

An Invitation to Statistics in Wasserstein Space

Shen Ziliang 2021213218

Shanghai University of Finance and Economics

February 26, 2022

- 1 Geodesics, the Log Map and the Exponential Map
- 2 Curvature and Compatibility of Measures

- 1 Geodesics, the Log Map and the Exponential Map
- 2 Curvature and Compatibility of Measures

Although the Wasserstein space $\mathcal{W}_p(\mathcal{X})$ is non-linear in terms of measures, it is linear in terms of maps. Indeed, if $\mu \in \mathcal{W}_p(\mathcal{X})$ and $T_i: \mathcal{X} \rightarrow \mathcal{X}$ are such that $\|T_i\| \rightarrow L_p(\mu)$ then $(\alpha T_1 + \beta T_2)\#\mu \in \mathcal{W}_p(\mathcal{X})$ for all $\alpha, \beta \in \mathbb{R}$. We assume here that \mathcal{X} is a Hilbert space and that $p = 2$; the results extend to any $p > 1$. Absolutely continuous measures are assumed to be so with respect to Lebesgue measure if $\mathcal{X} = \mathbb{R}^d$ and otherwise refer to Definition 1.6.4.

Definition 1.6.4: Gaussian Null Sets and Absolutely Continuous Measures

A subset $\mathcal{N} \subseteq \mathcal{X}$ is a Gaussian null set if whenever ν is a nondegenerate Gaussian measure, $\nu(\mathcal{N}) = 0$. A probability measure $\mu \in P(\mathcal{X})$ is absolutely continuous if μ vanishes on all Gaussian null sets.

Definition: Wasserstein space

Let X be a separable Banach space. The p -Wasserstein space on X is defined as:

$$\mathcal{W}_p(\mathcal{X}) = \left\{ \mu \in P(\mathcal{X}) : \int_{\mathcal{X}} \|x\|^p d\mu(x) < \infty \right\}, p \geq 1$$

We will sometimes abbreviate and write simply \mathcal{W}_p instead of $\mathcal{W}_p(\mathcal{X})$.

Geodesics in $\mathcal{W}_2(\mathcal{X})$

Let $\gamma \in \mathcal{W}_2(\mathcal{X})$ be absolutely continuous and $\mu \in \mathcal{W}_2(\mathcal{X})$ arbitrary. From Sect.1.6.1 (Theorem 1.6.2), we know that there exists a unique solution to the Monge–Kantorovich problem, and that solution is given by a transport map that we denote by t_γ^μ . Recalling that $\mathbf{i} : \mathcal{X} \rightarrow \mathcal{X}$ is the identity map, we can define a curve:

$$\gamma_t = [\mathbf{i} + t(t_\gamma^\mu - \mathbf{i})] \# \mu, t \in [0, 1]$$

This curve is known as McCann's [93] interpolant. As hinted in the introduction to this section, it is constructed via classical linear interpolation of the transport maps t_γ^μ and the identity. Clearly $\gamma_1 = \gamma, \gamma_0 = \mu$.

Definition: Constant speed geodesics

A curve $\gamma : [0, 1] \rightarrow \mathcal{S}$ is a (constant speed) geodesic if

$$d(\gamma_t, \gamma_s) = d(\gamma_0, \gamma_1)(t - s), \forall 0 \leq s \leq t \leq 1$$

Geodesics in $\mathcal{W}_2(\mathcal{X})$

McCann's interpolant is a constant-speed geodesic in $\mathcal{W}_2(\mathcal{X})$, that is

$$dW_2(\gamma_t, \gamma_s) = W_2(\gamma_0, \gamma_1)(t - s), \forall 0 \leq s \leq t \leq 1$$

Tangent space of $\mathcal{W}_2(\mathcal{X})$ at μ

$$Tan_\mu = \overline{\{t(\mathbf{t} - \mathbf{i}) : \mathbf{t} = t_\mu^\nu, \exists \nu \in \mathcal{W}_2(\mathcal{X}), t > 0\}}^{L^2(\mu)}$$

Strictly speaking, Tan_μ is a subset of the space of functions $f: \mathcal{X} \rightarrow \mathcal{X}$ such that $\|f\| \in L^2(\mu)$ rather than $L^2(\mu)$ itself, as in Definition, but we will write $Tan_\mu \subseteq L^2(\mu)$ for simplicity.

