

Convergence of a variable metric scheme with application to optimal transport

Adrien Vacher

ORSAY – 20th November, 2025

Joint work with François-Xavier Vialard – accepted at ICML 2023

Introduction

A variable metric scheme

Application to the minimization of the semi-dual

Conclusion

Introduction

The (quadratic) optimal transport problem

Given two probability measures μ, ν on \mathbb{R}^d , the squared Wasserstein distance is defined as

$$W_2^2(\mu, \nu) = \inf_{\pi \geq 0} \int \|x - y\|^2 d\pi(x, y) + \iota(\pi_1 = \mu) + \iota(\pi_2 = \nu),$$

which, under several mild conditions (Santambrogio, 2015), also admits the following dual formulation

$$W_2^2(\mu, \nu) = \sup_{\varphi(x) + \psi(y) \leq \|x - y\|^2} \int \varphi(x) d\mu(x) + \int \psi(y) d\nu(y).$$

Theorem (Brenier, 90s)

If μ has compact, convex support and has a C^1 density bounded from above and below, then φ is C^2 and verifies

$$T_{\#}(\mu) = \nu.$$

where $T(x) = x - \nabla\varphi(x)$ is the optimal transport map.

- For ML applications, one seeks to estimate T by \hat{T} from empirical measures $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$, $\hat{\nu} = \frac{1}{n} \sum_{i=1}^n \delta_{y_i}$ where the (x_i) 's (resp. the (y_i) 's) are iid samples from μ (resp. ν).
- Ideally, \hat{T} should be "not too expensive" to compute and get closer to the original OT map T as n grows.

A popular approach is to regularize the primal with the addition of an entropic term $+\varepsilon \text{KL}(\pi|\mu \otimes \nu)$ (Léonard, 2014). One can then recover approximate potentials $(\hat{\varphi}_\varepsilon, \hat{\psi}_\varepsilon)$ by solving the following dual problem

$$\sup_{\varphi, \psi} \int \varphi(x) d\hat{\mu}(x) + \int \psi(y) d\hat{\nu}(y) - \varepsilon \int e^{\frac{\varphi(x) + \psi(y) - \|x-y\|^2}{\varepsilon}} d\hat{\mu}(x) d\hat{\nu}(y).$$

This problem admits a finite reparametrization and can be approximately solved by the well-known Sinkhorn algorithm (Cuturi, 2013).

Drawbacks of entropic approach

1. Slow of convergence of Sinkhorn as regularization ε decreases: after k iterations, the error scales as $\min((1 - e^{-1/\varepsilon})^{2k}, 1/(k\varepsilon))$ (Peyré and Cuturi, 2019; Léger, 2021).
2. Approximation of the original potential degrades exponentially with the dimension:
 $\|\nabla\varphi - \nabla\hat{\varphi}_\varepsilon\| \sim \varepsilon^{-d/2}/\sqrt{n}$ (Pooladian and Niles-Weed, 2021).

These two drawbacks find the same origin: as ε goes to 0 we get closer to the original linear dual formulation. Hence

1. The problem lacks strong convexity, leading to slow rates.
2. We end up by simply discretizing the cost constraint $\varphi(\mathbf{x}) + \psi(\mathbf{y}) \leq \|\mathbf{x} - \mathbf{y}\|^2$ on the sample points, which is a way too loose relaxation (Vacher et al., 2021).

We pre-optimize the dual w.r.t the potential ψ so that the OT map is given by the gradient of:

$$\arg \min_f J(f) := \int f(x) d\mu(x) + \int f^*(y) d\nu(y),$$

where f^* given by: $f^*(y) = \sup_x x^\top y - f(x)$.

- The new objective gains convexity: if f has M -Lipschitz gradient $J(f) - J(f_0) \geq \frac{1}{2M} \|\nabla f - T\|_{L^2(\mu)}^2$ where $T = \nabla f_0$ (Hütter and Rigollet, 2021).
- The cost constraint is removed: more appealing statistical rates (Divol et al., 2025).

