

SNEkhorn : Dimension Reduction with Symmetric Entropic Affinities

Hugues Van Assel

PhD student ENS Lyon



Titouan Vayer



Rémi Flamary



Nicolas Courty

Overview of the talk

Part I: Symmetric Entropic Affinities

Part II: Application to dimensionality reduction

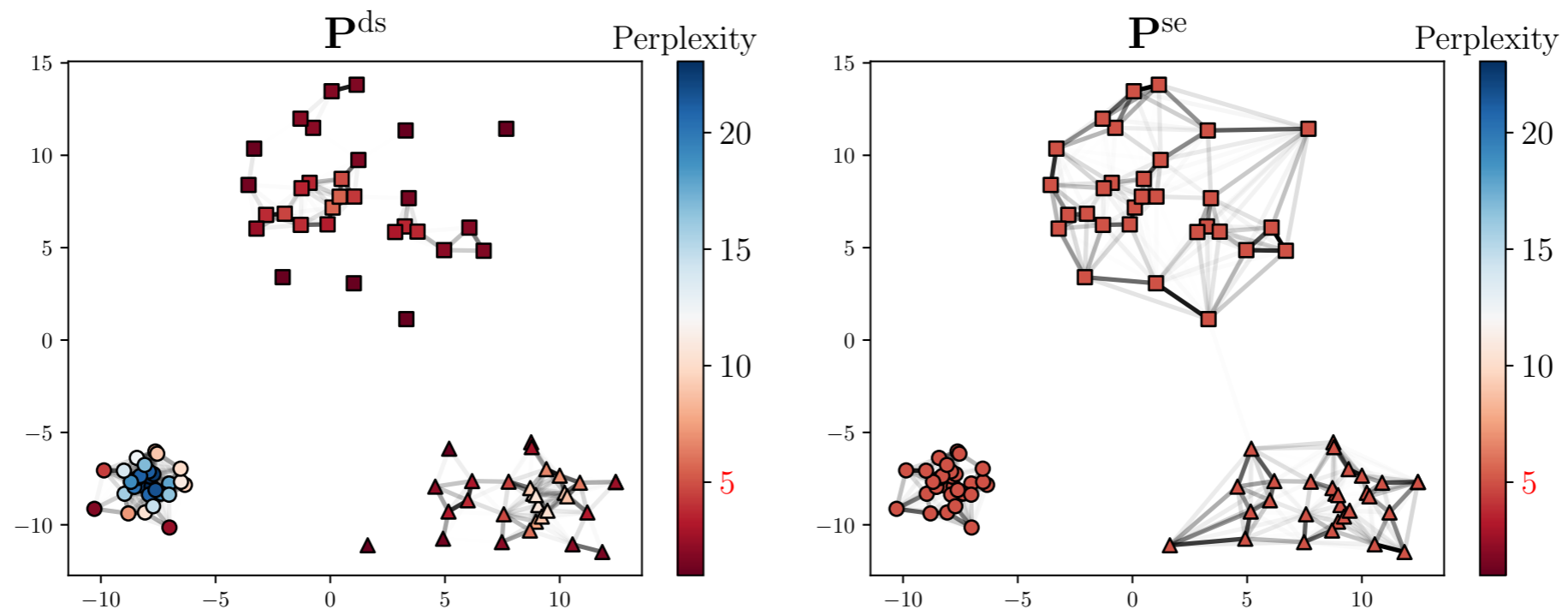
Part III: Future works

Overview of the talk

Part I: Symmetric Entropic Affinities

Part II: Application to dimensionality reduction

Part III: Future works

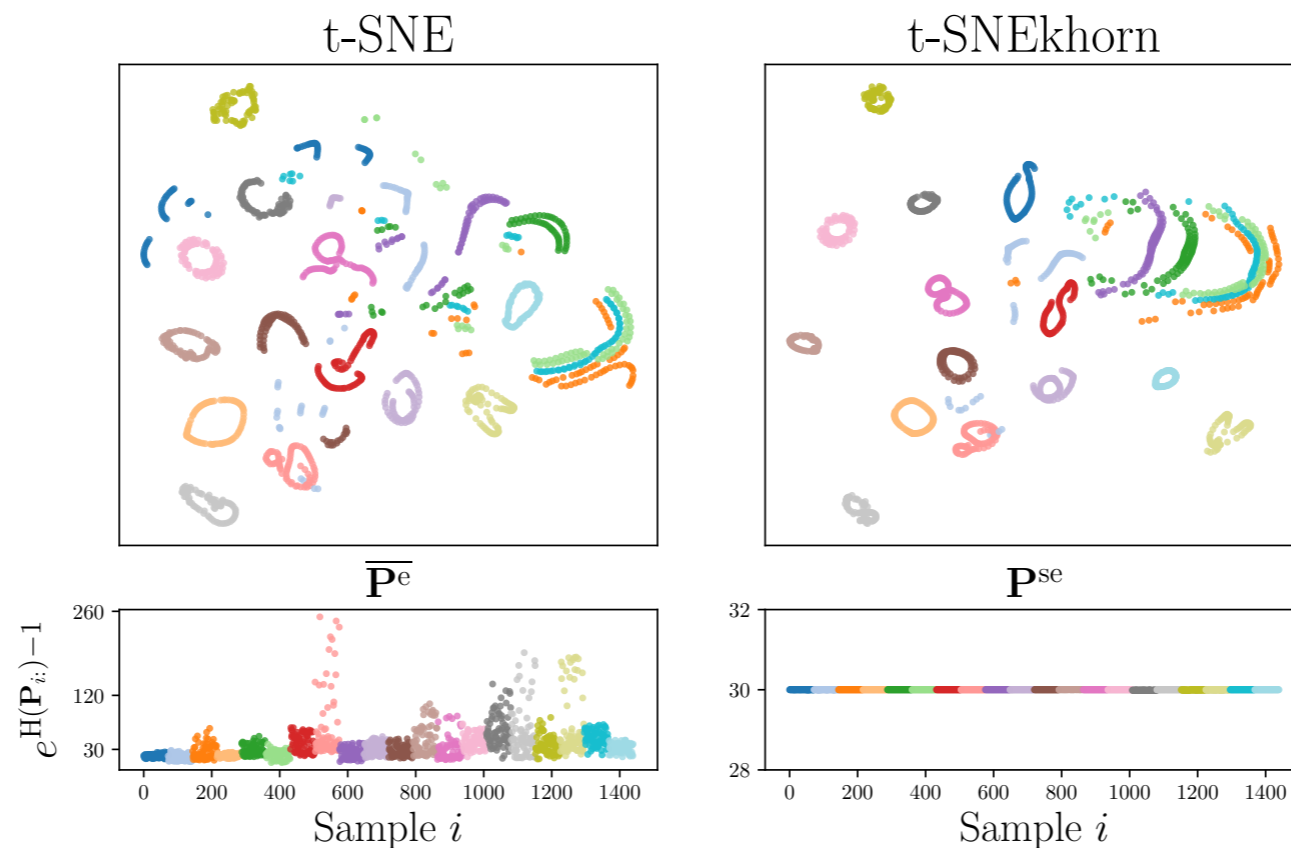


Overview of the talk

Part I: Symmetric Entropic Affinities

