

# The Information Geometry of Optimization

---

Gianluca Covini

Presentation for the exam of *Optimization*

## References

- **Reference for the main result:** G. Raskutti, S. Mukherjee, *The Information Geometry of Mirror Descent*, IEEE Transactions on Information Theory, 2015
- **Reference for general explanation of information geometry:** F. Nielsen, *An Elementary Introduction to Information Geometry*, Entropy, 2020

## Motivations

---

# Setting

## Setting

$$L : \Theta \rightarrow \mathbb{R}; \quad \text{find} \quad \theta^* = \arg \min_{\theta \in \Theta} L(\theta)$$

# Setting

## Setting

$$L : \Theta \rightarrow \mathbb{R}; \quad \text{find} \quad \theta^* = \arg \min_{\theta \in \Theta} L(\theta)$$

## Possible Approach

$$\text{Gradient Descent (GD):} \quad \theta_{t+1} = \theta_t - \alpha_t \nabla L(\theta_t)$$

# Setting

## Setting

$$L : \Theta \rightarrow \mathbb{R}; \quad \text{find} \quad \theta^* = \arg \min_{\theta \in \Theta} L(\theta)$$

## Possible Approach

$$\text{Gradient Descent (GD):} \quad \theta_{t+1} = \theta_t - \alpha_t \nabla L(\theta_t)$$

## Problem

GD is **coordinate-dependent**

$$\theta_{t+1} = \theta_t - \alpha_t \nabla L_\theta(\theta_t) \quad \xrightarrow[L_\eta(\eta) = L_\theta(\theta(\eta))]{\theta = \theta(\eta)} \quad \eta_{t+1} = \eta_t - \alpha_t \nabla L_\eta(\eta_t)$$

→ In general,  $\{\theta_t\}_t \neq \{\eta_t\}_t$ .

# Setting

## Setting

$$L : \Theta \rightarrow \mathbb{R}; \quad \text{find} \quad \theta^* = \arg \min_{\theta \in \Theta} L(\theta)$$

## Possible Approach

$$\text{Gradient Descent (GD):} \quad \theta_{t+1} = \theta_t - \alpha_t \nabla L(\theta_t)$$

## Problem

GD is **coordinate-dependent**

$$\theta_{t+1} = \theta_t - \alpha_t \nabla L_\theta(\theta_t) \quad \xrightarrow[L_\eta(\eta) = L_\theta(\theta(\eta))]{\theta = \theta(\eta)} \quad \eta_{t+1} = \eta_t - \alpha_t \nabla L_\eta(\eta_t)$$

- In general,  $\{\theta_t\}_t \neq \{\eta_t\}_t$ .
- GD does not take into account the **geometry** of the problem.

## Goal of the Presentation

Show a “geometry-wise” approach to optimization for a wide class of spaces of applicative interest.

## Goal of the Presentation

Show a “geometry-wise” approach to optimization for a wide class of spaces of applicative interest.

Which space?

find  $p^* = \arg \min_{p \in M} L(p)$  when  $M$  is an **information manifold**.

## Goal of the Presentation

Show a “geometry-wise” approach to optimization for a wide class of spaces of applicative interest.

Which space?

find  $p^* = \arg \min_{p \in M} L(p)$  when  $M$  is an **information manifold**.

- Riemannian manifolds with additional structure.
- Aim to geometrically represent how information passes from data to models.

# Motivating Example

## Example

$$M := \{p_\theta = \mathcal{N}(\theta) \mid \theta = (\Sigma^{-1}\mu, -\Sigma^{-1}/2) \in \Theta\}$$

with appropriate manifold structure is an **information manifold**.

# Motivating Example

## Example

$$M := \{p_\theta = \mathcal{N}(\theta) \mid \theta = (\Sigma^{-1}\mu, -\Sigma^{-1}/2) \in \Theta\}$$

with appropriate manifold structure is an **information manifold**.

## Optimization Problem

Given  $x$  data, we may want to find  $\theta^*$  the **maximum likelihood estimator**

find  $\theta^*$

solving  $\min_{\theta \in \Theta} \underbrace{\ell(\theta; x)}_{\text{negative log-likelihood}} = \min_{p_\theta \in M} -\log(p_\theta(x))$

# Table of contents

---

1. Motivations
2. Information Geometry
3. Optimization on Information Manifolds
4. Conclusion

# Information Geometry

---

## Manifold

A  **$D$ -dimensional manifold** is a topological space (locally) homeomorphic to an open set of  $\mathbb{R}^D$ .