\implies we aim at developing an algorithm to solve the semi-dual formulation when $\mu, \nu = \hat{\mu}, \hat{\nu}$.

A variable metric scheme

Euclidean gradient descent

The simplest way to minimize a Euclidean function $\theta : \mathbb{R}^d \rightarrow \mathbb{R}$ is to apply gradient descent

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \tau \nabla \theta(\mathbf{x}_k).$$

However, the gradient is not an intrinsic quantity, it is implicitly dependent of the metric:

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} d\theta_{\mathbf{x}_k}(\mathbf{x}) + \frac{1}{2\tau} \|\mathbf{x} - \mathbf{x}_k\|^2,$$

with $d\theta_y$ the differential of θ at y . Can be go beyond a euclidean metric?

Mirror descent: relative smoothness/strong convexity (Lu et al., 2018)

Let h be a differentiable convex function. Instead of using the euclidean metric, one can use the Bregman divergence associated to h and solve

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} d\theta_{\mathbf{x}_k}(\mathbf{x}) + \frac{1}{2\tau} \Delta_h(\mathbf{x}, \mathbf{x}_k),$$

with $\Delta_h(\mathbf{x}, \mathbf{x}_k) = h(\mathbf{x}_k) - h(\mathbf{x}) - dh_{\mathbf{x}_k}(\mathbf{x} - \mathbf{x}_k)$. If θ is convex with minimum θ^* and:

- if there exists $\beta \geq 0$ s.t. $\Delta_\theta(\mathbf{x}, \mathbf{y}) \leq \beta \Delta_h(\mathbf{x}, \mathbf{y})$ and $1/\tau = \beta$, then $\theta(\mathbf{x}_k) - \theta^* \leq O(\beta/k)$ (relative smoothness).
- if there further exists $\alpha > 0$ s.t. $\Delta_\theta(\mathbf{x}, \mathbf{y}) \geq \alpha \Delta_h(\mathbf{x}, \mathbf{y})$, then $\theta(\mathbf{x}_k) - \theta^* \leq O(\frac{\alpha}{(1+\alpha/(\beta-\alpha))^{k-1}})$ (relative strong convexity).

Taking $h(\mathbf{x}) = \|\mathbf{x}\|^2/2$ recovers standard euclidean g.d. results.

We wish to apply a mirror descent type algorithm for the semi-dual: which divergence should we choose?

Proposition (Stability of the semi-dual)

If (f, g) are γ -strongly convex potentials then

$$\Delta_J(f, g) \leq \frac{1}{2\gamma} \|\nabla f - \nabla g\|_{L^2((\nabla g^*)_{\#}(\nu))}^2.$$

If (f, g) have M -Lipschitz gradients and ∇g^* is well-defined over the support of ν then

$$\Delta_J(f, g) \geq \frac{1}{2M} \|\nabla f - \nabla g\|_{L^2((\nabla g^*)_{\#}(\nu))}^2.$$

- These results generalize previous stability results recovered for $g = f_0$ the optimal transport potential.
- Unfortunately, we cannot directly apply the previous results on mirror descent since $\|\nabla f - \nabla g\|_{(\nabla g^*)_{\#}(\nu)}^2$ cannot be taken as the Bregman divergence of a fixed function h : such a function h would solve for all (f, g)

$$\lim_{\lambda \rightarrow \infty} \frac{h(\lambda f)}{\lambda^2} = \|\nabla f\|_{L^2((\nabla g^*)_{\#}(\nu))}^2.$$

The left hand side does not depend on g while the right hand side does.

- We thus study the convergence of the following type of schemes

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x} \in C} d\theta_{\mathbf{x}_k}(\mathbf{x}) + \frac{1}{2\tau} \|\mathbf{x} - \mathbf{x}_k\|_{\mathbf{x}_k}^2 ,$$

with C a convex set and $\|\cdot\|_{\mathbf{x}}$ a pseudo-metric depending of the point \mathbf{x} . We need the constraint C since we need the potentials over which we optimize to remain strongly convex/and or smooth in order to get stability.