Part II: Application to dimensionality reduction

Part III: Future works

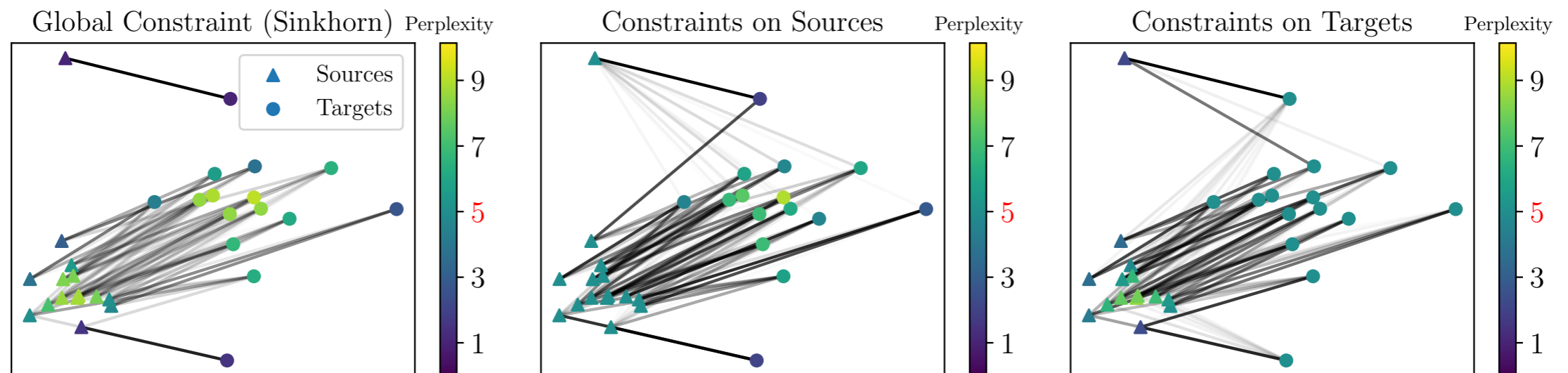


Overview of the talk

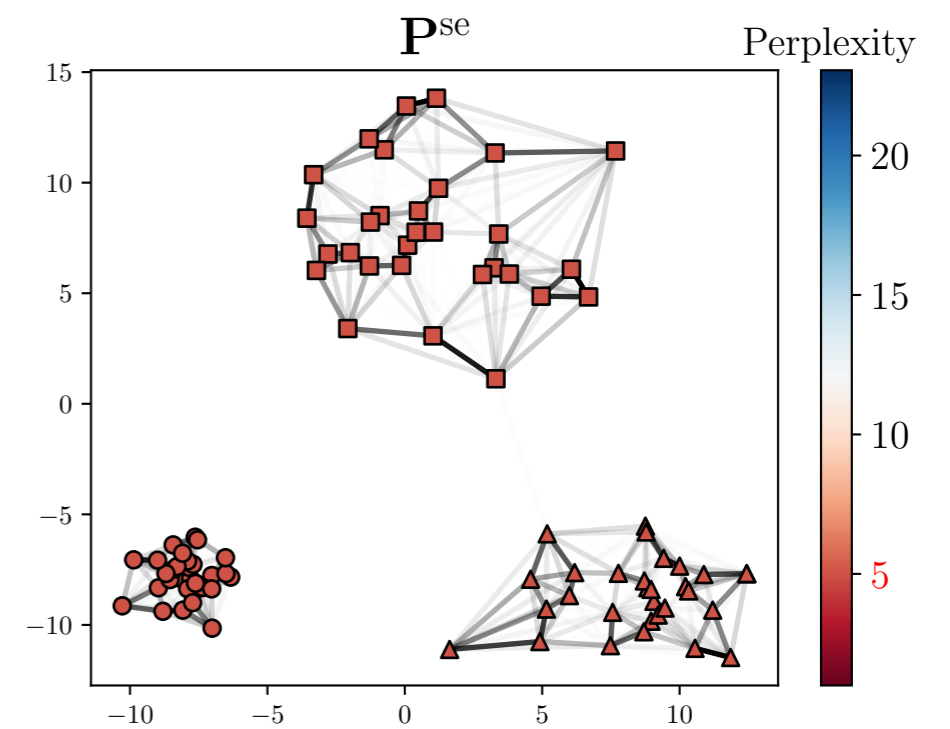
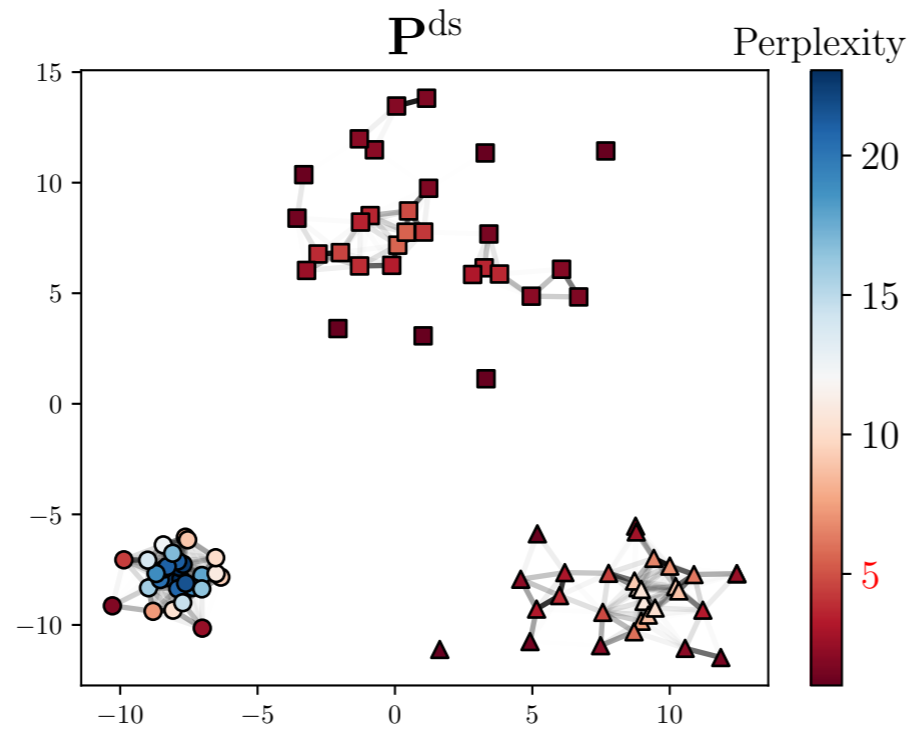
Part I: Symmetric Entropic Affinities

