

# Mirror and Preconditioned Gradient Descent in Wasserstein Space

**Clément Bonet**<sup>1</sup>, Théo Uscidda<sup>1</sup>, Adam David<sup>2</sup>,  
Pierre-Cyril Aubin-Frankowski<sup>3</sup>, Anna Korba<sup>1</sup>

<sup>1</sup>ENSAE, CREST, Institut Polytechnique de Paris

<sup>2</sup>TU Berlin

<sup>3</sup>TU Wien

NeurIPS in Paris 2024

05/12/2024



# Motivations

Let  $\mathcal{P}_2(\mathbb{R}^d) = \{\mu \in \mathcal{P}(\mathbb{R}^d), \int \|x\|_2^2 d\mu(x) < \infty\}$ ,  $\mathcal{F} : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ .

**Goal:**

$$\min_{\mu \in \mathcal{P}_2(\mathbb{R}^d)} \mathcal{F}(\mu)$$

**Applications:**

- Sampling from  $\nu \propto e^{-V}$  (Wibisono, 2018)
- Modeling dynamic of population of cells (Schiebinger et al., 2019)
- Learning neural networks (Mei et al., 2018; Chizat and Bach, 2018)

## Example of functionals

- Free energies:  $\mathcal{F}(\mu) = \int V d\mu + \iint W(x, y) d\mu(x)d\mu(y) + \mathcal{H}(\mu)$  where  $\mathcal{H}(\mu) = \int \log(\mu(x)) d\mu(x)$  for  $\mu \ll \text{Leb}$
- $\mathcal{F}(\mu) = \text{KL}(\mu||\nu) = \int V d\mu + \mathcal{H}(\mu) + \text{cst}$  for sampling from  $\nu \propto e^{-V(x)}$
- $\mathcal{F}(\mu) = D(\mu, \nu)$  for sampling from  $\nu$

# Detour by $\mathbb{R}^d$

Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ .

**Goal:**  $\min_{x \in \mathbb{R}^d} f(x)$ .

- Gradient descent:

$$\forall k \geq 0, x_{k+1} = x_k - \tau \nabla f(x_k)$$

- Non-increasing if  $f$   $\beta$ -smooth
- Converge if  $f$   $\beta$ -smooth and  $\alpha$ -strongly convex (i.e.  $f - \alpha \frac{\|\cdot\|_2^2}{2}$  convex)
- Mirror descent (Lu et al., 2018):

$$\forall k \geq 0, x_{k+1} = \nabla \phi^* (\nabla \phi(x_k) - \tau \nabla f(x_k))$$

- Non-increasing if  $f$   $\beta$ -smooth relative to  $\phi$  (i.e.  $\beta\phi - f$  convex)
- Converge if  $f$   $\beta$ -smooth and  $\alpha$ -convex relative to  $\phi$  (i.e.  $f - \alpha\phi$  convex)
- Preconditioned gradient descent (Maddison et al., 2021):

$$\forall k \geq 0, x_{k+1} = x_k - \tau \nabla h^* (\nabla f(x_k))$$

- Non-increasing if  $h^*$   $\beta$ -smooth relative to  $f^*$  (with  $f^*$  the Legendre transform)
- Converge if  $h^*$   $\beta$ -smooth and  $\alpha$ -convex relative to  $f^*$

# Wasserstein Geometry (Ambrosio et al., 2005)

## Definition (Wasserstein distance)

Let  $\mu, \nu \in \mathcal{P}_2(\mathbb{R}^d)$  and denote by  $\Pi(\mu, \nu)$  the set of coupling between  $\mu, \nu$ . Then, the Wasserstein distance is

$$W_2^2(\mu, \nu) = \inf_{\gamma \in \Pi(\mu, \nu)} \int \|x - y\|_2^2 d\gamma(x, y).$$

