

Optimal Transport for Data Analysis

Bernhard Schmitzer

2017-05-16

1 Introduction

1.1 Reminders on Measure Theory

Reference: Ambrosio, Fusco, Pallara: Functions of Bounded Variation and Free Discontinuity Problems, Chapters 1 & 2, [Ambrosio et al., 2000].

Definition 1.1 (σ -algebra). A collection \mathcal{E} of subsets of a set X is called σ -algebra if

- (i) $\emptyset \in \mathcal{E}$; $[A \in \mathcal{E}] \Rightarrow [X \setminus A \in \mathcal{E}]$;
- (ii) for a sequence $A_n \in \mathcal{E} \Rightarrow \bigcup_{n=0}^{\infty} A_n \in \mathcal{E}$.

Comment: Closed under finite unions, intersections and countable intersections. $A \cap B = X \setminus ((X \setminus A) \cup (X \setminus B))$.

Comment: Elements of \mathcal{E} : ‘measurable sets’. Pair (X, \mathcal{E}) : ‘measure space’.

Example 1.2. Borel algebra: smallest σ algebra containing all open sets of a topological space.

Comment: Intersection of two σ -algebras is again σ -algebra. ‘smallest’ is well-defined.

Definition 1.3 (Positive measure and vector measure). For measure space (X, \mathcal{E}) a function $\mu : \mathcal{E} \mapsto [0, +\infty]$ is called ‘positive measure’ if

- (i) $\mu(\emptyset) = 0$;
- (ii) for pairwise disjoint sequence $A_n \in \mathcal{E} \Rightarrow \mu(\bigcup_{n=0}^{\infty} A_n) = \sum_{n=0}^{\infty} \mu(A_n)$

For measure space (X, \mathcal{E}) and \mathbb{R}^m , $m \geq 1$, a function $\mu : \mathcal{E} \mapsto \mathbb{R}^m$ is called ‘measure’ if μ satisfies (i) and (ii) with absolute convergence.

Comment: Measures are vector space, measures are finite, positive measures may be infinite.

Example 1.4. Examples for measures:

1. counting measure: $\#(A) = |A|$ if A finite, $+\infty$ else.
2. Dirac measure: $\delta_x(A) = 1$ if $x \in A$, 0 else.
3. Lebesgue measure $\mathcal{L}([a, b]) = b - a$ for $b \geq a$.
4. Scaled measures: positive measure μ , function $f \in L^1(\mu)$, new measure $\nu = f \cdot \mu$. $\nu(A) \stackrel{\text{def.}}{=} \int_A f(x) d\mu(x)$.

5. Weak gradient of discontinuous function f , $\mu = Df$.

$$\int_{\Omega} \varphi(x) \cdot d\mu(x) = - \int_{\Omega} \operatorname{div} \varphi(x) f(x) dx$$

for $\varphi \in C^1(\Omega)$.

Definition 1.5 (Total variation). For measure μ on (X, \mathcal{E}) the total variation $|\mu|$ of $A \in \mathcal{E}$ is

$$|\mu|(A) = \sup \left\{ \sum_{n=0}^{\infty} |\mu(A_n)| \mid A_n \in \mathcal{E}, \text{ pairwise disjoint, } \bigcup_{n=0}^{\infty} A_n = A \right\}.$$

$|\mu|$ is finite, positive measure on (X, \mathcal{E}) .

Comment: Careful with nomenclature in image analysis.

Definition 1.6. A set $N \subset X$ is μ -negligible if $\exists A \in \mathcal{E}$ with $N \subset A$ and $\mu(A) = 0$. Two functions $f, g : X \rightarrow Y$ are identical ‘ μ -almost everywhere’ when $\{x \in X \mid f(x) \neq g(x)\}$ is μ -negligible.

Example 1.7. Null sets are Lebesgue-negligible sets.

Definition 1.8 (Measurable functions, push-forward). Let $(X, \mathcal{E}), (Y, \mathcal{F})$ be measurable spaces. A function $f : X \rightarrow Y$ is ‘measurable’ if $f^{-1}(A) \in \mathcal{E}$ for $A \in \mathcal{F}$.

For measure μ on (X, \mathcal{E}) the ‘push-forward’ of μ under f to (Y, \mathcal{F}) , we write $f_{\#}\mu$, is defined by $f_{\#}\mu(A) = \mu(f^{-1}(A))$ for $A \in \mathcal{F}$.

Change of variables formula:

$$\int_X g(f(x)) d\mu(x) = \int_Y g(y) df_{\#}\mu(y)$$

Sketch: Varying densities.

Example 1.9 (Marginal). Let $\operatorname{proj}_i : X \times X \rightarrow X$, $\operatorname{proj}_i(x_0, x_1) = x_i$. Marginals of measure γ on $X \times X$:

$$\operatorname{proj}_{0\#}\gamma(A) = \gamma(A \times X), \quad \operatorname{proj}_{1\#}\gamma(A) = \gamma(X \times A).$$

Sketch: Discuss pre-images of proj_i .

Definition 1.10 (Absolute continuity, singularity). Let μ be positive measure, ν measure on measurable space (X, \mathcal{E}) . ν is ‘absolutely continuous’ w.r.t. μ , we write $\nu \ll \mu$, if $[\mu(A) = 0] \Rightarrow [\nu(A) = 0]$.

Sketch: Density \ll Lebesgue, density $\not\ll$ density when support different, Dirac measures $\not\ll$ Lebesgue, mixed measures $\not\ll$ density, mixed measures \ll mixed measures when Diracs coincide.

Positive measures μ, ν are ‘mutually singular’, we write $\mu \perp \nu$, if $\exists A \in \mathcal{E}$ such that $\mu(A) = 0$, $\mu(X \setminus A) = 0$. For general measures replace μ, ν by $|\mu|, |\nu|$.

Definition 1.11 (σ -finite). A positive measure μ is called σ -finite if $X = \bigcup_{n=0}^{\infty} A_n$ for sequence $A_n \in \mathcal{E}$ with $\mu(A_n) < +\infty$.

Example 1.12. Lebesgue measure is not finite but σ -finite.

Theorem 1.13 (Radon–Nikodym, Lebesgue decomposition [[Ambrosio et al., 2000](#), Theorem 1.28]). Let μ be σ -finite positive measure. ν general measure.

Radon–Nikodym: For $\nu \ll \mu$ there is a function $f \in L^1(\mu)$ such that $\nu = f \cdot \mu$.

Lebesgue decomposition: there exist unique measures ν_a, ν_s such that

$$\nu = \nu_a + \nu_s, \quad \nu_a \ll \mu, \quad \nu_s \perp \mu.$$

Note: $\nu_a = f \cdot \mu$ for some $f \in L^1(\mu)$.

Corollary 1.14. A real-valued measure ν can be decomposed into $\nu = \nu_+ - \nu_-$ with ν_+, ν_- mutually singular positive measures.

Proof. Since $\nu \ll |\nu|$ there exists $f \in L^1(|\nu|)$ with $\nu = f \cdot |\nu|$. Set $A_+ = f^{-1}((0, +\infty))$, $A_- = f^{-1}((-\infty, 0))$ and set $\nu_{\pm}(B) = |\nu(B \cap A_{\pm})|$. \square

Comment: f is only unique $|\nu|$ -almost everywhere.

1.2 Duality

References: Kurdila, Zabrankin: Convex functional analysis [[Kurdila and Zabrankin, 2005](#)].
For Hilbert spaces: Bauschke, Combettes: Convex Analysis and Monotone Operator Theory in Hilbert Spaces [[Bauschke and Combettes, 2011](#)]

Definition 1.15 (Dual space). For normed vector space $(X, \|\cdot\|_X)$ its topological dual space is given by

$$X^* = \{y : X \rightarrow \mathbb{R} \mid y \text{ linear, continuous, i.e. } \exists C < \infty, |y(x)| \leq C \|x\|_X \forall x \in X\}.$$

Norm on X^* :

$$\|y\|_{X^*} = \sup \{|y(x)| \mid x \in X, \|x\|_X \leq 1\}$$

$(X^*, \|\cdot\|_{X^*})$ is Banach space. For $y(x)$ one often writes $\langle y, x \rangle$ or $\langle y, x \rangle_{X^*, X}$.

Comment: Linear not necessarily continuous in infinite dimensions. Dual norm is operator norm.

Definition 1.16 (Weak convergence). A sequence x_n in X converges weakly to $x \in X$ if $y(x_n) \rightarrow y(x)$ for all $y \in X^*$. We write $x_n \rightharpoonup x$.

Definition 1.17 (Weak* convergence). A sequence y_n in X^* converges weakly to $y \in X^*$ if $y_n(x) \rightarrow y(x)$ for all $x \in X$. We write $y_n \xrightarrow{*} y$.

Application to measures:

Definition 1.18 (Radon measures). Let (X, d) be compact metric space, let \mathcal{E} be Borel- σ -algebra. A finite measure (positive or vector valued) is called a ‘Radon measure’. Write:

- $\mathcal{M}_+(X)$: positive Radon measures,
- $\mathcal{P}(X) \subset \mathcal{M}_+(X)$: Radon probability measures (total mass = 1),
- $\mathcal{M}(X)^m$: (vector valued) Radon measures.

