# Numerical Optimal Transport

http://optimaltransport.github.io

# *Entropic Regularization*

## Gabriel Peyré

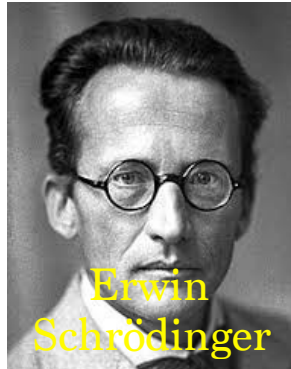www.numerical-tours.com

# Overview

- **Entropic Regularization and Sinkhorn**

- Convergence Analysis

- Sinkhorn Divergences
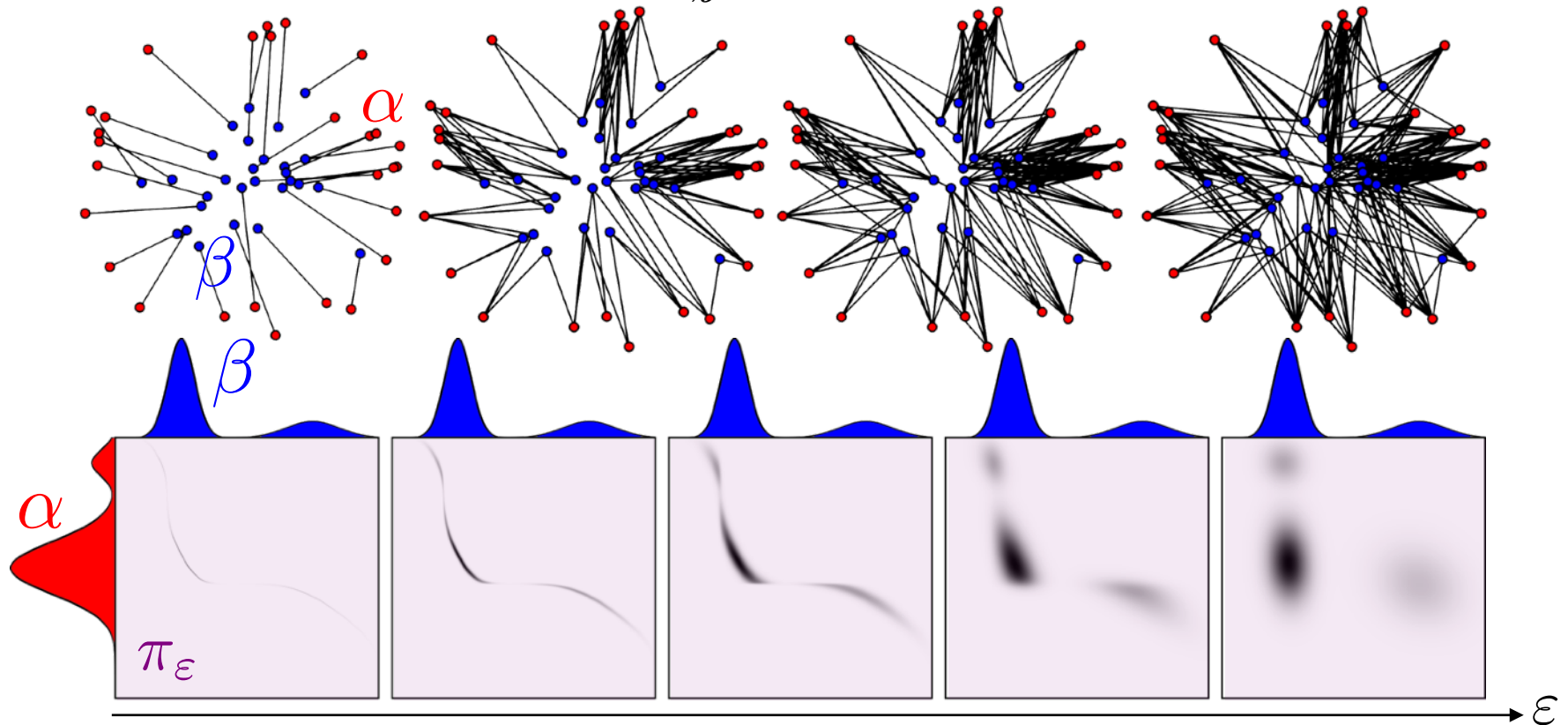
- Generative Model Fitting

# Entropic Regularization

*Schrödinger's problem:* [1931]

$$\min_{\mathbf{P} \in \mathbf{U}(\mathbf{a},\mathbf{b})} \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log\left(\frac{\mathbf{P}_{i,j}}{\mathbf{a}_i \mathbf{b}_j}\right)$$

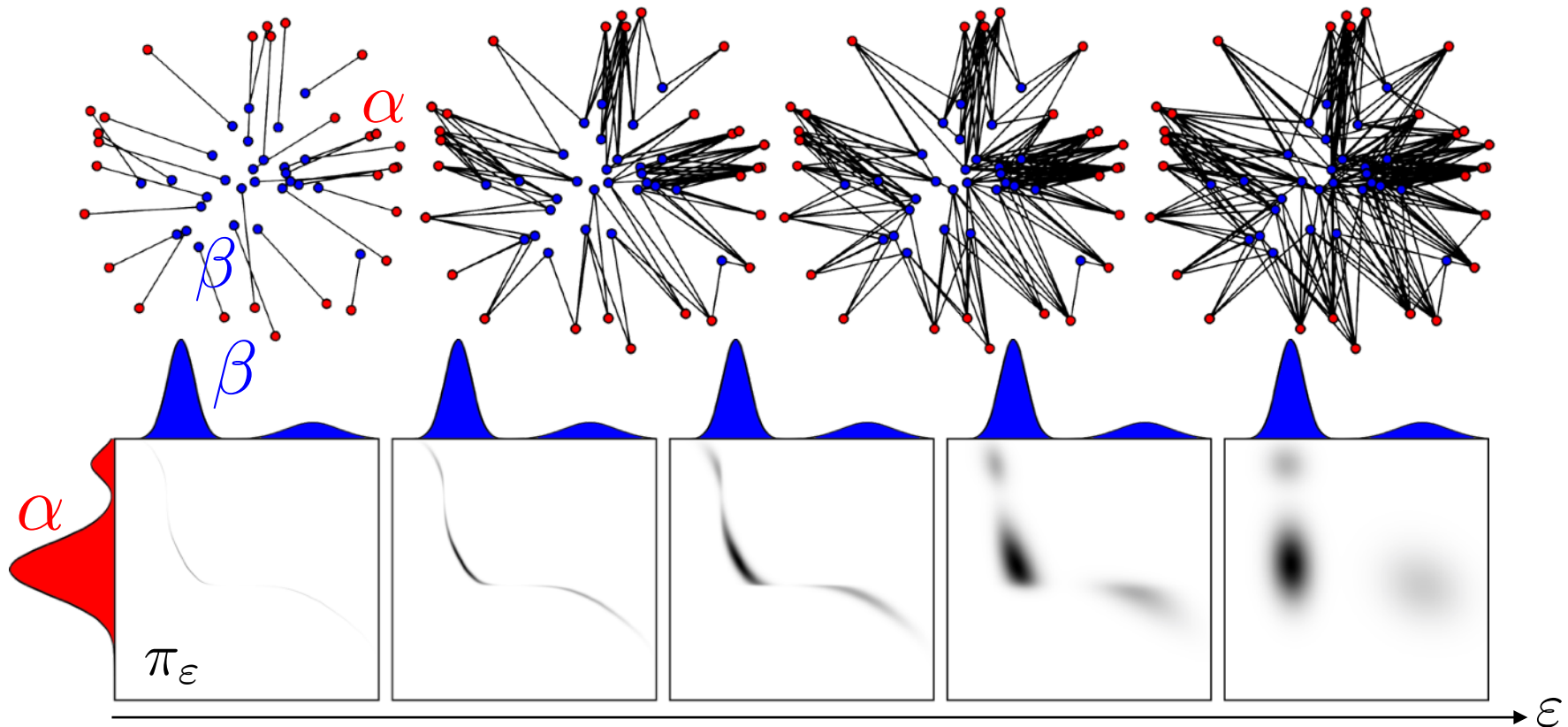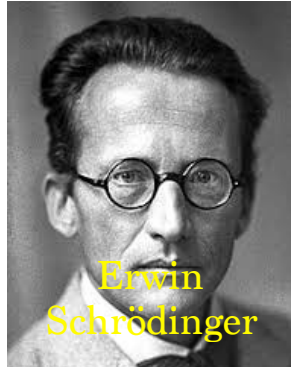Erwin Schrödinger

$$\pi = \sum_{i,j} \mathbf{P}_{i,j} \delta_{x_i, y_j}$$

# Entropic Regularization

Relative-entropy:  $\mathrm{KL}(\pi | \alpha \otimes \beta) \overset{\text{def.}}{=} \int_{\mathcal{X}^2} \log\left(\frac{\mathrm{d}\pi}{\mathrm{d}\alpha \mathrm{d}\beta}(x,y)\right) \mathrm{d}\pi(x,y)$

*Schrödinger's problem:* [1931]

$$\mathrm{W}^p_{\varepsilon,p}(\alpha, \beta) \overset{\text{def.}}{=} \min_{\pi_1=\alpha, \pi_2=\beta} \int_{\mathcal{X}^2} d^p(x,y)\mathrm{d}\pi(x,y) + \varepsilon\mathrm{KL}(\pi | \alpha \otimes \beta)$$


Erwin Schrödinger

# Probabilistic Interpretation

Relative-entropy:  $\mathrm{KL}(\pi | \alpha \otimes \beta) \stackrel{\mathrm{def.}}{=} \int_{\mathcal{X}^2} \log \left( \frac{\mathrm{d}\pi}{\mathrm{d}\alpha \mathrm{d}\beta}(x, y) \right) \mathrm{d}\pi(x, y)$

*Schrödinger's problem:*                    [1931]

$$\mathrm{W}^p_{\varepsilon, p}(\alpha, \beta) \stackrel{\mathrm{def.}}{=} \min_{\pi_1 = \alpha, \pi_2 = \beta} \int_{\mathcal{X}^2} d^p(x, y)\mathrm{d}\pi(x, y) + \varepsilon \mathrm{KL}(\pi | \alpha \otimes \beta)$$

$$\min_{(X,Y)} \left\{ \mathbb{E}(c(X, Y)) + \varepsilon \mathrm{I}(X, Y) \; ; \; X \sim \alpha, Y \sim \beta \right\}$$
Mutual information

Erwin Schrödinger

Christian Léonard

# Probabilistic Interpretation

Relative-entropy: $\mathrm{KL}(\pi | \alpha \otimes \beta) \overset{\text{def.}}{=} \int_{\mathcal{X}^2} \log\left( \frac{\mathrm{d}\pi}{\mathrm{d}\alpha \mathrm{d}\beta}(x,y) \right) \mathrm{d}\pi(x,y)$
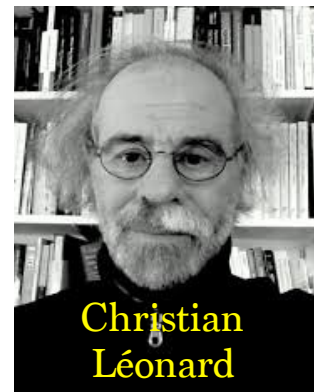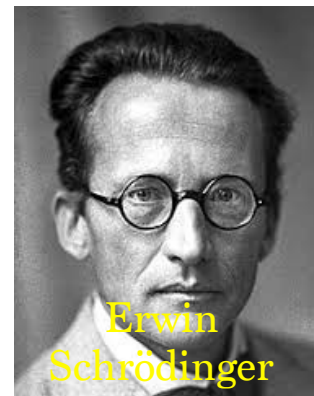
*Schrödinger's problem:* [1931]

$$\mathrm{W}^p_{\varepsilon,p}(\alpha, \beta) \overset{\text{def.}}{=} \min_{\pi_1=\alpha, \pi_2=\beta} \int_{\mathcal{X}^2} d^p(x,y)\mathrm{d}\pi(x,y) + \varepsilon\mathrm{KL}(\pi|\alpha \otimes \beta)$$

$$\min_{(X,Y)} \left\{ \mathbb{E}(c(X,Y)) + \varepsilon\mathrm{I}(X,Y) \;;\; X \sim \alpha, Y \sim \beta \right\}$$
Mutual information

Erwin Schrödinger

Christian Léonard



$\varepsilon = 0$     $\varepsilon = .05$     $\varepsilon = 0.2$     $\varepsilon = 1$