Part II: Application to dimensionality reduction

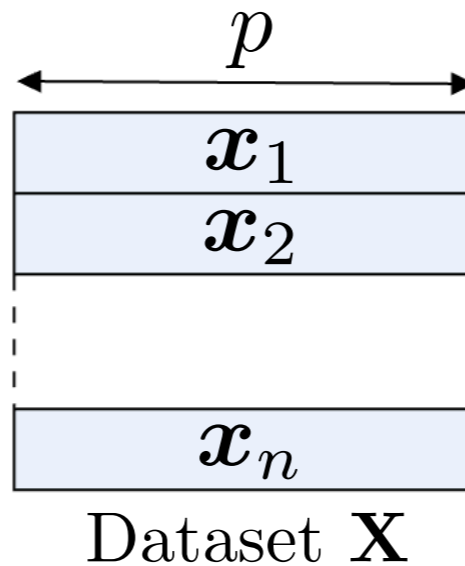
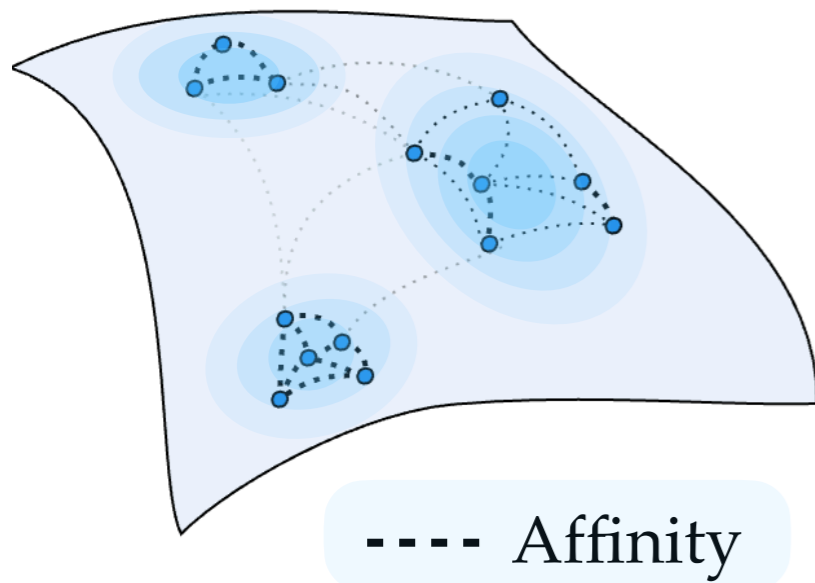
Part III: Future works



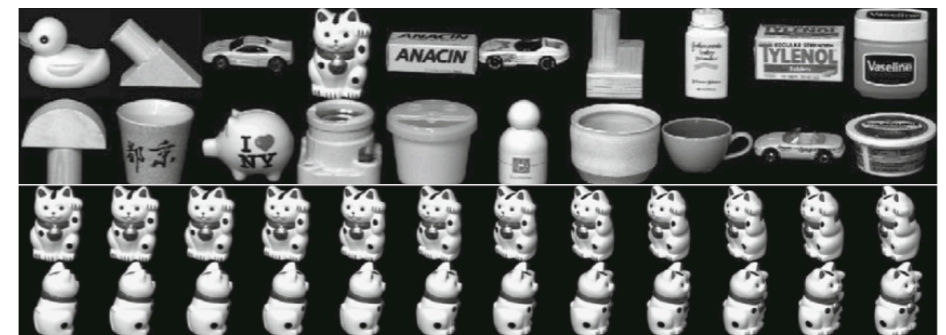
Part I: Symmetric Entropic Affinities



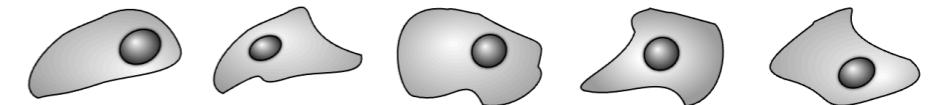
Affinity Matrices



Images



Cells

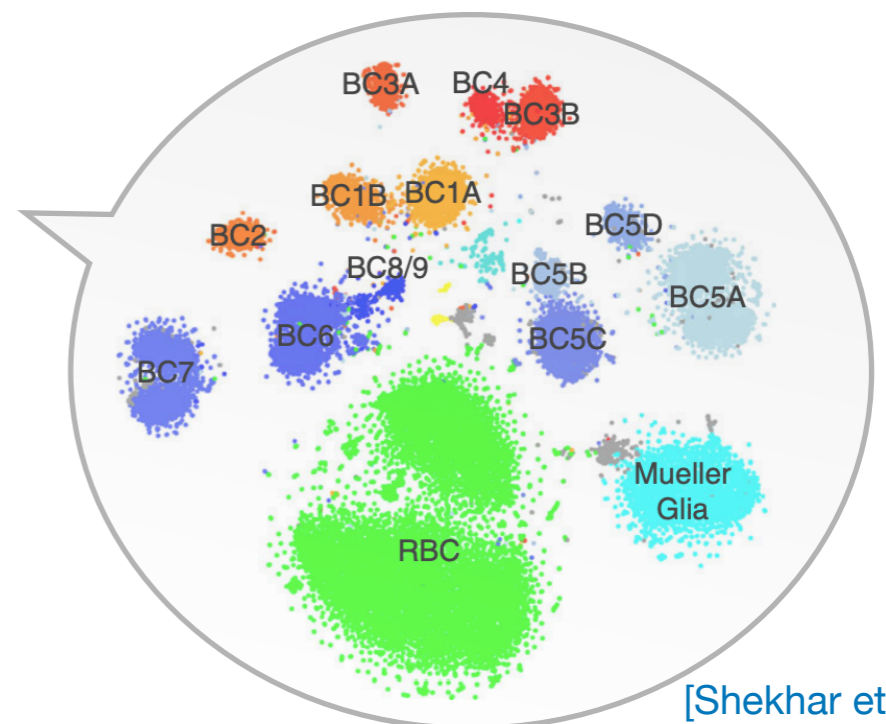


Symmetric matrix with non-negative coefficients.

