

EAI Math Reading Group

Information Geometry

$$g_{ij} = \int_{\Omega} (\partial_i \ln \varphi) (\partial_j \ln \varphi) \varphi \, d\omega$$

Motivation

Information Geometry = Differential Geometry applied to statistics

- Statistical Manifold : Manifold of (parametrized) distributions
- Fisher Information Metric : Riemannian metric

Applications :

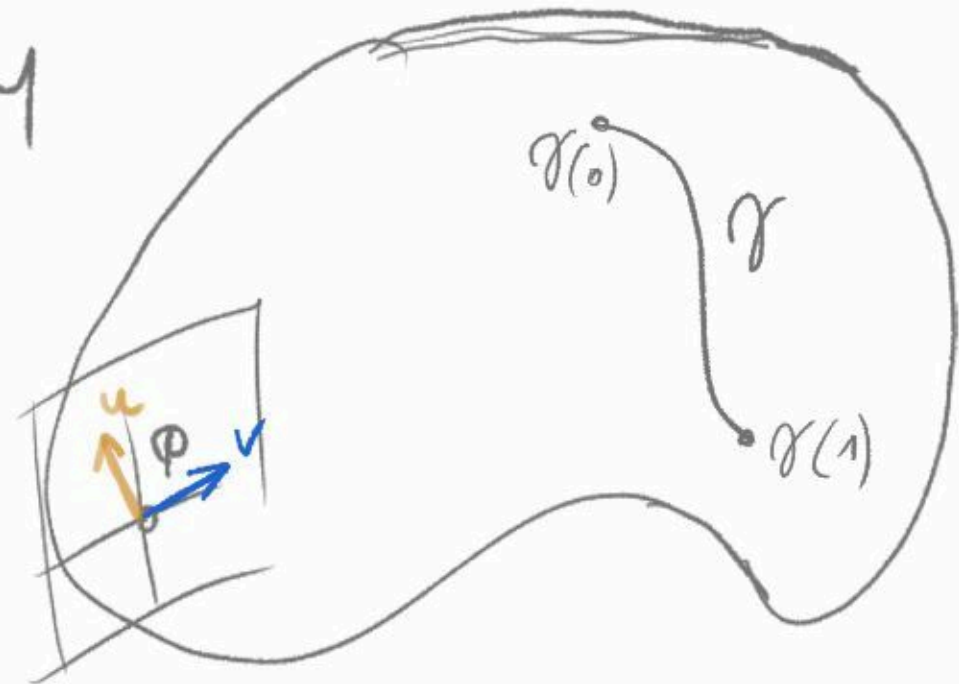
- Cramér-Rao bound
- "distances" between probability distributions
- "shortest-paths" between probability distributions

Outline

- 1) Recap on Differential Geometry
- 2) The Fisher Information Metric
- 3) Statistical Manifolds
- 4) Divergences