→ We will consider global homeomorphism.

We can define a set of coordinates:

$$\begin{aligned}\theta : M &\rightarrow \Theta \subseteq \mathbb{R}^D \\ p &\mapsto \theta(p) = (\theta_1(p), \dots, \theta_D(p))\end{aligned}$$

where  $\theta$  is a homeomorphism (continuous bijection).

>>

We can define functions  $L : M \rightarrow \mathbb{R}$  and  $\frac{\partial L}{\partial \theta_i}(p) := \frac{\partial(L \circ \theta^{-1})}{\partial x_i}(\theta(p))$

## Tangent Space

$\forall p \in M$ , we associate a **tangent space**  $T_p M$ , which can be seen as the space of directional derivatives

$$T_p M := \{v : C^\infty(M) \rightarrow \mathbb{R} \mid v \text{ linear}; v(fg) = v(f)g + fv(g)\}$$

→  $T_p M$  is a  $D$ -dimensional vector space with natural basis corresponding to partial derivatives  $\mathcal{B} := \{e_i, \quad i = 1, \dots, D\}$ .

→ We can define a **vector field** as a function

$$X : p \in M \mapsto v \in T_p M$$

>>

# Riemannian Manifolds

---

We want to define an **inner product**  $g_p$  on  $T_p M$ .

We want to define an **inner product**  $g_p$  on  $T_p M$ .

Two examples of information manifolds we will see:

## Statistical Manifolds

In general, we can consider the manifolds of parametric families induced by the **Fisher information**  $M := \{p_\theta \mid \theta \in \Theta\}$ ,  $\mathcal{I}(\theta) = \left( \mathbb{E} \left[ \frac{\partial}{\partial \theta_i} \ell(\theta; x) \frac{\partial}{\partial \theta_j} \ell(\theta; x) \right] \right)$

$$\rightarrow g_p(u, v) = u^T \mathcal{I}(\theta(p)) v \quad \forall u, v \in T_p M$$

# Riemannian Manifolds

We want to define an **inner product**  $g_p$  on  $T_p M$ .

Two examples of information manifolds we will see:

## Statistical Manifolds

In general, we can consider the manifolds of parametric families induced by the **Fisher information**  $M := \{p_\theta \mid \theta \in \Theta\}$ ,  $\mathcal{I}(\theta) = \left( \mathbb{E} \left[ \frac{\partial}{\partial \theta_i} \ell(\theta; x) \frac{\partial}{\partial \theta_j} \ell(\theta; x) \right] \right)$

$$\rightarrow g_p(u, v) = u^T \mathcal{I}(\theta(p)) v \quad \forall u, v \in T_p M$$

## Bregman Manifolds

The inner product can also be induced by **Bregman divergence**.

$$F : \Theta \rightarrow \mathbb{R} \quad \text{mirror map} \quad \rightsquigarrow \quad B_F(\theta \mid \theta') := F(\theta) - F(\theta') - (\theta - \theta')^T \nabla F(\theta')$$

$$\rightarrow g_p(u, v) = u^T \nabla^2 F(\theta(p)) v \quad \forall u, v \in T_p M$$

$\rightarrow$  Equivalent for exponential parametric families (e.g. Gaussians).

# Riemannian Manifolds

We want to define an **inner product**  $g_p$  on  $T_p M$ .

## Statistical Manifolds

In general, we can consider the manifolds of parametric families induced by the **Fisher information**  $M := \{p_\theta \mid \theta \in \Theta\}$ ,  $\mathcal{I}(\theta) = \left( \mathbb{E} \left[ \frac{\partial}{\partial \theta_i} \ell(\theta; x) \frac{\partial}{\partial \theta_j} \ell(\theta; x) \right] \right)$

$$\rightarrow g_p(u, v) = u^T \mathcal{I}(\theta(p)) v \quad \forall u, v \in T_p M$$

## Bregman Manifolds

The inner product can also be induced by **Bregman divergence**.

$$F : \Theta \rightarrow \mathbb{R} \quad \text{mirror map} \quad \rightsquigarrow \quad B_F(\theta \mid \theta') := F(\theta) - F(\theta') - (\theta - \theta')^T \nabla F(\theta')$$

$$\rightarrow g_p(u, v) = u^T \nabla^2 F(\theta(p)) v \quad \forall u, v \in T_p M$$

## Riemannian Manifolds

$(M, g)$  as defined is a Riemannian manifold.

