



# Asymptotic Guarantees for Learning Generative Models with the Sliced-Wasserstein Distance

Kimia Nadjahi<sup>1</sup>   Alain Durmus<sup>2</sup>   Umut Şimşekli<sup>1,3</sup>   Roland Badeau<sup>1</sup>

<sup>1</sup> Télécom Paris   <sup>2</sup> ENS Paris-Saclay   <sup>3</sup> University of Oxford

# Minimum Distance Estimation

$$\hat{\theta}_n = \operatorname{argmin}_{\theta \in \Theta} \mathbf{D}(\hat{\mu}_n, \mu_\theta)$$

$\mathbf{D}$ : distance between distributions

$\hat{\mu}_n$ : empirical distribution of **data points**  $Y_1, \dots, Y_n$  i.i.d from  $\mu_\star$

$\mu_\theta$ : distribution parametrized by  $\theta \in \Theta$

# Minimum Distance Estimation

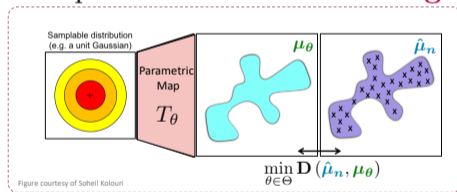
$$\hat{\theta}_n = \operatorname{argmin}_{\theta \in \Theta} \mathbf{D}(\hat{\mu}_n, \mu_\theta)$$

**D**: distance between distributions

$\hat{\mu}_n$ : empirical distribution of **data points**  $Y_1, \dots, Y_n$  i.i.d from  $\mu_\star$

$\mu_\theta$ : distribution parametrized by  $\theta \in \Theta$

Example: **Generative Modeling**



Asymptotic Guarantees for Learning Generative Models with the Sliced-Wasserstein Distance.

K. Nadjahi, A. Durmus, U. Şimşekli, R. Badeau

# Minimum **Expected** Distance Estimation

Directly optimizing  $\mu_\theta$  is often **not possible** (e.g. GANs)

$$\hat{\theta}_{n,m} = \operatorname{argmin}_{\theta \in \Theta} \mathbb{E} [\mathbf{D}(\hat{\mu}_n, \hat{\mu}_{\theta,m}) \mid Y_{1:n}]$$

$\hat{\mu}_{\theta,m}$ : empirical distribution of a sample  $Z_1, \dots, Z_m$  i.i.d. from  $\mu_\theta$

# Minimum Wasserstein Estimation

Choose  $\mathbf{D} = \mathbf{W}_p$  (Wasserstein distance of order  $p \geq 1$ )

- ✓ Robust and increasingly popular estimators: Wasserstein GAN [1], Wasserstein auto-encoders [2]
- ✓ Asymptotic guarantees [3]

[1] Arjovsky et al., 2017   [2] Tolstikhin et al., 2018   [3] Bernton et al., 2019

# Minimum Wasserstein Estimation

Choose  $\mathbf{D} = \mathbf{W}_p$  (Wasserstein distance of order  $p \geq 1$ )

- ✓ Robust and increasingly popular estimators: Wasserstein GAN [1], Wasserstein auto-encoders [2]
- ✓ Asymptotic guarantees [3]

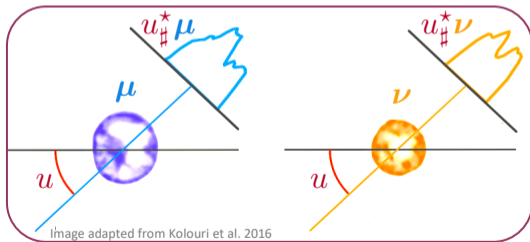
[1] Arjovsky et al., 2017   [2] Tolstikhin et al., 2018   [3] Bernton et al., 2019

- ✗  $\mathbf{W}_p$ : expensive + curse of dimensionality
- ✗ Central limit theorem in [3] valid in 1D

# Sliced-Wasserstein distance

In 1D,  $W_p$  has an analytical form  $\Rightarrow$  Motivates a practical alternative:

$$\text{SW}_p^p(\mu, \nu) = \int_{\mathbb{S}^{d-1}} W_p^p(u_{\#}^* \mu, u_{\#}^* \nu) d\sigma(u)$$



# Minimum Sliced-Wasserstein Estimation

$$\hat{\theta}_n = \operatorname{argmin}_{\theta \in \Theta} \mathbf{SW}_p(\hat{\mu}_n, \mu_\theta)$$

$$\hat{\theta}_{n,m} = \operatorname{argmin}_{\theta \in \Theta} \mathbb{E} [\mathbf{SW}_p(\hat{\mu}_n, \hat{\mu}_{\theta,m}) \mid Y_{1:n}]$$

Successful in generative modeling applications (e.g., SW-GAN, Deshpande et al., 2018)