## Properties:

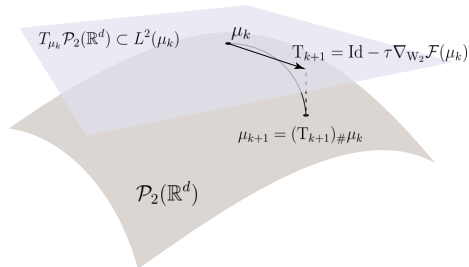
- $W_2$  distance,  $(\mathcal{P}_2(\mathbb{R}^d), W_2)$ : Wasserstein space
- Riemannian structure (with geodesics and tangent space  $\mathcal{T}_\mu \mathcal{P}_2(\mathbb{R}^d) \subset L^2(\mu)$ )
- Wasserstein gradient  $\nabla_{W_2} \mathcal{F}(\mu) \in \mathcal{T}_\mu \mathcal{P}_2(\mathbb{R}^d)$  of  $\mathcal{F}$  at  $\mu$  satisfies for all  $T \in L^2(\mu)$ ,

$$\mathcal{F}(T \# \mu) = \mathcal{F}(\mu) + \langle \nabla_{W_2} \mathcal{F}(\mu), T - \text{Id} \rangle_{L^2(\mu)} + o(\|T - \text{Id}\|_{L^2(\mu)})$$

## Example

- $\mathcal{V}(\mu) = \int V d\mu, \nabla_{W_2} \mathcal{V}(\mu) = \nabla V$
- $\mathcal{W}(\mu) = \iint W(x - y) d\mu(x) d\mu(y), \nabla_{W_2} \mathcal{W}(\mu) = \nabla W \star \mu$

# Wasserstein Gradient Descent



## Wasserstein Gradient Descent:

$$\begin{cases} T_{k+1} = \operatorname{argmin}_{T \in L^2(\mu_k)} \frac{1}{2} \|T - \text{Id}\|_{L^2(\mu_k)}^2 + \tau \langle \nabla_{W_2} \mathcal{F}(\mu_k), T - \text{Id} \rangle_{L^2(\mu_k)} \\ \mu_{k+1} = (T_{k+1})\# \mu_k \end{cases}$$

Taking the FOC:  $T_{k+1} = \text{Id} - \tau \nabla_{W_2} \mathcal{F}(\mu_k)$

**Particle approximation:**  $\hat{\mu}_k^n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i^k}$ ,  $x_i^{k+1} = T_{k+1}(x_i^k)$  for all  $i \in \{1, \dots, n\}$ .

# Contributions

Study schemes of the form

$$\begin{cases} \mathbb{T}_{k+1} = \operatorname{argmin}_{\mathbb{T} \in L^2(\mu_k)} d(\mathbb{T}, \operatorname{Id}) + \tau \langle \nabla_{W_2} \mathcal{F}(\mu_k), \mathbb{T} - \operatorname{Id} \rangle_{L^2(\mu_k)} \\ \mu_{k+1} = (\mathbb{T}_{k+1})_{\#} \mu_k, \end{cases}$$

and provide **convergence conditions**.

Considered divergences:

- For  $d(\mathbb{T}, \operatorname{Id}) = \frac{1}{2} \|\mathbb{T} - \operatorname{Id}\|_{L^2(\mu)}^2$ : **Wasserstein gradient descent**
- For  $d_{\phi_{\mu}}(\mathbb{T}, \operatorname{Id}) = \phi_{\mu}(\mathbb{T}) - \phi_{\mu}(\operatorname{Id}) - \langle \nabla \phi_{\mu}(\operatorname{Id}), \mathbb{T} - \operatorname{Id} \rangle_{L^2(\mu)}$  (**Bregman divergence** on  $L^2(\mu)$ ): extends **Mirror Descent** ([Beck and Teboulle, 2003](#)) to  $\mathcal{P}_2(\mathbb{R}^d)$ .
- For  $d(\mathbb{T}, \operatorname{Id}) = \int h(\mathbb{T}(x) - x) d\mu(x)$ : extends **Preconditioned Gradient Descent** ([Maddison et al., 2021](#)) to  $\mathcal{P}_2(\mathbb{R}^d)$ .

# Relative Convexity and Smoothness

Let  $\phi_\mu, \psi_\mu : L^2(\mu) \rightarrow \mathbb{R}$  convex,  $\mathcal{F}, \mathcal{G} : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ .

Define  $\tilde{\mathcal{F}}_\mu(T) = \mathcal{F}(T \# \mu)$ ,  $\tilde{\mathcal{G}}_\mu(T) = \mathcal{G}(T \# \mu)$ .

Relative smoothness/convexity on  $L^2(\mu)$ :

- $\phi_\mu$  is  $\beta$ -smooth relative to  $\psi_\mu$  if for all  $T, S \in L^2(\mu)$ ,  $d_{\phi_\mu}(T, S) \leq \beta d_{\psi_\mu}(T, S)$ .
- $\phi_\mu$  is  $\alpha$ -convex relative to  $\psi_\mu$  if for all  $T, S \in L^2(\mu)$ ,  $d_{\phi_\mu}(T, S) \geq \alpha d_{\psi_\mu}(T, S)$ .

