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

**Kimia Nadjahi**[1]    Alain Durmus[2]    Umut Şimşekli[1,3]    Roland Badeau[1]

[1] Télécom Paris    [2] ENS Paris-Saclay    [3] University of Oxford

# Minimum Distance Estimation

$$\hat{\theta}_n = \underset{\theta \in \Theta}{\arg\min} \; D(\hat{\mu}_n, \mu_\theta)$$

$D$ : distance between distributions

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

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

# Minimum Distance Estimation

$$\hat{\theta}_n = \arg\min_{\theta} \, D(\hat{\mu}_n; \mu_\theta)$$

$D$ : distance between distributions

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

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

Example: Generative Modeling

# Minimum Expected Distance Estimation

Directly optimizing     is often not possible   (e.g. GANs)

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

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

# Minimum Wasserstein Estimation

Choose $D = W_p$ (Wasserstein distance of order $p \geq 1$)

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

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

## Minimum Wasserstein Estimation

Choose $D = W_p$ (Wasserstein distance of order $p \geq 1$)

X Robust and increasingly popular estimators: Wasserstein GAN [1], Wasserstein auto-encoders [2]

X Asymptotic guarantees [3]

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

$W_p$: expensive + curse of dimensionality

Central limit theorem in [3] valid in 1D

# Sliced-Wasserstein distance

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

$$SW_p^p(\;\;;\;\;) = \int_{\mathbb{S}^{d-1}} W_p^p(u_\sharp^?\;\;;\;u_\sharp^?\;\;) d\;(u)$$

Asymptotic Guarantees for Learning Generative Models with the Sliced-Wasserstein Distance.
K. Nadjahi, A. Durmus, U. Şimsekli, R. Badeau

# Minimum Sliced-Wasserstein Estimation

$$\hat{\theta}_n = \text{argmin}_\theta \ \text{SW}_p(\hat{\mu}_n; \mu_\theta)$$

$$\hat{\theta}_{n;m} = \text{argmin}_\theta \ E\left[\text{SW}_p(\hat{\mu}_n; \hat{\mu}_{\theta;m}) \mid Y_{1:n}\right]$$

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

# Minimum Sliced-Wasserstein Estimation

$$\hat{\mu}_n = \arg\min_{\mu \in \mathcal{M}_2} SW_p(\hat{\mu}_n; \mu)$$

$$\hat{\mu}_{n,m} = \arg\min_{\mu \in \mathcal{M}_2} E[SW_p(\hat{\mu}_n; \hat{\mu}_{\mu,m}) \mid Y_{1:n}]$$

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

———————————— Our contributions: ————————————

Convergence in $SW_p$ ⟹ weak convergence of probability measures

Existence and consistency of $\hat{\mu}_n$, $\hat{\mu}_{n,m}$

Central limit theorem for $\hat{\mu}_n$: $\sqrt{n}$ convergence rate for any dimension

# Thank you!

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



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