma_{\nu\sigma}{}^\lambda \Big|_P \Gamma_{\mu\rho}{}^\sigma \Big|_P - \partial_\nu \Gamma_{\mu\rho}{}^\lambda \Big|_P + \partial_\mu \Gamma_{\nu\rho}{}^\lambda \Big|_P \right) \delta_1^\mu \delta_2^\nu X_P^\rho \hat{e}_\lambda.$$

The expression between brackets is called the Riemann tensor  $R_{\mu\nu\rho}{}^\lambda$ . It expresses the  $\lambda$ -th component of the difference between the result of parallelly transporting  $\hat{e}_\rho$  first in the  $\nu$ -th direction and then in the  $\mu$ -th direction, and the result of doing it in the reverse order.

Note for the attentive reader: in calculating  $X_Q$  to second order in  $\delta_1$  and  $\delta_2$ , we should of course really have taken second order terms into account when calculating  $X_R$ . However, careful inspection shows that any  $[\delta_1]^2$  and  $[\delta_2]^2$  terms cancel in the final result.

### A.6.2 Formal definition

For any connection  $\Gamma_{\mu\nu}{}^\rho$  the curvature tensor, or *Riemann tensor* is defined by

$$R_{\mu\nu\rho}{}^\lambda = \partial_\mu \Gamma_{\nu\rho}{}^\lambda - \partial_\nu \Gamma_{\mu\rho}{}^\lambda + \Gamma_{\mu\sigma}{}^\lambda \Gamma_{\nu\rho}{}^\sigma - \Gamma_{\nu\sigma}{}^\lambda \Gamma_{\mu\rho}{}^\sigma. \quad (\text{A.11})$$

There are several further observations one can make about the Riemann tensor:

- From the definition we see that it is antisymmetric in its first two indices:

$$R_{\mu\nu\rho}{}^\lambda = -R_{\nu\mu\rho}{}^\lambda. \quad (\text{A.12})$$

- For a symmetric (or torsion free) connection the last two terms cancel. The Riemann tensor for a torsion free connection then satisfies

$$R_{\mu\nu\rho}{}^\lambda + R_{\nu\rho\mu}{}^\lambda + R_{\rho\mu\nu}{}^\lambda = 0 \quad (\text{for torsion free connections}).$$

- For the metric connection we can also establish the following identity:

$$R_{\mu\nu\rho\lambda} = R_{\rho\lambda\mu\nu} \quad (\text{for the metric connection}).$$

Combining this with (A.12) we find  $R_{\mu\nu\rho\lambda} = -R_{\mu\nu\lambda\rho}$ . Finally the curvature for the metric connection also satisfies the Bianchi identities:

$$\nabla_\mu R_{\sigma\nu\rho}{}^\lambda + \nabla_\nu R_{\sigma\rho\mu}{}^\lambda + \nabla_\rho R_{\sigma\mu\nu}{}^\lambda = 0 \quad (\text{for the metric connection}).$$

(Both of these relations are easy to prove in geodesic coordinates, see § 6.6 of [2].)

### A.6.3 Affine flatness

A manifold on which the Riemann tensor vanishes everywhere is called (affine) flat. On an affine flat manifold, it is possible to introduce a coordinate system in which the connection vanishes everywhere<sup>1)</sup>. Such a coordinate system is called an affine coordinate system.

<sup>1)</sup>A proof of this fact can be found in § 6.7 of [2].

## A.7 Submanifolds

Consider an  $n$ -dimensional subspace  $\mathcal{S}$  of an  $m$ -dimensional manifold  $\mathcal{M}$  with coordinates  $\theta^\mu$ ,  $\mu = 1 \dots m$ . On  $\mathcal{S}$  we may introduce coordinates  $\vartheta^\alpha$ ,  $\alpha = 1 \dots n$ . We consider  $\mathcal{S}$  to be embedded in  $\mathcal{M}$ : any point in  $\mathcal{S}$  is also a point in  $\mathcal{M}$ . The  $\mathcal{M}$ -coordinates of such a point are a function of the  $\mathcal{S}$ -coordinates:  $\theta = \theta(\vartheta)$ . When this function is smooth and the coordinate vectors

$$\hat{e}_\alpha \equiv \frac{\partial \theta^\mu}{\partial \vartheta^\alpha} \hat{e}_\mu \quad (A.13)$$

are linearly independent, that is, the Jacobian

$$B_\alpha^\mu \equiv \frac{\partial \theta^\mu}{\partial \vartheta^\alpha}$$

is non-zero definite<sup>2)</sup>,  $\mathcal{S}$  is called a submanifold of  $\mathcal{M}$ .