Relative smoothness/convexity along a curve  $\mu_t = (T_t) \# \mu$  with  $T_t = (1-t)S + tT$  for all  $t \in [0, 1]$ ,  $T, S \in L^2(\mu)$ .

- $\mathcal{F}$   $\beta$ -smooth relative to  $\mathcal{G}$  along  $t \mapsto \mu_t$  if  $\forall s, t \in [0, 1]$ ,

$$d_{\tilde{\mathcal{F}}_\mu}(T_s, T_t) \leq \beta d_{\tilde{\mathcal{G}}_\mu}(T_s, T_t)$$

- $\mathcal{F}$   $\alpha$ -convex relative to  $\mathcal{G}$  along  $t \mapsto \mu_t$  if  $\forall s, t \in [0, 1]$ ,

$$d_{\tilde{\mathcal{F}}_\mu}(T_s, T_t) \geq \alpha d_{\tilde{\mathcal{G}}_\mu}(T_s, T_t)$$

# Mirror Descent on the Wasserstein Space

Let  $\phi_\mu : L^2(\mu) \rightarrow \mathbb{R}$  be strictly convex, proper and differentiable.

**Mirror Descent scheme:**

$$\begin{cases} \mathbb{T}_{k+1} = \operatorname{argmin}_{\mathbb{T} \in L^2(\mu_k)} d_{\phi_{\mu_k}}(\mathbb{T}, \operatorname{Id}) + \tau \langle \nabla_{W_2} \mathcal{F}(\mu_k), \mathbb{T} - \operatorname{Id} \rangle_{L^2(\mu_k)} \\ \mu_{k+1} = (\mathbb{T}_{k+1}) \# \mu_k. \end{cases}$$

By FOC:  $\nabla \phi_{\mu_k}(\mathbb{T}_{k+1}) = \nabla \phi_{\mu_k}(\operatorname{Id}) - \tau \nabla_{W_2} \mathcal{F}(\mu_k)$

**Computing the scheme:**

- For  $\phi_\mu(\mathbb{T}) = \int V \circ \mathbb{T} \, d\mu$ ,  $\mathbb{T}_{k+1} = \nabla V^* \circ (\nabla V - \tau \nabla_{W_2} \mathcal{F}(\mu_k))$
- For  $\phi_\mu$  pushforward compatible (i.e.  $\phi_\mu(\mathbb{T}) = \phi(\mathbb{T} \# \mu)$  with  $\phi : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ ):

$$\nabla_{W_2} \phi(\mu_{k+1}) \circ \mathbb{T}_{k+1} = \nabla_{W_2} \phi(\mu_k) - \tau \nabla_{W_2} \mathcal{F}(\mu_k)$$

In general: implicit in  $\mathbb{T}_{k+1} \rightarrow$  Newton method



# Descent Lemma

Let  $\phi_\mu : L^2(\mu) \rightarrow \mathbb{R}$  be strictly convex, proper and differentiable.

**Mirror Descent scheme:**

$$\begin{cases} \mathbb{T}_{k+1} = \operatorname{argmin}_{\mathbb{T} \in L^2(\mu_k)} d_{\phi_{\mu_k}}(\mathbb{T}, \operatorname{Id}) + \tau \langle \nabla_{W_2} \mathcal{F}(\mu_k), \mathbb{T} - \operatorname{Id} \rangle_{L^2(\mu_k)} \\ \mu_{k+1} = (\mathbb{T}_{k+1})_{\#} \mu_k. \end{cases}$$

## Proposition (Descent Lemma)

*Assumptions:*

- For all  $k \geq 0$ ,  $\mathcal{F}$  is  $\beta$ -smooth relative to  $\phi$  along  $t \mapsto ((1-t)\operatorname{Id} + t\mathbb{T}_{k+1})_{\#} \mu_k$

Then, for all  $k \geq 0$ ,

$$\mathcal{F}(\mu_{k+1}) \leq \mathcal{F}(\mu_k) - \beta d_{\phi_{\mu_k}}(\operatorname{Id}, \mathbb{T}_{k+1}).$$

# Convergence

## Proposition

*Assumptions:* Let  $\beta > 0, \alpha \geq 0$  and  $T_{\phi_{\mu_k}}^{\mu_k, \mu^*} = \operatorname{argmin}_{T_{\# \mu_k = \mu^*}} d_{\phi_{\mu_k}}(T, \operatorname{Id})$ .