The Space $\mathcal{L}_p(\mu)$

Let \mathcal{X} be a Banach space and μ a Borel measure on \mathcal{X} . Then the space $\mathcal{L}_p(\mu)$ is the space of measurable functions $f: \mathcal{X} \rightarrow \mathcal{X}$ such that

$$\|f\|_{\mathcal{L}_p(\mu)} = \left(\int_{\mathcal{X}} \|f(x)\|_{\mathcal{X}}^p d\mu(x) \right)^{1/p} < \infty$$

The Tangent Bundle

Although not obvious from the definition, this is a linear space. The reason is that, in \mathbb{R}^d , Lipschitz functions are dense in $L_2(\mu)$, and for t Lipschitz the negative of a tangent element

$$-t(\mathbf{t} - \mathbf{i}) = s(\mathbf{s} - \mathbf{i}), s > t\|\mathbf{t}\|_{Lip}, \mathbf{s} = \mathbf{i} + \frac{t}{s}(\mathbf{i} - \mathbf{t})$$

lies in the tangent space, since \mathbf{s} can be seen to belong to the subgradient of a convex function by definition of \mathbf{s} . This also shows that Tan_μ can be seen to be the $L_2(\mu)$ -closure of all gradients of C_c^∞ functions. The tangent space definition is valid for arbitrary measures in $\mathcal{W}_2(\mathcal{X})$. The exponential map at $\gamma \in \mathcal{W}_2(\mathcal{X})$ is the restriction to Tan_γ of the mapping that sends $\mathbf{r} \in L_2(\gamma)$ to $[\mathbf{r} + \mathbf{i}] \# \gamma \in \mathcal{W}_2(\mathcal{X})$.

The Log Map and The Exponential Map

The Exponential Map

$\exp_\gamma : \text{Tan}_\gamma \rightarrow \mathcal{W}_2(\mathcal{X})$ takes the form

$$\exp_\gamma(t(\mathbf{t} - \mathbf{i})) = \exp_\gamma([\mathbf{t}\mathbf{t} + (1 - t)\mathbf{i}] - \mathbf{i}) = [\mathbf{t}\mathbf{t} + (1 - t)\mathbf{i}] \#_\gamma$$

Thus, when γ is absolutely continuous, $\exp \exp_\gamma$ is surjective, as can be seen from its right inverse, the log map

The Log Map

$\log_\gamma : \mathcal{W}_2(\mathcal{X}) \rightarrow \text{Tan}_\gamma$ takes the form

$$\log_\gamma(\mu) = \mathbf{t}_\gamma^\mu - \mathbf{i}$$

because convex combinations of optimal maps are optimal maps as well. In particular, McCann's interpolant $[\mathbf{t}\mathbf{t} + (1 - t)\mathbf{i}] \#_\gamma$ is mapped bijectively to the line segment $t(\mathbf{t}_\gamma^\mu - \mathbf{i}) \in \text{Tan}_\gamma$ through the log map.

Other Definition McCann's interpolant

McCann's interpolant can also be defined as

$$[tp_2 + (1 - t)p_1] \# \pi, p_1(x, y) = x, p_2(x, y) = y$$

where $p_1, p_2 : \mathcal{X}^2 \rightarrow \mathcal{X}$ are projections and π is any optimal transport plan between γ and μ . This is defined for arbitrary measures $\gamma, \mu \in \mathcal{W}_2(\mathcal{X})$, and reduces to the previous definition if γ is absolutely continuous.