Coefficient $(i, j) =$ similarity between x_i and x_j .

Useful in many ML methods :

- Dimensionality reduction [Van der Maaten and Hinton, 2008]
- Spectral clustering [Von Luxburg, 2007]
- Kernel machines [Schölkopf and Smola, 2002]
- Semi-supervised learning [Zhou et al., 2003]
- Self-supervised learning [Zbontar et al., 2021]



[Shekhar et al., 2016]

Gaussian Affinity (or Gibbs kernel)

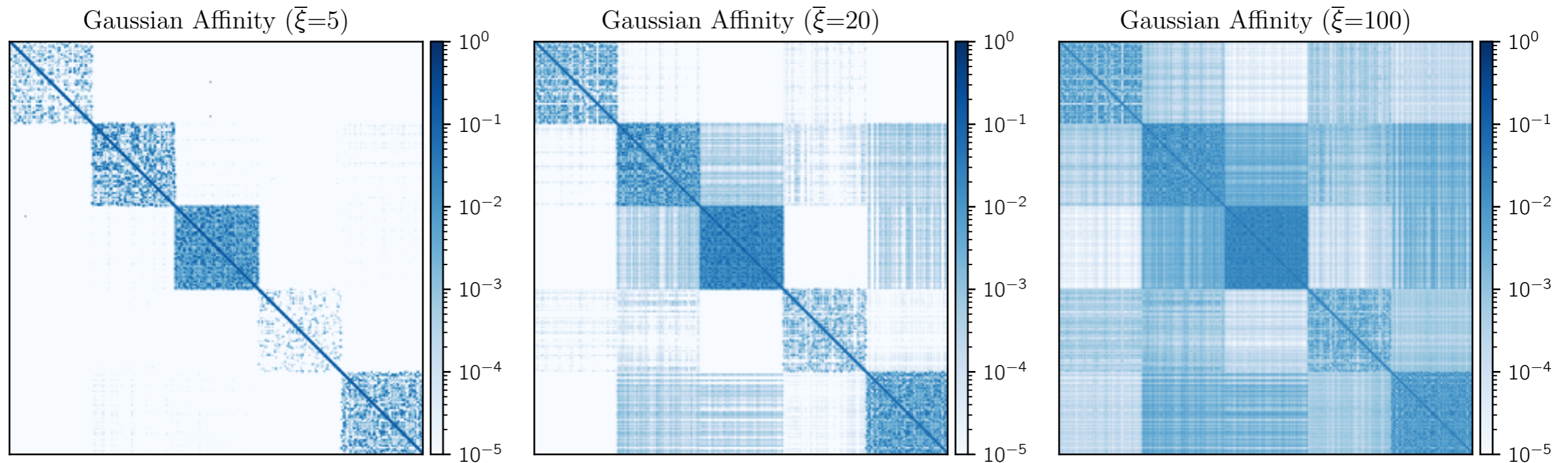


Fig : Affinity on 5 classes of the COIL Dataset [Nene et al., 1996]

Cost matrix: $\mathbf{C} \in \mathbb{R}_+^{n \times n}$ such that $\mathbf{C} = \mathbf{C}^\top$ and $C_{ij} = 0 \iff i = j$.

Example: $C_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$.

Gibbs kernel : $\mathbf{K} = \exp(-\mathbf{C}/\sigma)$.

Gaussian Affinity (or Gibbs kernel)