**Asymptotic Guarantees for Learning Generative Models with the Sliced-Wasserstein Distance.**

K. Nadjahi, A. Durmus, U. Şimşekli, R. Badeau



# Minimum Sliced-Wasserstein Estimation

$$\hat{\theta}_n = \operatorname{argmin}_{\theta \in \Theta} \mathbf{SW}_p(\hat{\mu}_n, \mu_\theta)$$

$$\hat{\theta}_{n,m} = \operatorname{argmin}_{\theta \in \Theta} \mathbb{E} [\mathbf{SW}_p(\hat{\mu}_n, \hat{\mu}_{\theta,m}) \mid Y_{1:n}]$$

Successful in generative modeling applications (e.g., SW-GAN, Deshpande et al., 2018)

---


Our contributions: 

---

- Convergence in  $\mathbf{SW}_p \Rightarrow$  weak convergence of probability measures
- Existence and consistency of  $\hat{\theta}_n, \hat{\theta}_{n,m}$
- Central limit theorem for  $\hat{\theta}_n$ :  $\sqrt{n}$  convergence rate for any dimension

# Thank you!

## Our Poster: East Exhibition Hall B + C #226




IP PARIS

### Asymptotic Guarantees for Learning Generative Models with the Sliced-Wasserstein Distance

Kimia Nadjahi<sup>1</sup>, Alain Durmus<sup>2</sup>, Umut Şimşekli<sup>1,3</sup>, Roland Badeau<sup>1</sup>

(kimia.nadjahi, umut.simsekli, roland.badeau)@telecom-paris.fr, alain.durus@mla.ens-cachan.fr

1: LTCI, Télécom Paris, Institut Polytechnique de Paris 2: CMLA, ENS Paris-Saclay 3: Department of Statistics, University of Oxford



---

#### Minimum Distance Estimation

- Observations  $Y_{1:n} = (Y_1, \dots, Y_n)$ ,  $Y_i \in Y \subset \mathbb{R}^d$ , i.i.d. from  $\mu_n \in \mathcal{P}(Y)$ , with  $\mathcal{P}(Y)$ : set of probability measures on  $Y$ .
- A family of distributions on  $Y$  parametrized by  $\theta \in \Theta \subset \mathbb{R}^p$ :  $\mathcal{M} = \{\mu_\theta \in \mathcal{P}(Y), \theta \in \Theta\}$ .
- Purely generative model: We can generate  $n \in \mathbb{N}$  i.i.d. samples from  $\mu_\theta$ , but the likelihood is intractable.  $\mu_{\theta, n}$  is the empirical distribution.

Given  $Y_{1:n}$ , its empirical distribution  $\hat{\mu}_n$  and a distance  $D$  on  $\mathcal{P}(Y)$ , we perform **Minimum Distance Estimation (MDE)**:

$$\hat{\theta}_n = \operatorname{argmin}_{\theta \in \Theta} D(\mu_\theta, \hat{\mu}_n) \quad (1)$$

or **Minimum Expected Distance Estimation (MEDE)**:

$$\hat{\theta}_{n, \infty} = \operatorname{argmin}_{\theta \in \Theta} \mathbb{E}[D(\mu_\theta, \mu_n) | Y_{1:n}] \quad (2)$$

#### Theoretical Results

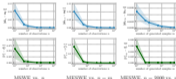
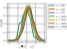
The convergence to **SW**, implies the weak convergence in  $\mathcal{P}(\mathbb{R}^d)$ .

#### Key assumptions.

- Continuity:** For any  $(\theta_n)_{n \in \mathbb{N}}$  in  $\Theta$  such that  $\lim_{n \rightarrow \infty} \mu_n(\theta_n, \theta) = 0$ ,  $\mathbb{A1}$ .  $\{\mu_n\}_{n \in \mathbb{N}}$  converges weakly  $\left(\frac{\cdot}{n}\right)$  to  $\mu_n$ .
- $\mathbb{A2}$ .  $\lim_{n \rightarrow \infty} \mathbb{E}[\operatorname{SW}_p(\mu_n, \mu_n)] = 0$ .
- Data-generating process:**  $\mathbb{A3}$ .  $\lim_{n \rightarrow \infty} \mathbb{E}[\operatorname{SW}_p(\mu_n, \mu_n)] = 0$ ,  $\mathbb{P}$ -almost surely.
- Bounded sets:** For some  $\epsilon > 0$ ,  $\mathbb{A4}$ .  $\Theta_\epsilon = \{\theta \in \Theta : \operatorname{SW}_p(\mu_n, \mu_n) \leq \epsilon_n + \epsilon\}$ , with  $\epsilon_n = \inf_{\theta \in \Theta} \operatorname{SW}_p(\mu_n, \mu_n)$ , is bounded.
- $\mathbb{A5}$ .  $\Theta_{\epsilon_n} = \{\theta \in \Theta : \operatorname{SW}_p(\mu_n, \mu_n) \leq \epsilon_n + \epsilon\}$ , with  $\epsilon_n = \inf_{\theta \in \Theta} \operatorname{SW}_p(\mu_n, \mu_n)$ , is bounded almost surely.