# Impact of Regularization



$$\pi_\varepsilon = \mathrm{argmin}_\pi \left\{ \int_{\mathbb{R}^2} \left( \|x - y\|^2 + \varepsilon \log \left( \frac{\mathrm{d}\pi}{\mathrm{d}\alpha \mathrm{d}\beta}(x,y) \right) \right) \mathrm{d}\pi(x,y) + \; ; \; \pi_1 = \textcolor{red}{\alpha}, \pi_2 = \textcolor{blue}{\beta} \right\}$$

$$\textit{Theorem:} \qquad \pi_\varepsilon \xrightarrow{\varepsilon \to +\infty} \textcolor{red}{\alpha} \otimes \textcolor{blue}{\beta} \qquad \pi_\varepsilon \xrightarrow{\varepsilon \to 0} \pi_{\mathrm{OT}}$$
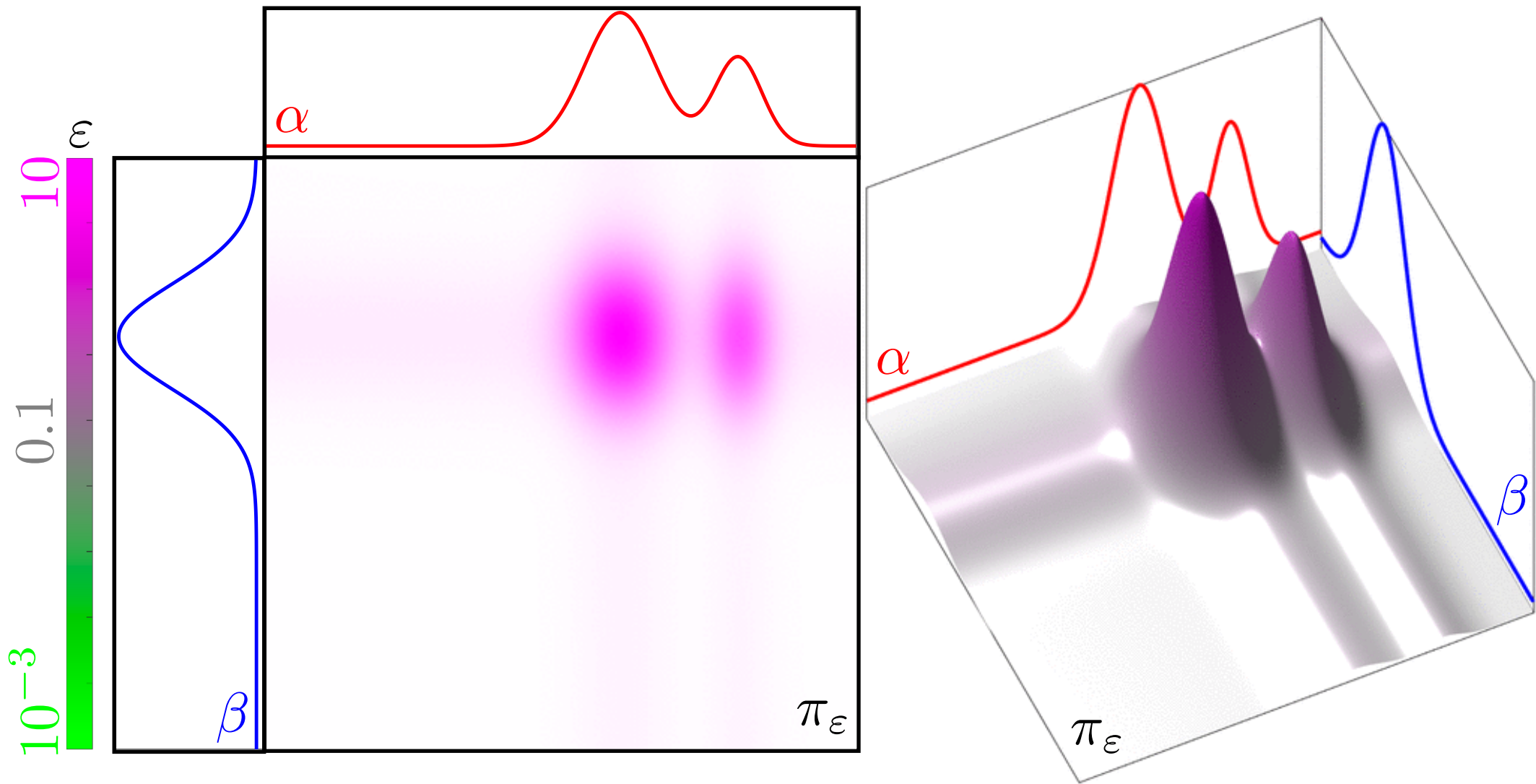
# Impact of Regularization



$$\pi_\varepsilon = \operatorname{argmin}_\pi \left\{ \int_{\mathbb{R}^2} \left( \|x - y\|^2 + \varepsilon \log \left( \frac{\mathrm{d}\pi}{\mathrm{d}\alpha\mathrm{d}\beta}(x,y) \right) \right) \mathrm{d}\pi(x,y) + \ ; \ \pi_1 = \textcolor{red}{\alpha}, \pi_2 = \textcolor{blue}{\beta} \right\}$$

$$\textit{Theorem:} \qquad \pi_\varepsilon \xrightarrow{\varepsilon \to +\infty} \textcolor{red}{\alpha} \otimes \textcolor{blue}{\beta} \qquad \pi_\varepsilon \xrightarrow{\varepsilon \to 0} \pi_{\mathrm{OT}}$$

# Impact of Regularization

$$Cumulative: \quad C_\pi(x,y) \overset{\text{def.}}{=} \int_{-\infty}^{x} \int_{-\infty}^{y} \mathrm{d}\pi(x,y)$$

$$Copula: \quad \chi_\pi(s,t) \overset{\text{def.}}{=} C_\pi(C_{\color{red}\alpha}^{-1}(s), C_{\color{blue}\beta}^{-1}(t))$$



$$Theorem: \quad \chi_{\pi_\varepsilon}(s,t) \quad \begin{array}{c} \xrightarrow{\varepsilon \to 0} \min(s,t) \quad \text{(dependence)} \\ \xrightarrow{\varepsilon \to +\infty} st \quad \text{(independence)} \end{array}$$

# Impact of Regularization

$$Cumulative: \quad C_\pi(x,y) \overset{\text{def.}}{=} \int_{-\infty}^{x} \int_{-\infty}^{y} \mathrm{d}\pi(x,y)$$

$$Copula: \quad \chi_\pi(s,t) \overset{\text{def.}}{=} C_\pi(C_{\color{red}\alpha}^{-1}(s), C_{\color{blue}\beta}^{-1}(t))$$



$$Theorem: \quad \chi_{\pi_\varepsilon}(s,t) \begin{array}{c} \overset{\varepsilon \to 0}{\nearrow} \min(s,t) \quad \text{(dependence)} \\ \underset{\varepsilon \to +\infty}{\searrow} st \quad \text{(independence)} \end{array}$$

# Sinkhorn's Algorithm

$$\min_{\mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b})} \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log\left(\frac{\mathbf{P}_{i,j}}{\mathbf{a}_i \mathbf{b}_j}\right)$$

*Proposition:* $\quad \mathbf{P}_{i,j} = \mathbf{u}_i \, \mathbf{K}_{i,j} \, \mathbf{v}_j \qquad \mathbf{K}_{i,j} \stackrel{\text{def.}}{=} e^{-\frac{d(x_i, y_j)^p}{\varepsilon}}$

# Sinkhorn's Algorithm

$$\min_{\mathbf{P} \in \mathbf{U}(\mathbf{a},\mathbf{b})} \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log\left(\frac{\mathbf{P}_{i,j}}{\mathbf{a}_i \mathbf{b}_j}\right)$$

*Proposition:* $\qquad \mathbf{P}_{i,j} = \mathbf{u}_i \mathbf{K}_{i,j} \mathbf{v}_j \qquad\qquad \mathbf{K}_{i,j} \overset{\text{def.}}{=} e^{-\frac{d(x_i, y_j)^p}{\varepsilon}}$

Row constraint: $\mathbf{u} \odot (\mathbf{K}\mathbf{v}) = \mathbf{a}$ $\qquad$ Col. constraint: $\mathbf{v} \odot (\mathbf{K}^\top \mathbf{u}) = \mathbf{b}$

# Sinkhorn's Algorithm

$$\min_{\mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b})} \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log \left( \frac{\mathbf{P}_{i,j}}{\mathbf{a}_i \mathbf{b}_j} \right)$$

*Proposition:* $\qquad \mathbf{P}_{i,j} = \mathbf{u}_i \, \mathbf{K}_{i,j} \, \mathbf{v}_j \qquad\qquad \mathbf{K}_{i,j} \overset{\text{def.}}{=} e^{-\frac{d(x_i, y_j)^p}{\varepsilon}}$

Row constraint: $\mathbf{u} \odot (\mathbf{K}\mathbf{v}) = \mathbf{a}$ $\qquad$ Col. constraint: $\mathbf{v} \odot (\mathbf{K}^\top \mathbf{u}) = \mathbf{b}$

Sinkhorn iterations: $\qquad \mathbf{u} \leftarrow \dfrac{\mathbf{a}}{\mathbf{K}\mathbf{v}} \qquad\qquad \mathbf{v} \leftarrow \dfrac{\mathbf{b}}{\mathbf{K}^\top \mathbf{u}}$

*Theorem:* [Sinkhorn 1964] $\quad (\mathbf{u}, \mathbf{v})$ converges.

# Sinkhorn's Algorithm

$$\min_{\mathbf{P}\in\mathbf{U}(\mathbf{a},\mathbf{b})} \sum_{i,j} d(x_i,y_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log\left(\frac{\mathbf{P}_{i,j}}{\mathbf{a}_i \mathbf{b}_j}\right)$$

*Proposition:* $\quad \mathbf{P}_{i,j} = \mathbf{u}_i \, \mathbf{K}_{i,j} \, \mathbf{v}_j \qquad \mathbf{K}_{i,j} \overset{\text{def.}}{=} e^{-\frac{d(x_i,y_j)^p}{\varepsilon}}$

Row constraint: $\mathbf{u} \odot (\mathbf{K}\mathbf{v}) = \mathbf{a}$ $\qquad$ Col. constraint: $\mathbf{v} \odot (\mathbf{K}^\top \mathbf{u}) = \mathbf{b}$
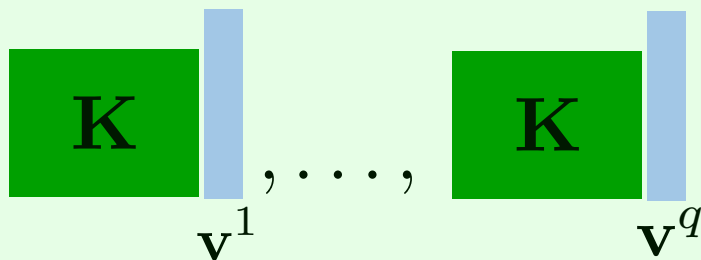
Sinkhorn iterations: $\qquad \mathbf{u} \leftarrow \dfrac{\mathbf{a}}{\mathbf{K}\mathbf{v}} \qquad\qquad \mathbf{v} \leftarrow \dfrac{\mathbf{b}}{\mathbf{K}^\top \mathbf{u}}$

*Theorem:* [Sinkhorn 1964] $\quad (\mathbf{u}, \mathbf{v})$ converges.

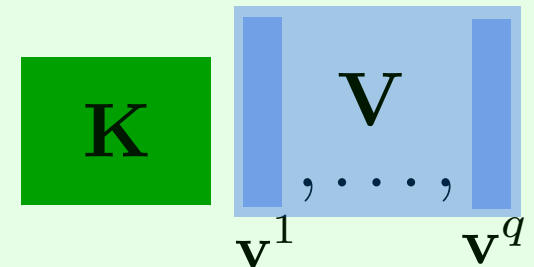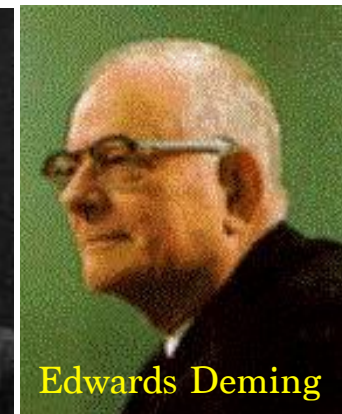Only matrix/vector multiplications.



Matrix-vectors $\quad\xrightarrow[\text{GPU}]{\text{parallelization}}\quad$ Matrix-matrix

$\rightarrow$ Convolution on regular grids, separable kernels.