Differential Geometry Recap : Riemannian metric

Manifold M



• inner product

$$\langle u, v \rangle_\phi = u^i v^j \underbrace{g_{ij}}_{\text{metric}}$$

tangent
space

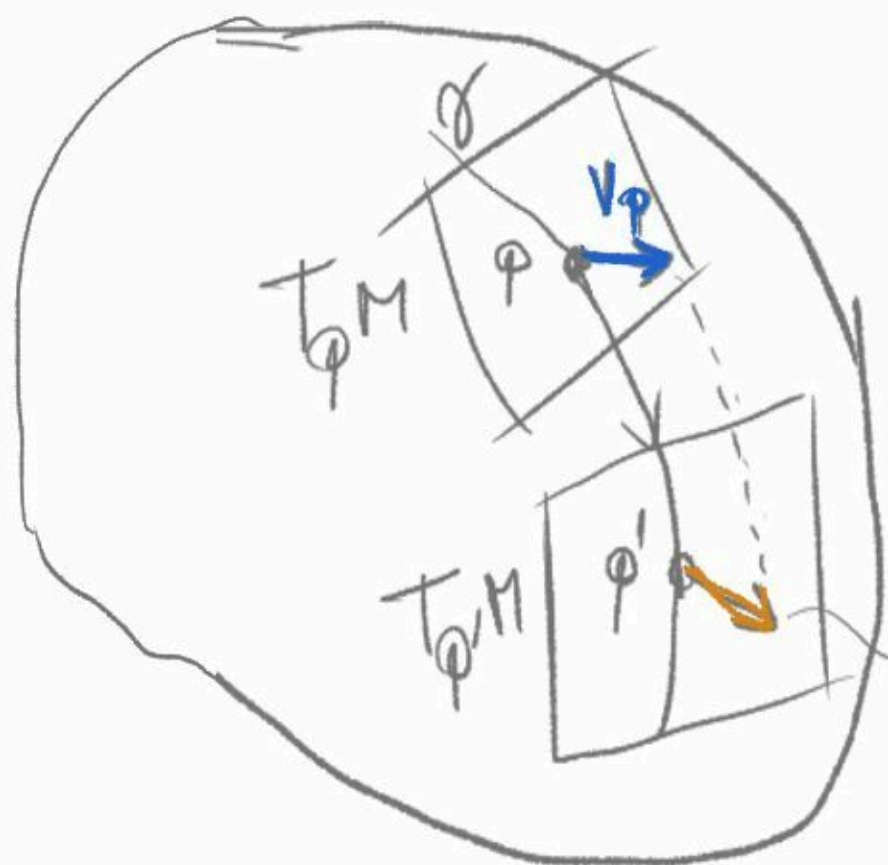
$T_\phi M$

Riemannian
manifold (M, g)

• length of a path

$$\| \gamma \| = \int_a^b \left\| \frac{d\gamma}{dt} \right\| dt = \int_a^b \sqrt{g_{ij} \dot{\gamma}^i \dot{\gamma}^j} dt$$

Recap: Affine Connection



$$v_{\phi'} \approx \bigtriangledown_{\gamma(0) \rightarrow \gamma(t)} v \in T_{\phi} M$$

$$\approx v_\phi - dv^i \Gamma_{ij}^k \partial_k$$

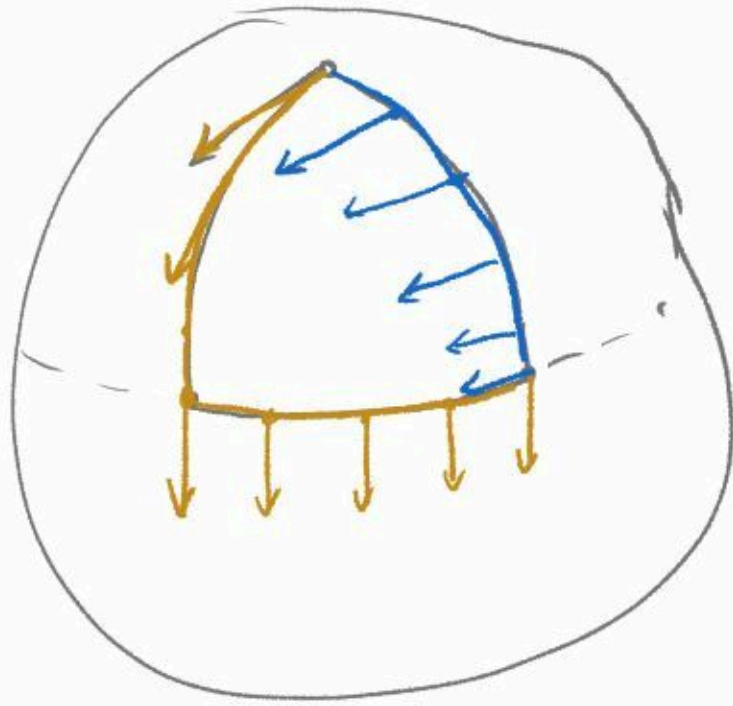
$$\bigtriangledown_{\gamma(0) \rightarrow \gamma(t)} ((\partial_j)_\phi) = (\partial_j)_{\phi'} - (d\theta_i (\Gamma_{ij}^k) (\partial_k)_\phi)$$