- $\mathcal{F}$   $\beta$ -smooth relative to  $\phi$  along  $t \mapsto ((1-t)\operatorname{Id} + tT_{k+1})_{\# \mu_k}$
- $\mathcal{F}$   $\alpha$ -convex relative to  $\phi$  along  $t \mapsto ((1-t)\operatorname{Id} + tT_{\phi_{\mu_k}}^{\mu_k, \mu^*})_{\# \mu_k}$
- Assume  $d_{\phi_{\mu_k}}(T_{\phi_{\mu_k}}^{\mu_k, \mu^*}, T_{k+1}) \geq d_{\phi_{\mu_{k+1}}}(T_{\phi_{\mu_{k+1}}}^{\mu_{k+1}, \mu^*}, \operatorname{Id})$

Then, for all  $k \geq 1$ ,  $\mathcal{F}(\mu_k) - \mathcal{F}(\mu^*) \leq \frac{\beta - \alpha}{k} d_{\phi_{\mu_0}}(T_{\phi_{\mu_0}}^{\mu_0, \mu^*}, \operatorname{Id})$ .

If  $\alpha > 0$ , for all  $k \geq 0$ ,  $d_{\phi_{\mu_k}}(T_{\phi_{\mu_k}}^{\mu_k, \mu^*}, \operatorname{Id}) \leq \left(1 - \frac{\alpha}{\beta}\right)^k d_{\phi_{\mu_0}}(T_{\phi_{\mu_0}}^{\mu_0, \mu^*}, \operatorname{Id})$ .

Let  $\phi_\mu$  be pushforward compatible. Define the OT problem:

$$\begin{aligned} W_{\phi}(\nu, \mu) &= \inf_{\gamma \in \Pi(\nu, \mu)} \phi(\nu) - \phi(\mu) - \int \langle \nabla_{W_2} \phi(\mu)(y), x - y \rangle d\gamma(x, y) \\ &\leq d_{\phi_\eta}(T, S) \quad \text{for } (T, S)_{\# \eta} \in \Pi(\nu, \mu) \end{aligned}$$

**Property:** If  $\mu \ll \operatorname{Leb}$  and  $\nabla_{W_2} \phi(\mu)$  is invertible, then  $\gamma^* = (T_{\phi_\mu}^{\mu, \nu}, \operatorname{Id})_{\# \mu}$ , and  $W_{\phi}(\nu, \mu) = d_{\phi_\mu}(T_{\phi_\mu}^{\mu, \nu}, \operatorname{Id})$ .

# Preconditioned GD

Let  $h : \mathbb{R}^d \rightarrow \mathbb{R}$  strictly convex, proper and differentiable.

**Preconditioned Gradient Descent scheme:** Let  $\phi_\mu^h(\mathbb{T}) = \int h \circ \mathbb{T} \, d\mu$ ,

$$\begin{cases} \mathbb{T}_{k+1} = \operatorname{argmin}_{\mathbb{T} \in L^2(\mu_k)} \phi_{\mu_k}^h \left( \frac{\operatorname{Id} - \mathbb{T}}{\tau} \right) \tau + \langle \nabla_{W_2} \mathcal{F}(\mu_k), \mathbb{T} - \operatorname{Id} \rangle_{L^2(\mu_k)} \\ \mu_{k+1} = (\mathbb{T}_{k+1}) \# \mu_k \end{cases}$$

By FOC:  $\mathbb{T}_{k+1} = \operatorname{Id} - \tau \nabla h^* \circ \nabla_{W_2} \mathcal{F}(\mu_k)$

Under relative smoothness and convexity of  $\phi_\mu^{h^*}$  relative to  $\mathcal{F}^*$ :

$$\forall k \geq 0, \phi_{\mu_{k+1}}^{h^*}(\nabla_{W_2} \mathcal{F}(\mu_{k+1})) \leq \phi_{\mu_k}^{h^*}(\nabla_{W_2} \mathcal{F}(\mu_k)) - \beta d_{\tilde{\mathcal{F}}_{\mu_k}}(\mathbb{T}_{k+1}, \operatorname{Id}),$$

$$\forall k \geq 1, \phi_{\mu_k}^{h^*}(\nabla_{W_2} \mathcal{F}(\mu_k)) - h^*(0) \leq \frac{\beta - \alpha}{k} (\mathcal{F}(\mu_0) - \mathcal{F}(\mu^*)).$$

# Showing Relative Smoothness and Convexity

Smoothness and convexity of  $\mathcal{F} : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$  relative to  $\phi : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ ?