# Sinkhorn, IPFP, RAS, ...



Richard Dennis Sinkhorn

Yule Udny    Edwards Deming    Frederick Stephan

Many names:

Sinkhorn algorithm

DAD scaling

Iterative proportional fitting
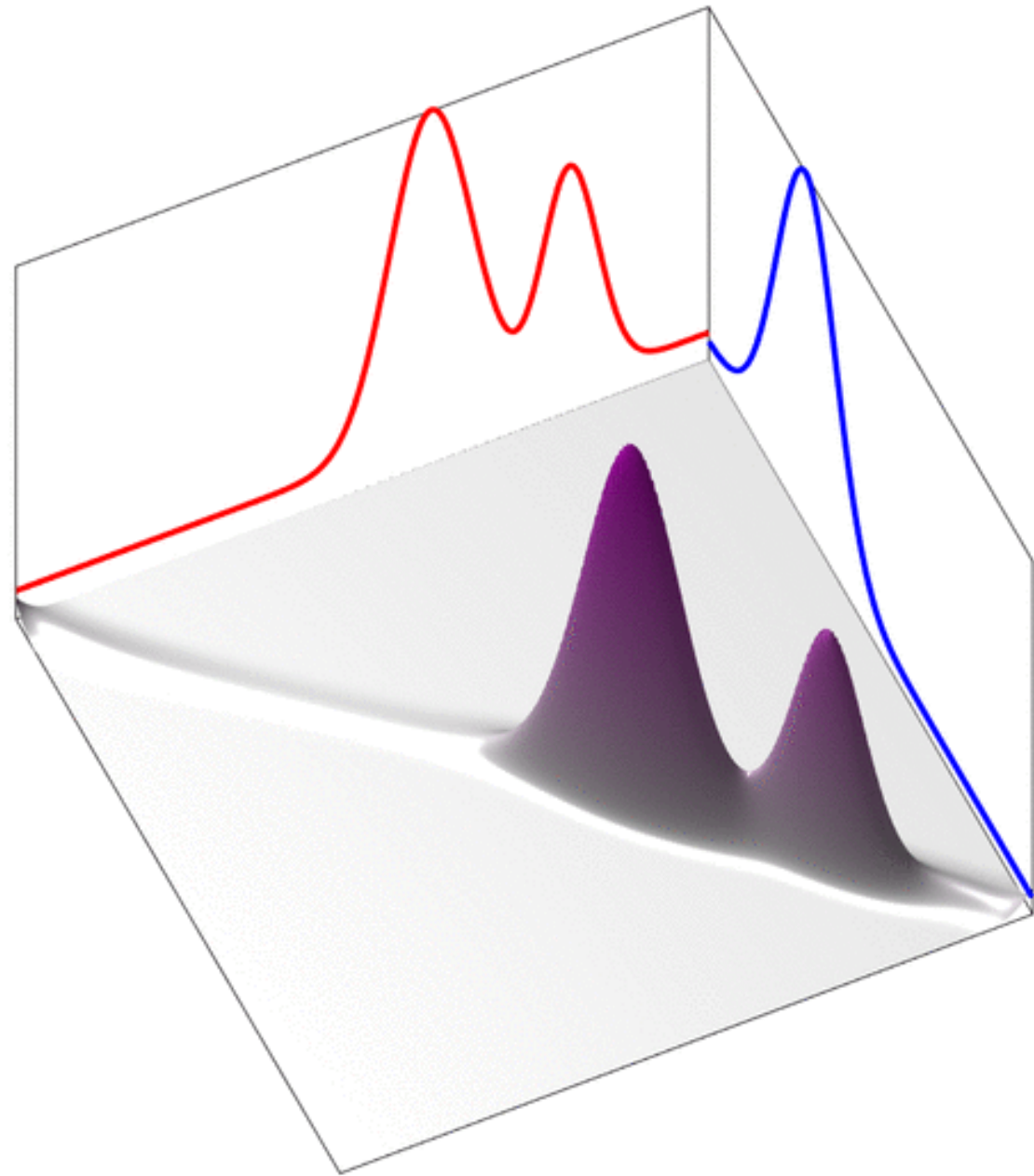
Biproportional fitting

RAS algorithm

Matrix scaling

Udny 1912

Kruithof, 1937

Deming and Stephan in 1940

Sinkhorn 1964

# Sinkhorn Evolution

# Other Regularizations

$$\min_{\pi}\left\{\int_{\mathbb{R}^2}\|x-y\|^2\mathrm{d}\pi(x,y)+\varepsilon R(\pi)\;;\;\pi_1=\textcolor{red}{\alpha},\pi_2=\textcolor{blue}{\beta}\right\}$$



Dykstra's algorithm

$$R(\pi)=\int\log\left(\frac{\mathrm{d}\pi}{\mathrm{d}x\mathrm{d}y}\right)\mathrm{d}\pi(x,y)$$

$$R(\pi)=\int\left(\frac{\mathrm{d}\pi}{\mathrm{d}x\mathrm{d}y}\right)^2\mathrm{d}x\mathrm{d}y$$

# Other Regularizations

$$\min_{\pi}\left\{\int_{\mathbb{R}^2}\|x-y\|^2\mathrm{d}\pi(x,y)+\varepsilon R(\pi)\ ;\ \pi_1=\textcolor{red}{\alpha},\pi_2=\textcolor{blue}{\beta}\right\}$$



$$R(\pi)=\int\log\left(\frac{\mathrm{d}\pi}{\mathrm{d}x\mathrm{d}y}\right)\mathrm{d}\pi(x,y)$$

Dykstra's algorithm

$$R(\pi)=\int\left(\frac{\mathrm{d}\pi}{\mathrm{d}x\mathrm{d}y}\right)^2\mathrm{d}x\mathrm{d}y$$

$$W_p^{\tau,p}(\textcolor{red}{\alpha}, \textcolor{blue}{\beta}) \overset{\text{def.}}{=} \min_{\pi} \int d^p \mathrm{d}\textcolor{purple}{\pi} + \tau \mathrm{KL}(\textcolor{purple}{\pi_1}|\textcolor{red}{\alpha}) + \tau \mathrm{KL}(\textcolor{purple}{\pi_2}|\textcolor{blue}{\beta})$$

[Liereo, Mielke, Savaré 2015]     See also:     [Chizat, Schmitzer, Peyré, Vialard 2015]

[Kondratyev, Monsaingeon, Vorotnikov 2015]

# Unbalanced OT

$$W_p^{\tau,p}(\textcolor{red}{\alpha}, \textcolor{blue}{\beta}) \stackrel{\text{def.}}{=} \min_{\pi} \int d^p \mathrm{d}\textcolor{purple}{\pi} + \tau \mathrm{KL}(\textcolor{purple}{\pi_1}|\textcolor{red}{\alpha}) + \tau \mathrm{KL}(\textcolor{purple}{\pi_2}|\textcolor{blue}{\beta})$$

[Liereo, Mielke, Savaré 2015]     See also:     [Chizat, Schmitzer, Peyré, Vialard 2015]
[Kondratyev, Monsaingeon, Vorotnikov 2015]

$$\int (\sqrt{\textcolor{red}{\alpha}} - \sqrt{\textcolor{blue}{\beta}})^2 \xleftarrow{\quad \tau \to 0 \quad} W_p^{\tau,p}(\textcolor{red}{\alpha}, \textcolor{blue}{\beta}) \xrightarrow{\quad \tau \to +\infty \quad} W_p^p(\textcolor{red}{\alpha}, \textcolor{blue}{\beta})$$

# Unbalanced OT

$$W_p^{\tau,p}(\textcolor{red}{\alpha}, \textcolor{blue}{\beta}) \overset{\text{def.}}{=} \min_{\pi} \int d^p \mathrm{d}\pi + \tau \mathrm{KL}(\pi_1 | \textcolor{red}{\alpha}) + \tau \mathrm{KL}(\pi_2 | \textcolor{blue}{\beta})$$
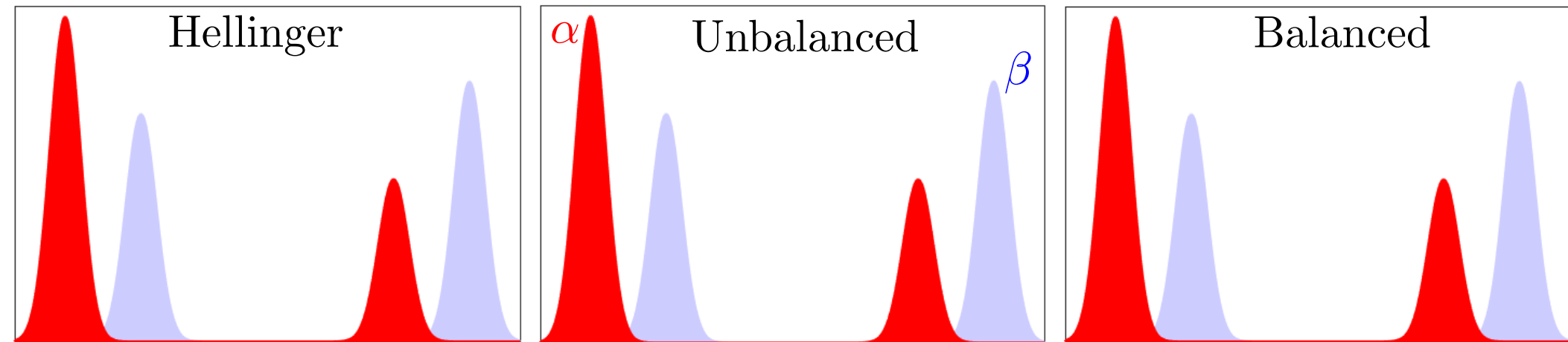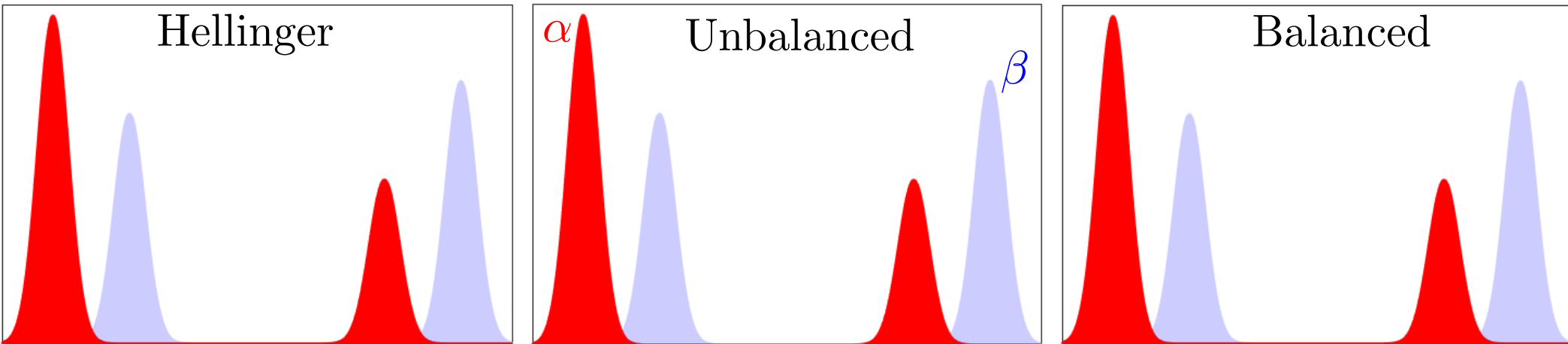
[Liereo, Mielke, Savaré 2015]     See also:     [Chizat, Schmitzer, Peyré, Vialard 2015]

[Kondratyev, Monsaingeon, Vorotnikov 2015]

$$\int(\sqrt{\textcolor{red}{\alpha}} - \sqrt{\textcolor{blue}{\beta}})^2 \xleftarrow{\ \tau \to 0\ } W_p^{\tau,p}(\textcolor{red}{\alpha}, \textcolor{blue}{\beta}) \xrightarrow{\ \tau \to +\infty\ } W_p^p(\textcolor{red}{\alpha}, \textcolor{blue}{\beta})$$



Hellinger

Unbalanced

Balanced

$$W_{\varepsilon,p}^{\tau,p}(\textcolor{red}{\alpha}, \textcolor{blue}{\beta}) \overset{\text{def.}}{=} \min_{\pi} \int d^p \mathrm{d}\pi + \tau \mathrm{KL}(\pi_1 | \textcolor{red}{\alpha}) + \tau \mathrm{KL}(\pi_2 | \textcolor{blue}{\beta}) + \varepsilon \mathrm{KL}(\pi | \textcolor{red}{\alpha} \otimes \textcolor{blue}{\beta})$$

Sinkhorn's algorithm:     $\mathbf{u} \leftarrow \left( \dfrac{\mathbf{a}}{\mathbf{Kv}} \right)^{1 + \frac{\varepsilon}{\tau}} \longleftrightarrow \mathbf{v} \leftarrow \left( \dfrac{\mathbf{b}}{\mathbf{K}^\top \mathbf{u}} \right)^{1 + \frac{\varepsilon}{\tau}}$

# Wasserstein Barycenters

Barycenters of measures $(\alpha_s)_s$:    $\sum_s \lambda_s = 1$

$$\beta^\star \in \underset{\beta}{\operatorname{argmin}} \sum_s \lambda_s \mathrm{W}_p^p(\alpha_s, \beta)$$