- This variable metric scheme was previously studied but it either assumes that the metric gets finer across each iterations $\|\cdot\|_{\mathbf{x}_k} \leq \|\cdot\|_{\mathbf{x}_{k+1}}$ either assume that the metrics are all equivalent (Combettes and Vũ, 2014) \implies this is typically not the case for the semi-dual.

Convergence result: convexity and relative smoothness

Theorem

Let E be a Banach space, let F be a real-valued convex function with Gateaux derivative dF satisfying for all $(x, y) \in E$, $\Delta_F(x, y) \leq \frac{\beta}{2} A^y(x - y)$ where for all $y \in E$, $A^y(\cdot)$ is a 2-homogeneous form over E depending on y and where β is a strictly positive constant and let $C \subset E$ be a closed convex subset of E . Assuming that $\sup_{(x,y) \in C^2} A^y(x - y) \leq K$, that a minimizer $\bar{x} \in C$ exists and that the iterates $x_0 \in C$, (x_k) generated as

$$x_{k+1} \in \arg \min_{x \in C} dF(x_k)(x - x_k) + \frac{\beta}{2} A^{x_k}(x - x_k), \quad (1)$$

exist, we have $F(x_k) - F(\bar{x}) \leq \frac{2\beta K}{k+1}$.

Theorem

Let E be a Banach space, let F be a real-valued convex function with Gateaux derivative dF , and let $C \subset E$ be a closed convex subset of E . If there exists $\alpha, \beta > 0$ and $A^y(\cdot)$ a 2-homogeneous form such that for all $(x, y) \in E$,
 $\frac{\alpha}{2} A^y(x - y) \leq \Delta_F(x, y) \leq \frac{\beta}{2} A^y(x - y)$. If a minimizer $\bar{x} \in C$ and the iterates $x_0 \in C$, (x_k) generated by (1) exist, we have

$$F(x_k) - F(\bar{x}) \leq \left(1 - \frac{\alpha}{\beta}\right)^k [F(x_0) - F(\bar{x})].$$

Application to the minimization of the
semi-dual

An exponentially convergent scheme

Theorem (Balanced case, smooth and strongly convex)

Let C be a convex set of γ -strongly convex, M -smooth functions. The minimum $\min_{f \in C} J(f)$ is attained at $\bar{f} \in C$. Furthermore, the iterates

$$f_{k+1} = \arg \min_{f \in C} dJ(f_k)(f - f_k) + \frac{1}{2\gamma} \|\nabla f - \nabla f_k\|_{L^2((\nabla f_k^*)_{\#}(\nu))}^2,$$

are well defined and verify

$$J(f_{k+1}) - J(\bar{f}) \leq \left(1 - \frac{\gamma}{M}\right)^k (J(f_1) - J(\bar{f})).$$

- By standard regularity of quadratic OT results, the ground truth optimal transport potential verifies the smoothness and strong convexity constraints when μ, ν have compact, convex support (Caffarelli, 2000). Hence in some cases, these extra constraints are benign.
- These results generalize to unbalanced quadratic OT.
- In Jacobs and Léger (2020), a similar type of scheme is considered yet instead the authors assume that μ, ν are supported on a fixed grid Ω and thus consider the fixed metric $\|\nabla f - \nabla g\|_{L^2(\Omega)}^2$. Yet the potentials are not constrained.

Practical implementation?

Recalling that $dJ(f)(g) = \int g(x)(d\mu(x) - d(\nabla f^*)_{\#}(\nu)(x))$, the iterations we need to solve for empirical measures $\hat{\mu}, \hat{\nu}$ are given by

$$\begin{aligned} f_{k+1} = \arg \min_{f \in C} & \frac{1}{n} \sum_{i=1}^n f(x_i) - f((\nabla f_k^*)(y_i)) \\ & + \frac{1}{2n\gamma} \sum_{i=1}^n \|\nabla f((\nabla f_k^*)(y_i)) - y_i\|^2. \end{aligned}$$

Even though we observe that f and ∇f are evaluated on a finite set of points, the set C remains infinite dimensional and the problem above may not admit a finite reparametrization.