Theorem 1.19 (Regularity [Ambrosio et al., 2000, Proposition 1.43]). For positive Radon measures on (X, \mathcal{E}) one has for $A \in \mathcal{E}$

$$\mu(A) = \sup \{ \mu(B) \mid B \in \mathcal{E}, B \subset A, B \text{ compact} \} = \inf \{ \mu(B) \mid B \in \mathcal{E}, A \subset B, B \text{ open} \} .$$

Theorem 1.20 (Duality [Ambrosio et al., 2000, Theorem 1.54]). Let (Ω, d) be compact metric space. Let $C(\Omega)^m$ be space of continuous functions from Ω to \mathbb{R}^m , equipped with sup-norm. The topological dual of $C(\Omega)^m$ can be identified with the space $\mathcal{M}(\Omega)^m$ equipped with the total variation norm $\|\mu\|_{\mathcal{M}} \stackrel{\text{def.}}{=} |\mu|(\Omega)$. Duality pairing for $\mu \in \mathcal{M}(\Omega)^m, f \in C(\Omega)^m$:

$$\mu(f) = \langle \mu, f \rangle_{\mathcal{M}, C} = \int_{\Omega} f(x) d\mu(x)$$

Corollary 1.21. Two measures $\mu, \nu \in \mathcal{M}(\Omega)^m$ with $\mu(f) = \nu(f)$ for all $f \in C(\Omega)^m$ coincide.

Theorem 1.22 (Banach–Alaoglu [Kurdila and Zabaranin, 2005, Theorem 2.4.4]). Let X be a separable normed space. Any bounded sequence in X^* has a weak* convergent subsequence.

Comment: Since $C(\Omega)$ is separable, any bounded sequence in $\mathcal{M}(\Omega)$ has a weak* convergent subsequence.

1.3 Monge formulation of optimal transport

Comment: Gaspard Monge: French mathematician and engineer, 18th century. Studied problem of optimal allocation of resources to minimize transport cost.

Sketch: Bakeries and cafes

Example 1.23 (According to Villani). Every morning in Paris bread must be transported from bakeries to cafes for consumption. Every bakery produces prescribed amount of bread, every cafe orders prescribed amount. Assume: total amounts identical. Look for most economical way to distribute bread.

Mathematical model:

- $\Omega \subset \mathbb{R}^2$: area of Paris
- $\mu \in \mathcal{P}(\Omega)$: distribution of bakeries and produced amount of bread,
- $\nu \in \mathcal{P}(\Omega)$: distribution of cafes and consumed amount of bread
- Cost function $c : \Omega \times \Omega \rightarrow \mathbb{R}_+$. $c(x, y)$ gives cost of transporting 1 unit of bread from bakery at x to cafe at y .
- Describe transport by map $T : \Omega \rightarrow \Omega$. Bakery at x will deliver bread to cafe at $T(x)$. Consistency condition: $T_{\#}\mu = \nu$.

Comment: Each cafe receives precisely ordered amount of bread.

- Total cost of transport map

$$C_M(T) = \int_{\Omega} c(x, T(x)) d\mu(x)$$

Comment: For bakery at location x pay $c(x, T(x)) \cdot \mu(x)$. Sum (i.e. integrate) over all bakeries.

Definition 1.24. Monge optimal transport problem: find T that minimizes C_M .

Problems:

- Do maps T with $T_{\#}\mu = \nu$ exist? Can not split mass.

Sketch: Splitting of mass.

- Does minimal T exist? Non-linear, non-convex constraint and objective.
-

Comment: \Rightarrow problem remained unsolved for long time.

1.4 Kantorovich formulation of optimal transport

Comment: Leonid Kantorovich: Russian mathematician, 20th century. Founding father of linear programming, proposed modern formulation of optimal transport. (Nobel prize in economics 1975.)

Do not describe transport by map T , but by positive measure $\pi \in \mathcal{M}_+(\Omega \times \Omega)$.

Definition 1.25 (Coupling / Transport Plan). Let $\mu, \nu \in \mathcal{P}(\Omega)$. Set of ‘couplings’ or ‘transport plans’ $\Pi(\mu, \nu)$ is given by

$$\Pi(\mu, \nu) = \left\{ \pi \in \mathcal{P}(\Omega \times \Omega) \mid \text{proj}_{0\#}\pi = \mu, \text{proj}_{1\#}\pi = \nu \right\}.$$

Example 1.26. $\Pi(\mu, \nu) \neq \emptyset$, contains at least product measure $\mu \otimes \nu \in \Pi(\mu, \nu)$. $(\mu \otimes \nu)(A \times B) = \mu(A) \cdot \nu(B)$ for measurable $A, B \subset \Omega$.

Definition 1.27. For compact metric space (Ω, d) , $\mu, \nu \in \mathcal{P}(\Omega)$, $c \in C(\Omega \times \Omega)$ the Kantorovich optimal transport problem is given by

$$\mathcal{C}(\mu, \nu) = \inf \left\{ \int_{\Omega \times \Omega} c(x, y) d\pi(x, y) \mid \pi \in \Pi(\mu, \nu) \right\} \quad (1)$$

Comment: Linear (continuous) objective, affine constraint set.

Comment: Language of measures covers finite dimensional and infinite dimensional case.

Theorem 1.28. Minimizers of (1) exist.

Proof. • Let π_n be minimizing sequence. Since $\pi_n \in \mathcal{P}(\Omega \times \Omega)$ have $\|\pi_n\|_{\mathcal{M}} = 1$. By Banach-Alaoglu (Theorem 1.22) \exists converging subsequence. After extraction of subsequence have convergent minimizing sequence $\pi_n \xrightarrow{*} \pi$.

- Positivity: π is a positive measure. Otherwise find function $\phi \in C(\Omega \times \Omega)$ with $\int_{\Omega \times \Omega} \phi d\pi < 0$ (use Corollary 1.14 and Theorem 1.19 for construction) which contradicts weak* convergence.
- Unit mass: $\pi(\Omega \times \Omega) = \int_{\Omega \times \Omega} 1 d\pi = \lim_{n \rightarrow \infty} \int_{\Omega \times \Omega} 1 d\pi_n = \pi_n(\Omega \times \Omega) = 1$
- Marginal constraint: For every $\phi \in C(\Omega)$

$$\begin{aligned} \int_{\Omega} \phi d\text{proj}_{0\#}\pi &= \int_{\Omega \times \Omega} \phi \circ \text{proj}_0 d\pi \\ &= \lim_{n \rightarrow \infty} \int_{\Omega \times \Omega} \phi \circ \text{proj}_0 d\pi_n = \lim_{n \rightarrow \infty} \int_{\Omega} \phi d\text{proj}_{0\#}\pi_n = \int_{\Omega} \phi d\mu \end{aligned}$$

So $\text{proj}_{0\#}\pi = \mu$. Analogous: $\text{proj}_{1\#}\pi = \nu$.

- So: $\pi \in \Pi(\mu, \nu)$.
- Since $c \in C(\Omega \times \Omega)$ and $\pi_n \xrightarrow{*} \pi$ have

$$\int_{\Omega \times \Omega} c \, d\pi = \lim_{n \rightarrow \infty} \int_{\Omega \times \Omega} c \, d\pi_n.$$

Therefore, π is minimizer. □

Comment: For proof under more general conditions see for instance [Villani, 2009, Chapter 4]

Proof two additional useful results to get some practice and intuition.

Proposition 1.29 (Restriction [Villani, 2009, Theorem 4.6]). Let $\mu, \nu \in \mathcal{P}(\Omega)$, $c \in C(\Omega \times \Omega)$, let π be optimizer for $\mathcal{C}(\mu, \nu)$. Let $\tilde{\pi} \in \mathcal{M}_+(\Omega \times \Omega)$, $\tilde{\pi}(\Omega \times \Omega) > 0$, $\tilde{\pi}(A) \leq \pi(A)$ for all measurable $A \subset \Omega \times \Omega$. Set $\pi' = \frac{\tilde{\pi}}{\tilde{\pi}(\Omega \times \Omega)}$, $\pi' \in \mathcal{P}(\Omega \times \Omega)$. Let $\mu' = \text{proj}_{0\#}\pi'$, $\nu' = \text{proj}_{1\#}\pi'$. Then π' is minimal for $\mathcal{C}(\mu', \nu')$.

Example 1.30. $\tilde{\pi}(A) \stackrel{\text{def.}}{=} \pi(A \cap (\Omega_0 \times \Omega_1))$ for $\Omega_0, \Omega_1 \subset \Omega$.

Sketch: Restriction to subset. More general restriction.

Proof. • Assume π' is not optimal. Then there is a measure $\pi'' \in \Pi(\mu', \nu')$ with strictly better cost.

- Consider the measure $\hat{\pi} = \pi - \tilde{\pi} + \tilde{\pi}(\Omega \times \Omega) \cdot \pi''$. $\hat{\pi}$ is a positive measure since $\tilde{\pi} \leq \pi$. $\hat{\pi} \in \mathcal{P}(\Omega \times \Omega)$ since $\pi'' \in \mathcal{P}(\Omega \times \Omega)$.

$$\begin{aligned} \text{proj}_{0\#}\hat{\pi} &= \text{proj}_{0\#}\pi - \text{proj}_{0\#}\tilde{\pi} + \tilde{\pi}(\Omega \times \Omega) \cdot \text{proj}_{0\#}\pi'' \\ &= \mu - \tilde{\pi}(\Omega \times \Omega) \cdot (\mu' - \mu') = \mu \end{aligned}$$

So $\hat{\pi} \in \Pi(\mu, \nu)$.

