

# Information Geometry from Divergences

---

**Armin van de Venn, Prof. Tomoi Koide**

June 30, 2026

Laboratory of Theoretical Physics

University of Tartu

Idea: Use Differential Geometry to study Statistics/Probability Theory

- Fisher information matrix:

$$g_{ij}(\theta) = \mathbb{E}_{p(x|\theta)} \left[ \frac{\partial}{\partial \theta_i} \log p(x|\theta) \frac{\partial}{\partial \theta_j} \log p(x|\theta) \right]$$

“How sensitive is the distribution to small changes of the parameters”

- C. R. Rao noticed ( $\approx$  1945) that this defines a Riemannian metric on a parametric family of probability distributions
- Shun'ichi Amari (pioneer of AI!) developed the modern framework of information geometry, including dual affine connections and divergences in the late 1970s and early 1980s

# Statistical Manifolds

A **statistical manifold** is a smooth family of probability distributions

$$\mathcal{M} = \{p(x|\theta) \mid \theta = (\theta^1, \dots, \theta^n) \in \Theta \subseteq \mathbb{R}^n\}$$

- Each element of  $\mathcal{M}$  is a probability distribution
- Topology usually inherited from  $\Theta$
- Smooth structure defined with the charts  $\varphi: \mathcal{M} \rightarrow \mathbb{R}^n$ ,  $\varphi(p(x|\theta)) = \theta$
- The parameters  $\theta^i$  act as local coordinates on  $\mathcal{M}$ , thus  $\mathcal{M}$  is  $n$ -dimensional

Example:

$$\mathcal{M} = \{\mathcal{N}(x|\mu, \sigma^2) \mid (\mu, \sigma) \in \mathbb{R} \times \mathbb{R}_{>0} \subseteq \mathbb{R}^2\}$$

is a two-dimensional statistical manifold.

# Divergences

## Metric

A metric is a map

$$d: X \times X \rightarrow \mathbb{R}$$

satisfying

- (i)  $d(x, y) \geq 0$
- (ii)  $d(x, y) = 0 \iff x = y$
- (iii)  $d(x, y) = d(y, x)$
- (iv)  $d(x, z) \leq d(x, y) + d(y, z)$

## Divergence

A divergence is a smooth map

$$D: \mathcal{M} \times \mathcal{M} \rightarrow \mathbb{R}$$

satisfying

- (i)  $D(p, q) \geq 0$
- (ii)  $D(p, q) = 0 \iff p = q$

In general  $D$  is not symmetric and does not satisfy the triangle inequality.

**Observation:** A divergence has a minimum along the diagonal

- Near the diagonal  $\theta' = \theta$ , the linear terms vanish and

$$D(\theta, \theta + d\theta) = \frac{1}{2}H_{ij}(\theta)d\theta^i d\theta^j + O(d\theta^3),$$

where

$$H_{ij}(\theta) := \partial'_i \partial'_j D|_{\theta=\theta'}, \quad \partial'_i D := \frac{\partial}{\partial \theta'^i} D(\theta, \theta').$$

### Additional Requirement

We require that  $H_{ij}$  is positive definite, i.e.

$$H_{ij}v^i v^j > 0 \quad \text{for all } v \neq 0.$$

For two vector fields  $X, Y \in \mathfrak{X}(\mathcal{M})$  we define

$$D[X|Y] := X^{(1)}Y^{(2)}D \Big|_{\theta=\theta'},$$

where  $X^{(1)}$  and  $Y^{(2)}$  are the canonical lifts to the first and second factor of  $\mathcal{M} \times \mathcal{M}$ :  
 $X_{(p,q)}^{(1)} = (X_p, 0)$ ,  $Y_{(p,q)}^{(2)} = (0, Y_q)$ . In other words  $X^{(1)} = X^i(\theta) \frac{\partial}{\partial \theta^i}$ ,  $Y^{(2)} = Y^i(\theta') \frac{\partial}{\partial \theta'^i}$ .