A metric on  $\mathcal{M}$  naturally induces a metric on  $\mathcal{S}$ :

$$g_{\alpha\beta} \equiv \langle \hat{e}_\alpha, \hat{e}_\beta \rangle = \langle B_\alpha^\mu \hat{e}_\mu, B_\beta^\nu \hat{e}_\nu \rangle = B_\alpha^\mu B_\beta^\nu g_{\mu\nu}.$$

We may also study the covariant derivative of vector fields in  $T(\mathcal{S})$ :

$$\begin{aligned} \nabla_\alpha \hat{e}_\beta &= \nabla_\alpha (B_\beta^\nu \hat{e}_\nu) = (\partial_\alpha B_\beta^\nu) \hat{e}_\nu + B_\beta^\nu \nabla_{\hat{e}_\alpha} \hat{e}_\nu \\ &= (\partial_\alpha B_\beta^\nu) \hat{e}_\nu + B_\beta^\mu B_\alpha^\nu \nabla_\mu \hat{e}_\nu = \left\{ \partial_\alpha B_\beta^\rho + B_\alpha^\mu B_\beta^\nu \Gamma_{\mu\nu}^\rho \right\} \hat{e}_\rho \end{aligned} \quad (A.14)$$

When  $\nabla_\alpha \hat{e}_\beta$  lies entirely in  $T(\mathcal{S})$  for any  $\alpha$  and  $\beta$ ,  $\mathcal{S}$  is called a *flat* submanifold of  $\mathcal{M}$ . In general, however,  $\nabla_\alpha \hat{e}_\beta$  has components orthogonal to  $T(\mathcal{S})$ , and thus cannot be viewed as a covariant derivative on  $\mathcal{S}$ . This may be cured by projection: the projection on  $T(\mathcal{S})$  of a vector field in  $T(\mathcal{M})$  is given by

$$\hat{X} = \langle X, \hat{e}_\alpha \rangle g^{\alpha\beta} \hat{e}_\beta.$$

By applying this projection to (A.14) we obtain a proper covariant derivative on  $\mathcal{S}$ :

$$\hat{\nabla}_\alpha \hat{e}_\beta = \left\{ \partial_\alpha B_\beta^\rho + B_\alpha^\mu B_\beta^\nu \Gamma_{\mu\nu}^\rho \right\} \langle \hat{e}_\rho, \hat{e}_\gamma \rangle g^{\gamma\delta} \hat{e}_\delta. \quad (A.15)$$

Noting that, given a metric, the components of an affine connection may be computed from the covariant derivative of the coordinate vectors:

$$\Gamma_{\mu\nu}^\rho = g^{\rho\lambda} \langle \nabla_\mu \hat{e}_\nu, \hat{e}_\lambda \rangle,$$

(A.15) yields the following connection on  $\mathcal{S}$ :

$$\Gamma_{\alpha\beta}^\gamma \equiv g^{\gamma\delta} \langle \hat{\nabla}_\alpha \hat{e}_\beta, \hat{e}_\delta \rangle = g^{\gamma\delta} \langle \nabla_\alpha \hat{e}_\beta, \hat{e}_\delta \rangle = g^{\gamma\delta} \left\{ \partial_\alpha B_\beta^\rho + B_\alpha^\mu B_\beta^\nu \Gamma_{\mu\nu}^\rho \right\} B_\delta^\lambda g_{\rho\lambda}.$$

Note that a flat manifold may very well contain non-flat submanifolds, as is exemplified by the possibility of having a sphere as a submanifold of  $\mathbb{R}^3$ . Also note that affine flatness of  $\mathcal{S}$  does *not* imply that  $\mathcal{S}$  is a flat submanifold of  $\mathcal{M}$ , even if  $\mathcal{M}$  is itself affine flat<sup>3)</sup>. On the other hand, a flat submanifold of an affine flat manifold *is* affine flat.

