

Regularized Optimal Transport

From Computational and Statistical Perspectives

Tao Li, Yuan Ni, Michael Stanley, Nikhil Supekar

New York University

December 23, 2020



Deepfakes and Generative Models

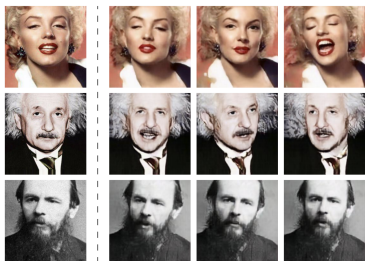


Figure 1: Deepfake images (right) generated from single, original (left)

- **Deepfakes** generates new facial expressions on existing images.
- **Generative adversarial networks (GANs)** learn a distribution from training data

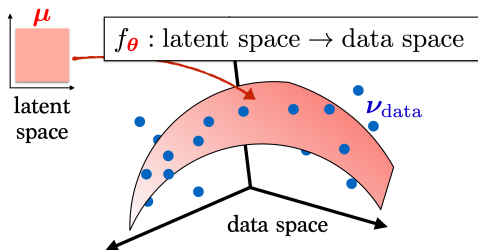


Figure 2: A generative model is a mapping from the latent space to the data space (Peyré and Cuturi 2019)

- ν_{data} : real faces
- f_{θ} : a parameterized mapping, e.g., deep neural networks
- $f_{\theta}(\mu)$: deepfakes

Discrepancy Function

Measuring how close $f_{\theta}(\mu)$ is to ν_{data}

- ▶ Total Variation
- ▶ Kullback-Leibler (KL) divergence, Jensen-Shannon (JS) divergence
- ▶ Wasserstein Distance (Arjovsky, Chintala, and Bottou 2017)

$$W(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y) \quad (\text{OT})$$

Training problem:

$$\min_{\theta} W(f_{\theta}(\mu), \nu)$$

Regularized Optimal Transport

- ▶ $W(f_\theta(\mu), \nu)$ is not smooth with respect to θ .
- ▶ Though $\nabla_\theta W(f_\theta(\mu), \nu)$ can be approximated, such methods fail to converge sometimes. (Gulrajani et al. 2017) (Bousquet et al. 2017)
- ▶ Adding regularization to make it smooth (Sanjabi et al. 2018)

$$W_\varepsilon(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \underbrace{\int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y)}_{\text{original eq. (OT)}} + \overbrace{\varepsilon \text{KL}(\pi \mid \mu \otimes \nu)}^{\text{entropy regularizer}}$$

(ROT)

Optimal Transport

$$\min_{\pi \in \Pi(\boldsymbol{\mu}, \boldsymbol{\nu})} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y) \quad (\text{Primal})$$

$$\begin{aligned} & \max_{u \in \mathcal{C}(\mathcal{X}), v \in \mathcal{C}(\mathcal{Y})} \int_{\mathcal{X}} u(x) d\boldsymbol{\mu} + \int_{\mathcal{Y}} v(y) d\boldsymbol{\nu} \quad (\text{Dual}) \\ & \text{subject to } u(x) + v(y) \leq c(x, y) \end{aligned}$$

Regularized Optimal Transport

$$\min_{\pi \in \Pi(\boldsymbol{\mu}, \boldsymbol{\nu})} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi \mid \boldsymbol{\mu} \otimes \boldsymbol{\nu}) \quad (\text{Primal})$$

$$\begin{aligned} & \max_{u \in \mathcal{C}(\mathcal{X}), v \in \mathcal{C}(\mathcal{Y})} \int_{\mathcal{X}} u(x) d\boldsymbol{\mu} + \int_{\mathcal{Y}} v(y) d\boldsymbol{\nu} \quad (\text{Dual}) \\ & - \varepsilon \int_{\mathcal{X} \times \mathcal{Y}} \exp \left(\frac{u(x) + v(y) - c(x, y)}{\varepsilon} \right) d\boldsymbol{\mu} d\boldsymbol{\nu} \end{aligned}$$

Define

$$f_\varepsilon(x, y, u, v) := u(x) + v(y) - \varepsilon \exp\left(\frac{u(x) + v(y) - c(x, y)}{\varepsilon}\right)$$

Dual Formulation

$$\max_{u \in \mathcal{C}(\mathcal{X}), v \in \mathcal{C}(\mathcal{Y})} \mathbb{E}_{X \sim \mu, Y \sim \nu} [f_\varepsilon(X, Y, u, v)]$$

- ▶ Alternating maximization: Sinkhorn algorithm (Cuturi 2013)
- ▶ Stochastic Gradient Descent (Aude et al. 2016)
- ▶ More methods from stochastic programming

Wasserstein distance on empirical measures:

$$W(\hat{\mu}, \hat{\nu}), \quad \hat{\mu} = 1/n \sum_i^n \delta_{x_i}, \hat{\nu} = 1/m \sum_j^m \delta_{y_j},$$

Curse of dimensionality

For $\mathbb{R}^d, d \geq 3$, $\mathbb{E}[|W(\hat{\mu}, \hat{\nu}) - W(\mu, \nu)|] = O(n^{-1/d})$

Regularization for breaking the curse

$$\mathbb{E} \left| W_\varepsilon(\alpha, \beta) - W_\varepsilon(\hat{\alpha}_n, \hat{\beta}_n) \right| = O \left(\frac{e^{\frac{\kappa}{\varepsilon}}}{\varepsilon^{\lfloor d/2 \rfloor} \sqrt{n}} \right) \text{ as } \varepsilon \rightarrow 0$$

$$\mathbb{E} \left| W_\varepsilon(\alpha, \beta) - W_\varepsilon(\hat{\alpha}_n, \hat{\beta}_n) \right| = O \left(\frac{1}{\sqrt{n}} \right) \text{ as } \varepsilon \rightarrow +\infty$$

[Klatt, Tameling, and Munk 2019]

As $m, n \rightarrow \infty$

$$\sqrt{\frac{nm}{n+m}} \{W_{p,\varepsilon}(\hat{\mu}_n, \hat{\nu}_m) - W_{p,\varepsilon}(\mu, \nu)\} \xrightarrow{D} \mathcal{N}_1(0, \sigma_{p,\varepsilon}^2(\mu, \nu))$$

- ▶ Asymptotic to a Gaussian
- ▶ Empirical Sinkhorn Divergence

[Sommerfeld and Munk 2017]

Under the null hypothesis $\mu = \nu$, as $m, n \rightarrow \infty$

$$\sqrt{\frac{mn}{m+n}}^{\frac{1}{p}} W_p(\hat{\mu}_n, \hat{\nu}_m) \xrightarrow{D} \{\max_u \langle G, u \rangle\}^{1/p}$$