## Metric Tensor

$$g(X, Y) := -D[X|Y]$$

$$g_{ij} = -\partial_i \partial'_j D \Big|_{\theta=\theta'} = \partial'_i \partial'_j D \Big|_{\theta=\theta'}$$

## Connection

$$g(\nabla_X Y, Z) := -D[XY|Z]$$

$$\Gamma^l_{ji} g_{kl} = -\partial_i \partial_j \partial'_k D \Big|_{\theta=\theta'}$$

## Dual Divergence

Existence of  $D$  implies the existence of the **dual divergence**

$$D^*(p, q) := D(q, p).$$

It satisfies the properties of a divergence as well.  $D^*$  induces

- $g^*_{ij} = -\partial_i \partial'_j D(\theta', \theta)|_{\theta'=\theta} = -\partial_j \partial'_i D(\theta, \theta')|_{\theta=\theta'} = g_{ji}$
- $\Gamma^* \neq \Gamma$

One can show:

$$Xg(Y, Z) = g(\nabla_X Y, Z) + g(Y, \nabla^*_X Z)$$

Two connections  $\nabla, \nabla^*$  that satisfy this relation are said to be **dual** wrt the metric  $g$ .

## Dual Divergence

Notice that

$$\Gamma^l_{ji}g_{kl} = -\partial_i\partial_j\dots \quad \Gamma^{*l}_{ji}g_{kl} = -\partial'_i\partial'_j\dots$$

hence  $\nabla$  and  $\nabla^*$  are both **torsion-free**. However:

$$\nabla_i g_{jk} = -T_{jik} \quad \nabla^*_i g_{jk} = T_{kij},$$

where

$$T^l_{ik} := \Gamma^l_{ik} - \Gamma^{*l}_{ik}$$

is the so-called **Amari–Chentsov Tensor** (or **cubic tensor**) which we know as the nonmetricity.

Nonmetricity naturally emerges within the setup of information geometry and is not added by hand! In fact, nonmetricity arises due to the asymmetry of the divergence!

## Case Study: Exponential Families

Probability distributions of form

$$p(x|\theta) = h(x) \exp(\theta^i F_i(x) - \psi(\theta)), \quad \psi(\theta) = \ln \left( \int h(x) \exp(\theta^i F_i(x)) dx \right).$$

- Example: Gaussian distributions

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

with choice  $h(x) = 1$ ,  $F(x) = (x, x^2)$ ,  $(\theta^1, \theta^2) = (\mu/\sigma^2, -1/(2\sigma^2))$ , and

$$\psi(\theta) = -\frac{(\theta^1)^2}{4\theta^2} + \frac{1}{2} \ln(\pi) - \frac{1}{2} \ln(-\theta^2).$$

## Case Study: Exponential Families

Standard choice of divergence: **Kullback–Leibler (KL) divergence**

$$D_{\text{KL}}(p, q) = \int p(x) \ln \left( \frac{p(x)}{q(x)} \right) dx.$$

- Related to maximum likelihood estimation
- Induces Fisher information metric
- Behaves naturally for exponential families

Induced geometry:

$$g_{ij}(\theta) = \partial_i \partial_j \psi(\theta)$$

$$\Gamma^l_{ij}(\theta) g_{kl}(\theta) = 0$$

$$\Gamma^{*l}_{ij}(\theta) g_{kl}(\theta) = \partial_i \partial_j \partial_k \psi(\theta)$$

## Dual Coordinates

Since  $g_{ij}(\theta) = \partial_i \partial_j \psi(\theta)$  is positive definite,  $\psi(\theta)$  is strictly convex. Therefore we may Legendre-transform

$$\phi(\eta) = \theta^i \eta_i - \psi(\theta)$$

with

$$\eta_i = \frac{\partial \psi(\theta)}{\partial \theta^i}.$$

The coordinates  $\eta$  are called **dual coordinates** (dual to  $\theta$ ).

- Example: Gaussian distributions

$$\theta^1 = \frac{\eta_1}{\eta_2 - \eta_1^2}, \quad \theta^2 = \frac{1}{2(\eta_1^2 - \eta_2)}$$

and

$$\phi(\eta) = -\frac{1}{2} \ln(2\pi e(\eta_2 - \eta_1^2)).$$

## Application: Brownian Bridge

A **Brownian bridge** (or **pinned Brownian motion**) describes Brownian motion that has a fixed start  $(x_a = 0, t_a = 0)$  and a fixed end  $(x_b, t_b)$ . Its probability distribution at some  $t \in (0, t_b)$  is given by a Gaussian distribution:

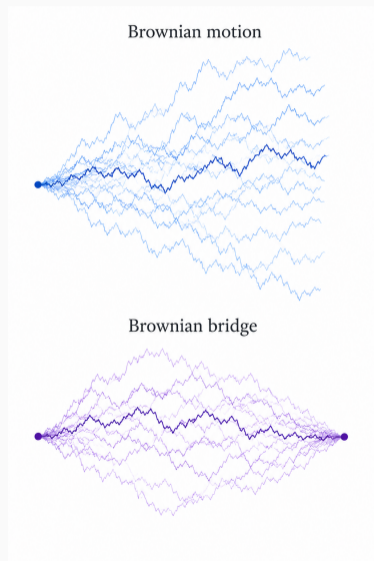
$$\mathcal{N}(x \mid \mu_{\text{BB}}(t), \sigma_{\text{BB}}^2(t)),$$

where

$$\mu_{\text{BB}}(t) = \frac{t}{t_b} x_b,$$

$$\sigma_{\text{BB}}^2(t) = 2\mathfrak{D} \frac{t(t_b - t)}{t_b}.$$

Here,  $\mathfrak{D}$  is the diffusion coefficient.