- 1 Geodesics, the Log Map and the Exponential Map
- 2 Curvature and Compatibility of Measures

Curvature

Let $\gamma, \mu, \nu \in \mathcal{W}_2(\mathcal{X})$ be absolutely continuous measures.

$$W_2^2(\mu, \nu) \leq \int_{\mathcal{X}} \|t_\gamma^\mu(x) - t_\gamma^\nu(x)\|^2 d\gamma(x) = \|\log_\gamma(\mu) - \log_\gamma(\nu)\|^2$$

In other words, the distance between μ and ν is smaller in $\mathcal{W}_2(\mathcal{X})$ than the distance between the corresponding vectors $\log_\gamma(\mu)$ and $\log_\gamma(\nu)$ in the tangent space Tan_γ . In the terminology of differential geometry, this means that the Wasserstein space has nonnegative sectional curvature at any absolutely continuous γ . It is instructive to see when equality holds. As $t_\nu^\gamma = (t_\gamma^\nu)^{-1}$, a change of variables gives

$$W_2^2(\mu, \nu) \leq \int_{\mathcal{X}} \|t_\gamma^\mu(t_\nu^\gamma(x)) - x\|^2 d\nu(x)$$

Since the map $t_\gamma^\mu \circ t_\nu^\gamma$ pushes forward ν to μ , equality holds if and only if $t_\gamma^\mu \circ t_\nu^\gamma = t_\nu^\mu$. This motivates the following definition.

Compatibility

Definition :Compatible Measures

A collection of absolutely continuous measures $\mathcal{C} \subseteq \mathcal{W}_2(\mathcal{X})$ is compatible if for all $\gamma, \mu, \nu \in \mathcal{C}$ we have $t_\gamma^\mu \circ t_\nu^\gamma = t_\nu^\mu$ (in $L_2(\nu)$).

Remark

The absolute continuity is not necessary and was introduced for notational simplicity. A more general definition that applies to general measures is the following: every finite subcollection of \mathcal{C} admits an optimal multicoupling whose relevant projections are simultaneously pairwise optimal.

A collection of two (absolutely continuous) measures is always compatible.

Example: Absolutely continuous measures in \mathbb{R}

if $\mathcal{X} = \mathbb{R}$, then the entire collection of absolutely continuous (or even just continuous) measures is compatible. This is because of the simple geometry of convex functions in \mathbb{R} : gradients of convex functions are nondecreasing, and this property is stable under composition. In a more probabilistic way of thinking, one can always push-forward μ to ν via the uniform distribution $Leb|_{[0,1]}$ (see Sect.1.5). Letting F_μ^{-1} and F_ν^{-1} denote the quantile functions, we have seen that

$$W_2^2(\mu, \nu) = \|F_\mu^{-1} - F_\nu^{-1}\|_{L_2(0,1)}$$

In other words, $\mu \rightarrow F_\mu^{-1}$ is an isometry from $\mathscr{W}_2(\mathbb{R})$ to the subset of $L_2(0,1)$ formed by (equivalence classes of) left-continuous nondecreasing functions on $(0,1)$.

If $\gamma = \text{Leb}|[0, 1]$, then $F_\mu^{-1} = t_\gamma^\mu$ for all μ , so we can write the above equality as

$$W_2^2(\mu, \nu) = \|F_\mu^{-1} - F_\nu^{-1}\|_{L_2(0,1)} = \|\log_\gamma(\mu) - \log_\gamma(\nu)\|^2$$

so that if $\mathcal{X} = \mathbb{R}$, the Wasserstein space is essentially flat (has zero sectional curvature). The importance of compatibility can be seen as mimicking the simple one-dimensional case in terms of a Hilbert space embedding. Let $\mathcal{C} \subseteq \mathcal{W}_2(\mathcal{X})$ be compatible and fix $\gamma \in \mathcal{C}$. Then for all $\mu, \nu \in \mathcal{C}$

$$W_2^2(\mu, \nu) = \int_{\mathcal{X}} \|t_\gamma^\mu(x) - t_\gamma^\nu(x)\|^2 d\gamma(x) = \|\log_\gamma(\mu) - \log_\gamma(\nu)\|_{L^2(\gamma)}^2$$

Consequently, once again, $\mu \rightarrow t_\gamma^\mu$ is an isometric embedding of \mathcal{C} into $L_2(\gamma)$.