Guillaume Carlier

Martial Agueh

$\alpha_2$

$\beta^\star$

$\alpha_1$

$\alpha_3$

[Solomon et al, SIGGRAPH 2015]

# Wasserstein Barycenters

Barycenters of measures $(\alpha_s)_s$: $\qquad \sum_s \lambda_s = 1$

$$\beta^\star \in \underset{\beta}{\mathrm{argmin}} \ \sum_s \lambda_s W_p^p(\alpha_s, \beta)$$

$\alpha_2$

$\beta^\star$

$\alpha_1$

$\alpha_3$
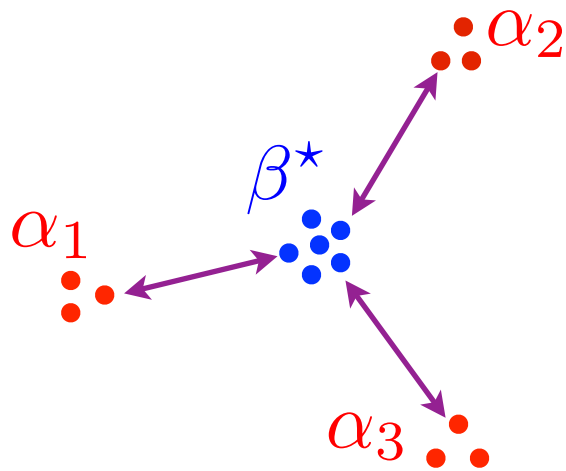
Guillaume Carlier

Martial Agueh

[Solomon et al, SIGGRAPH 2015]

Sinkhorn's algorithm: $\qquad \left( \mathbf{u}_s \leftarrow \dfrac{\mathbf{a}_s}{\mathbf{K}\mathbf{v}_s} \right)_s \qquad \longleftarrow \qquad \left( \mathbf{v}_s \leftarrow \dfrac{\mathbf{b}}{\mathbf{K}^\top \mathbf{u}_s} \right)_s$

$$\mathbf{b} \leftarrow \prod_s (\mathbf{K}^\top \mathbf{u}_s)^{\lambda_s}$$

# Overview

- Entropic Regularization and Sinkhorn

- **Convergence Analysis**

- Sinkhorn Divergences

- Generative Model Fitting

KL divergence:

$$\mathbf{KL}(\mathbf{P}|\mathbf{K}) \stackrel{\text{def.}}{=} \sum_{i,j} \mathbf{P}_{i,j} \log\left(\frac{\mathbf{P}_{i,j}}{\mathbf{K}_{i,j}}\right) - \mathbf{P}_{i,j} + \mathbf{K}_{i,j}$$

$$\mathrm{KL}(\mathbf{P}|\mathbf{K}) = D_\varphi(\mathbf{P}|\mathbf{K}) \quad \text{for} \quad \varphi(\mathbf{P}) = \sum_{i,j} \mathbf{P}_{i,j} \log(\mathbf{P}_{i,j})$$

$$\min_{\mathbf{P} \in \mathbf{U}(\mathbf{a},\mathbf{b})} \langle \mathbf{P}, \mathbf{C} \rangle + \varepsilon \, \mathrm{KL}(\mathbf{P}|\mathbf{a} \otimes \mathbf{b}) \Leftrightarrow \min_{\mathbf{P} \in \mathbf{U}(\mathbf{a},\mathbf{b})} \varepsilon \, \mathrm{KL}(\mathbf{P}|\mathbf{K}) \qquad \mathbf{K}_{i,j} \stackrel{\text{def.}}{=} e^{-\frac{\mathbf{C}_{i,j}}{\varepsilon}}$$

Iterative projections: $\mathbf{P}^{(\ell+1)} \stackrel{\text{def.}}{=} \mathrm{Proj}^{\mathbf{KL}}_{\mathcal{C}^1_{\mathbf{a}}}(\mathbf{P}^{(\ell)}) \quad \text{and} \quad \mathbf{P}^{(\ell+2)} \stackrel{\text{def.}}{=} \mathrm{Proj}^{\mathbf{KL}}_{\mathcal{C}^2_{\mathbf{b}}}(\mathbf{P}^{(\ell+1)})$

*Theorem:* $\quad \mathbf{P}^{(\ell)} \to \mathbf{P}^\star = \underset{\mathbf{P} \in \mathcal{C}^1_{\mathbf{a}} \cap \mathcal{C}^1_{\mathbf{b}}}{\mathrm{argmin}} \ \mathrm{KL}(\mathbf{P}|\mathbf{K})$

For affine $(\mathcal{C}^1_{\mathbf{a}}, \mathcal{C}^2_{\mathbf{b}})$,

# Bregman Iterative Projections

$$\langle \mathbf{P}, \mathbf{C} \rangle + \varepsilon \operatorname{KL}(\mathbf{P}|\mathbf{a} \otimes \mathbf{b}) = \varepsilon \operatorname{KL}(\mathbf{P}|\mathbf{K}) + \operatorname{cst} \quad \text{where} \quad \mathbf{K}_{i,j} = e^{-\frac{\mathbf{C}_{i,j}}{\varepsilon}} \mathbf{a}_i \mathbf{b}_j$$

$$\textit{Shrödinger problem:} \quad \min_{\mathbf{P} \in \mathbf{U}(\mathbf{a},\mathbf{b})} \operatorname{KL}(\mathbf{P}|\mathbf{K})$$

$$\mathbf{U}(\mathbf{a}, \mathbf{b}) = \mathcal{C}_{\mathbf{a}}^1 \cup \mathcal{C}_{\mathbf{b}}^2 \qquad \mathcal{C}_{\mathbf{a}}^1 \overset{\text{def.}}{=} \{\mathbf{P} \ : \ \mathbf{P}\mathbb{1}_m = \mathbf{a}\} \qquad \mathcal{C}_{\mathbf{b}}^2 \overset{\text{def.}}{=} \{\mathbf{P} \ : \ \mathbf{P}^{\mathrm{T}}\mathbb{1}_m = \mathbf{b}\}$$

# Bregman Iterative Projections

$$\langle \mathbf{P}, \mathbf{C} \rangle + \varepsilon \, \mathrm{KL}(\mathbf{P} | \mathbf{a} \otimes \mathbf{b}) = \varepsilon \, \mathrm{KL}(\mathbf{P} | \mathbf{K}) + \mathrm{cst} \quad \text{where} \quad \mathbf{K}_{i,j} = e^{-\frac{\mathbf{C}_{i,j}}{\varepsilon}} \mathbf{a}_i \mathbf{b}_j$$

$$\textit{Shrödinger problem:} \quad \min_{\mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b})} \mathrm{KL}(\mathbf{P} | \mathbf{K})$$

$$\mathbf{U}(\mathbf{a}, \mathbf{b}) = \mathcal{C}_{\mathbf{a}}^1 \cup \mathcal{C}_{\mathbf{b}}^2 \qquad \mathcal{C}_{\mathbf{a}}^1 \overset{\text{def.}}{=} \{ \mathbf{P} \; : \; \mathbf{P} \mathbb{1}_m = \mathbf{a} \} \qquad \mathcal{C}_{\mathbf{b}}^2 \overset{\text{def.}}{=} \{ \mathbf{P} \; : \; \mathbf{P}^{\mathrm{T}} \mathbb{1}_m = \mathbf{b} \}$$

Iterative projections: $\quad \mathbf{P}^{(\ell+1)} \overset{\text{def.}}{=} \mathrm{Proj}_{\mathcal{C}_{\mathbf{a}}^1}^{\mathbf{KL}}(\mathbf{P}^{(\ell)}) \quad \text{and} \quad \mathbf{P}^{(\ell+2)} \overset{\text{def.}}{=} \mathrm{Proj}_{\mathcal{C}_{\mathbf{b}}^2}^{\mathbf{KL}}(\mathbf{P}^{(\ell+1)})$

$\textit{Theorem:} \qquad \mathbf{P}^{(\ell)} \to \mathbf{P}^\star = \underset{\mathbf{P} \in \mathcal{C}_{\mathbf{a}}^1 \cap \mathcal{C}_{\mathbf{b}}^1}{\mathrm{argmin}} \; \mathrm{KL}(\mathbf{P} | \mathbf{K})$

For affine $(\mathcal{C}_{\mathbf{a}}^1, \mathcal{C}_{\mathbf{b}}^2)$,

[Bregman, 1967]

# Bregman Iterative Projections

$$\langle \mathbf{P}, \mathbf{C} \rangle + \varepsilon \, \text{KL}(\mathbf{P}|\mathbf{a} \otimes \mathbf{b}) = \varepsilon \, \text{KL}(\mathbf{P}|\mathbf{K}) + \text{cst} \quad \text{where} \quad \mathbf{K}_{i,j} = e^{-\frac{\mathbf{C}_{i,j}}{\varepsilon}} \mathbf{a}_i \mathbf{b}_j$$

*Shrödinger problem:* $\displaystyle\min_{\mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b})} \text{KL}(\mathbf{P}|\mathbf{K})$
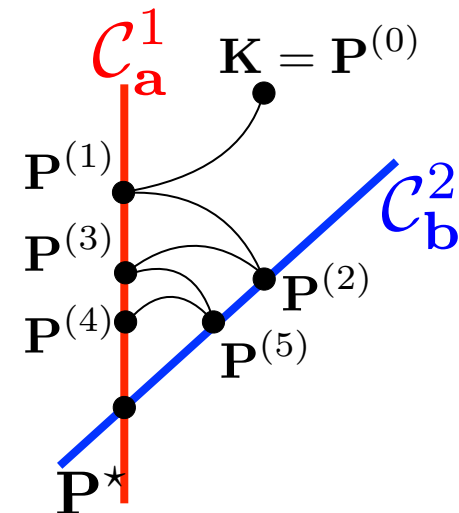
$$\mathbf{U}(\mathbf{a}, \mathbf{b}) = \mathcal{C}_{\mathbf{a}}^1 \cup \mathcal{C}_{\mathbf{b}}^2 \qquad \mathcal{C}_{\mathbf{a}}^1 \overset{\text{def.}}{=} \{\mathbf{P} \, : \, \mathbf{P}\mathbb{1}_m = \mathbf{a}\} \qquad \mathcal{C}_{\mathbf{b}}^2 \overset{\text{def.}}{=} \left\{\mathbf{P} \, : \, \mathbf{P}^{\mathrm{T}}\mathbb{1}_m = \mathbf{b}\right\}$$

Iterative projections: $\mathbf{P}^{(\ell+1)} \overset{\text{def.}}{=} \text{Proj}_{\mathcal{C}_{\mathbf{a}}^1}^{\mathbf{KL}}(\mathbf{P}^{(\ell)})$ and $\mathbf{P}^{(\ell+2)} \overset{\text{def.}}{=} \text{Proj}_{\mathcal{C}_{\mathbf{b}}^2}^{\mathbf{KL}}(\mathbf{P}^{(\ell+1)})$

*Theorem:* $\mathbf{P}^{(\ell)} \to \mathbf{P}^\star = \underset{\mathbf{P} \in \mathcal{C}_{\mathbf{a}}^1 \cap \mathcal{C}_{\mathbf{b}}^1}{\text{argmin}} \text{KL}(\mathbf{P}|\mathbf{K})$
For affine $(\mathcal{C}_{\mathbf{a}}^1, \mathcal{C}_{\mathbf{b}}^2)$,

Sinkhorn $\iff$ iterative projections.