<sup>1)</sup>We shall write  $\hat{e}_\alpha$  for the basis vectors of  $T(\mathcal{S})$ , relying on the choice of indices  $\alpha, \beta, \dots$  to distinguish them from the basis vectors  $\hat{e}_\mu$  of  $T(\mathcal{M})$ . Note again that (A.13) is obvious if we identify  $\hat{e}_\mu \equiv \frac{\partial}{\partial \theta^\mu}$ , and  $\hat{e}_\alpha \equiv \frac{\partial}{\partial \vartheta^\alpha}$ .

<sup>2)</sup>or – equivalently – has full rank.

<sup>3)</sup>Consider for example a torus (which is flat with respect to the metric connection) as a submanifold of  $\mathbb{R}^3$ .

## SOME SPECIAL FAMILIES OF PROBABILITY DISTRIBUTIONS

Two families of probability distributions that appear time and time again are the exponential family and the mixture family. We shall introduce them briefly and discuss the link with  $\alpha$ -connections.

### B.1 Exponential family

An exponential family of probability distributions is a family of distributions that can be written as

$$\mathcal{M} = \left\{ p_\theta \mid p_\theta(\mathbf{x}) = e^{\sum_\mu \theta^\mu x_\mu + c_0(\theta)} P_0(\mathbf{x}) \right\},$$

where  $P_0(\mathbf{x})$  is some fixed ‘carrier’ measure (which may be discrete or continuous) and  $c_0$  ensures that the distribution is normalized.

An important subset of the exponential family is the set of Gaussian distributions: while in the normal way of writing there is an  $x^2$  term in the exponent which seems to disqualify them, we may write

$$P_{(\mu, \sigma)}(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}\right) = \exp\left(\frac{\mu}{\sigma^2} x - \frac{1}{\sigma^2} y - \frac{1}{2} \frac{\mu^2}{\sigma^2}\right) \frac{1}{\sqrt{2\pi}\sigma} \delta(y - x^2),$$

showing that in fact the normal distributions *do* form an exponential family.

We shall now compute the metric and  $\alpha$ -connections for exponential families. We have:

$$\ell = \log p = \theta^k x_k + c_0(\theta) + \log P(\mathbf{x}).$$

Therefore:

$$\partial_\mu \ell = x_\mu + \partial_\mu c_0(\theta),$$

(where, as usual,  $\partial_\mu \equiv \frac{\partial}{\partial \theta^\mu}$ ).

We thus have:

$$g_{\mu\nu} \equiv \int d\mathbf{x} p \partial_\mu \ell \partial_\nu \ell = \int d\mathbf{x} p [x_\mu x_\nu + x_\mu \partial_\nu c_0 + \partial_\mu c_0 x_\nu + \partial_\mu c_0 \partial_\nu c_0].^{1)}$$

We see that

$$x_\mu p = \left( \partial_\mu e^{\theta^k x_k} \right) e^{c_0} P_0$$

<sup>1)</sup>Here, and in the following, the integration should be replaced by summation if  $P_0$  is a discrete distribution.

and

$$x_\mu x_\nu p = \left( \partial_\mu \partial_\nu e^{\theta^k x_k} \right) e^{c_0} P_0.$$

Noting that  $\int d\mathbf{x} p = 1$  implies

$$\int d\mathbf{x} e^{\theta^k x_k} P_0(\mathbf{x}) = e^{-c_0(\boldsymbol{\theta})},$$

we find:

$$\begin{aligned} g_{\mu\nu}(\boldsymbol{\theta}) &= \int d\mathbf{x} \left[ \partial_\mu \partial_\nu e^{\theta^k x_k} + \partial_\mu e^{\theta^k x_k} \partial_\nu c_0 + \partial_\mu c_0 \partial_\nu e^{\theta^k x_k} + \partial_\mu \partial_\nu c_0 \right] e^{c_0(\boldsymbol{\theta})} P_0(\mathbf{x}) \\ &= \left[ \partial_\mu \partial_\nu e^{-c_0} + \partial_\mu e^{-c_0} \partial_\nu c_0 + \partial_\mu c_0 \partial_\nu e^{-c_0} + \partial_\mu c_0 \partial_\nu c_0 \right] e^{c_0} \\ &= -\partial_\mu \partial_\nu c_0(\boldsymbol{\theta}). \end{aligned}$$