• Affine Connection
→ Covariant derivative

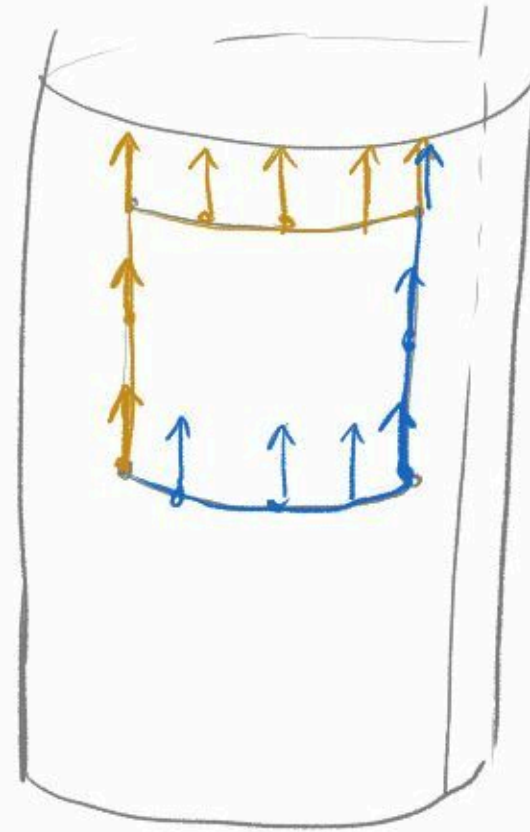
$$(\nabla_X Y)^k = X^i (\nabla_i Y)^k = X^i \left(\frac{\partial Y^k}{\partial x^i} + \Gamma_{ij}^k Y^j \right)$$

- parallel transport
- geodesics $\nabla_{\dot{\gamma}} \dot{\gamma} = 0$
- curvature, torsion...

Curvature [the hand wavy version]



constant positive
curvature



zero curvature

Metric Connection and Levi-Civita connection

a connection ∇ is compatible with the metric g if

$$Z\langle X, Y \rangle = \langle \nabla_Z X, Y \rangle + \langle X, \nabla_Z Y \rangle$$

• Parallel transport via a metric connection preserves inner products:

$$\left\langle \nabla_r^\nabla u, \nabla_r^\nabla v \right\rangle_{\varphi'} = \langle u, v \rangle_\varphi$$

• There is a unique symmetric (torsion-free) connection $\Gamma_{ij}^k = \Gamma_{ji}^k$ compatible with g

called the Levi-Civita connection ∇_{ij}^k

Fisher Information Matrix (FIM)

Let $\{P_\theta\}_{\theta \in \Theta}$ be a parametric family of distributions,

then the matrix

$$g_{ij}(\theta) = \mathbb{E}_\theta [\partial_i \log p_\theta \partial_j \log p_\theta] = \int_{\Omega} \partial_i \log p_\theta(x) \partial_j \log p_\theta(x) \underbrace{p_\theta(x)}_{p_\theta(x)} dx$$

with $\log p_\theta(x) = \log p_\theta(x)$, the log-likelihood,

is a symmetric definite positive matrix that can be used as a Riemannian metric

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

Example: (Normal distribution)

$$g_{ij}(\mu, \sigma) = \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{2\sigma^4} \end{pmatrix}$$

Application: Sufficient Statistics

For some function $Y = F(X)$ w/ $X \sim \phi(x, \theta)$ we can factor

$$\phi(x; \theta) = q(F(x); \theta) \underbrace{r(x; \theta)}$$

If $r(x; \theta) = r(x)$ \leadsto $Y = F(X)$ is a sufficient statistic for X

Theorem

Let $G(\theta)$ be the FIM of $S = \phi(x; \theta)$ and $G_F(\theta)$ the FIM of $S_F = q(y; \theta)$

then $G_F(\theta) \leq G(\theta)$, i.e. $\underbrace{G(\theta) - G_F(\theta)}$ is definite positive

If F is exhaustive $\Rightarrow G(\theta) - G_F(\theta) = 0$

Proof

$$G_F(\theta)_{ij} = \mathbb{E}_\theta \left[\partial_i \log q(y; \theta) \partial_j \log q(y; \theta) \right]$$

• Since $\phi(x; \theta) = q(F(x); \theta) r(x; \theta)$,

$$\partial_i \ell(x; \theta) = \partial_i \log q(F(x); \theta) + \partial_i \log r(x, \theta)$$