$\mathbf{P}^{(2\ell)} \overset{\text{def.}}{=} \text{diag}(\mathbf{u}^{(\ell)})\mathbf{K}\,\text{diag}(\mathbf{v}^{(\ell)}), \quad \mathbf{P}^{(2\ell+1)} \overset{\text{def.}}{=} \text{diag}(\mathbf{u}^{(\ell+1)})\mathbf{K}\,\text{diag}(\mathbf{v}^{(\ell)})$
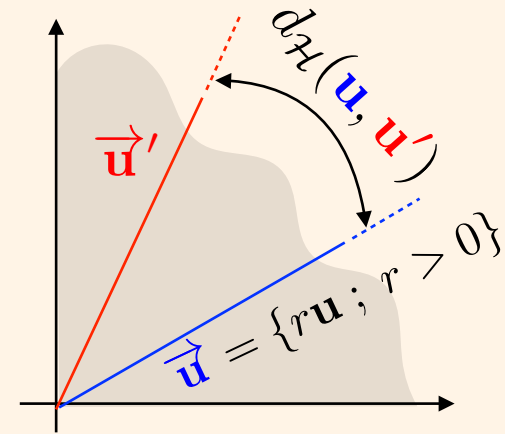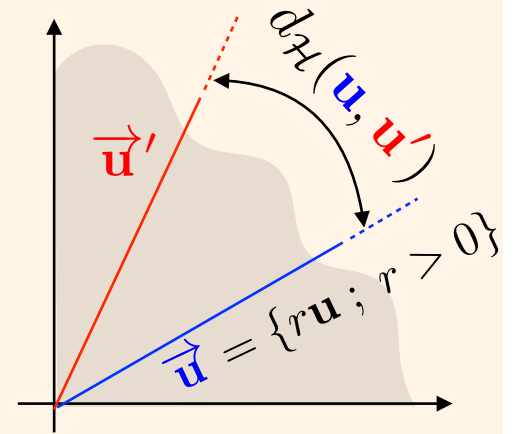
[Bregman, 1967]

# Hilbert Projective Metric

*Hilbert's projective metric:* $\quad \forall (\mathbf{u}, \mathbf{u}') \in (\mathbb{R}_{+,*}^n)^2$

$$d_{\mathcal{H}}(\mathbf{u}, \mathbf{u}') \stackrel{\text{def.}}{=} \|\log(\mathbf{u}) - \log(\mathbf{u}')\|_V$$

$$\|f\|_V \stackrel{\text{def.}}{=} \max(f) - \min(f)$$
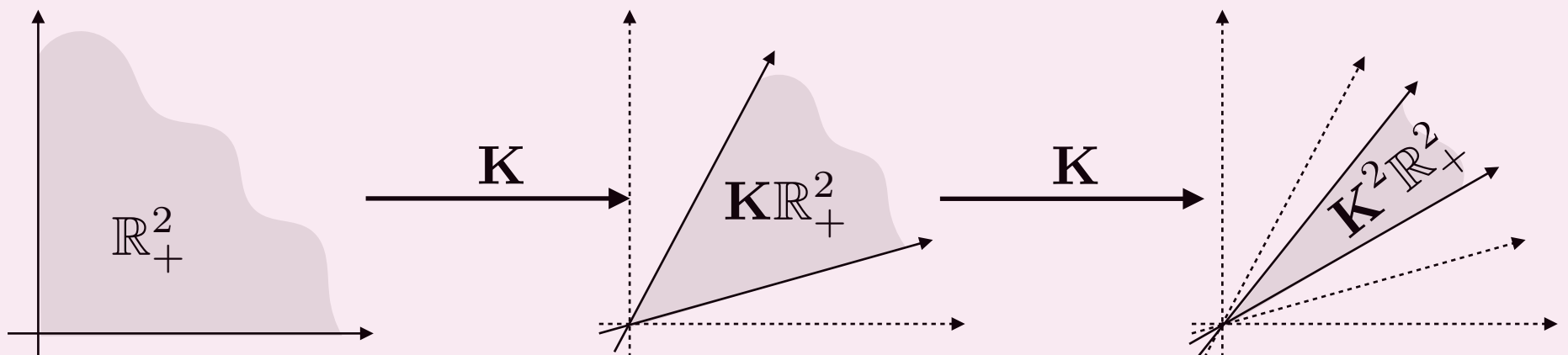
$d_{\mathcal{H}}$ is a distance on the set of rays $\overrightarrow{\mathbf{u}}$.

# Hilbert Projective Metric

*Hilbert's projective metric:* $\quad \forall (\mathbf{u}, \mathbf{u}') \in (\mathbb{R}^n_{+,*})^2$

$$d_{\mathcal{H}}(\mathbf{u}, \mathbf{u}') \overset{\text{def.}}{=} \|\log(\mathbf{u}) - \log(\mathbf{u}')\|_V$$

$$\|f\|_V \overset{\text{def.}}{=} \max(f) - \min(f)$$

$d_{\mathcal{H}}$ is a distance on the set of rays $\overrightarrow{\mathbf{u}}$.

*Birkhoff's contraction theorem:*

**Theorem 1.1.** Let $\mathbf{K} \in \mathbb{R}^{n \times m}_{+,*}$, then for $(\mathbf{v}, \mathbf{v}') \in (\mathbb{R}^m_{+,*})^2$

$$d_{\mathcal{H}}(\mathbf{K}\mathbf{v}, \mathbf{K}\mathbf{v}') \le \lambda(\mathbf{K}) d_{\mathcal{H}}(\mathbf{v}, \mathbf{v}') \text{ where } \begin{cases} \lambda(\mathbf{K}) \overset{\text{def.}}{=} \frac{\sqrt{\eta(\mathbf{K})}-1}{\sqrt{\eta(\mathbf{K})}+1} < 1 \\ \eta(\mathbf{K}) \overset{\text{def.}}{=} \max_{i,j,k,\ell} \frac{\mathbf{K}_{i,k}\mathbf{K}_{j,\ell}}{\mathbf{K}_{j,k}\mathbf{K}_{i,\ell}}. \end{cases}$$

# Perron Frobenius

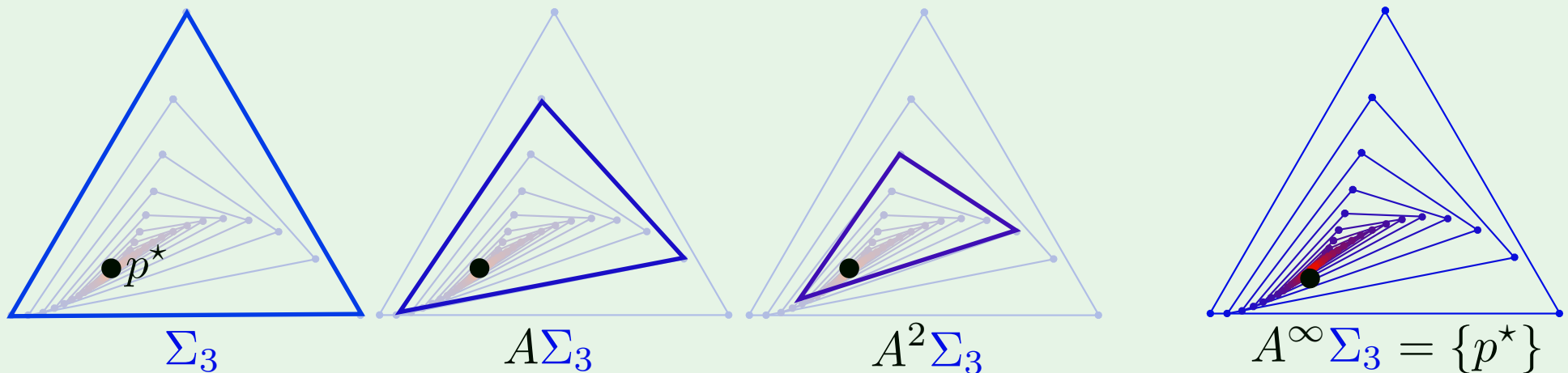Simplex: $\Sigma_k = \left\{ p \in \mathbb{R}_+^k \ ; \ \sum_i p_i = 1 \right\}$

Stochastic matrix: $A \in \mathbb{R}_+^n, \ A^\top \mathbb{1}_k = \mathbb{1}_k$

$A : \Sigma_k \to \Sigma_k$

*Theorem:* [Perron-Frobenius]

If $A > 0, \ \exists! p^\star, \ Ap^\star = p^\star.$

$\exists \rho \in [0, 1[, \ \|A^k p - p^\star\| \leqslant \rho^k$



$\Sigma_3$     $A\Sigma_3$     $A^2\Sigma_3$     $A^\infty \Sigma_3 = \{p^\star\}$

Sinkhorn iterations:

$$\mathbf{u}^{(\ell+1)} \stackrel{\text{def.}}{=} \frac{\mathbf{a}}{\mathbf{K}\mathbf{v}^{(\ell)}} \quad \text{and} \quad \mathbf{v}^{(\ell+1)} \stackrel{\text{def.}}{=} \frac{\mathbf{b}}{\mathbf{K}^{\mathsf{T}}\mathbf{u}^{(\ell+1)}}$$

*Theorem:* One has $(\mathbf{u}^{(\ell)}, \mathbf{v}^{(\ell)}) \to (\mathbf{u}^{\star}, \mathbf{v}^{\star})$

$$d_{\mathcal{H}}(\mathbf{u}^{(\ell)}, \mathbf{u}^{\star}) = O(\lambda(\mathbf{K})^{2\ell}), \quad d_{\mathcal{H}}(\mathbf{v}^{(\ell)}, \mathbf{v}^{\star}) = O(\lambda(\mathbf{K})^{2\ell}).$$

$$d_{\mathcal{H}}(\mathbf{u}^{(\ell)}, \mathbf{u}^{\star}) \leq \frac{d_{\mathcal{H}}(\mathbf{P}^{(\ell)}\mathbb{1}_m, \mathbf{a})}{1 - \lambda(\mathbf{K})^2} \quad d_{\mathcal{H}}(\mathbf{v}^{(\ell)}, \mathbf{v}^{\star}) \leq \frac{d_{\mathcal{H}}(\mathbf{P}^{(\ell),\top}\mathbb{1}_n, \mathbf{b})}{1 - \lambda(\mathbf{K})^2}$$

$$\|\log(\mathbf{P}^{(\ell)}) - \log(\mathbf{P}^{\star})\|_{\infty} \leq d_{\mathcal{H}}(\mathbf{u}^{(\ell)}, \mathbf{u}^{\star}) + d_{\mathcal{H}}(\mathbf{v}^{(\ell)}, \mathbf{v}^{\star})$$

[Franklin and Lorenz, 1989]

# Local Analysis of Sinkhorn

Sinkhorn fixed point: $\mathbf{f}^{(\ell+1)} = \Phi(\mathbf{f}^{(\ell)})$

$$\Phi \overset{\text{def.}}{=} \Phi_2 \odot \Phi_1 \quad \text{where} \quad \begin{cases} \Phi_1(\mathbf{f}) = \varepsilon \log \mathbf{K}^{\mathrm{T}}(e^{\mathbf{f}/\varepsilon}) - \log(\mathbf{b}), \\ \Phi_2(\mathbf{g}) = \varepsilon \log \mathbf{K}(e^{\mathbf{g}/\varepsilon}) - \log(\mathbf{a}). \end{cases}$$

*Proposition:* $\quad \partial\Phi(\mathbf{f}) = \mathrm{diag}(\mathbf{a})^{-1} \odot \mathbf{P} \odot \mathrm{diag}(\mathbf{b})^{-1} \odot \mathbf{P}^{\mathrm{T}}$.

For $\ell$ large enough, $\quad \|\mathbf{f}^{(\ell)} - \mathbf{f}\| = O((1-\kappa)^{\ell})$

# Local Analysis of Sinkhorn

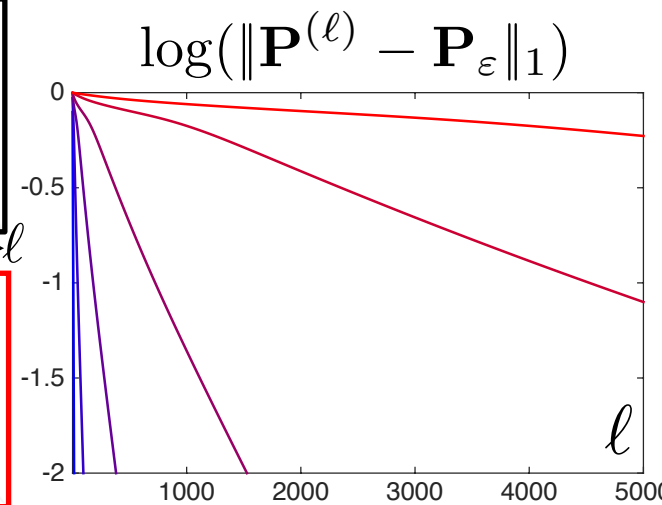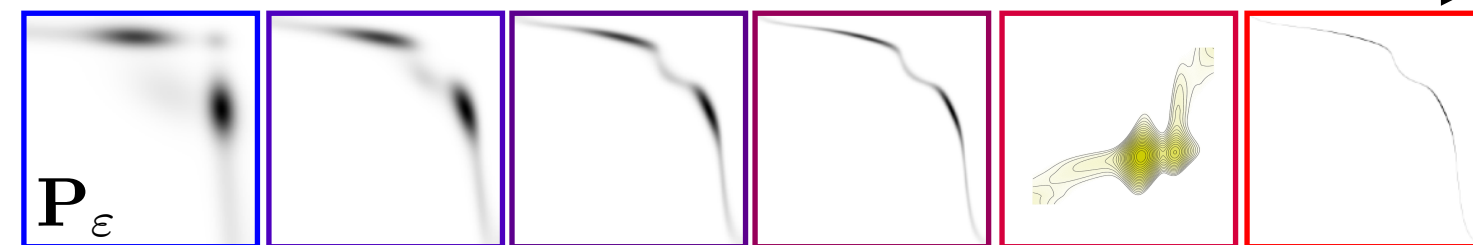Sinkhorn fixed point: $\mathbf{f}^{(\ell+1)} = \Phi(\mathbf{f}^{(\ell)})$

$$\Phi \stackrel{\text{def.}}{=} \Phi_2 \odot \Phi_1 \quad \text{where} \quad \begin{cases} \Phi_1(\mathbf{f}) = \varepsilon \log \mathbf{K}^{\mathrm{T}}(e^{\mathbf{f}/\varepsilon}) - \log(\mathbf{b}), \\ \Phi_2(\mathbf{g}) = \varepsilon \log \mathbf{K}(e^{\mathbf{g}/\varepsilon}) - \log(\mathbf{a}). \end{cases}$$

*Proposition:* $\quad \partial\Phi(\mathbf{f}) = \mathrm{diag}(\mathbf{a})^{-1} \odot \mathbf{P} \odot \mathrm{diag}(\mathbf{b})^{-1} \odot \mathbf{P}^{\mathrm{T}}$.