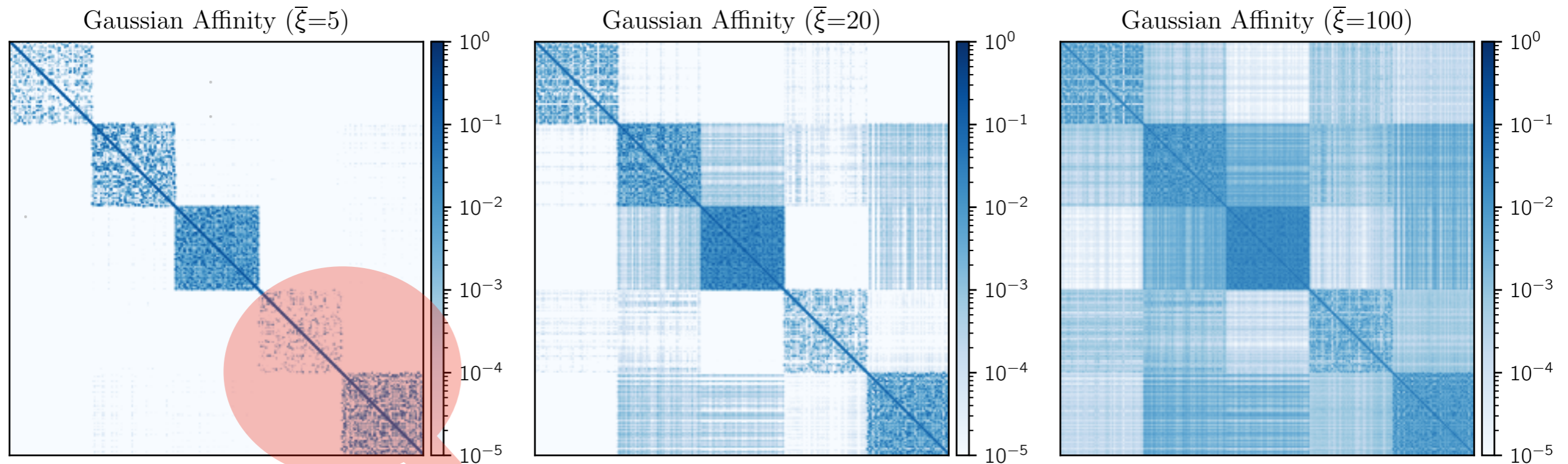


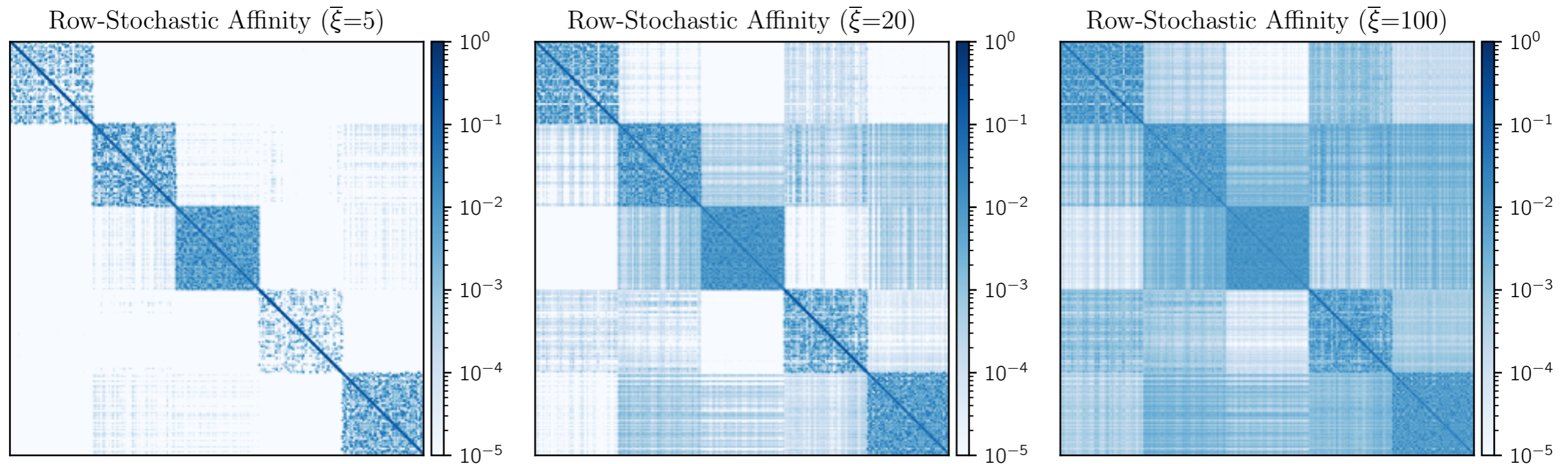
Fig : Affinity on 5 classes of the COIL Dataset [Nene et al., 1996]

Cost matrix: $\mathbf{C} \in \mathbb{R}_+^{n \times n}$ such that $\mathbf{C} = \mathbf{C}^\top$ and $C_{ij} = 0 \iff i = j$.

Example: $C_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$. **How to better reveal classes ?**

Gibbs kernel : $\mathbf{K} = \exp(-\mathbf{C}/\sigma)$.