**Note:** from now on, we will consider information manifolds as induced by Bregman divergence.

To define an information manifold we also need more structure

## Affine Connection

We define the **affine connection** a  $\nabla : (X, Y) \mapsto \nabla_X Y$  vector field.

→ For Bregman manifolds,  $\nabla_{e_i} e_j = 0$  for every  $e_i, e_j \in \mathcal{B}$  natural basis of  $T_p M$ .

**Note:** from now on, we will consider information manifolds as induced by Bregman divergence.

To define an information manifold we also need more structure

## Affine Connection

We define the **affine connection** a  $\nabla : (X, Y) \mapsto \nabla_X Y$  vector field.

→ For Bregman manifolds,  $\nabla_{e_i} e_j = 0$  for every  $e_i, e_j \in \mathcal{B}$  natural basis of  $T_p M$ .

We also need a **dual structure**,

$$\nabla^* : (X, Y) \mapsto \nabla_X^* Y$$

## Three ways to see duality

$$\nabla \longrightarrow \nabla^* \text{ s.t. } X\langle Y, Z \rangle = \langle \nabla_X Y, Z \rangle + \langle Y, \nabla_X^* Z \rangle$$

$\nabla$  induced by  $B_F$   $\longrightarrow$   $\nabla^*$  induced by  $B_F^*(\theta \mid \theta') := B_F(\theta' \mid \theta)$

$\nabla$  induced by  $B_F$   $\longrightarrow$   $\nabla^*$  induced by  $B_{F^*}$

$$\text{where } F^*(\eta) := \sup_{\theta \in \Theta} \{\theta^T \eta - F(\theta)\}$$

## Three ways to see duality

$$\nabla \longrightarrow \nabla^* \text{ s.t. } X\langle Y, Z \rangle = \langle \nabla_X Y, Z \rangle + \langle Y, \nabla_X^* Z \rangle$$

$\nabla$  induced by  $B_F$   $\longrightarrow$   $\nabla^*$  induced by  $B_F^*(\theta \mid \theta') := B_F(\theta' \mid \theta)$

$\nabla$  induced by  $B_F$   $\longrightarrow$   $\nabla^*$  induced by  $B_{F^*}$

$$\text{where } F^*(\eta) := \sup_{\theta \in \Theta} \{\theta^T \eta - F(\theta)\}$$

## Information manifold

$(M, g_P, \nabla, \nabla^*)$  induced by  $B_F$  is an **information manifold**.

## Three ways to see duality

$$\nabla \longrightarrow \nabla^* \text{ s.t. } X\langle Y, Z \rangle = \langle \nabla_X Y, Z \rangle + \langle Y, \nabla_X^* Z \rangle$$

$\nabla$  induced by  $B_F$   $\longrightarrow$   $\nabla^*$  induced by  $B_F^*(\theta \mid \theta') := B_F(\theta' \mid \theta)$

$\nabla$  induced by  $B_F$   $\longrightarrow$   $\nabla^*$  induced by  $B_{F^*}$

$$\text{where } F^*(\eta) := \sup_{\theta \in \Theta} \{\theta^T \eta - F(\theta)\}$$

## Information manifold

$(M, g_P, \nabla, \nabla^*)$  induced by  $B_F$  is an **information manifold**.

What does duality imply? It defines two sets of **coordinates**

$$\theta = \nabla F^*(\eta) \longleftrightarrow \eta = \nabla F(\theta)$$

# Optimization on Information Manifolds

---

## Riemannian Gradient Descent

Given  $(M, g)$  Riemannian manifold.