For $\ell$ large enough, $\quad \|\mathbf{f}^{(\ell)} - \mathbf{f}\| = O((1-\kappa)^{\ell})$

Global rate: $\kappa \sim e^{-\frac{1}{\varepsilon}}$
[Franklin and Lorenz, 1989]

$\longleftrightarrow$

Local rate: $\kappa \sim \varepsilon$
[Robert Berman 2017]

# Local Analysis of Sinkhorn

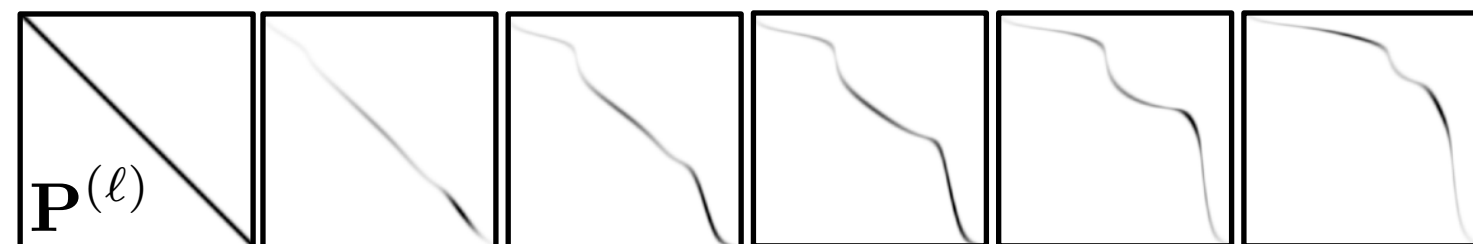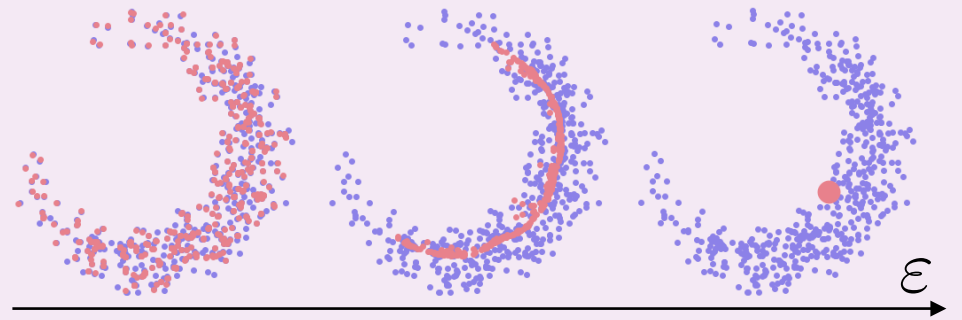Sinkhorn fixed point: $\mathbf{f}^{(\ell+1)} = \Phi(\mathbf{f}^{(\ell)})$

$$\Phi \overset{\text{def.}}{=} \Phi_2 \odot \Phi_1 \quad \text{where} \quad \begin{cases} \Phi_1(\mathbf{f}) = \varepsilon \log \mathbf{K}^{\mathrm{T}}(e^{\mathbf{f}/\varepsilon}) - \log(\mathbf{b}), \\ \Phi_2(\mathbf{g}) = \varepsilon \log \mathbf{K}(e^{\mathbf{g}/\varepsilon}) - \log(\mathbf{a}). \end{cases}$$

*Proposition:* $\partial\Phi(\mathbf{f}) = \mathrm{diag}(\mathbf{a})^{-1} \odot \mathbf{P} \odot \mathrm{diag}(\mathbf{b})^{-1} \odot \mathbf{P}^{\mathrm{T}}$.

For $\ell$ large enough, $\|\mathbf{f}^{(\ell)} - \mathbf{f}\| = O((1-\kappa)^{\ell})$

Global rate: $\kappa \sim e^{-\frac{1}{\varepsilon}}$
[Franklin and Lorenz, 1989]

$\longleftrightarrow$

Local rate: $\kappa \sim \varepsilon$
[Robert Berman 2017]



$\mathbf{P}^{(\ell)}$

$\mathbf{P}_\varepsilon$

$\log(\|\mathbf{P}^{(\ell)} - \mathbf{P}_\varepsilon\|_1)$

# Overview

- Entropic Regularization and Sinkhorn

- Convergence Analysis

- **<span style="color:red">Sinkhorn Divergences</span>**

- Generative Model Fitting

# Sinkhorn Divergences

$$\mathrm{W}_{\varepsilon,p}^p(\textcolor{red}{\alpha},\textcolor{blue}{\beta}) \stackrel{\mathrm{def.}}{=} \min_{\pi_1=\textcolor{red}{\alpha},\,\pi_2=\textcolor{blue}{\beta}} \int_{\mathcal{X}^2} d^p(\textcolor{red}{x},\textcolor{blue}{y})\mathrm{d}\textcolor{purple}{\pi}(\textcolor{red}{x},\textcolor{blue}{y}) + \varepsilon\mathrm{KL}(\pi|\xi)$$

*Problem:* $\mathrm{W}_{\varepsilon}(\textcolor{red}{\alpha},\textcolor{red}{\alpha}) \neq 0$

$$\min_{\textcolor{red}{\alpha}} \mathrm{W}_{\varepsilon,p}^p(\textcolor{red}{\alpha},\textcolor{blue}{\beta})$$



$\varepsilon$

# Sinkhorn Divergences

$$\mathrm{W}^p_{\varepsilon,p}(\textcolor{red}{\alpha},\textcolor{blue}{\beta}) \overset{\text{def.}}{=} \min_{\pi_1=\textcolor{red}{\alpha},\pi_2=\textcolor{blue}{\beta}} \int_{\mathcal{X}^2} d^p(\textcolor{red}{x},\textcolor{blue}{y})\mathrm{d}\textcolor{purple}{\pi}(\textcolor{red}{x},\textcolor{blue}{y}) + \varepsilon\mathrm{KL}(\pi|\xi)$$

*Problem:* $\mathrm{W}_\varepsilon(\textcolor{red}{\alpha},\textcolor{red}{\alpha}) \neq 0$

$$\min_{\textcolor{red}{\alpha}} \mathrm{W}^p_{\varepsilon,p}(\textcolor{red}{\alpha},\textcolor{blue}{\beta})$$



$\varepsilon$

$$\overline{\mathrm{W}}^p_{p,\varepsilon}(\textcolor{red}{\alpha},\textcolor{blue}{\beta}) \overset{\text{def.}}{=} \mathrm{W}^p_{p,\varepsilon}(\textcolor{red}{\alpha},\textcolor{blue}{\beta}) - \tfrac{1}{2}\mathrm{W}^p_{p,\varepsilon}(\textcolor{red}{\alpha},\textcolor{red}{\alpha}) - \tfrac{1}{2}\mathrm{W}^p_{p,\varepsilon}(\textcolor{blue}{\beta},\textcolor{blue}{\beta})$$

[Ramdas, García Trillos, Cuturi, 2017]

# Sinkhorn Divergences

$$\mathrm{W}_{\varepsilon,p}^{p}(\alpha,\beta) \overset{\text{def.}}{=} \min_{\pi_1=\alpha,\pi_2=\beta} \int_{\mathcal{X}^2} d^p(x,y)\mathrm{d}\pi(x,y) + \varepsilon\mathrm{KL}(\pi|\xi)$$

*Problem:* $\mathrm{W}_\varepsilon(\alpha,\alpha) \neq 0$

$$\min_\alpha \mathrm{W}_{\varepsilon,p}^p(\alpha,\beta)$$



$\varepsilon$

$$\overline{\mathrm{W}}_{p,\varepsilon}^{p}(\alpha,\beta) \overset{\text{def.}}{=} \mathrm{W}_{p,\varepsilon}^{p}(\alpha,\beta) - \tfrac{1}{2}\mathrm{W}_{p,\varepsilon}^{p}(\alpha,\alpha) - \tfrac{1}{2}\mathrm{W}_{p,\varepsilon}^{p}(\beta,\beta)$$

[Ramdas, García Trillos, Cuturi, 2017]

*Theorem:* $\quad \mathrm{W}_p^p(\alpha,\beta) \xleftarrow{\ \varepsilon\to 0\ } \overline{\mathrm{W}}_{\varepsilon,p}^{p}(\alpha,\beta) \xrightarrow{\ \varepsilon\to +\infty\ } \|\alpha-\beta\|_{-d^p}^2$

[Léonard 2012]
[Carlier et al 2017]

[Ramdas, García Trillos, Cuturi, 2017]

*Kernel norms (MMD):* $\quad \|\xi\|_{-d^p}^2 \overset{\text{def.}}{=} -\int_{\mathcal{X}^2} d(x,y)^p \mathrm{d}\xi(x)\mathrm{d}\xi(y)$

*Proposition:* $\|\cdot\|_{-\|\cdot\|^p}$ is a norm for $0 < p < 2$.

Arthur
Gretton

# Sinkhorn Divergences

$$\overline{\mathrm{W}}_{p,\varepsilon}^{p}(\textcolor{red}{\alpha},\textcolor{blue}{\beta}) \overset{\text{def.}}{=} \mathrm{W}_{p,\varepsilon}^{p}(\textcolor{red}{\alpha},\textcolor{blue}{\beta}) - \tfrac{1}{2}\mathrm{W}_{p,\varepsilon}^{p}(\textcolor{red}{\alpha},\textcolor{red}{\alpha}) - \tfrac{1}{2}\mathrm{W}_{p,\varepsilon}^{p}(\textcolor{blue}{\beta},\textcolor{blue}{\beta})$$

$\textcolor{red}{\text{concave}}$ $\qquad$ $\textcolor{blue}{\text{concave}}$

*Theorem:* [Feydy, Séjourné, P, Vialard, Trouvé, Amari 2018]

If $e^{-\frac{d^{p}}{\varepsilon}}$ is positive:

$\overline{\mathrm{W}}_{\varepsilon,p} \geqslant 0$ and $\overline{\mathrm{W}}_{\varepsilon,p}^{p}(\cdot,\textcolor{blue}{\beta})$ is convex.

$\overline{\mathrm{W}}_{\varepsilon,p}(\textcolor{red}{\alpha_n},\textcolor{blue}{\beta}) \to 0 \iff \textcolor{red}{\alpha_n} \overset{\text{weak}*}{\longrightarrow} \textcolor{blue}{\beta}$

$$\min_{\textcolor{red}{\alpha}} \mathrm{W}_{\varepsilon,p}^{p}(\textcolor{red}{\alpha},\textcolor{blue}{\beta}) \qquad\qquad \min_{\textcolor{red}{\alpha}} \overline{\mathrm{W}}_{\varepsilon,p}^{p}(\textcolor{red}{\alpha},\textcolor{blue}{\beta})$$

# Sample Complexity



$$Theorem: \quad \mathbb{E}(|\mathrm{W}_p(\hat{\alpha}, \hat{\beta}) - \mathrm{W}_p(\alpha, \beta)|) = O(n^{-\frac{1}{d}})$$

$$\mathbb{E}(|\|\hat{\alpha} - \hat{\beta}\|_k - \|\alpha - \beta\|_k|) = O(n^{-\frac{1}{2}})$$

*Optimal transport:* suffers from curse of dimensionality.