Similarly boring calculations show that

$$E(\partial_\mu \partial_\nu \ell \partial_\rho \ell) = 0,$$

while

$$E(\partial_\mu \ell \partial_\nu \ell \partial_\rho \ell) = -\partial_\mu \partial_\nu \partial_\rho c_0(\boldsymbol{\theta}).$$

Inserting these two into the definition (1.16), we find

$$\Gamma_{\mu\nu\rho}^{(\alpha)}(\boldsymbol{\theta}) = \frac{\alpha - 1}{2} \partial_\mu \partial_\nu \partial_\rho c_0(\boldsymbol{\theta}).$$

In particular, we see that the exponential family is +1-flat.

## B.2 Mixture family

Any family of probability distributions that can be written as

$$\mathcal{M} = \{ p_\theta(\mathbf{x}) \mid p_\theta(\mathbf{x}) = \theta^\mu P_\mu(\mathbf{x}) + c_0(\boldsymbol{\theta}) P_0(\mathbf{x}) \}$$

is called a mixture family. In this definition,  $P_0$  and each  $P_\mu$  should be properly normalized probability distributions, each of the  $\theta^\mu$  should be in  $[0, 1]$ , and  $c_0$  is a normalization constant:  $c_0 = 1 - \sum_\mu \theta^\mu$ , which should also be in  $[0, 1]$ .

We shall again try to compute the metric and  $\alpha$ -connections. To this end, note that

$$\partial_\mu \ell(\boldsymbol{\theta}) = \frac{1}{p_\theta} \partial_\mu p_\theta = \frac{1}{p_\theta} [P_\mu - P_0].$$

Therefore:

$$g_{\mu\nu}(\boldsymbol{\theta}) = \int d\mathbf{x} \frac{1}{p_\theta(\mathbf{x})} [P_\mu(\mathbf{x}) - P_0(\mathbf{x})] [P_\nu(\mathbf{x}) - P_0(\mathbf{x})].<sup>1)</sup>$$

---

<sup>1)</sup>Again, integration should be replaced by summation if the  $P_\mu$  are discrete distributions.

It does not seem possible to calculate this integral for general  $\{P_\mu\}$ . However, in one special case we may: namely when we limit ourselves to distributions over a countable set of atoms:

$$p_\theta(x) = \sum_{\mu=1}^N \theta^\mu \delta_{x,a_\mu} + c_0(\theta) \delta_{x,a_0}, \quad (\text{B.1})$$

where the same conditions on the  $\theta^\mu$  apply as before. As long as these conditions are met,  $N$  may be infinity.

In this case we find

$$\partial_\mu \ell(\theta) = \frac{1}{p_\theta} \partial_\mu p_\theta = \frac{1}{p_\theta} [\delta_{x,a_\mu} - \delta_{x,a_0}].$$

Therefore:

$$g_{\mu\nu}(\theta) = \sum_x \frac{1}{p_\theta} \partial_\mu p_\theta \partial_\nu p_\theta = \sum_{\rho=0}^N \frac{(\delta_{\rho\mu} - \delta_{\rho 0})(\delta_{\rho\nu} - \delta_{\rho 0})}{\sum_{\lambda=1}^N \theta^\lambda \delta_{\rho\lambda} + c_0(\theta) \delta_{\rho 0}}.$$

The numerator is always zero except when  $\rho = 0$ , or when  $\rho = \mu = \nu$ . We thus find:

$$g_{\mu\nu}(\theta) = \frac{\delta_{\mu\nu}}{\theta^\mu} + \frac{1}{1 - \sum_\lambda \theta^\lambda}, \quad (\text{B.2})$$

a result first obtained by Čencov in 1972 [8].