- **Generalization of GD** in the sense of optimizing by moving in the direction of steepest descent;
- We define the map  $\exp_p(v)$  which gives the arrival point after a unit of time of the shortest curve starting from  $p$  with velocity  $v$ .

$$\text{RGD: } p_{t+1} = \exp_{p_t}(-\alpha_t \nabla_M L(p_t))$$

>>

## Problem

$\exp_p(v)$  is computationally intractable.

## Problem

$\exp_p(v)$  is computationally intractable.

On the **Bregman manifold**  $(M, F)$ ,

## Natural Gradient Descent (NGD)

We can replace  $\exp_p(v)$  with its first-order Taylor approximation  $\exp_p(v) \approx p + v$ .

$$\text{NGD: } \theta_{t+1} = \theta_t - \alpha_t \underbrace{(\nabla_{\theta}^2 F(\theta_t))^{-1} \nabla_{\theta} (L_{\theta}(\theta_t))}_{\nabla^{(NG)} L_{\theta}(\theta_t) \text{ natural gradient}}$$

# NGD and MD

## Problem

$\exp_p(v)$  is computationally intractable.

On the **Bregman manifold**  $(M, F)$ ,

## Natural Gradient Descent (NGD)

We can replace  $\exp_p(v)$  with its first-order Taylor approximation  $\exp_p(v) \approx p + v$ .

$$\text{NGD: } \theta_{t+1} = \theta_t - \alpha_t \underbrace{(\nabla_\theta^2 F(\theta_t))^{-1} \nabla_\theta(L_\theta(\theta_t))}_{\nabla^{(NG)} L_\theta(\theta_t) \text{ natural gradient}}$$

## Mirror Descent (MD)

$$\text{MD: } \theta_{t+1} = \arg \min_{\theta \in \Theta} \left\{ \theta^T \nabla_\theta L_\theta(\theta_t) + \frac{1}{\alpha_t} B_F(\theta \mid \theta_t) \right\}$$

# Equivalence Result

## Theorem [Raskutti, Mukherjee]

Given an information manifold  $(M, g, \nabla, \nabla^*)$  induced by a Bregman divergence  $B_F$ , **MD on  $(M, F)$  is equivalent to NGD in the dual space  $(M, F^*)$ .**

# Equivalence Result

## Theorem [Raskutti, Mukherjee]

Given an information manifold  $(M, g, \nabla, \nabla^*)$  induced by a Bregman divergence  $B_F$ , **MD on  $(M, F)$  is equivalent to NGD in the dual space  $(M, F^*)$ .**

## Proof

$$\text{MD: } \theta_{t+1} = \arg \min_{\theta \in \Theta} \left\{ \theta^T \nabla L_\theta(\theta_t) + \frac{1}{\alpha_t} B_F(\theta \mid \theta_t) \right\}$$

# Equivalence Result

## Theorem [Raskutti, Mukherjee]

Given an information manifold  $(M, g, \nabla, \nabla^*)$  induced by a Bregman divergence  $B_F$ , **MD on  $(M, F)$  is equivalent to NGD in the dual space  $(M, F^*)$ .**

## Proof

$$\text{MD: } \theta_{t+1} = \arg \min_{\theta \in \Theta} \left\{ \theta^T \nabla L_\theta(\theta_t) + \frac{1}{\alpha_t} B_F(\theta \mid \theta_t) \right\}$$

Finding the minimum by differentiation yields the step:

$$\nabla F(\theta_{t+1}) = \nabla F(\theta_t) - \alpha_t \nabla_\theta L(\theta_t)$$

# Equivalence Result

## Theorem [Raskutti, Mukherjee]

Given an information manifold  $(M, g, \nabla, \nabla^*)$  induced by a Bregman divergence  $B_F$ , **MD on  $(M, F)$  is equivalent to NGD in the dual space  $(M, F^*)$ .**

## Proof

$$\text{MD: } \theta_{t+1} = \arg \min_{\theta \in \Theta} \left\{ \theta^T \nabla L_\theta(\theta_t) + \frac{1}{\alpha_t} B_F(\theta \mid \theta_t) \right\}$$

Finding the minimum by differentiation yields the step:

$$\nabla F(\theta_{t+1}) = \nabla F(\theta_t) - \alpha_t \nabla_\theta L(\theta_t)$$