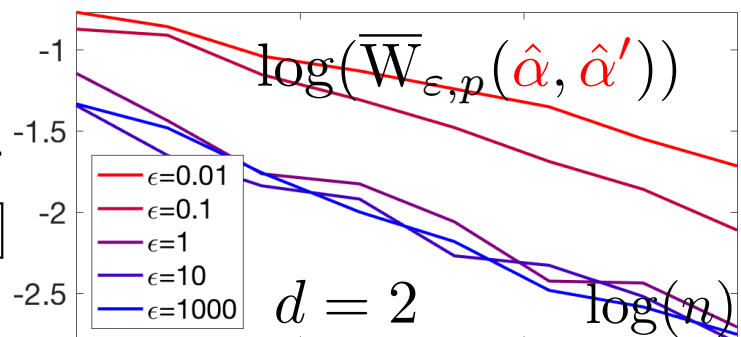$\to$ Adapt to support dimensionality [Weed, Bach 2017]

# Sample Complexity



*Theorem:*
$$\mathbb{E}(|W_p(\hat{\alpha}, \hat{\beta}) - W_p(\alpha, \beta)|) = O(n^{-\frac{1}{d}})$$
$$\mathbb{E}(|\|\hat{\alpha} - \hat{\beta}\|_k - \|\alpha - \beta\|_k|) = O(n^{-\frac{1}{2}})$$

*Optimal transport:* suffers from curse of dimensionality.
$\rightarrow$ Adapt to support dimensionality [Weed, Bach 2017]

*Theorem:* [Genevay, Bach, P, Cuturi]
$$\mathbb{E}(|\overline{W}_{\varepsilon,p}(\hat{\alpha}, \hat{\beta}) - \overline{W}_{\varepsilon,p}(\alpha, \beta)|) = O(\varepsilon^{-\frac{d}{2}} n^{-\frac{1}{2}})$$

# Overview

- Entropic Regularization and Sinkhorn

- Convergence Analysis

- Sinkhorn Divergences

- **Generative Model Fitting**

# Density Fitting and Generative Models

*Observations:* $\beta \stackrel{\text{def.}}{=} \frac{1}{n} \sum_{i=1}^{n} \delta_{x_i}$

*Parametric model:* $\theta \mapsto \alpha_\theta$

# Density Fitting and Generative Models

*Observations:* $\beta \overset{\text{def.}}{=} \frac{1}{n} \sum_{i=1}^{n} \delta_{x_i}$

*Parametric model:* $\theta \mapsto \alpha_\theta$



*Density fitting:* $\mathrm{d}\alpha_\theta(x) = \rho_\theta(x)\mathrm{d}x$

$$\min_\theta \widehat{\mathrm{KL}}(\alpha_\theta | \beta) \overset{\text{def.}}{=} -\sum_i \log(\rho_\theta(x_i))$$

Maximum
likelihood (MLE)

# Density Fitting and Generative Models

*Observations:* $\beta \stackrel{\text{def.}}{=} \frac{1}{n} \sum_{i=1}^{n} \delta_{x_i}$

*Parametric model:* $\theta \mapsto \alpha_\theta$

*Density fitting:* $\mathrm{d}\alpha_\theta(x) = \rho_\theta(x)\mathrm{d}x$

$$\min_\theta \widehat{\mathrm{KL}}(\alpha_\theta | \beta) \stackrel{\text{def.}}{=} -\sum_i \log(\rho_\theta(x_i))$$

Maximum likelihood (MLE)

*Generative model fit:* $\alpha_\theta = g_{\theta,\sharp}\zeta$

$$\widehat{\mathrm{KL}}(\alpha_\theta | \beta) = +\infty$$

$\to$ MLE undefined.

$\to$ Need a weaker metric.

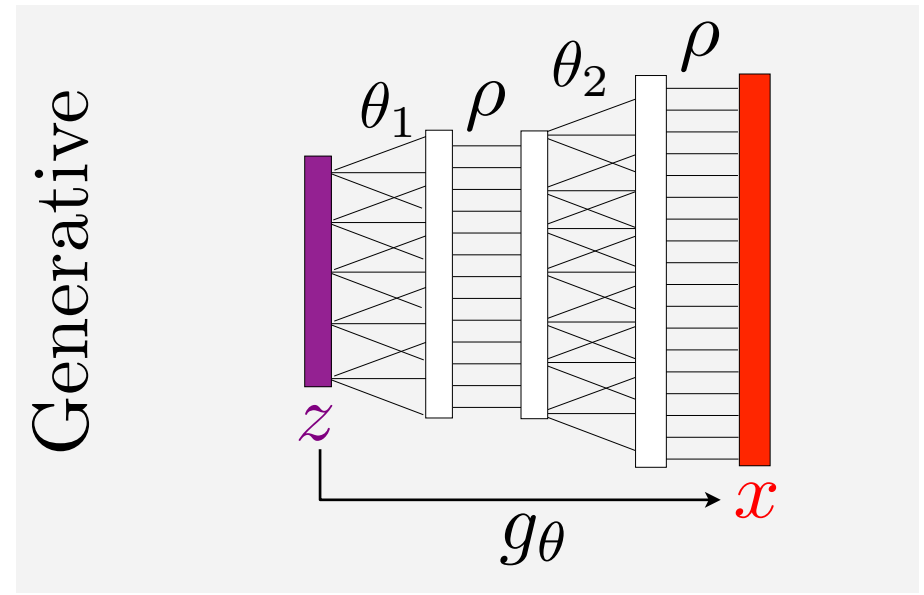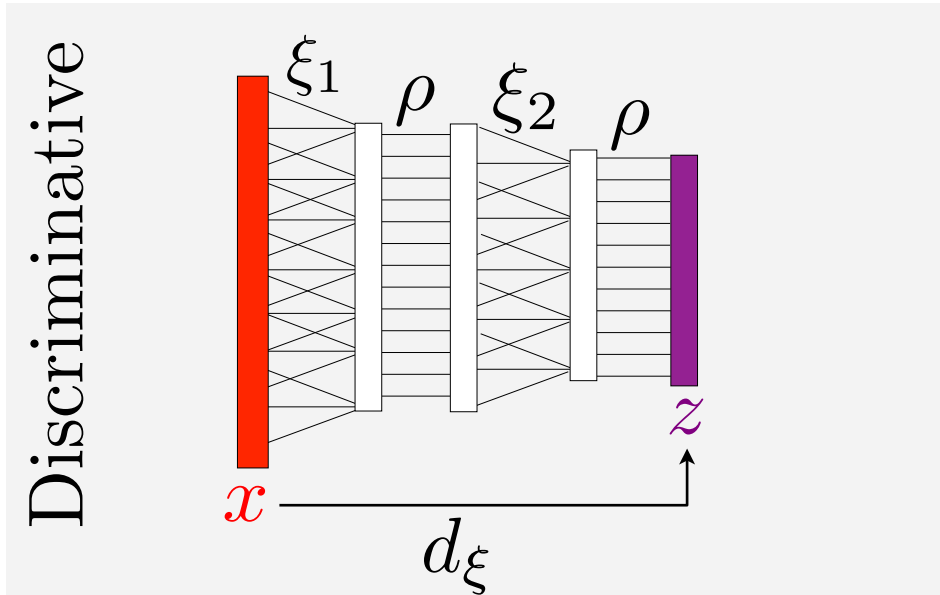$$\min_\theta \overline{\mathrm{W}}_{\varepsilon,p}^p(\alpha_\theta, \beta)$$

# Deep Discriminative vs Generative Models

Deep networks:

$$d_\xi(x) = \rho(\xi_K(\ldots \rho(\xi_2(\rho(\xi_1(x)\ldots)$$

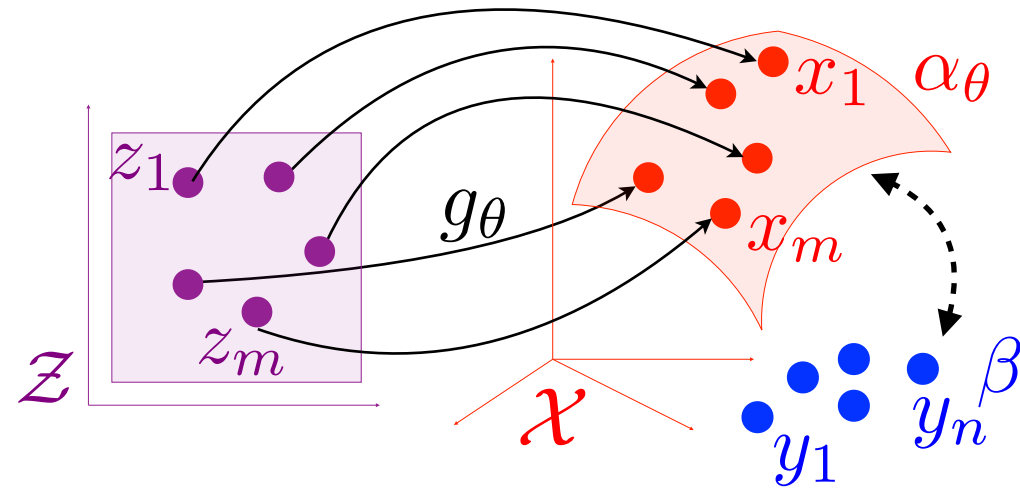$$g_\theta(z) = \rho(\theta_K(\ldots \rho(\theta_2(\rho(\theta_1(z)\ldots)$$

# Deep Discriminative vs Generative Models

Deep networks:
$$d_\xi(x) = \rho(\xi_K(\ldots \rho(\xi_2(\rho(\xi_1(x)\ldots)$$
$$g_\theta(z) = \rho(\theta_K(\ldots \rho(\theta_2(\rho(\theta_1(z)\ldots)$$
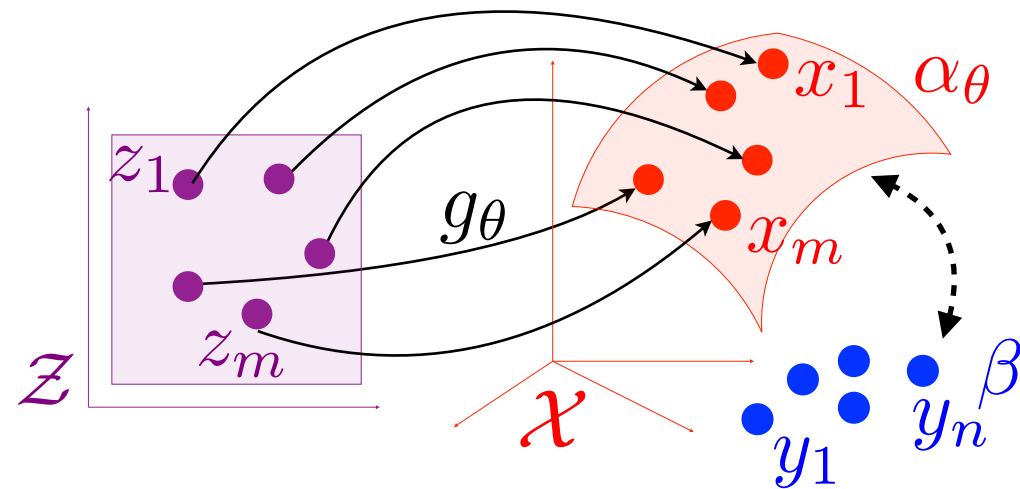
# Training Architecture



$$\min_\theta \mathcal{E}(\theta) \stackrel{\text{def.}}{=} \overline{\mathrm{W}}_{\varepsilon,p}^p(\textcolor{red}{\alpha_\theta}, \textcolor{blue}{\beta})$$

Stochastic gradient descent

$$\theta \leftarrow \theta - \tau\nabla\hat{\mathcal{E}}(\theta)$$

$$\hat{\mathcal{E}}(\theta) \stackrel{\text{def.}}{=} \overline{\mathrm{W}}_{\varepsilon,p}^p(\tfrac{1}{m}\sum_i \delta_{g_\theta(z_i)}, \textcolor{blue}{\beta})$$

# Training Architecture



$$\min_{\theta} \mathcal{E}(\theta) \stackrel{\text{def.}}{=} \overline{\mathrm{W}}_{\varepsilon,p}^{p}(\textcolor{red}{\alpha_\theta}, \textcolor{blue}{\beta})$$
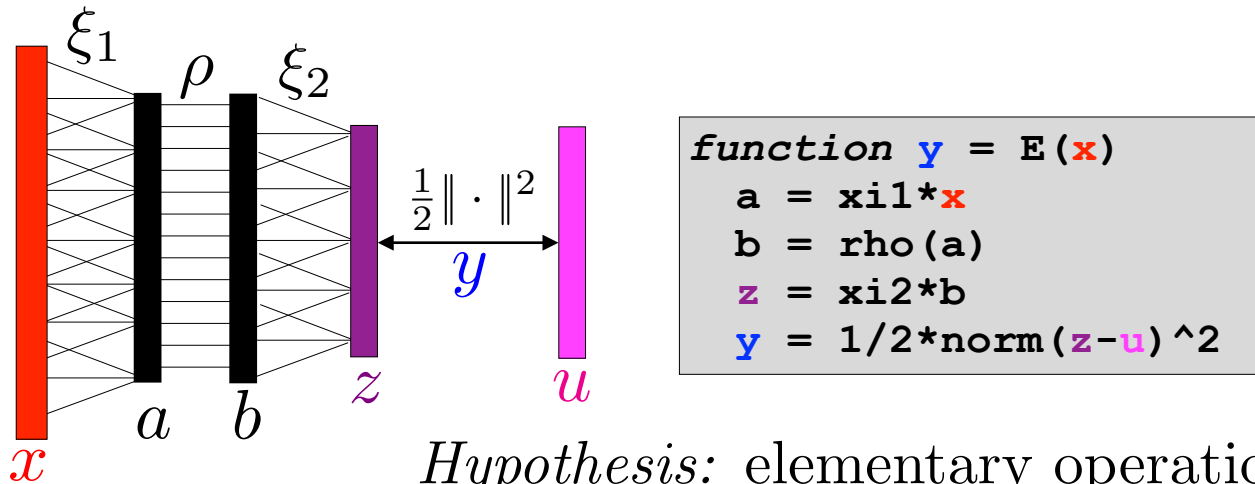
Stochastic gradient descent

$$\theta \leftarrow \theta - \tau \nabla \hat{\mathcal{E}}(\theta)$$

$$\hat{\mathcal{E}}(\theta) \stackrel{\text{def.}}{=} \overline{\mathrm{W}}_{\varepsilon,p}^{p}\left(\tfrac{1}{m}\sum_i \delta_{g_\theta(z_i)}, \textcolor{blue}{\beta}\right)$$