- $\hat{\pi}$ has lower transport cost than π :

$$\int_{\Omega \times \Omega} c \, d\hat{\pi} = \int_{\Omega \times \Omega} c \, d\pi - \tilde{\pi}(\Omega \times \Omega) \int_{\Omega \times \Omega} c \, d\pi' + \tilde{\pi}(\Omega \times \Omega) \int_{\Omega \times \Omega} c \, d\pi'' < \int_{\Omega \times \Omega} c \, d\pi$$

- So π is not optimal which is a contradiction. Therefore π' must be optimal. □

Proposition 1.31 (Convexity [Villani, 2009, Theorem 4.8]). The function $\mathcal{P}(\Omega)^2 \rightarrow \mathbb{R}$, $(\mu, \nu) \mapsto \mathcal{C}(\mu, \nu)$ is convex.

Proof. • Let $\mu_0, \mu_1, \nu_0, \nu_1 \in \mathcal{P}(\Omega)$. Let π_i be corresponding minimizers in $\mathcal{C}(\mu_i, \nu_i)$, $i \in \{0, 1\}$.

- For $\lambda \in (0, 1)$ set

$$\hat{\mu} = (1 - \lambda)\mu_0 + \lambda\mu_1, \quad \hat{\nu} = (1 - \lambda)\nu_0 + \lambda\nu_1, \quad \hat{\pi} = (1 - \lambda)\pi_0 + \lambda\pi_1.$$

- $\hat{\pi} \in \Pi(\hat{\mu}, \hat{\nu})$ since

$$\text{proj}_{0\#}\hat{\pi} = (1 - \lambda)\text{proj}_{0\#}\pi_0 + \lambda\text{proj}_{0\#}\pi_1 = (1 - \lambda)\mu_0 + \lambda\mu_1 = \hat{\mu}.$$

- Convexity:

$$\mathcal{C}(\hat{\mu}, \hat{\nu}) \leq \int_{\Omega \times \Omega} c \, d\hat{\pi} = (1 - \lambda) \int_{\Omega \times \Omega} c \, d\pi_0 + \lambda \int_{\Omega \times \Omega} c \, d\pi_1 = (1 - \lambda)\mathcal{C}(\mu_0, \nu_0) + \lambda\mathcal{C}(\mu_1, \nu_1)$$

□

2 Kantorovich duality

2.1 More duality

Definition 2.1 (Topologically paired spaces). Two vector spaces X, X^* with locally convex Hausdorff topology are called *topologically paired spaces* if all continuous linear functionals on one space can be identified with all elements of the other.

Example 2.2. Let (Ω, d) be a compact metric space. $C(\Omega)$ and $\mathcal{M}(\Omega)$ with the sup-norm topology and the weak-* topology are topologically paired spaces.

Any continuous linear functional on $C(\Omega)$ can be identified with an element in $\mathcal{M}(\Omega)$ by construction. If Φ is a weak-* continuous linear functional on $\mathcal{M}(\Omega)$ it can be identified with the continuous function $\varphi : x \mapsto \Phi(\delta_x)$.

Definition 2.3 (Fenchel–Legendre conjugates). Let X, X^* be topologically paired spaces. Let $f : X \rightarrow \mathbb{R} \cup \{\infty\}$. Its Fenchel–Legendre conjugate $f^* : X^* \rightarrow \mathbb{R} \cup \{\infty\}$ is given by

$$f^*(y) = \sup\{\langle y, x \rangle - f(x) \mid x \in X\}.$$

f^* is convex, lsc on X^* . Likewise, for $g : X^* \rightarrow \mathbb{R} \cup \{\infty\}$ define conjugate g^* . If f, g convex, lsc then $f = f^{**}, g = g^{**}$.

Comment: Lsc: lower semicontinuous, $[x_n \rightarrow x] \Rightarrow [f(x) \leq \liminf_{n \rightarrow \infty} f(x_n)]$

Theorem 2.4 (Fenchel–Rockafellar [[Rockafellar, 1967](#)]). Let $(X, X^*), (Y, Y^*)$ be two pairs of topologically paired spaces. Let $f : X \rightarrow \mathbb{R} \cup \{\infty\}, g : Y \rightarrow \mathbb{R} \cup \{\infty\}, f, g$ convex, $A : X \rightarrow Y$ linear, continuous. Assume $\exists x \in X$ such that f finite at x, g finite and continuous at Ax . Then

$$\inf\{f(x) + g(Ax) \mid x \in X\} = \max\{-f^*(-A^*z) - g^*(z) \mid z \in Y^*\}.$$

In particular a maximizer of the problem on the right exists. $A^* : Y^* \rightarrow X^*$ is adjoint of A defined by $\langle z, Ax \rangle_{Y^*, Y} = \langle A^*z, x \rangle_{X^*, X}$.

Comment: Can sometimes be used ‘in both directions’ to establish existence of both primal and dual problem.

2.2 Dual Kantorovich problem

Theorem 2.5. Given the setting of Definition 1.27 one finds

$$\mathcal{C}(\mu, \nu) = \sup \left\{ \int_{\Omega} \alpha \, d\mu + \int_{\Omega} \beta \, d\nu \mid \alpha, \beta \in C(\Omega), \alpha(x) + \beta(y) \leq c(x, y) \text{ for all } (x, y) \in \Omega^2 \right\} \quad (2)$$

Proof. • Problem (2) can be written as

$$\mathcal{C}(\mu, \nu) = - \inf \{ f(\alpha, \beta) + g(A(\alpha, \beta)) \mid (\alpha, \beta) \in C(\Omega)^2 \}$$

with

$$\begin{aligned} f : C(\Omega)^2 &\rightarrow \mathbb{R}, & (\alpha, \beta) &\mapsto - \int_{\Omega} \alpha \, d\mu - \int_{\Omega} \beta \, d\nu \\ g : C(\Omega^2) &\rightarrow \mathbb{R} \cup \{\infty\}, & \psi &\mapsto \begin{cases} 0 & \text{if } \psi(x, y) \leq c(x, y) \text{ for all } (x, y) \in \Omega^2 \\ +\infty & \text{else.} \end{cases} \\ A : C(\Omega)^2 &\rightarrow C(\Omega^2), & [A(\alpha, \beta)](x, y) &= \alpha(x) + \beta(y). \end{aligned}$$

- f, g are convex, lsc. A is bounded, linear.
- Let (α, β) be two constant, finite functions with $\alpha(x) + \beta(y) < \min\{c(x', y') | (x', y') \in \Omega^2\}$. Then $f(\alpha, \beta) < \infty$, $g(A(\alpha, \beta)) < \infty$ and g is continuous at $A(\alpha, \beta)$. Thus, with Theorem 2.4 (and Example 2.2)

$$\mathcal{C}(\mu, \nu) = \min \{f^*(-A^*\pi) + g^*(\pi) | \pi \in \mathcal{M}(\Omega^2)\}.$$

- One obtains:

$$\begin{aligned} f^*(-\rho, -\sigma) &= \sup \left\{ -\int_{\Omega} \alpha d\rho - \int_{\Omega} \beta d\sigma + \int_{\Omega} \alpha d\mu + \int_{\Omega} \beta d\nu \mid (\alpha, \beta) \in C(\Omega^2)^2 \right\} \\ &= \begin{cases} 0 & \text{if } \rho = \mu, \sigma = \nu, \\ +\infty & \text{else.} \end{cases} \end{aligned}$$

(Reasoning similar than for positivity of limit π in proof of Theorem 1.28.)

$$\begin{aligned} g^*(\pi) &= \sup \left\{ \int_{\Omega^2} \psi d\pi \mid \psi \in C(\Omega^2), \psi(x, y) \leq c(x, y) \text{ for all } (x, y) \in \Omega^2 \right\} \\ &= \begin{cases} \int_{\Omega^2} c d\pi & \text{if } \pi \in \mathcal{M}_+(\Omega^2), \\ +\infty & \text{else.} \end{cases} \end{aligned}$$

□

So far we have not yet proven existence of dual maximizers. For this we need some additional arguments. We follow the presentation in [Santambrogio, 2015, Section 1.2].

Definition 2.6 (c -transform). For $\psi \in C(\Omega)$ define its c -transform $\psi^c \in C(\Omega)$ by

$$\psi^c(y) = \inf \{c(x, y) - \psi(x) | x \in \Omega\}$$

and its \bar{c} -transform $\psi^{\bar{c}} \in C(\Omega)$ by

$$\psi^{\bar{c}}(x) = \inf \{c(x, y) - \psi(y) | y \in \Omega\}$$

A function ψ is called \bar{c} -concave if it can be written as $\psi = \phi^c$ for some $\phi \in C(\Omega)$. Analogously, ψ is c -concave if it can be written as $\psi = \phi^{\bar{c}}$.