Example: Gaussian compatible measures

The Gaussian case presented in Sect.1.6.3 is helpful in shedding light on the structure imposed by the compatibility condition. Let $\gamma \in \mathcal{W}_2(\mathbb{R}^d)$ be a standard Gaussian distribution with identity covariance matrix. Let Σ_μ denote the covariance matrix of a measure $\mu \in \mathcal{W}_2(\mathbb{R}^d)$. When μ and ν are centred nondegenerate Gaussian measures,

$$t_\gamma^\mu = \Sigma_\mu^{1/2}; \quad t_\gamma^\nu = \Sigma_\nu^{1/2}; \quad t_\mu^\nu = \Sigma_\mu^{-1/2} [\Sigma_\mu^{1/2} \Sigma_\nu \Sigma_\mu^{1/2}]^{1/2} \Sigma_\mu^{-1/2}$$

so that $\gamma, \mu,$ and ν are compatible if and only if

$$t_\mu^\nu = t_\gamma^\nu \circ t_\mu^\gamma = \Sigma_\nu^{1/2} \Sigma_\mu^{-1/2}$$

equivalently, Σ_μ and Σ_ν commute.

We see that a collection \mathcal{C} of Gaussian measures on \mathbb{R}^d that includes the standard Gaussian distribution is compatible if and only if all the covariance matrices of the measures in \mathcal{C} are *simultaneously diagonalisable*.

Gaussian compatible measures

there exists an orthogonal matrix U such that $D_\mu = U\Sigma_\mu U^T$ is diagonal for all $\mu \in \mathcal{C}$. In that case, formula

$$\begin{aligned} W_2^2(\mu, \nu) &= \text{tr}(\Sigma_\mu + \Sigma_\nu - 2(\Sigma_\mu^{1/2}\Sigma_\nu\Sigma_\mu^{1/2})^{1/2}) = \text{tr}(\Sigma_\mu + \Sigma_\nu - 2\Sigma_\mu^{1/2}\Sigma_\nu^{1/2}) \\ &= \text{tr}(D_\mu + D_\nu - 2D_\mu^{1/2}D_\nu^{1/2}) = \sum_{i=1}^d (\sqrt{\alpha_i} - \sqrt{\beta_i})^2, \alpha_i = [D_\mu]_{ii}, \beta_i = [D_\nu]_{ii} \end{aligned}$$

and identifying the (nonnegative) number $a \in \mathbb{R}$, with the map $x \rightarrow ax$ on \mathbb{R} , the optimal maps take the “orthogonal separable” form

$$t_\mu^\nu = \Sigma_\nu^{1/2}\Sigma_\mu^{-1/2} = UD_\nu^{1/2}D_\mu^{1/2}U^T = U \circ (\sqrt{\beta_1/\alpha_1}, \dots, \sqrt{\beta_d/\alpha_d}) \circ U^T$$

In other words, up to an orthogonal change of coordinates, the optimal maps take the form of d nondecreasing real-valued functions. This is yet another crystallisation of the one-dimensional-like structure of compatible measures.

With the intuition of the Gaussian case at our disposal, we can discuss a more general case. Suppose that the optimal maps are continuously differentiable. Then differentiating the equation $t'_\mu = t'_\gamma \circ t^\gamma_\mu$ gives

$$\nabla t'_\mu(x) = \nabla t'_\gamma(t^\gamma_\mu(x)) \nabla t^\gamma_\mu(x)$$

Since optimal maps are gradients of convex functions, their derivatives must be symmetric and positive semidefinite matrices. A product of such matrices stays symmetric if and only if they commute, so in this differentiable setting, compatibility is equivalent to commutativity of the matrices $\nabla t'_\gamma(t^\gamma_\mu(x))$ and $\nabla t^\gamma_\mu(x)$ for μ -almost all x . In the Gaussian case, the optimal maps are linear functions, so x does not appear in the matrices.

Radial transformations

Consider the collection of functions $\mathbf{t} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ of the form $\mathbf{t}(x) = xG(\|x\|)$ with $G : \mathbb{R}_+ \rightarrow \mathbb{R}$ differentiable. Then a straightforward calculation shows that

$$\nabla \mathbf{t}(x) = G(\|x\|)\mathbf{I} + [G'(\|x\|)/\|x\|]xx^T$$

Since both \mathbf{I} and xx^T are positive semidefinite, the above matrix is so if both G and G' are nonnegative. If $\mathbf{s}(x) = xH(\|x\|)$ is a function of the same form, then $\mathbf{s}(\mathbf{t}(x)) = xG(\|x\|)H(\|x\|)$ which belongs to that family of functions (since G is nonnegative). Clearly

$$\nabla \mathbf{s}(\mathbf{t}(x)) = H(\|x\|G(\|x\|))\mathbf{I} + [G(\|x\|)H'(\|x\|)/\|x\|]xx^T$$

commutes with $\nabla \mathbf{t}(x)$, since both matrices are of the form $a\mathbf{I} + bxx^T$ with a, b scalars (that depend on x).