#### Numerical Experiments

- Multivariate Gaussians.**  $\mathcal{M} = \{N(\mu, \Sigma) : \mu \in \mathbb{R}^2, \Sigma \succ 0\}$ , and  $(\mu_n, \Sigma_n) = (0, I)$ .



Comparison Wasserstein and SW    MSWE vs.  $n = 10$     MSWE vs.  $n = 200$  vs.  $n = 100$

---


#### Optimal Transport (OT) Metrics

For  $p \geq 1$ ,  $\mathcal{P}_p(Y)$ : set of probability measures on  $Y$  with finite  $p$ th moment. Let  $\mu, \nu \in \mathcal{P}_p(Y)$ .

**Wasserstein distance ( $W_p$ ).** Computationally expensive, except in  $\mathbb{R}^d$  ( $Y \subset \mathbb{R}$ ) - analytical form.

**Sliced-Wasserstein (SW) distances.**  $\mathbb{R}^{2d-1}$ :  $d$ -dimensional unit sphere,  $\sigma$ : uniform distribution on  $\mathbb{S}^{2d-1}$ .

Practical metric based on projection:  $W_p(\mu, \nu) = \int_{\mathbb{S}^{2d-1}} W_p(\mu|_{\sigma}, \nu|_{\sigma}) d\sigma$


$$\operatorname{SW}_p^p(\mu, \nu) = \int_{\mathbb{S}^{2d-1}} W_p^p(\mu|_{\sigma}, \nu|_{\sigma}) d\sigma$$

#### Existence and consistency of MSWE

Assume  $\mathbb{A1}$ ,  $\mathbb{A3}$ ,  $\mathbb{A4}$ . Then, there exists  $\mathbb{E}$  with  $\mathbb{P}(\mathbb{E}) = 1$  such that, for all  $\omega \in \mathbb{E}$ ,

$$\lim_{n \rightarrow \infty} \inf_{\theta \in \Theta} \operatorname{SW}_p(\mu_n(\omega), \mu_n) = \inf_{\theta \in \Theta} \operatorname{SW}_p(\mu_n, \mu_n)$$

$\lim_{n \rightarrow \infty} \operatorname{argmin}_{\theta \in \Theta} \operatorname{SW}_p(\mu_n(\omega), \mu_n) \subset \operatorname{argmin}_{\theta \in \Theta} \operatorname{SW}_p(\mu_n, \mu_n)$

Holds, for all  $\omega \in \mathbb{E}$ , there exists  $n(\omega)$  such that, for all  $n \geq n(\omega)$ ,  $\operatorname{argmin}_{\theta \in \Theta} \operatorname{SW}_p(\mu_n(\omega), \mu_n)$  is non-empty.

---

#### Combining MDE and OT

Minimum Wasserstein estimators, defined in (1) and (2) with  $D = W_p$ , have asymptotic guarantees [5] but are not practical.

$\Rightarrow$  With  $D = \operatorname{SW}_p$ , in (1) and (2), we get the **minimum (expected) SW estimators (MDE) (MSWE)** of order  $p$ .

Recent studies show the empirical success of SW-based estimators on generative modeling, but lack of theoretical guarantees.

$\Rightarrow$  We investigate the asymptotic properties of these estimators.

#### Central limit theorem for MSWE with $p = 1$

Consider  $\mathbb{A1}$ ,  $\mathbb{A3}$ ,  $\mathbb{A4}$ ,  $\mu_n = \mu_n$  (with  $\theta_n \in \Theta$  self-separated) and  $H: \theta \mapsto \int_{\mathbb{S}^{2d-1}} \int_{\mathbb{R}^d} H(x, y) d(\mu_n|_{\sigma} - \nu_n|_{\sigma})(x, y) d\sigma d\mu_n(x)$ , with

- $\sqrt{n}(\mu_n - \mu_n) \xrightarrow{d} G$ , where  $G_n$  and  $P_n$  contain the CDFs of the projected  $\mu_n$  and  $\nu_n$ .
- $P_n(\cdot, \cdot)$ : the "derivative" of  $F_p(\cdot, \cdot)$  in  $\theta$ .

Then,  $\sqrt{n} \inf_{\theta \in \Theta} \operatorname{SW}_p(\mu_n, \mu_n) \xrightarrow{d} \inf_{\theta \in \Theta} H(\theta)$ .

$$\sqrt{n}(\hat{\theta}_n - \theta_n) \xrightarrow{d} \operatorname{argmin}_{\theta \in \Theta} H(\theta), \text{ as } n \rightarrow \infty$$

$\Rightarrow$  Convergence rate of  $\sqrt{n}$  independent of the dimension

---

#### Main References

- [1] E. Bernton, P. E. Jacob, M. Gerber, C. P. Robert. On parameter estimation with the Wasserstein distance. *Information and Inference: A Journal of the IMA*, Jan 2019.
- [2] I. Drouot, Z. Zhang, A. G. Schwing. Generative modeling using the sliced Wasserstein distance. *CVPR* 2018.