Comment: Setting $\beta = \alpha^c$ (or $\alpha = \beta^{\bar{c}}$) in (2) corresponds to optimization over β for fixed α (and vice versa). In general alternating optimization of (2) in α and β does not yield an optimal solution.

Lemma 2.7. The set of c -concave and \bar{c} -concave functions are equicontinuous.

Proof. • Since $c \in C(\Omega \times \Omega)$ and (Ω, d) compact there is a continuous function $\omega : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ with $\omega(0) = 0$ such that $|c(x, y) - c(x, y')| \leq \omega(d(y, y'))$.

- Let $\psi = \phi^c$. Set $\phi_x : y \mapsto c(x, y) - \phi(x)$. For every $x \in \Omega$ have $|\phi_x(y) - \phi_x(y')| \leq \omega(d(y, y'))$. One finds

$$\psi(y) \leq \phi_x(y) \leq \phi_x(y') + \omega(d(y, y'))$$

for all $x, y, y' \in \Omega$. Taking the infimum over x one gets $\psi(y) \leq \psi(y') + \omega(d(y, y'))$ and by symmetry $|\psi(y) - \psi(y')| \leq \omega(d(y, y'))$. This implies equicontinuity of \bar{c} -concave functions.

- Argument for $\phi^{\bar{c}}$ analogous. □

Theorem 2.8 (Arzelà-Ascoli [Rudin, 1986, Thm. 11.28]). If (Ω, d) is a compact metric space and $(f_n)_n$ is a sequence of uniformly bounded, equicontinuous functions in $C(\Omega)$ then $(f_n)_n$ has a uniformly converging subsequence.

Theorem 2.9 ([Santambrogio, 2015, Prop. 1.11]). Maximizers of (2) exist.

Proof. • For feasible (α, β) with finite score in (2) one can always replace β by α^c and subsequently α by $(\alpha^c)^{\bar{c}}$ which are still feasible and do not decrease the functional value. Hence, we may impose the additional constraint that (α, β) in (2) are (c, \bar{c}) -concave.

- Replacing feasible (α, β) in (2) by

$$(x \mapsto \alpha(x) - C, y \mapsto \beta(y) + C) \quad \text{with } C = \min_{x' \in \Omega} \alpha(x')$$

does not change the functional value or affect the constraints.

- Arguing as in Lemma 2.7 one finds for c -concave α with $\min_x \alpha(x) = 0$ that $\alpha(x) \in [0, \omega(\text{diam } \Omega)]$ and for the corresponding $\beta = \alpha^c$ that $\beta(y) \in [\min c - \omega(\text{diam } \Omega), \max c]$.
- Hence, we may consider maximizing sequences of (2) that are uniformly bounded and equicontinuous. By the Arzelà-Ascoli Theorem there exists a uniformly converging subsequence. Since the objective (and the constraints) of (2) are upper semicontinuous (see proof of Theorem 2.5), the limit must be a maximizer. □

Corollary 2.10 (Primal-dual optimality condition). π solves (1) and (α, β) solve (2) if and only if $\alpha(x) + \beta(y) = c(x, y)$ π -almost everywhere.

Proof. • \Rightarrow : Assume $\pi, (\alpha, \beta)$ are primal and dual optimal then:

$$\int_{\Omega \times \Omega} c(x, y) d\pi(x, y) = \int_{\Omega} \alpha d\mu + \int_{\Omega} \beta d\nu = \int_{\Omega \times \Omega} [\alpha(x) + \beta(y)] d\pi(x, y)$$

And $\alpha(x) + \beta(y) \leq c(x, y)$ for all $(x, y) \in \Omega^2$. Therefore $\alpha(x) + \beta(y) = c(x, y)$ π -a.e..

- \Leftarrow : Assume $\alpha(x) + \beta(y) = c(x, y)$ π -a.e..

$$\int_{\Omega} \alpha d\mu + \int_{\Omega} \beta d\nu = \int_{\Omega \times \Omega} [\alpha(x) + \beta(y)] d\pi(x, y) = \int_{\Omega \times \Omega} c(x, y) d\pi(x, y)$$

□

Remark 2.11 (Economic Interpretation of Kantorovich Duality). Bakeries and cafes hire a third-party company to do the transportation and agree to split the transport cost. When picking up bread at bakery x in the morning, the company charges an advance payment $\alpha(x)$ per unit of bread for the transport. Upon delivery at a cafe at y it charges a final payment $\beta(y)$ per unit of bread from the cafe.

The total payment to the company will be $\int_{\Omega} \alpha d\mu + \int_{\Omega} \beta d\nu$. It is left to the company to decide which bread to deliver where. And they will want to minimize the total transport cost, i.e. to find the *global minimum* of $\int_{\Omega \times \Omega} c d\pi$.

But it can never charge more than $c(x, y) - \alpha(x)$ when dropping of bread from x at y , otherwise the cafe y may complain and try to hire another company to get bread from bakery x at a lower price. When every cafe receives bread from its ‘subjectively cheapest’ bakery (and similarly each bakery delivers to its ‘subjectively cheapest’ cafe), the transport plan is said to be at equilibrium: no party will attempt to change its partner in a local attempt to reduce its costs.

Kantorovich duality states that for the optimal transport model the global minimum and equilibrium coincide.

A useful application of duality is the following result which is also the foundation for the numerical approximation of the Kantorovich problem.

Proposition 2.12 (Stability of optimal plans). Let $(\mu_n)_n, (\nu_n)_n$ be sequences in $\mathcal{P}(\Omega)$ converging weak* to μ and ν respectively. Let $(\pi_n)_n$ be a corresponding sequence of optimal plans. Then any cluster point of $(\pi_n)_n$ is an optimal coupling for $\mathcal{C}(\mu, \nu)$.

Comment: $(\pi_n)_n$ will always have cluster points due to Theorem 1.22.

Proof. • Let π be a cluster point of $(\pi_n)_n$. Without changing notation let $(\pi_n)_n$ be a subsequence converging weak* to π . Then $\pi \in \Pi(\mu, \nu)$ as for any $\phi \in C(\Omega)$:

$$\int_{\Omega} \phi d(\text{proj}_{0\#}\pi) = \int_{\Omega \times \Omega} \phi \circ \text{proj}_0 d\pi = \lim_{n \rightarrow \infty} \int_{\Omega \times \Omega} \phi \circ \text{proj}_0 d\pi_n = \lim_{n \rightarrow \infty} \int_{\Omega} \phi d\mu_n = \int_{\Omega} \phi d\mu$$

- Since π is feasible for $\mathcal{C}(\mu, \nu)$, for this converging subsequence:

$$\mathcal{C}(\mu, \nu) \leq \int_{\Omega \times \Omega} c d\pi = \lim_{n \rightarrow \infty} \int_{\Omega \times \Omega} c d\pi_n = \lim_{n \rightarrow \infty} \mathcal{C}(\mu_n, \nu_n)$$

- Since this is true for any converging sub-sequence, get:

$$\mathcal{C}(\mu, \nu) \leq \liminf_{n \rightarrow \infty} \mathcal{C}(\mu_n, \nu_n)$$

- Let $\varepsilon > 0$.
- Let $(\alpha_n, \beta_n)_n$ be a sequence of dual optimizers for $\mathcal{C}(\mu_n, \nu_n)$. Arguing as in the proof of Theorem 2.9 we can extract a converging subsequence $(\alpha_n, \beta_n)_n$, converging uniformly to some (α, β) . Note that $\alpha(x) + \beta(y) \leq c(x, y)$ for all $(x, y) \in \Omega \times \Omega$.
- There is some $N \in \mathbb{N}$ such that $|\alpha - \alpha_n| \leq \varepsilon/2, |\beta - \beta_n| \leq \varepsilon/2$ for all $n \geq N$. So:

$$\int_{\Omega} \alpha d\mu_n + \int_{\Omega} \beta d\nu_n \geq \int_{\Omega} \alpha_n d\mu_n + \int_{\Omega} \beta_n d\nu_n - \varepsilon$$

- This is true for all converging subsequences $(\alpha_n, \beta_n)_n$. Taking the supremum over the limit superior for all such subsequences, we get

$$\mathcal{C}(\mu, \nu) \geq \int_{\Omega} \alpha d\mu + \int_{\Omega} \beta d\nu \geq \limsup_{n \rightarrow \infty} \mathcal{C}(\mu, \nu) - \varepsilon$$

- Since this is true for any $\varepsilon > 0$ we find

$$\mathcal{C}(\mu, \nu) \geq \limsup_{n \rightarrow \infty} \mathcal{C}(\mu, \nu)$$

and thus π is optimal for $\mathcal{C}(\mu, \nu)$. (And $\mathcal{C}(\mu_n, \nu_n)$ converges.) □

Comment: The proof can be extended to cover a sequence of changing cost functions $(c_n)_n$ in $C(\Omega \times \Omega)$, where c_n is used for $\mathcal{C}(\mu_n, \nu_n)$ if $(c_n)_n$ converges uniformly to a limit $c \in C(\Omega \times \Omega)$.

Comment: For treatment of duality in more general regularity setting see for instance [Villani, 2009, Chapter 5]. A preview of the required concepts is given in subsection below.