• $\partial_i \log q(y; \theta) = \mathbb{E}_\theta \left[\partial_i \ell(x; \theta) | Y \right] :$

$$\int_{\mathcal{B}} \partial_i \log q(y; \theta) q(y; \theta) dy = \int_{F^{-1}(\mathcal{B})} \partial_i \ell(x; \theta) \phi(x; \theta) dx$$

• $\Rightarrow \mathbb{E}_\theta \left[\partial_i \log r(x; \theta) | F(x) \right] = 0$

$\Rightarrow \partial_i \log r(x; \theta) \perp \phi(F(x))$ for the inner product

$$\langle \phi, \psi \rangle_\theta = \mathbb{E}_\theta [\phi(x) \psi(x)]$$

\Rightarrow The information loss is given by

$$\left(\Delta G(\theta) \right)_{ij} = \mathbb{E}_\theta [\partial_i \log r(x; \theta) \partial_j \log r(x; \theta)] = \mathbb{E}_\theta [\underbrace{\text{Cov}_\theta [\partial_i \ell, \partial_j \ell | Y]}_{\geq 0}]$$

$\Rightarrow G_F \leq G$ and $G_F = G$ if $\partial_i \log r(x, \theta) = 0$

Cramer-Rao Bound

$$\mathbb{E}[\hat{\theta}] = \theta$$

The variance of an unbiased estimator $\hat{\theta}$ of θ is at least

$$V_{\theta}(\hat{\theta}) \geq I(\theta)^{-1}$$

Proof

• $V_{\theta}(A) = \underbrace{\|dE[A]_{\theta}\|_{\theta}^2}_{\text{norm induced by the Fisher metric.}} \quad \text{where } A \text{ is some random variable}$

• For some submanifold $\mathcal{P} \subset \mathcal{P}_X$ $df \in T_{\theta}^*$

$$V_{\theta}[A] \geq \|dE[A]|_{\mathcal{P}}\|_{\theta}^2 \quad \text{so taking } A = \hat{\theta}$$

$$V_{\theta}[\hat{\theta}] \geq \sum_{i,j} (\partial_i \mathbb{E}[\hat{\theta}] \partial_j \mathbb{E}[\hat{\theta}]) G_{ij}^{-1}$$

"Fisher connection"

Let S be a parametric model. Define

$$(\Gamma_{ij}^k)_\theta = \mathbb{E}_\theta \left[\left(\partial_i \partial_j l_\theta + \frac{1}{2} \partial_i l_\theta \partial_j l_\theta \right) (\partial_k l_\theta) \right]$$

• Γ_{ij}^k are the Christoffel coefficients of the metric connection associated to the Fisher metric! (\Leftarrow)

• α -connections for $\alpha \in \mathbb{R}$ construct the connection ∇^α from

$$(\Gamma_{ij}^{\alpha k})_\theta = \mathbb{E}_\theta \left[\left(\partial_i \partial_j l_\theta + \frac{1-\alpha}{2} \partial_i l_\theta \partial_j l_\theta \right) (\partial_k l_\theta) \right]$$

• $\nabla^0 = \nabla$ (\Leftarrow)

Chentsov Theorem

For a model \mathcal{S} , $F(x)$ a sufficient statistic and \mathcal{S}_F the induced model, $\int_i \log p_\theta(x) = \int_i \log q_\theta(F(x)) \Rightarrow g_{ij}, \Gamma_{ij}^{(\alpha)k}$ coincide

Theorem: Chentsov (1972)

Suppose (g, ∇) invariant for sufficient statistics.

Then there exists $c \in \mathbb{R}$ and $\alpha \in \mathbb{R}$ s.t.
 $c > 0$

e.g. is the Fisher metric and

$$\nabla = \nabla^\alpha$$

Conjugate Connections

Let (M, g, ∇) be a manifold w/ connection ∇ , a connection ∇^* is conjugate to ∇ w/ respect to g if

$$X \langle Y, Z \rangle = \langle \nabla_X Y, Z \rangle + \langle Y, \nabla_X^* Z \rangle \quad \text{for } X, Y, Z \in \chi(M)$$

$$\Leftrightarrow X g(X, Z) = g(\nabla_X Y, Z) + g(Y, \nabla_X^* Z)$$

$$\bullet (\nabla^*)^* = \nabla$$

• Dual Parallel transport preserves the metric:

$$\left\langle \prod_{c(0) \rightarrow c(t)}^{\nabla} u, \prod_{c(0) \rightarrow c(t)}^{\nabla^*} v \right\rangle_{c(t)} = \langle u, v \rangle_{c(0)}$$