Dual change of variable:  $\eta = \nabla F(\theta)$ ,  $\theta = \nabla F^*(\eta)$ ,

$$\eta_{t+1} = \eta_t - \alpha_t \nabla_\theta L(\nabla F^*(\eta_t))$$

# Equivalence Result

## Proof

$$\text{MD: } \theta_{t+1} = \arg \min_{\theta \in \Theta} \left\{ \theta^T \nabla L_\theta(\theta_t) + \frac{1}{\alpha_t} B_F(\theta \mid \theta_t) \right\}$$

Finding the minimum by differentiation yields the step:

$$\nabla F(\theta_{t+1}) = \nabla F(\theta_t) - \alpha_t \nabla_\theta L(\theta_t)$$

Dual change of variable:  $\eta = \nabla F(\theta)$ ,  $\theta = \nabla F^*(\eta)$ ,

$$\eta_{t+1} = \eta_t - \alpha_t \nabla_\theta L(\nabla F^*(\eta_t))$$

Chain rule:  $\nabla_\eta L(\nabla F^*(\eta)) = \nabla_\eta^2 F^*(\eta) \nabla_\theta L(\nabla F^*(\eta))$

# Equivalence Result

## Proof

$$\text{MD: } \theta_{t+1} = \arg \min_{\theta \in \Theta} \left\{ \theta^T \nabla L_\theta(\theta_t) + \frac{1}{\alpha_t} B_F(\theta \mid \theta_t) \right\}$$

Finding the minimum by differentiation yields the step:

$$\nabla F(\theta_{t+1}) = \nabla F(\theta_t) - \alpha_t \nabla_\theta L(\theta_t)$$

Dual change of variable:  $\eta = \nabla F(\theta)$ ,  $\theta = \nabla F^*(\eta)$ ,

$$\eta_{t+1} = \eta_t - \alpha_t \nabla_\theta L(\nabla F^*(\eta_t))$$

Chain rule:  $\nabla_\eta L(\nabla F^*(\eta)) = \nabla_\eta^2 F^*(\eta) \nabla_\theta L(\nabla F^*(\eta))$

Therefore,  $\eta_{t+1} = \eta_t - \alpha_t (\nabla^2 F^*(\eta_t))^{-1} \nabla_\eta L(\nabla F^*(\eta_t))$

which corresponds to the natural gradient descent step. □

# Application

Back to our motivating example,

$$M := \{p_\theta = \mathcal{N}(\theta) \mid \theta = (\Sigma^{-1}\mu, -\Sigma^{-1}/2) \in \Theta\} \text{ with } F(\theta) := \frac{1}{2}\|\theta\|_2^2$$

We want to find the MLE  $\theta^*$  given  $x$  data:

$$\text{find } \theta^* = \arg \min_{\theta \in \Theta} \ell(\theta; x) = \arg \min_{p_\theta \in M} -\log(p_\theta(x))$$

## Application

- NGD moves in the direction of steepest descent of  $\ell$  and asymptotically achieves the minimum possible asymptotic variance (CR bound) **but** it requires  $\nabla^2 F$ ;
- for MD we don't have guarantees of moving in direction of steepest descent and of achieving asymptotical CR bound, **but** it is a 1-order method.

→ The equivalence result guarantees that we have a 1-order method that achieves CR bound.

## Conclusion

---

# Conclusion

## Takeaways

- Optimization is strongly influenced by space geometry;

## Takeaways

- Optimization is strongly influenced by space geometry;
- we saw information manifolds induced by parametric families and by Bregman divergence;

## Takeaways

- Optimization is strongly influenced by space geometry;
- we saw information manifolds induced by parametric families and by Bregman divergence;
- on information manifolds optimization can be performed through NGD and MD;

## Takeaways

- Optimization is strongly influenced by space geometry;
- we saw information manifolds induced by parametric families and by Bregman divergence;
- on information manifolds optimization can be performed through NGD and MD;
- the two are equivalent on Bregman manifolds.

## Takeaways

- Optimization is strongly influenced by space geometry;
- we saw information manifolds induced by parametric families and by Bregman divergence;
- on information manifolds optimization can be performed through NGD and MD;
- the two are equivalent on Bregman manifolds.

**Thank you for your attention!**