2.3 c -cyclical monotonicity and duality

The proof for Proposition 2.12 relies on the uniform convergence of the dual potentials (α_n, β_n) . This is not available in less regular settings and a fundamentally different argument has to be used relying on the following property:

Definition 2.13 (c -cyclical monotonicity [Santambrogio, 2015, Def. 1.36]). Let $c \in C(\Omega \times \Omega)$. A set $\Gamma \subset \Omega \times \Omega$ is c -cyclical monotone (short: c -CM) if for every $n \in \mathbb{N}$ and every tuple of points $((x_1, y_1), \dots, (x_n, y_n)) \in \Gamma^n$ one has

$$\sum_{i=1}^n c(x_i, y_i) \leq \sum_{i=1}^n c(x_i, y_{i-1})$$

with convention $y_0 := y_n$.

Comment: Intuition behind this: assume bakery x_i delivers to cafe y_i , $(x_i, y_i) \in \Gamma \subset \Omega \times \Omega$ and Γ is c -CM. Now assume bakery x_1 decides to deliver to cafe y_n instead, which then rejects bread from bakery x_n . This bakery has now reroute its bread to cafe y_{n-1} and so forth, until eventually a cycle occurs and bakery x_2 reroutes its bread to cafe y_1 . The fact that Γ is c -CM implies that such a cyclic rerouting can never improve the transport cost.

Definition 2.14 (Support of measure). Let (Ω, d) be a compact metric space with its Borel σ -algebra and $\mu \in \mathcal{M}_+(\Omega)$. The support of μ , denoted $\text{spt } \mu$ is the smallest closed set $A \subset \Omega$ such that $\mu(A) = \mu(\Omega)$. For $x \in \text{spt } \mu$ one has $\mu(B_r(x)) > 0$ for any $r > 0$.

From the above comment we deduce intuitively: if π is optimal transport plan one must have $\text{spt } \pi$ is c -CM. Otherwise a cyclic rerouting of mass, as above, could yield an improved plan. Formal statement:

Proposition 2.15 ([Santambrogio, 2015, Thm. 1.38]). If π is an optimal transport plan for $\mathcal{C}(\mu, \nu)$ then $\text{spt } \pi$ is c -CM.

Proof. • Assume $\text{spt } \pi$ is not c -CM. Then there is $n \in \mathbb{N}$, and a tuple $((x_1, y_1), \dots, (x_n, y_n)) \in (\text{spt } \pi)^n$ with

$$\sum_{i=1}^n c(x_i, y_i) - \sum_{i=1}^n c(x_i, y_{i-1}) = \varepsilon > 0.$$

(and all x_i, y_i different).

- Select small open environments $U_i \subset \Omega$ of x_i , and $V_i \subset \Omega$ of y_i such that

$$\begin{aligned} |c(x, y) - c(x_i, y)| &\leq \frac{\varepsilon}{4n} \quad \text{for all } y \in \Omega, x \in U_i, \\ |c(x, y) - c(x, y_i)| &\leq \frac{\varepsilon}{4n} \quad \text{for all } x \in \Omega, y \in V_i. \end{aligned}$$

(and U_i, V_i pairwise disjoint).

- Set $\delta = \min_i \pi(U_i \times V_i) > 0$ since $(x_i, y_i) \in \text{spt } \pi$.

- Set:

$$\pi_i = \frac{\pi \llcorner (U_i \times V_i)}{\pi(U_i \times V_i)}, \quad \mu_i = \text{proj}_{0\#} \pi_i, \quad \nu_i = \text{proj}_{1\#} \pi_i, \quad \hat{\pi}_i = \mu_i \otimes \nu_{i-1}.$$

- Let:

$$\hat{\pi} = \pi - \delta \sum_{i=1}^n \pi_i + \delta \sum_{i=1}^n \hat{\pi}_i$$

Note: $\hat{\pi} \geq 0$, $\hat{\pi} \in \Pi(\mu, \nu)$.

- New cost:

$$\begin{aligned} \int c \, d\hat{\pi} &= \int c \, d\pi + \delta \sum_{i=1}^n \left(\underbrace{\int c \, d\hat{\pi}_i}_{\leq c(x_i, y_{i-1}) + \frac{\varepsilon}{4n}} - \underbrace{\int c \, d\pi_i}_{\geq c(x_i, y_i) - \frac{\varepsilon}{4n}} \right) \\ &\leq \int c \, d\pi + \delta \left[\underbrace{\sum_{i=1}^n (c(x_i, y_{i-1}) - c(x_i, y_i))}_{=-\varepsilon} + \frac{\varepsilon}{2} \right] \\ &= \int c \, d\pi - \frac{\delta \varepsilon}{2}. \end{aligned}$$

So π cannot be optimal. □

The converse implication is much less clear: if $\text{spt } \pi$ is c -CM, is π an optimal transport plan? While there are no cyclical rearrangements, possibly there is a more complicated way to improve the plan. With duality we can show that c -CM is indeed sufficient for optimality.

Proposition 2.16 ([Santambrogio, 2015, Thm. 1.37]). Let $c \in C(\Omega \times \Omega)$, $\Gamma \subset \Omega^2$, $\Gamma \neq \emptyset$, Γ is c -CM. Then there exists a c -concave function $\alpha \in C(\Omega)$ such that

$$\alpha(x) + \alpha^c(y) = c(x, y) \quad \text{for all } (x, y) \in \Gamma.$$

For proof use small auxiliary Lemma.

Lemma 2.17. Let $c \in C(\Omega \times \Omega)$, $\beta : \Omega \rightarrow \mathbb{R} \cup \{-\infty\}$, β bounded from above, β not identical $-\infty$. Set

$$\alpha(x) := \inf\{c(x, y) - \beta(y) \mid y \in \Omega\}.$$

Then $\alpha \in C(\Omega)$, $(\alpha^c)^{\bar{c}} = \alpha$ which also implies that α is c -concave.

Proof of Proposition 2.16. • β is bounded from above, and finite at least at one point: family of functions $(x \mapsto c(x, y) - \beta(y))_{y \in \Omega: \beta(y) > -\infty}$ is non-empty and uniformly bounded from below. $\Rightarrow \alpha$ is pointwise infimum over non-empty family of equicontinuous functions uniformly bounded from below. $\Rightarrow \alpha \in C(\Omega)$.

- Now:

$$\begin{aligned}\alpha(x) &= \inf_y \{c(x, y) - \beta(y)\} \\ \alpha^c(y') &= \inf_x \sup_y \{c(x, y') - c(x, y) + \beta(y)\} \\ (\alpha^c)^{\bar{c}}(x') &= \inf_{y'} \sup_x \inf_y \{c(x', y') - c(x, y') + c(x, y) - \beta(y)\}\end{aligned}$$

- By setting $x = x'$ in supremum get: $(\alpha^c)^{\bar{c}}(x') \geq \alpha(x')$.
- By setting $y = y'$ in inner infimum get: $(\alpha^c)^{\bar{c}}(x') \leq \alpha(x')$.

□

Proof. • Pick $(x_1, y_1) \in \Gamma$. For $y \in \Omega$ set

$$\beta(y) = \sup \left\{ \sum_{i=1}^n c(x_i, y_i) - \sum_{i=2}^n c(x_i, y_{i+1}) \mid n \in \mathbb{N}, (x_i, y_i) \in \Gamma \text{ for } i = 1, \dots, n, y_n = y \right\}$$

- For $y \notin \text{proj}_1(\Gamma)$ find $\beta(y) = -\infty$ (supremum over empty set).
- For $y \in \text{proj}_1(\Gamma)$ use c -CM of Γ :

$$\beta(y) = \sup \left\{ \underbrace{\sum_{i=1}^n c(x_i, y_i) - \sum_{i=2}^n c(x_i, y_{i+1})}_{\leq c(x_1, y_n)} \mid n \in \mathbb{N}, (x_i, y_i) \in \Gamma \text{ for } i = 1, \dots, n, y_n = y \right\}$$

So β is bounded from above.

- For $y = y_1$ get by setting $n = 1$:

$$\beta(y_1) \geq c(x_1, y_1)$$

So β is not identical to $-\infty$.

- Now for $x \in \Omega$ set

$$\alpha(x) = \inf \{c(x, y) - \beta(y) \mid y \in \Omega\} .$$

By Lemma 2.17: $\alpha \in C(\Omega)$, $(\alpha^c)^{\bar{c}} = \alpha$.

- Now let $(x, y) \in \Gamma$. We need: $\alpha(x) + \alpha^c(y) = c(x, y)$. Since $\alpha^c \geq \beta$ and $\alpha(x) + \alpha^c(y) \leq c(x, y)$ a sufficient condition is $\alpha(x) + \beta(y) \geq c(x, y)$.
- For every $\varepsilon > 0$ there is some \hat{y} such that

$$\alpha(x) \geq c(x, \hat{y}) - \beta(\hat{y}) - \varepsilon.$$

And since $\alpha(x) \in \mathbb{R}$ have $\hat{y} \in \text{proj}_1(\Gamma)$.