Radial transformations

In order to be able to change the base measure γ , we need to check that the inverses belong to the family. But if $y = \mathbf{t}(x)$, then $x = ay$ for some scalar a that solves the equation

$$aG(a\|y\|) = 1$$

Such a is guaranteed to be unique if $a \mapsto aG(a)$ is strictly increasing and it will exist (for y in the range of \mathbf{t}) if it is continuous. As a matter of fact, since the eigenvalues of $\nabla \mathbf{t}(x)$ are $G(a)$ and

$$G(a) + G'(a)a = (aG(a))', a = \|x\|$$

the condition that $a \mapsto aG(a)$ is strictly increasing is sufficient (this is weaker than G itself increasing). Finally, differentiability of G is not required, so it is enough if G is continuous and $aG(a)$ is strictly increasing.

Separable variables

Consider the collection of functions $\mathbf{t} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ of the form.

$$\mathbf{t}(x_1, \dots, x_d) = (T_1(x), \dots, T_d(x)), T_i : \mathbb{R} \rightarrow \mathbb{R}, 2.8)$$

with T_i continuous and strictly increasing. This is a generalisation of the compatible Gaussian case discussed above in which all the T_i 's were linear. Here, it is obvious that elements in this family are optimal maps and that the family is closed under inverses and composition, so that compatibility follows immediately.

This family is characterised by measures having a *common dependence structure*.

Separable variables

copula

we say that $C : [0, 1]^d \rightarrow [0, 1]$ is a copula if C is (the restriction of) a distribution function of a random vector having uniform margins.

In other words, if there is a random vector $V = (V_1, \dots, V_d)$ with $\mathbb{P}(V_i \leq a) = a$ for all $a \in [0, 1]$ and all $j = 1, \dots, d$, and

$$\mathbb{P}(V_1 < v_1, \dots, V_d < v_d) = C(v_1, \dots, v_d), v_i \in [0, 1]$$

To any d -dimensional probability measure μ , one can assign a copula $C = C_\mu$ in terms of the distribution function G of μ and its marginals G_j as

$$G(a_1, \dots, a_d) = \mu((-\infty, a_1] \times \dots \times (-\infty, a_d]) = C(G_1(a_1), \dots, G_d(a_d)).$$

If each G_j is surjective on $(0, 1)$, which is equivalent to it being continuous, then this equation defines C uniquely on $(0, 1)^d$, and consequently on $[0, 1]^d$. If some marginal G_j is not continuous, then uniqueness is lost, but C still exists. The connection of copulae to compatibility becomes clear in the following lemma.

Lemma 2.3.3 (Compatibility and Copulae)

The copulae associated with absolutely continuous measures $\mu, \nu \in \mathcal{W}_2(\mathbb{R}^d)$ are equal if and only if t_{μ}^{ν} takes the separable form (2.8).

Composition with linear functions

If $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex with gradient \mathbf{t} and A is a dd matrix, then the gradient of the convex function $x \mapsto \phi(Ax)$ at x is $\mathbf{t}_A = A^T \mathbf{t}(Ax)$. Suppose ψ is another convex function with gradient \mathbf{s} and that compatibility holds, i.e., $\nabla \mathbf{s}(\mathbf{t}(x))$ commutes with $\nabla \mathbf{t}(x)$ for all x . Then in order for

$$\nabla \mathbf{s}_A(\mathbf{t}_A(x)) = A^T \nabla \mathbf{s}(AA^T \mathbf{t}(Ax))A, \nabla \mathbf{t}_A(x) = A^T \nabla \mathbf{t}(Ax)A$$

to commute, it suffices that $AA^T = \mathbf{I}$, i.e., that A be orthogonal. Consequently, if $\{\mathbf{t}_{\#} \mu_{t \in T}\}$ are compatible, then so are $\{\mathbf{t}_U \# \mu_{t \in T}\}$ for any orthogonal matrix U .