ℓ_1 Norm - Row Stochastic

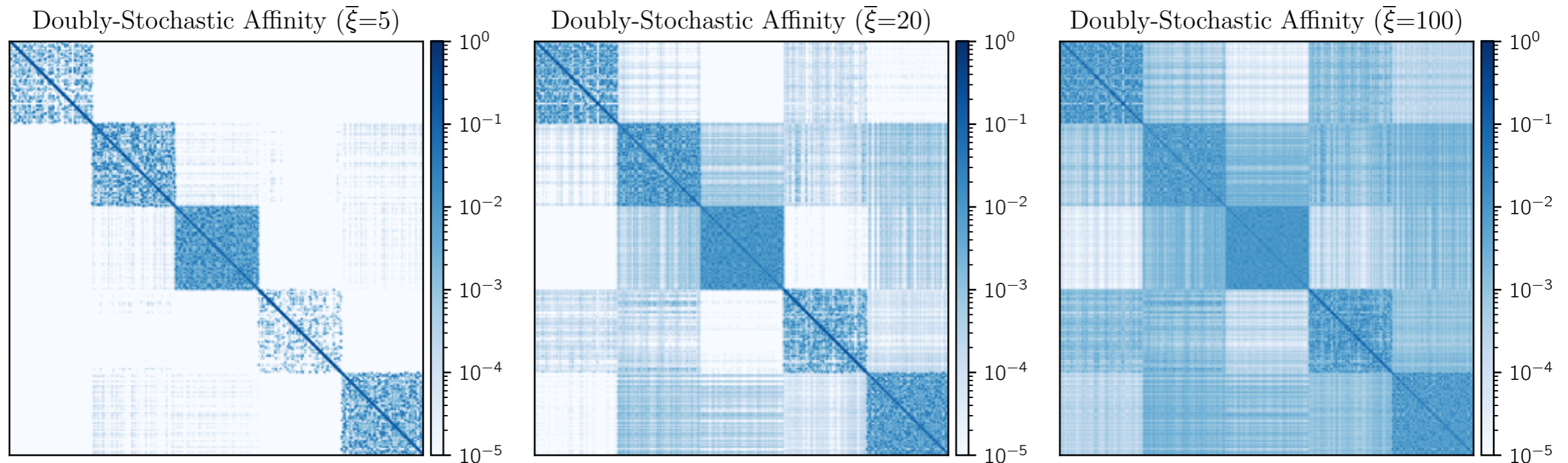


init $\mathbf{K} = \exp(-\mathbf{C}/\sigma)$

$$\mathbf{K} \leftarrow \text{diag}(\mathbf{K}\mathbf{1})^{-1}\mathbf{K}$$

normalize rows

ℓ_1 Norm - Doubly Stochastic



Sinkhorn Algorithm

While not converged:

$$\mathbf{K} \leftarrow \text{diag}(\mathbf{K}\mathbf{1})^{-1} \mathbf{K}$$

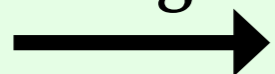
normalize rows

$$\mathbf{K} \leftarrow \mathbf{K} \text{diag}(\mathbf{1}\mathbf{K})^{-1}$$

normalize columns

$$\text{init } \mathbf{K} = \exp(-\mathbf{C}/\sigma)$$

converges to



$$\mathcal{DS} = \{ \mathbf{P} \text{ s.t. } \mathbf{P}\mathbf{1} = \mathbf{P}^\top \mathbf{1} = \mathbf{1} \} .$$

Doubly Stochastic Affinity

Sinkhorn Algorithm

$$\text{init } \mathbf{K} = \exp(-\mathbf{C}/\sigma)$$

While not converged:

$$\mathbf{K} \leftarrow \text{diag}(\mathbf{K}\mathbf{1})^{-1}\mathbf{K} \quad \# \text{ normalize rows}$$

$$\mathbf{K} \leftarrow \mathbf{K} \text{diag}(\mathbf{1}\mathbf{K})^{-1} \quad \# \text{ normalize columns}$$

converges to $\mathcal{DS} = \{\mathbf{P} \text{ s.t. } \mathbf{P}\mathbf{1} = \mathbf{P}^\top \mathbf{1} = \mathbf{1}\}.$