- Let  $\mathcal{F}(\mu) = \int V d\mu$  and  $\phi(\mu) = \int U d\mu$ :

$V$   $\beta$ -smooth relative to  $U \implies \mathcal{F}$   $\beta$ -smooth relative to  $\phi$

$V$   $\alpha$ -convex relative to  $U \implies \mathcal{F}$   $\alpha$ -convex relative to  $\phi$

- Let  $\mathcal{F}(\mu) = \iint W(x - y) d\mu(x)d\mu(y)$  and  $\phi(\mu) = \iint K(x - y) d\mu(x)d\mu(y)$ :

$W$   $\beta$ -smooth relative to  $K \implies \mathcal{F}$   $\beta$ -smooth relative to  $\phi$

$W$   $\alpha$ -convex relative to  $K \implies \mathcal{F}$   $\alpha$ -convex relative to  $\phi$

- For  $\mathcal{F} = \mathcal{G} + \mathcal{H}$ ,  $d_{\tilde{\mathcal{F}}_\mu} = d_{\tilde{\mathcal{G}}_\mu} + d_{\tilde{\mathcal{H}}_\mu}$  and  $\mathcal{F}$  1-convex relative to  $\mathcal{G}$  and  $\mathcal{H}$
- In general: look at the Hessian

# Mirror Descent on Interaction Energy

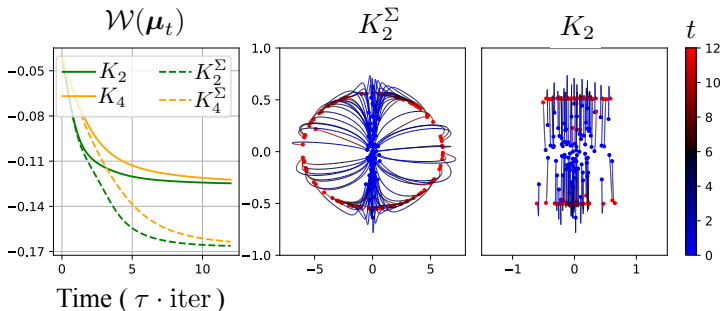
**Goal:** Let  $\Sigma \in S_d^{++}(\mathbb{R})$  possibly ill-conditioned,

$$\min_{\mu} \mathcal{W}(\mu) = \iint W(x - y) d\mu(x)d\mu(y) \quad \text{with} \quad W(z) = \frac{1}{4}\|z\|_{\Sigma^{-1}}^4 - \frac{1}{2}\|z\|_{\Sigma^{-1}}^2$$

Bregman potential:  $\phi_{\mu}(T) = \iint K(T(x) - T(y)) d\mu(x)d\mu(y)$  with

$$K_2(z) = \frac{1}{2}\|z\|_2^2, \quad K_2^{\Sigma}(z) = \frac{1}{2}\|z\|_{\Sigma^{-1}}^2,$$

$$K_4(z) = \frac{1}{4}\|z\|_2^4 + \frac{1}{2}\|z\|_2^2, \quad K_4^{\Sigma}(z) = \frac{1}{4}\|z\|_{\Sigma^{-1}}^4 + \frac{1}{2}\|z\|_{\Sigma^{-1}}^2.$$



# Mirror Descent on Gaussian

Goal:

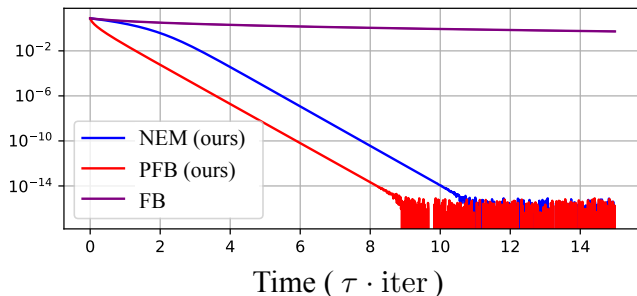
$$\min_{\mu} \mathcal{F}(\mu) = \int V d\mu + \mathcal{H}(\mu) \quad \text{with} \quad V(x) = \frac{1}{2} x^T \Sigma^{-1} x$$

→ minimum  $\mu^* = \mathcal{N}(0, \Sigma)$ .