- For β get recursive formula:

$$\beta(y) = \sup \{c(x, y) - c(x, \hat{y}) + \beta(\hat{y}) \mid (x, y) \in \Gamma, \hat{y} \in \text{proj}_1(\Gamma)\}$$

So:

$$\beta(y) \geq c(x, y) - c(x, \hat{y}) + \beta(\hat{y}).$$

- Combining lower bounds for $\alpha(x)$ and $\beta(y)$ we get:

$$\alpha(x) + \beta(y) \geq c(x, \hat{y}) - \beta(\hat{y}) - \varepsilon + c(x, y) - c(x, \hat{y}) + \beta(\hat{y}) = c(x, y) - \varepsilon.$$

Since this is true for any $\varepsilon > 0$ it is true for $\varepsilon = 0$. □

Application: sufficient condition for optimality of transport plans:

Corollary 2.18. If $\pi \in \Pi(\mu, \nu)$ and $\text{spt } \pi$ is c -CM, then π is an optimal coupling for $\mathcal{C}(\mu, \nu)$.

Proof. Use $(\alpha, \beta = \alpha^c)$ for function α provided by Proposition 2.16 as dual feasible candidates. □

Another application: alternative proof for stability result Proposition 2.12. Need a few ingredients.

Definition 2.19. For metric space (Ω, d) define Hausdorff distance for subsets A, B of Ω :

$$d_H(A, B) = \max\{\max\{d(x, B) \mid x \in A\}, \max\{d(A, y) \mid y \in B\}\}$$

Theorem 2.20 (Blaschke [Ambrosio and Tilli, 2004, Thm. 4.4.15]). For a compact metric space (Ω, d) the set of compact subsets of Ω with the distance d_H is a compact metric space.

Lemma 2.21. Let A_n be a sequence of compact subsets of Ω , $A \subset \Omega$ compact, let $A_n \rightarrow A$ in the Hausdorff distance. Then for every $x \in A$ there is a sequence $(x_n)_n$, $x_n \in A_n$ such that $x_n \rightarrow x$.

Proof. • Let $x \in A$. For any $\varepsilon > 0$ there is some N such that $d_H(A, A_n) \leq \varepsilon$ for $n \geq N$. Then:

$$d(x, A_n) \leq d_H(A, A_n) \leq \varepsilon$$

So there is some $x_n \in A_n$ such that $d(x, x_n) \leq \varepsilon$. □

Lemma 2.22. Compact metric space (Ω, d) . Hausdorff convergent sequence of compact subsets A_n to A . Weak* convergent sequence of measures $(\mu_n)_n$ in $\mathcal{P}(\Omega)$ to μ . $\text{spt } \mu_n \subset A_n$. Then $\text{spt } \mu \subset A$.

Proof. • Consider sequence of functions $f_n \in C(\Omega)$, $f_n : x \mapsto d(x, A_n)$ and set $f : x \mapsto d(x, A)$. Then $f_n \rightarrow f$ uniformly (with Lemma 2.21).

- So for every $\varepsilon > 0$ there is some $N < \infty$ such that $|f_n - f| \leq \varepsilon$ for $n \geq N$.
- Then:

$$\int f_n \, d\mu_n = 0$$

$$\int f \, d\mu = \lim_{n \rightarrow \infty} \int f \, d\mu_n \leq \lim_{n \rightarrow \infty} \int f_n \, d\mu_n + \varepsilon = \varepsilon$$

□

Alternative proof for Proposition 2.12. • As in first proof: let $(\pi_n)_n$ be sequence of optimizers. For any converging subsequence the limit π is in $\Pi(\mu, \nu)$.

- With Theorem 2.20 the sequence $(\text{spt } \pi_n)_n$ has a convergent subsequence in the Hausdorff metric. Denote limit set by Γ .
- By Lemma 2.22 have $\text{spt } \pi \subset \Gamma$ for any cluster point π of $(\pi_n)_n$.
- Every $\text{spt } \pi_n$ is c -CM by Proposition 2.15. Therefore, so is Γ . Indeed: let $n \in \mathbb{N}$, $((x_1, y_1), \dots, (x_n, y_n)) \in \Gamma^n$. For $i = 1, \dots, n$ let $((x_{i,k}, y_{i,k}))_k$ be sequence with $(x_{i,k}, y_{i,k}) \in \text{spt } \pi_k$, with $(x_{i,k}, y_{i,k}) \rightarrow (x_i, y_i)$. For every k find by c -CM of $\text{spt } \pi_k$:

$$\sum_{i=1}^n c(x_{i,k}, y_{i,k}) \leq \sum_{i=1}^n c(x_{i,k}, y_{i-1,k})$$

Hence this is also true in limit.

- So $\text{spt } \pi \subset \Gamma$ is c -CM. With Corollary 2.18 π is optimal for $C(\mu, \nu)$.

□

2.4 Solution to the Monge problem

We now consider a special case for which the Monge problem has a solution. Duality will be an important ingredient in the proof.

First we briefly discuss that the Kantorovich formulation of optimal transport, Definition 1.27, can be interpreted as a relaxation of the Monge formulation, Definition 1.24.

Proposition 2.23 (Kantorovich is a relaxation of the Monge problem). Assume $T : \Omega \rightarrow \Omega$ is a feasible transport map for the Monge problem between μ and ν , Definition 1.24. In particular $T_{\#}\mu = \nu$.

Let

$$(\text{id}, T) : \Omega \rightarrow \Omega \times \Omega, \quad x \mapsto (x, T(x)).$$

Then $\pi = (\text{id}, T)_{\#}\mu \in \Pi(\mu, \nu)$ and

$$\int_{\Omega \times \Omega} c \, d\pi = \int_{\Omega} c(x, T(x)) \, d\mu(x).$$

Proof. • Clearly $\pi \in \mathcal{P}(\Omega \times \Omega)$.

- $\text{proj}_0 \circ T = \text{id}$. Hence

$$\text{proj}_0 \# \pi = \text{proj}_0 \# T \# \mu = \mu.$$

- Similarly $\text{proj}_1 \circ T = T$. Hence

$$\text{proj}_1 \# \pi = \text{proj}_1 \# T \# \mu = T \# \mu = \nu.$$

- Equality of cost follows from change of variables under (id, T) . □

This implies in particular that $\mathcal{C}(\mu, \nu) \leq C_M(\mu, \nu)$ since every feasible Monge map induces a Kantorovich coupling of equal cost.

The converse inequality is in general not true but we will now prove it for a special case.

Theorem 2.24 (Solution to the Monge problem). Let $\Omega \subset \mathbb{R}^d$ be compact, let the cost function c be given by $c(x, y) = h(x - y)$ for a strictly convex function $h : \mathbb{R}^d \rightarrow \mathbb{R}$. Let μ be Lebesgue-absolutely continuous and let $\partial\Omega$ be μ -negligible.

Then the optimal transport plan π is supported on the graph of a transport map $T : \Omega \rightarrow \Omega$.

Proof. • h is convex and finite. Hence it is continuous and locally Lipschitz. Therefore, c is continuous and Lipschitz (since Ω is compact).

- Therefore Theorem 1.28 and Theorem 2.9 apply and provide existence of primal and dual optimizers π and (α, β) .
- From the proof of Theorem 2.9 know: $\beta = \alpha^c$, $\alpha = \beta^c$. Analogous to Lemma 2.7 this implies that α and β are Lipschitz.
- By Rademacher's theorem (see e.g. [Ziemer, 1989, Theorem 2.2.1]) α is Lebesgue-almost everywhere differentiable in $\text{int } \Omega$. And consequently μ -almost everywhere on Ω (since $\partial\Omega$ is μ -negligible).
- From Corollary 2.10: $\alpha(x) + \beta(y) = c(x, y)$ π -almost everywhere. For (x_0, y_0) with $\alpha(x_0) + \beta(y_0) = c(x_0, y_0)$ we find

$$x \mapsto c(x, y_0) - \alpha(x)$$

is minimal at x_0 (since $\beta(y_0) = \inf_x \{c(x, y_0) - \alpha(x)\} = c(x_0, y_0) - \alpha(x_0)$). If α is differentiable at x_0 (which it is μ -a.e., i.e. for (x, y) π -a.e.), then $\nabla\alpha(x_0) \in \partial h(x_0 - y_0)$.

- For a strictly convex function ∂h is 'invertible'. That is, for every $v \in \mathbb{R}^d$ there is a unique $w \in \mathbb{R}^d$ such that $v \in \partial h(w)$. We denote this map by ∂h^{-1} and find

$$x_0 - y_0 = \partial h^{-1}(\nabla\alpha(x_0)).$$

- This relation is still true π -almost everywhere. Set $T(x) = x - \partial h^{-1}(\nabla\alpha(x))$. Then $y = T(x)$ π -almost everywhere.