Finite reparametrization: the case of smooth and strongly convex functions.

Proposition (Taylor et al. (2017))

For $C := C_{\gamma, M}$, the set of M -smooth, γ -strongly convex functions, the previous problem can be reformulated as

$$\inf_{\substack{u \in \mathbb{R}^{2n} \\ g \in \mathbb{R}^{(2n) \times d}}} \sum_{i=1}^n u_i - \sum_{i=n+1}^{2n} u_i + \sum_{i=n}^{n+1} \|g_i - z_i\|^2 / (2\gamma),$$
$$u_i \geq u_j + g_j^\top (z_i - z_j) + \frac{\|(g_i - g_j) / \sqrt{M} - \sqrt{\lambda}(z_i - z_j)\|^2}{2(1 - \lambda/M)}$$

where $z_i = x_i$ if $i \in \{1, \dots, n\}$ and $z_i = y_i$ if $i \in \{n+1, \dots, 2n\}$. Furthermore, the cost to solve this problem with an Interior Point Method (Nesterov and Nemirovskii, 1994) with precision τ requires $O(n^3 d^3 \log(\log(\tau)))$ operations.

Proposition

Denoting (u, g) a solution of the previous problem, then point-wise, the corresponding potential $f(x)$ as well as its gradient $\nabla f(x)$ can be recovered as the solution (v, p) of the following program

$$\begin{aligned} & \min_{v, g} v \\ & v \geq u_i + g_i(x - z_i) + \frac{\|(g_i - p)/\sqrt{M} - \sqrt{\lambda}(z_i - x)\|^2}{2(1 - \lambda/M)}. \end{aligned}$$

In particular, one can compute $\nabla f^*(y)$ point-wise via standard gradient descent. Note that since explicit regularity bounds are known for the recovered potential f , we can achieve optimal rates.

A word about the statistical/computational trade-off

- Even though the number of iterations is $O(1)$, the overall complexity is dominated by the cost per iteration and scales as $\tilde{O}(n^3)$ which may seem disappointing. However, it yields an approximation guarantee of the original OT map of order $n^{-2/d}$ hence the overall approximation cost of the original OT map scales as $\tau^{-\min(7, 3d/2+1)}$
- For Sinkhorn, once the regularization parameter is optimally chosen, we obtain $\tau^{-2(d+1)-7}$ (Vacher and Vialard, 2023).

For $M = +\infty$, the problem we seek to solve is now

$$\inf_{f \in C_{\lambda, \infty}} J(f).$$

This problem was recently studied in Gallouët et al. (2024). Strikingly, it admits a weak optimal transport reformulation. Especially we may obtain a better complexity than $\tilde{O}(n^3)$.

We apply this method to compute the OT map between μ , a discrete uniform probability measure across 50-regularly spaced points on $[-1, 1]$ and $\nu = (\nabla f_0)_\#(\mu)$ with f_0 convex so that an OT map T is given by $T = \nabla f_0$. We choose for this experiment, $f_0 = |x| + 0.25x^2/2$.

Frank-Wolfe as a baseline

In the introduction, we recalled that gradient descent was implicitly dependent of a metric. One way to avoid choosing such a metric is simply to compute the iterates as

$$\tilde{f} = \arg \min_{f \in C_{\lambda, M}} dF(f_k)(f),$$

and take $f_{k+1} = f_k + \frac{2}{k+1}(\tilde{f} - f_k)$ which is exactly the Frank-Wolfe algorithm for the semi-dual.

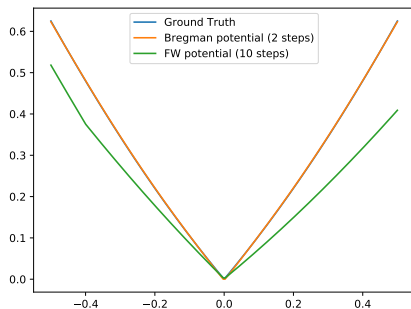


Figure 1: Potential generated by Frank-Wolfe after 10 iterations (in green) vs generated by our algorithm after 2 iterations (in orange) vs. ground truth (in blue).