Comparison between:

- Forward-Backward (FB) on the Bures-Wasserstein space (Diao et al., 2023)
- Preconditioned Forward-Backward (PFB) scheme with  $\phi(\mu) = \int V d\mu$
- NEM: MD with  $\phi(\mu) = \mathcal{H}(\mu)$  and restriction to Gaussian

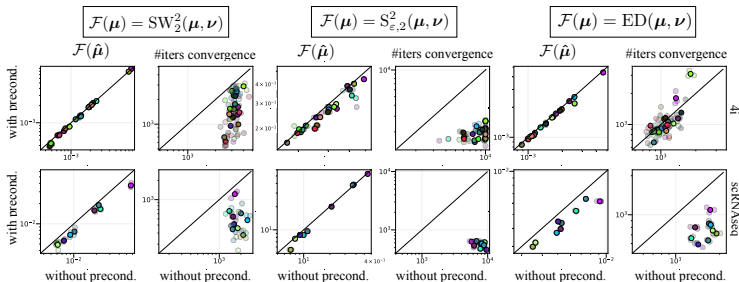
$$\text{KL}(\mu_t || \mu^*)$$



# Preconditioned GD on Single-Cells

**Goal:**  $\min_{\mu} \mathcal{F}(\mu) = D(\mu, \nu)$  with  $\mu_0$  untreated cell and  $\nu$  perturbed cell

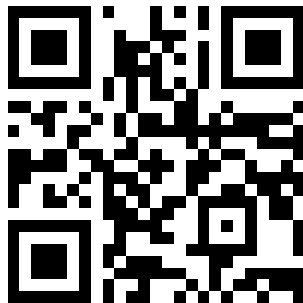
Use PGD with  $h^*(x) = (\|x\|_2^a + 1)^{1/a} - 1$  with  $a \in \{1.25, 1.5, 1.75\}$ , which is well suited to minimize functions growing in  $\|x - x^*\|^{a/(a-1)}$  near  $x^*$ .



- Rows: 2 profiling technologies
  - Columns/subcolumns: Different objectives  $\mathcal{F}$ /measure of convergence and number of iterations to converge
  - Points: For treatment  $i$ ,  $z_i = (x_i, y_i)$  with  $x_i$  value of  $\mathcal{F}(\hat{\mu}) = D(\hat{\mu}, \nu)$  (1st subcolumn) or number of iterations (2nd subcolumn) without preconditioning and  $y_i$  with preconditioning
  - Colors: treatments
- **Points below the diagonal: PGD provides a better minimum or converges faster**

# Thank you!

Paper: <https://arxiv.org/abs/2406.08938>





# References I

- Luigi Ambrosio, Nicola Gigli, and Giuseppe Savaré. *Gradient Flows: in Metric Spaces and in the Space of Probability Measures*. Springer Science & Business Media, 2005.
- Amir Beck and Marc Teboulle. Mirror Descent and Nonlinear Projected Subgradient Methods for Convex Optimization. *Operations Research Letters*, 31(3):167–175, 2003.
- Lenaïc Chizat and Francis Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. *Advances in neural information processing systems*, 31, 2018.
- Michael Ziyang Diao, Krishna Balasubramanian, Sinho Chewi, and Adil Salim. Forward-backward Gaussian variational inference via JKO in the Bures-Wasserstein Space. In *International Conference on Machine Learning*, pages 7960–7991. PMLR, 2023.
- Haihao Lu, Robert M Freund, and Yurii Nesterov. Relatively Smooth Convex Optimization by First-Order Methods, and Applications. *SIAM Journal on Optimization*, 28(1):333–354, 2018.

## References II

- Chris J Maddison, Daniel Paulin, Yee Whye Teh, and Arnaud Doucet. Dual Space Preconditioning for Gradient Descent. *SIAM Journal on Optimization*, 31(1): 991–1016, 2021.
- Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the landscape of two-layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33):E7665–E7671, 2018.
- Geoffrey Schiebinger, Jian Shu, Marcin Tabaka, Brian Cleary, Vidya Subramanian, Aryeh Solomon, Joshua Gould, Siyan Liu, Stacie Lin, Peter Berube, et al. Optimal-transport analysis of single-cell gene expression identifies developmental trajectories in reprogramming. *Cell*, 176(4):928–943, 2019.
- Andre Wibisono. Sampling as optimization in the space of measures: The langevin dynamics as a composite optimization problem. In *Conference on Learning Theory*, pages 2093–3027. PMLR, 2018.