- Equality of cost:

$$\int_{\Omega \times \Omega} c(x, y) d\pi(x, y) = \int_{\Omega \times \Omega} c(x, T(x)) d\pi(x, y) = \int_{\Omega} c(x, T(x)) d\mu(x)$$

- Push-forward condition:

$$\begin{aligned} \int_{\Omega} \phi(y) dT_{\#}\mu(y) &= \int_{\Omega} \phi(T(x)) d\mu(x) = \int_{\Omega \times \Omega} \phi(T(x)) d\pi(x, y) \\ &= \int_{\Omega \times \Omega} \phi(y) d\pi(x, y) = \int_{\Omega} \phi(y) d\nu(y) \end{aligned}$$

□

Example 2.25 (Quadratic case: $c(x, y) = \frac{1}{2}\|x - y\|^2$). This corresponds to $h(x) = \frac{1}{2}\|x\|^2$. Consequently, $\partial h(x) = \{x\}$ and $\partial h^{-1}(x) = x$. Moreover since $\alpha = \beta^c$:

$$\begin{aligned} \alpha(x) &= \inf \left\{ \frac{1}{2}\|x - y\|^2 - \beta(y) \mid y \in \Omega \right\} = \frac{1}{2}\|x\|^2 + \inf \left\{ -\langle x, y \rangle + \frac{1}{2}\|y\|^2 - \beta(y) \mid y \in \Omega \right\} \\ &= \frac{1}{2}\|x\|^2 - \underbrace{\sup \left\{ \langle x, y \rangle - g(y) \mid y \in \mathbb{R}^d \right\}}_{:=\phi(x) : \text{convex}} \end{aligned}$$

Therefore,

$$T(x) = x - \nabla\alpha(x) = x - (x - \nabla\phi(x)) = \nabla\phi(x).$$

So T is almost everywhere the gradient of a convex function. This is part of the famous polar factorization theorem by [Brenier, 1991].

3 Wasserstein spaces

Definition 3.1 (Wasserstein distance). Let (Ω, d) be a compact metric space. For $p \in [1, \infty)$ let $W_p : \mathcal{P}(\Omega) \times \mathcal{P}(\Omega) \rightarrow \mathbb{R}$,

$$W_p(\mu, \nu) = \left(\inf \left\{ \int_{\Omega \times \Omega} d(x, y)^p d\pi(x, y) \mid \pi \in \Pi(\mu, \nu) \right\} \right)^{1/p}$$

Comment: On non-compact spaces one usually restricts the Wasserstein space to measures with finite moment of order p , i.e., $\int_{\Omega} d(x, x_0)^p d\mu < +\infty$ for some arbitrary reference point $x_0 \in \Omega$. This is a sufficient condition to keep W_p finite.

Example 3.2. Dirac measures are isometric embedding of Ω into $\mathcal{P}(\Omega)$: $W_p(\delta_x, \delta_y) = d(x, y)$, since $\Pi(\delta_x, \delta_y) = \{\delta_{(x,y)}\}$.

To prove that W_p is indeed a distance we will rely on the following powerful theorem which is often useful to dissect and reassemble measures with certain sought-after properties.

Theorem 3.3 (Disintegration [Ambrosio et al., 2005, Theorem 5.3.1]). Let $\tilde{\Omega}, \Omega$ be compact metric spaces, let $f : \tilde{\Omega} \rightarrow \Omega$ be measurable and $\pi \in \mathcal{P}(\tilde{\Omega})$. Set $\mu = f_{\#}\pi \in \mathcal{P}(\Omega)$. Then there is a family $(\pi_y)_{y \in \Omega}$ in $\mathcal{P}(\tilde{\Omega})$, unique μ -a.e., such that $\pi_y(f^{-1}(\{y\})) = 1$ and for $\phi \in C(\tilde{\Omega})$ one has

$$\int_{\tilde{\Omega}} \phi d\pi = \int_{\Omega} \left(\int_{\tilde{\Omega}} \phi d\pi_y \right) d\mu(y).$$

Sketch: Table, disintegration.

Comment: Disintegration formalizes the notion of conditional probability. It is easiest to visualize in a discrete case when $\tilde{\Omega} = \Omega \times \Omega$ and $f = \text{proj}_0$. Then π can be interpreted as table and any π_y will be the restriction of π to row y , renormalized to mass 1 (if the row is non-empty). π_y gives the probabilities of picking a given column under the condition that row y has already been selected.

Example 3.4 (Disintegration of transport plan). Let $\pi \in \Pi(\mu, \nu)$. Let $(\gamma_x)_{x \in \Omega}$ be the disintegration of π with respect to proj_0 . That is, for any $\phi \in C(\Omega \times \Omega)$ have

$$\int_{\Omega \times \Omega} \phi(x, y) d\pi(x, y) = \int_{\Omega} \left(\int_{\Omega} \phi(x, y) d\gamma_x(y) \right) d\mu(x).$$

γ_x can be interpreted as describing where mass particles starting in x are going. Note that it is only uniquely defined μ -a.e..

Comment: By the disintegration theorem γ_x would be in $\mathcal{P}(\Omega \times \Omega)$. But since $\gamma_x(\text{proj}_0^{-1}(\{x\})) = \gamma_x(\{x\} \times \Omega) = 1$ we can interpret γ_x as element of $\mathcal{P}(\Omega)$.

Theorem 3.5. W_p is a metric on $\mathcal{P}(\Omega)$.

Proof. • W_p is non-negative (since $d(x, y)^p \geq 0$), symmetric (since $d(x, y)^p$ is symmetric) and finite (since Ω is compact, i.e., d is bounded).

- Let $T : \Omega \rightarrow \Omega \times \Omega$, $T(x) = (x, x)$ be the ‘diagonal’ embedding of Ω into $\Omega \times \Omega$. $W_p(\mu, \mu) = 0$, since $\pi = T_{\#}\mu \in \Pi(\mu, \mu)$ and $\int d^p d\pi = 0$: Note that $(\text{proj}_i \circ T)(x) = x$ and that $f_{\#}(g_{\#}\rho) = (f \circ g)_{\#}\rho$. Hence, $\text{proj}_{i\#}T_{\#}\mu = \mu$. Further,

$$\int_{\Omega \times \Omega} d^p d\pi = \int_{\Omega \times \Omega} d^p d(T_{\#}\mu) = \int_{\Omega} d^p \circ T d\mu = 0.$$

- Let $W_p(\mu, \nu) = 0$. Then there must be some $\pi \in \Pi(\mu, \nu)$ with $\int_{\Omega \times \Omega} d(x, y)^p d\pi(x, y) = 0$, which implies $d(x, y) = 0$ π -a.e., i.e., $x = y$ π -a.e.. So for $\phi \in C(\Omega)$

$$\int_{\Omega \times \Omega} \phi(x) d\pi(x, y) = \int_{\Omega \times \Omega} \phi(y) d\pi(x, y)$$

and thus $\text{proj}_{0\#}\pi = \text{proj}_{1\#}\pi$ which implies $\mu = \nu$.

- Towards triangle inequality: Let $\mu, \nu, \rho \in \mathcal{P}(\Omega)$, let π_{01}, π_{12} be optimal couplings for $W_p(\mu, \nu)$ and $W_p(\nu, \rho)$. Let $(\gamma_{01,y})_{y \in \Omega}$ be the disintegration of π_{01} with respect to proj_1 . That is, for any $\phi \in C(\Omega \times \Omega)$ have

$$\int_{\Omega \times \Omega} \phi(x, y) d\pi_{01}(x, y) = \int_{\Omega} \left(\int_{\Omega} \phi(x, y) d\gamma_{01,y}(x) \right) d\nu(y).$$

Similarly, let $(\gamma_{12,y})_{y \in \Omega}$ be the disintegration of π_{12} with respect to proj_0 .

- Define a new measure $\pi \in \mathcal{P}(\Omega \times \Omega)$ via

$$\int_{\Omega \times \Omega} \phi(x, z) d\pi(x, z) = \int_{\Omega} \left(\int_{\Omega \times \Omega} \phi(x, z) d\gamma_{01,y}(x) d\gamma_{12,y}(z) \right) d\nu(y).$$

Sketch: Some intuition for π .

- Claim: $\pi \in \Pi(\mu, \rho)$. For $\phi \in C(\Omega)$ get

$$\begin{aligned} \int_{\Omega \times \Omega} \phi(x) d\pi(x, z) &= \int_{\Omega} \left(\int_{\Omega \times \Omega} \phi(x) d\gamma_{01,y}(x) d\gamma_{12,y}(z) \right) d\nu(y) \\ &= \int_{\Omega} \left(\int_{\Omega} \phi(x) d\gamma_{01,y}(x) \right) d\nu(y) = \int_{\Omega \times \Omega} \phi(x) d\pi_{01}(x, y) = \int_{\Omega} \phi d\mu \end{aligned}$$

- Triangle inequality:

$$\begin{aligned} W_p(\mu, \rho) &\leq \left(\int_{\Omega \times \Omega} d(x, z)^p d\pi(x, z) \right)^{1/p} = \left(\int_{\Omega} \left(\int_{\Omega \times \Omega} d(x, z)^p d\gamma_{01,y}(x) d\gamma_{12,y}(z) \right) d\nu(y) \right)^{1/p} \\ &\leq \left(\int_{\Omega} \left(\int_{\Omega \times \Omega} (d(x, y) + d(y, z))^p d\gamma_{01,y}(x) d\gamma_{12,y}(z) \right) d\nu(y) \right)^{1/p} \\ &\stackrel{\text{Minkowski ineq.}}{\leq} \left(\int_{\Omega} \left(\int_{\Omega \times \Omega} d(x, y)^p d\gamma_{01,y}(x) d\gamma_{12,y}(z) \right) d\nu(y) \right)^{1/p} \\ &\quad + \left(\int_{\Omega} \left(\int_{\Omega \times \Omega} d(y, z)^p d\gamma_{01,y}(x) d\gamma_{12,y}(z) \right) d\nu(y) \right)^{1/p} \\ &= \left(\int_{\Omega \times \Omega} d(x, y)^p d\pi_{01}(x, y) \right)^{1/p} + \left(\int_{\Omega \times \Omega} d(y, z)^p d\pi_{12}(x, y) \right)^{1/p} \\ &= W_p(\mu, \nu) + W_p(\nu, \rho). \end{aligned}$$