$\leadsto (M, g, \nabla, \nabla^*)$ is called a Conjugate Connection Manifold

N.B. $\bar{\nabla} = \frac{\nabla + \nabla^*}{2}$ recovers the Levi-Civita connection

Statistical Manifold

- Amani-Chentsov tensor :

$$C_{ijk} = \Gamma_{ij}^k - \Gamma_{ij}^{*k}$$

↳ totally symmetric $C_{ijk} = C_{\sigma(i)\sigma(j)\sigma(k)}$ for σ permutation
(0,3) tensor

9. Statistical manifold (M, g, C)

↳ totally symmetric

- α - ECMs: For any pair (∇, ∇^*) of conjugate connections, $\alpha \in \mathbb{R}$, let

$$\begin{cases} T_{ij}^{\alpha,k} = \Gamma_{ij}^{0k} - \frac{\alpha}{2} C_{ijk} \\ T_{ij}^{-\alpha,k} = \Gamma_{ij}^{0,k} + \frac{\alpha}{2} C_{ijk} \end{cases}$$

where $T_{ij}^{0k} = {}^{LC}T_{ij}^k$

$$\nabla \leq \nabla^1, \quad \nabla^* \leq \nabla^{-1}$$

$\leadsto (M, g, \nabla^{-\alpha}, \nabla^{\alpha})$ is a conjugate connection manifold

Fundamental Theorem of Information Theory

If a torsion-free (symmetric) connection has constant curvature K then so does its conjugate ∇^*

$$\Rightarrow (M, g, \nabla^{-\alpha}, \nabla^{\alpha}) \text{ is } \nabla^{\alpha}\text{-flat} \Leftrightarrow \nabla^{-\alpha}\text{-flat}$$
$$\nabla\text{-flat} \Leftrightarrow \nabla^*\text{-flat}$$

Why do we care?

$$\hookrightarrow \Gamma_{ij}^k = 0$$

- Geodesics are easy to compute in flat space (affine curves)
- Adjust α to get a flat space

Divergences

$$\mathbb{D}: M \times M \rightarrow [0, \infty) \quad , \quad \mathbb{D} \in C^3$$

1. $\mathbb{D}(\theta: \theta') \geq 0$ and $\mathbb{D}(\theta: \theta') = 0 \Leftrightarrow \theta = \theta'$
2. $\left. \begin{aligned} \partial_{i,\cdot} \mathbb{D}(\theta: \theta')|_{\theta=\theta'} &= 0, \partial_{\cdot,j} \mathbb{D}(\theta: \theta')|_{\theta=\theta'} = 0 \end{aligned} \right\}$
3. $-\partial_{\cdot,i} \partial_{\cdot,j} \mathbb{D}(\theta: \theta')|_{\theta'=\theta}$ is positive definite

Example: Kullback-Liebler divergence

$$\mathbb{D}(p \parallel q) = \int_{\Omega} \frac{p}{q} \ln \frac{p}{q} q \, d\omega$$

$$g_{ij}(\theta) = \int_{\Omega} (\partial_i \ln p) (\partial_j \ln p) p \, d\omega \quad \Bigg) \text{ Fisher Matrix!}$$

Divergence to manifold

given a divergence \mathbb{D} , let

$${}^{\mathbb{D}}g = -\mathcal{I}_{ij} \mathbb{D}(\theta; \theta')|_{\theta=\theta'}$$

$$\Gamma_{ijk} = -\mathcal{I}_{ij,k} \mathbb{D}(\theta; \theta')|_{\theta=\theta'}$$

$$\Gamma_{ijk}^* = -\mathcal{I}_{k,ij} \mathbb{D}(\theta; \theta')|_{\theta=\theta'}$$

then $(\eta, {}^{\mathbb{D}}g, {}^{\mathbb{D}}\nabla, {}^{\mathbb{D}}\nabla^*)$ is a ecm