# Application: Brownian Bridge

**Idea:** Use information geometry in order to describe the Brownian bridge

At any  $t \in (0, t_b)$  the distribution of the Brownian Bridge is Gaussian. Consider a path within the Gaussian statistical manifold:

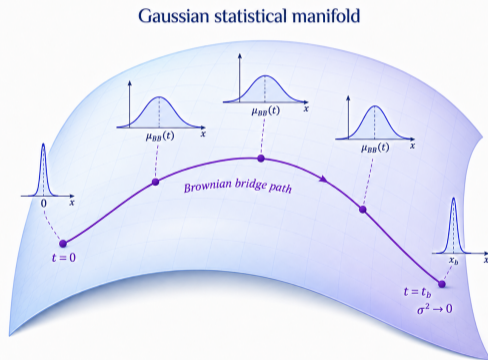
$$(\theta^1, \theta^2) = \left( \frac{\mu_{\text{geo}}(t)}{\sigma_{\text{geo}}^2(t)}, \frac{-1}{2\sigma_{\text{geo}}^2(t)} \right)$$

and

$$\psi(\theta) = -\frac{(\theta^1)^2}{4\theta^2} + \frac{1}{2} \ln(\pi) - \frac{1}{2} \ln(-\theta^2).$$

Parameters  $\theta^i$  are a bit awkward:  $\sigma_{\text{BB}} \rightarrow 0$  as  $t \rightarrow t_b$ , therefore  $\theta^i \rightarrow \infty$ .

Consider instead the dual coordinates.



## Application: Brownian Bridge

The dual coordinates are

$$\eta_1 = \frac{\partial\psi(\theta)}{\partial\theta^1} = -\frac{\theta^1}{2\theta^2} = \mu_{\text{geo}}(t)$$
$$\eta_2 = \frac{\partial\psi(\theta)}{\partial\theta^2} = \frac{(\theta^1)^2}{4(\theta^2)^2} - \frac{1}{2\theta^2} = \mu_{\text{geo}}^2(t) + \sigma_{\text{geo}}^2(t).$$

Dual connection (using KL divergence) in dual coordinates:

$$\Gamma^{*k}_{ij}(\eta) = 0.$$

Hence, geodesics in dual coordinates can be described by straight lines. Let us therefore make the ansatz

$$\eta_i(t) = c_i t + d_i.$$

## Application: Brownian Bridge

We already know the mean and variance of a Brownian Bridge at the start and endpoint, hence we can find the coefficients of the straight lines. Then:

$$\mu_{\text{geo}}(t) = \frac{t}{t_b} x_b,$$

$$\sigma_{\text{geo}}^2(t) = \frac{x_b^2}{t_b^2} t(t_b - t).$$

Comparing, we indeed find  $\mu_{\text{BB}}(t) = \mu_{\text{geo}}(t)$ . Furthermore  $\sigma_{\text{BB}}^2(t) = \sigma_{\text{geo}}^2(t)$  only if the diffusion coefficient is given by

$$\mathfrak{D} = \frac{1}{2} \frac{x_b^2}{t_b}.$$

Time evolution of a **canonical** Brownian bridge is identical to a geodesic on the associated statistical manifold!

## Analogies with Gravity:

- spacetime  $\leftrightarrow$  statistical manifold
- spacetime interval  $\leftrightarrow$  statistical distinguishability
- nonmetricity  $\leftrightarrow$  Amari-Chentsov tensor
- free fall along geodesics  $\leftrightarrow$  geodesic motion of stochastic processes?

## **(Bold) Conjecture**

Stochastic processes are described by geodesics on an appropriate statistical manifold

## Applications in AI:

- Optimization: Natural gradient descent
- Classification: Pattern recognition and clustering
- Generative AI: diffusion models

... see e.g. Amari, "Information Geometry and Its Applications";  
Karczewski et al. arXiv:2505.17517;

**Thank you!**