□

Theorem 3.6 (W_p metrizes weak* convergence). Let (Ω, d) be a compact metric space. W_p metrizes the weak* convergence on $\mathcal{P}(\Omega)$. That is, for a sequence $(\mu_n)_n$ and some μ in $\mathcal{P}(\Omega)$ one has:

$$[W_p(\mu_n, \mu) \rightarrow 0] \quad \Leftrightarrow \quad [\mu_n \xrightarrow{*} \mu]$$

Proof. • \Rightarrow : assume $W_p(\mu_n, \mu) \rightarrow 0$. Let $(\pi_n)_n$ be a corresponding sequence of optimal transport plans. Let $\tilde{\mu}$ be a cluster point of $(\mu_n)_n$ and let $\pi \in \Pi(\tilde{\mu}, \mu)$ a corresponding cluster point of $(\pi_n)_n$. As before, denote the converging subsequence also by $(\pi_n)_n$. One has:

$$W_p(\tilde{\mu}, \mu) \leq \lim_{n \rightarrow \infty} \int_{\Omega \times \Omega} d^p d\pi_n = \lim_{n \rightarrow \infty} W_p(\mu_n, \mu) = 0$$

Since W_p is a metric, $\tilde{\mu} = \mu$. Hence, $\mu_n \xrightarrow{*} \mu$.

- \Leftarrow : assume $\mu_n \xrightarrow{*} \mu$. Let π_n be optimal plans for $W_p(\mu_n, \mu)$. Extract a converging subsequence, again denoted by $(\pi_n)_n$. By Proposition 2.12 (stability of optimal plans) any cluster point π of $(\pi_n)_n$ is an optimal coupling for $W_p(\mu, \mu)$. So:

$$0 = W_p(\mu, \mu) = \int_{\Omega \times \Omega} d^p d\pi = \lim_{n \rightarrow \infty} \int_{\Omega \times \Omega} d^p d\pi_n = \lim_{n \rightarrow \infty} W_p(\mu_n, \mu)$$

□

3.1 Displacement interpolation

An intriguing property of the Wasserstein space $(\mathcal{P}(\Omega), W_p)$ is that it is a length space if (Ω, d) is a length space.

Definition 3.7 (Length space). A metric space (Ω, d) is a length space if for every pair $(x, y) \in \Omega$ there is a continuous map $\gamma_{x,y} \in C([0, 1], \Omega)$ with

$$\gamma_{x,y}(0) = x, \quad \gamma_{x,y}(1) = y, \quad d(\gamma_{x,y}(s), \gamma_{x,y}(t)) = d(x, y) \cdot |s - t|$$

for $s, t \in [0, 1]$.

Theorem 3.8. If (Ω, d) is a length space and the map $(x, y) \mapsto \gamma_{x,y}$ that takes start and endpoint to a shortest path between them is measurable then $(\mathcal{P}(\Omega), W_p)$ is a length space.

Comment: Sufficient conditions for the measurability of $(x, y) \mapsto \gamma_{x,y}$ can be found for instance in [Villani, 2009, Proposition 7.16].

Proof. • Let $(\gamma_{x,y})_{(x,y) \in \Omega^2}$ be the family of maps for (Ω, d) as given by Definition 3.7. For fixed $s, t \in [0, 1]$ let

$$\begin{aligned} \Gamma_s &: \Omega \times \Omega \rightarrow \Omega, & (x, y) &\mapsto \gamma_{x,y}(s), \\ \Gamma_{s,t} &: \Omega \times \Omega \rightarrow \Omega \times \Omega, & (x, y) &\mapsto (\gamma_{x,y}(s), \gamma_{x,y}(t)). \end{aligned}$$

Comment: Between γ and Γ the roles of ‘index’ and ‘arguments’ of the functions are exchanged. This is formally helpful to use the push-forward of Γ .

- For given $\mu, \nu \in \mathcal{P}(\Omega)$ let π be an optimal coupling for $W_p(\mu, \nu)$. Denote $\rho_s = \Gamma_{s\#}\pi$.

Sketch: Interpretation of ρ_s .

- Claim: $s \mapsto \rho_s$ is a geodesic in $(\mathcal{P}(\Omega), W_p)$ between μ and ν . A ‘length space map’ for $(\mathcal{P}(\Omega), W_p)$ between μ and ν , $\gamma_{\mu,\nu} : [0, 1] \rightarrow \mathcal{P}(\Omega)$ is given by $\gamma_{\mu,\nu}(s) = \rho_s$. We will now show this.
- Measurability of Γ_s : By assumption $S : (x, y) \mapsto \gamma_{x,y}$ is measurable. For fixed $t \in [0, 1]$ the map $e_t : C([0, 1], \Omega) \rightarrow \Omega$, $\gamma \mapsto \gamma(t)$ is continuous and thus measurable. We find $\Gamma_s = e_s \circ S$. Similarly, $\Gamma_{s,t} = (\Gamma_s, \Gamma_t) = (e_s, e_t) \circ S$ is measurable.
- Claim: $\Gamma_{s,t\#}\pi \in \Pi(\rho_s, \rho_t)$.

$$\text{proj}_{0\#}\Gamma_{s,t\#}\pi = (\text{proj}_0 \circ \Gamma_{s,t})\# \pi = \Gamma_{s\#}\pi = \rho_s$$

- Claim: $W_p(\rho_s, \rho_t) = |s - t| \cdot W_p(\mu, \nu)$.

$$\begin{aligned} W_p(\rho_s, \rho_t)^p &\leq \int_{\Omega \times \Omega} d(x, y)^p d(\Gamma_{s,t\#}\pi)(x, y) = \int_{\Omega \times \Omega} ((d \circ \Gamma_{s,t})(x, y))^p d\pi(x, y) \\ &= \int_{\Omega \times \Omega} (d(\gamma_{x,y}(s), \gamma_{x,y}(t)))^p d\pi(x, y) = |s - t|^p \int_{\Omega \times \Omega} (d(x, y))^p d\pi(x, y) \\ W_p(\rho_s, \rho_t) &\leq |s - t| \cdot W_p(\mu, \nu) \end{aligned}$$

So for $0 \leq s \leq t \leq 1$ have

$$W_p(\mu, \rho_s) \leq s \cdot W_p(\mu, \nu), \quad W_p(\rho_s, \rho_t) \leq (t - s) \cdot W_p(\mu, \nu), \quad W_p(\rho_t, \nu) \leq (1 - t) \cdot W_p(\mu, \nu)$$

So

$$W_p(\mu, \rho_s) + W_p(\rho_s, \rho_t) + W_p(\rho_t, \nu) \leq W_p(\mu, \nu)$$

and by the triangle inequality

$$W_p(\mu, \rho_s) + W_p(\rho_s, \rho_t) + W_p(\rho_t, \nu) \geq W_p(\mu, \nu).$$

Hence we must have equality and in particular $W_p(\rho_s, \rho_t) = |s - t| \cdot W_p(\mu, \nu)$.

□

References

- L. Ambrosio, N. Fusco, and D. Pallara. *Functions of Bounded Variation and Free Discontinuity Problems*. Oxford mathematical monographs. Oxford University Press, 2000.
- L. Ambrosio and P. Tilli. *Topics on Analysis in Metric Spaces*. Number 25 in Oxford Lecture Series in Mathematics and its Applications. Oxford University Press, 2004.
- L. Ambrosio, N. Gigli, and G. Savaré. *Gradient Flows in Metric Spaces and in the Space of Probability Measures*. Lectures in Mathematics. Birkhäuser Boston, 2005.
- H. H. Bauschke and P. L. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. CMS Books in Mathematics. Springer, 1st edition, 2011.
- Y. Brenier. Polar factorization and monotone rearrangement of vector-valued functions. *Comm. Pure Appl. Math.*, 44(4):375–417, 1991.
- A. J. Kurdila and M. Zabrankin. *Convex functional analysis*, volume 1 of *Systems and Control: Foundations and Applications*. Birkhäuser, 2005.
- R. T. Rockafellar. Duality and stability in extremum problems involving convex functions. *Pacific J. Math*, 21(1):167–187, 1967.
- W. Rudin. *Real and complex analysis*. McGraw-Hill Book Company, 3rd edition, 1986.
- F. Santambrogio. *Optimal Transport for Applied Mathematicians*, volume 87 of *Progress in Nonlinear Differential Equations and Their Applications*. Birkhäuser Boston, 2015.
- C. Villani. *Optimal Transport: Old and New*, volume 338 of *Grundlehren der mathematischen Wissenschaften*. Springer, 2009.
- W. P. Ziemer. *Weakly Differentiable Functions*, volume 120 of *Graduate Texts in Mathematics*. Springer New York, 1989.