Calculating the  $\alpha$ -connections one stumbles across similar difficulties as for the metric, however for one special value of  $\alpha$  these difficulties collapse:

$$\begin{aligned} \Gamma_{\mu\nu\rho}^{(\alpha)} &= \int d\mathbf{x} p \left[ \partial_\mu \partial_\nu \ell \partial_\rho \ell + \frac{1-\alpha}{2} \partial_\mu \ell \partial_\nu \ell \partial_\rho \ell \right] \\ &= \int d\mathbf{x} p \left[ \left( \frac{1}{p} \partial_\mu \partial_\nu p - \frac{1}{p^2} \partial_\mu p \partial_\nu p \right) \frac{1}{p} \partial_\rho p + \frac{1-\alpha}{2} \frac{1}{p^3} \partial_\mu p \partial_\nu p \partial_\rho p \right]. \end{aligned}$$

Noting that  $\partial_\mu \partial_\nu p_\theta(\mathbf{x}) = \frac{\partial}{\partial \theta^\mu} (P_\nu(\mathbf{x}) - P_0(\mathbf{x})) = 0$ , this reduces to

$$\Gamma_{\mu\nu\rho}^{(\alpha)}(\theta) = -\frac{1+\alpha}{2} \int d\mathbf{x} \frac{1}{p_\theta^2} \partial_\mu p_\theta \partial_\nu p_\theta \partial_\rho p_\theta,$$

and in particular  $\Gamma_{\mu\nu\rho}^{(-1)} = 0$ . The mixture family is thus found to be  $-1$ -flat.

As an aside we note that the calculation may be performed for general  $\alpha$  if we limit ourselves to distributions over a countable set of atoms again: in the special case of (B.1) we find:

$$\begin{aligned} \Gamma_{\mu\nu\rho}^{(\alpha)}(\theta) &= \frac{1+\alpha}{2} \sum_{\lambda=0}^N \frac{1}{p_\theta^2} (\delta_{\lambda\mu} - \delta_{\lambda 0})(\delta_{\lambda\nu} - \delta_{\lambda 0})(\delta_{\lambda\rho} - \delta_{\lambda 0}) \\ &= \frac{1+\alpha}{2} \left[ \delta_{\mu\nu} \delta_{\mu\rho} \left( \frac{1}{\theta^\mu} \right)^2 + \left( \frac{1}{1 - \sum_\lambda \theta^\lambda} \right)^2 \right]. \end{aligned}$$



---

## BIBLIOGRAPHY

- [1] S. Amari, *Differential-geometrical methods in statistics*, Lecture notes in statistics, Springer-Verlag, Berlin, 1985.
- [2] R. d'Inverno, *Introducing Einstein's relativity*, Clarendon Press, Oxford, 1992.
- [3] J. M. Corcuera and F. Giummolè, *A characterization of monotone and regular divergences*, Accepted by Annals of the Institute of Statistical Mathematics, 1998.
- [4] S. Amari, *Natural gradient works efficiently in learning*, Neural Computation **10** (1998) 251.
- [5] Ma, Ji, and Farmer, *An efficient EM-based training algorithm for feedforward neural networks*, Neural Networks **10** (1997) 243.
- [6] M. Rattray and D. Saad, *Transients and asymptotics of natural gradient learning*, Submitted to ICANN 98, 1998.
- [7] A. Fujiwara and S. Amari, *Gradient systems in view of information geometry*, Physica **D 80** (1995) 317.
- [8] N. Čencov, *Statistical decision rules and optimal inference*, Transl. Math. Monographs, vol. 53, Amer. Math. Soc., Providence USA, 1981, (translated from the Russian original, published by Nauka, Moscow in 1972).
- [9] L. L. Campbell, *An extended Čencov characterization of the information metric*, Proc. AMS **98** (1986) 135.
- [10] K. R. Parthasarathy, *Introduction to probability and measure*, MacMillan, Delhi, 1977.
- [11] C. W. Misner, K. S. Thorne, and J. A. Wheeler, *Gravitation*, W. H. Freeman and company, San Fransisco, 1973.
- [12] W. A. van Leeuwen, *Structuur van Ruimte-tijd, syllabus naar het college van G. G. A. Bäuerle en W. A. van Leeuwen, uitgewerkt door F. Hoogeveen*, in Dutch; obtainable from the University of Amsterdam, Institute for Theoretical Physics, Valckenierstraat 65, Amsterdam.
- [13] R. H. Dijkgraaf, *Meetkunde, topologie en fysica, syllabus voor het college Wiskunde voor Fysici II*, in Dutch; obtainable from the University of Amsterdam, Dept. of Mathematics, Plantage Muidergracht 24, Amsterdam.