Towards $O(n)$ optimal transport?

The choice $C_{\lambda,M}$ is way too expressive \implies not so good
statistical rates + computationally quite expensive. How about
a parametric model? For instance take

$$C_{z,\varepsilon} := \{f(x) = \sum_{i=1}^p w_i \sqrt{\|x - z_i\|^2 + \varepsilon} \mid w_i \geq 0\},$$

where the centroids z_i and the regularization ε are fixed. In this
case the iterates cost $O(p^2dn)$ and if the ground truth OT
potential belongs to $C_{z,\varepsilon}$, we can recover an approximation rate
in $O(1/n)$.

Conclusion

How to choose the set C such that:

1. The iterations admit a finite reparametrization and are cheap to compute (typically $O(n)$).
2. C is not "too complex" in order to recover appealing statistical rates.
3. C can approximate reasonably well a large class of OT potentials \implies need additional structure on OT, probably very hard problem.

Mirror descent is roughly Riemannian gradient descent with hessian metric when the step size goes to zero. What is the space $(C_{\lambda,M}, g_f(p, p) = \|\nabla p\|_{L^2((\nabla f^*)_{\#}(\nu))})$?

References

- Caffarelli, L. A. (2000). Monotonicity properties of optimal transportation and the fkg and related inequalities. *Communications in Mathematical Physics*.
- Combettes, P. L. and Vũ, B. C. (2014). Variable metric forward–backward splitting with applications to monotone inclusions in duality. *Optimization*.
- Cuturi, M. (2013). Sinkhorn distances: Lightspeed computation of optimal transport. In *NeurIPS*.
- Divol, V., Niles-Weed, J., and Pooladian, A.-A. (2025). Optimal transport map estimation in general function spaces. *The Annals of Statistics*.

- Gallouët, T., Natale, A., and Todeschi, G. (2024). From geodesic extrapolation to a variational bdf2 scheme for wasserstein gradient flows. *Mathematics of Computation*.
- Hütter, J.-C. and Rigollet, P. (2021). Minimax estimation of smooth optimal transport maps. *The Annals of Statistics*.
- Jacobs, M. and Léger, F. (2020). A fast approach to optimal transport: the back-and-forth method. *Numerische Mathematik*.
- Léger, F. (2021). A gradient descent perspective on sinkhorn. *Applied Mathematics and Optimization*.
- Lu, H., Freund, R. M., and Nesterov, Y. (2018). Relatively smooth convex optimization by first-order methods, and applications. *SIAM Journal on Optimization*.

- Léonard, C. (2014). A survey of the schrodinger problem and some of its connections with optimal transport. Discrete and Continuous Dynamical Systems.
- Nesterov, Y. E. and Nemirovskii, A. S. (1994). Interior point methods in convex programming: theory and applications. Society for Industrial and Applied Mathematics, Philadelphia.
- Peyré, G. and Cuturi, M. (2019). Computational optimal transport: With applications to data science. Foundations and Trends in Machine Learning.
- Pooladian, A.-A. and Niles-Weed, J. (2021). Entropic estimation of optimal transport maps. In NeurIPS, OTML workshop.
- Santambrogio, F. (2015). Optimal transport for applied mathematicians. Birkäuser, NY.

- Taylor, A. B., Hendrickx, J. M., and Glineur, F. (2017). Exact worst-case performance of first-order methods for composite convex optimization. *SIAM Journal on Optimization*.
- Vacher, A., Muzellec, B., Rudi, A., Bach, F., and Vialard, F.-X. (2021). A dimension-free computational upper-bound for smooth optimal transport estimation. In *COLT*.
- Vacher, A. and Vialard, F.-X. (2023). Semi-dual unbalanced quadratic optimal transport: fast statistical rates and convergent algorithm. In *ICML*.