# Automatic Differentiation

**Setup:** $\mathcal{E} : \mathbb{R}^n \to \mathbb{R}$ computable in $K$ operations.



```
function y = E(x)
  a = xi1*x
  b = rho(a)
  z = xi2*b
  y = 1/2*norm(z-u)^2
```
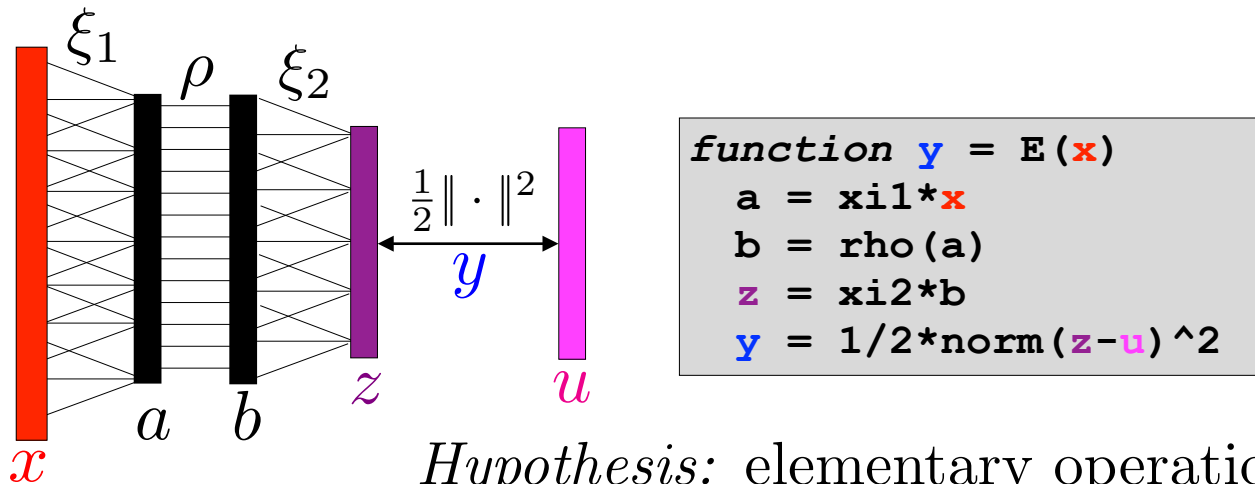
*Hypothesis:* elementary operations $(a \times b, \log(a), \sqrt{a} \ldots)$
and their derivatives cost $O(1)$.

**Question:** What is the complexity of computing $\nabla \mathcal{E} : \mathbb{R}^n \to \mathbb{R}^n$?

# Automatic Differentiation

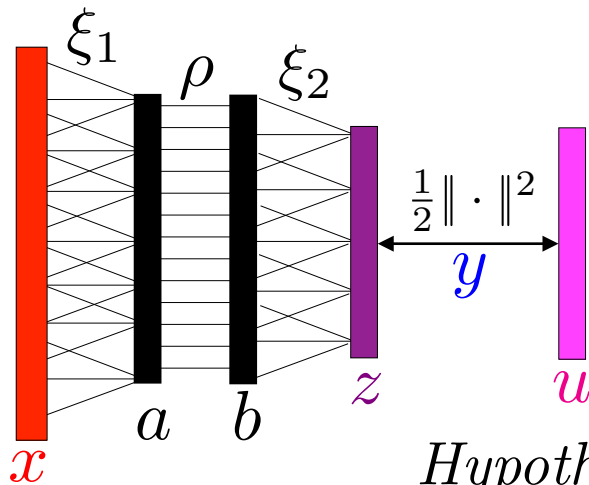**Setup:** $\mathcal{E} : \mathbb{R}^n \to \mathbb{R}$ computable in $K$ operations.



```
function y = E(x)
  a = xi1*x
  b = rho(a)
  z = xi2*b
  y = 1/2*norm(z-u)^2
```

*Hypothesis:* elementary operations $(a \times b, \log(a), \sqrt{a} \ldots)$ and their derivatives cost $O(1)$.

**Question:** What is the complexity of computing $\nabla \mathcal{E} : \mathbb{R}^n \to \mathbb{R}^n$?

Finite differences: $\quad \nabla \mathcal{E}(\theta) \approx \dfrac{1}{\varepsilon} (\mathcal{E}(\theta + \varepsilon \delta_1) - \mathcal{E}(\theta), \ldots \mathcal{E}(\theta + \varepsilon \delta_n) - \mathcal{E}(\theta))$

$K(n+1)$ operations, intractable for large $n$.

# Automatic Differentiation

**Setup:** $\mathcal{E} : \mathbb{R}^n \to \mathbb{R}$ computable in $K$ operations.



```
function y = E(x)
  a = xi1*x
  b = rho(a)
  z = xi2*b
  y = 1/2*norm(z-u)^2
```

```
function dx = nablaE(x)
  dz = z-u
  db = xi2'*dz
  da = diag(dphi(a)) * db
  dx = xi1'*da
```

*Hypothesis:* elementary operations $(a \times b, \log(a), \sqrt{a} \dots)$
and their derivatives cost $O(1)$.

**Question:** What is the complexity of computing $\nabla \mathcal{E} : \mathbb{R}^n \to \mathbb{R}^n$?

Finite differences:
$$\nabla \mathcal{E}(\theta) \approx \frac{1}{\varepsilon}(\mathcal{E}(\theta + \varepsilon\delta_1) - \mathcal{E}(\theta), \dots \mathcal{E}(\theta + \varepsilon\delta_n) - \mathcal{E}(\theta))$$
$K(n+1)$ operations, intractable for large $n$.

*Theorem:* there is an algorithm to compute $\nabla \mathcal{E}$
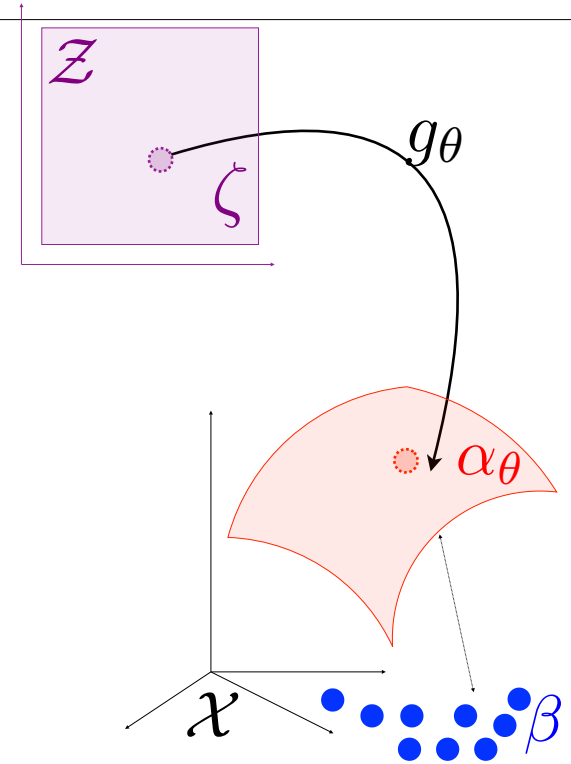in $O(K)$ operations.     [Seppo Linnainmaa, 1970]

This algorithm is reverse mode automatic differentiation

# Examples of Images Generation

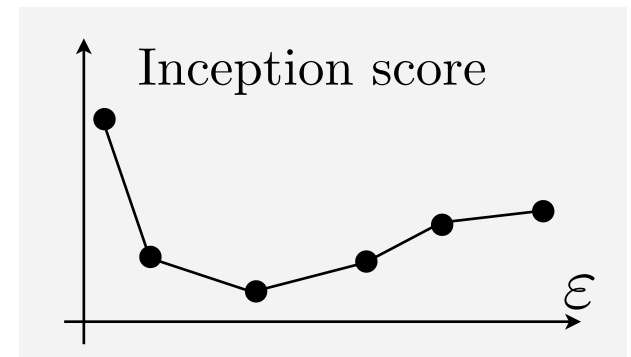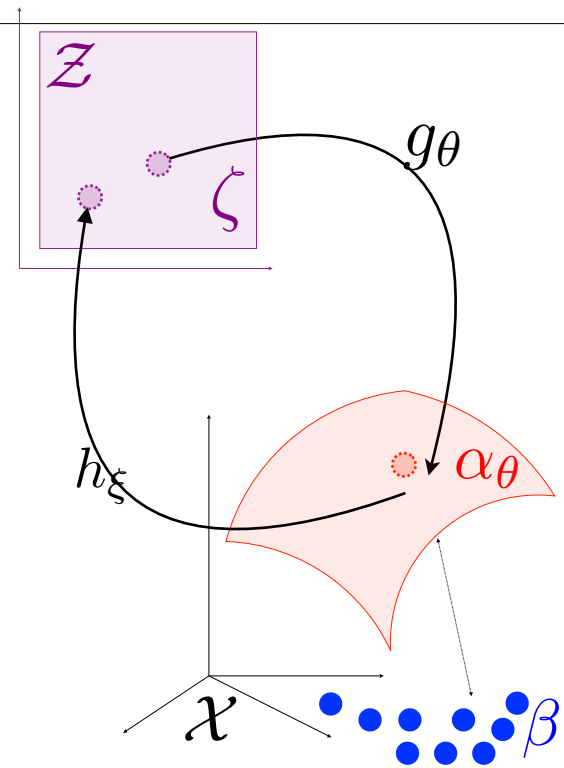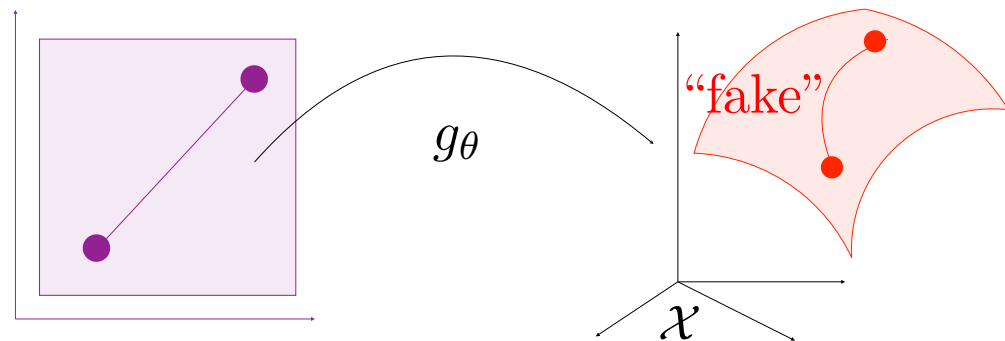# Examples of Images Generation



Inputs $\beta$     Generated $\alpha_\theta$

$\mathcal{Z}$   $\zeta$   $g_\theta$   $h_\xi$   $\alpha_\theta$   $\mathcal{X}$   $\beta$

Inception score

$\varepsilon$

Ian Goodfellow

$\rightarrow$ Need to learn the metric $d(x, y) = \|h_\xi(x) - h_\xi(y)\|$ (GANs)

$\rightarrow$ Influence of $\varepsilon$?

$\rightarrow$ Performance evaluation of generative models is an open problem.

*Progressive Growing of GANs for Improved Quality, Stability, and Variation*
Tero Karras, Timo Aila, Samuli Laine, Jaakko Lehtinen, ICLR 2018

$g_\theta$

"fake"

$\mathcal{X}$

*Progressive Growing of GANs for Improved Quality, Stability, and Variation*
Tero Karras, Timo Aila, Samuli Laine, Jaakko Lehtinen, ICLR 2018