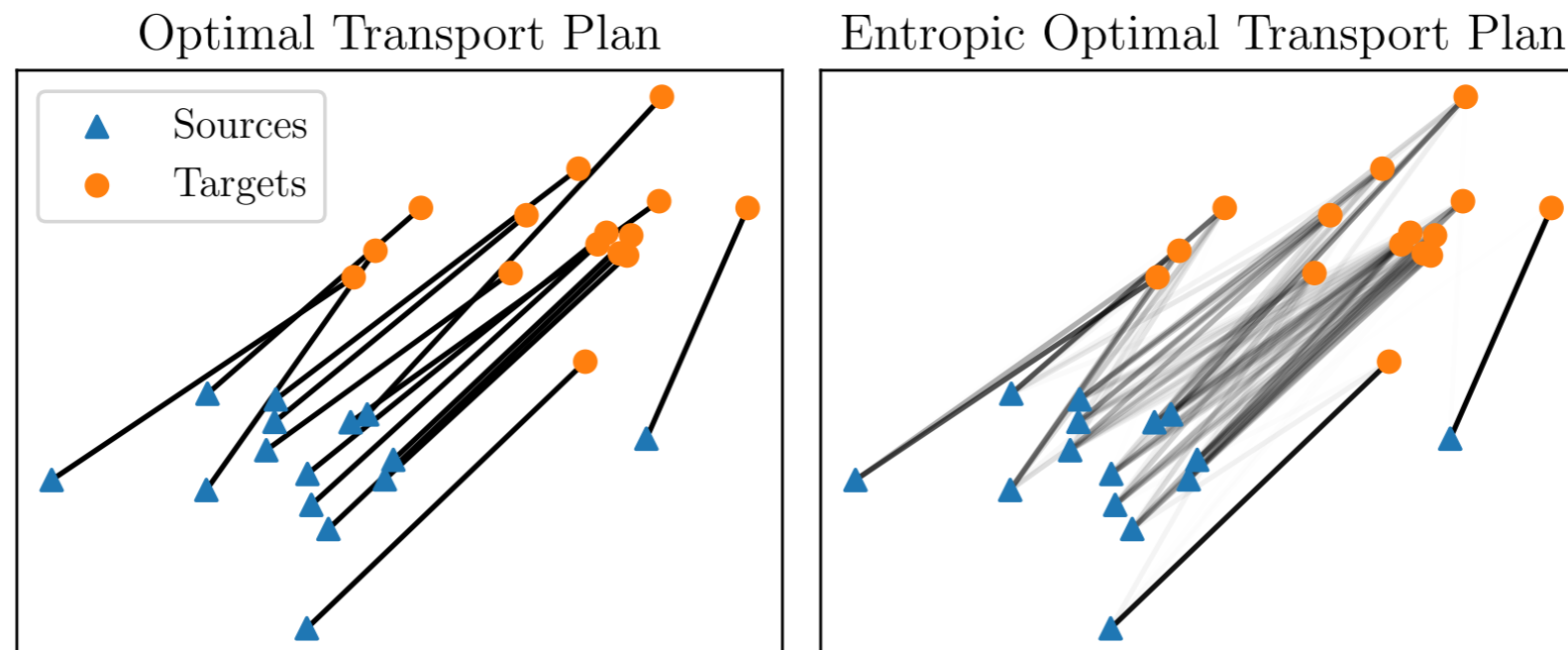
For any set \mathcal{E} : $\text{Proj}_{\mathcal{E}}^{\text{KL}}(\mathbf{K}) = \arg \min_{\mathbf{P} \in \mathcal{E}} \langle \mathbf{P}, \log(\mathbf{P} \oslash \mathbf{K}) \rangle.$

Sinkhorn iterations compute $\mathbf{P}^{\text{ds}} := \text{Proj}_{\mathcal{DS}}^{\text{KL}}(\mathbf{K}).$

\mathbf{P}^{ds} solves the optimal transport problem:

$$\min_{\mathbf{P} \in \mathbb{R}_+^{n \times n}} \langle \mathbf{P}, \mathbf{C} \rangle - \sigma \sum_i H(\mathbf{P}_{i:}) \quad \text{s.t. } \mathbf{P} \in \mathcal{DS}.$$

Entropic Optimal Transport



OT **OT plan** **Pairwise cost between sources and targets**

$$\min_{\mathbf{P} \in \mathbb{R}_+^{n \times n}} \langle \mathbf{P}, \mathbf{C} \rangle \quad \text{s.t.} \quad \mathbf{P}\mathbf{1} = \boldsymbol{\alpha}, \quad \mathbf{P}^\top \mathbf{1} = \boldsymbol{\beta}.$$

Marginals

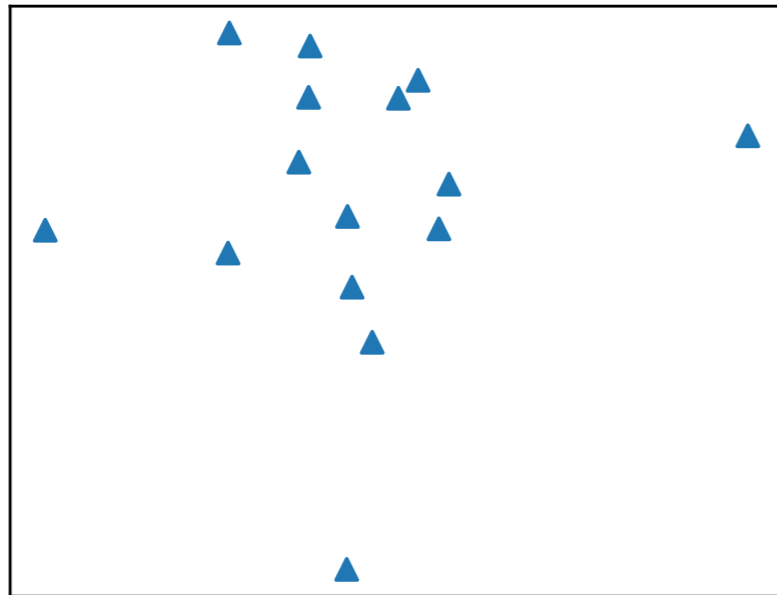
$H(\mathbf{p}) = -\langle \mathbf{p}, \log \mathbf{p} - \mathbf{1} \rangle$ is the Shannon entropy.

Entropic OT [Cuturi, 2013] **Entropic regularizer**

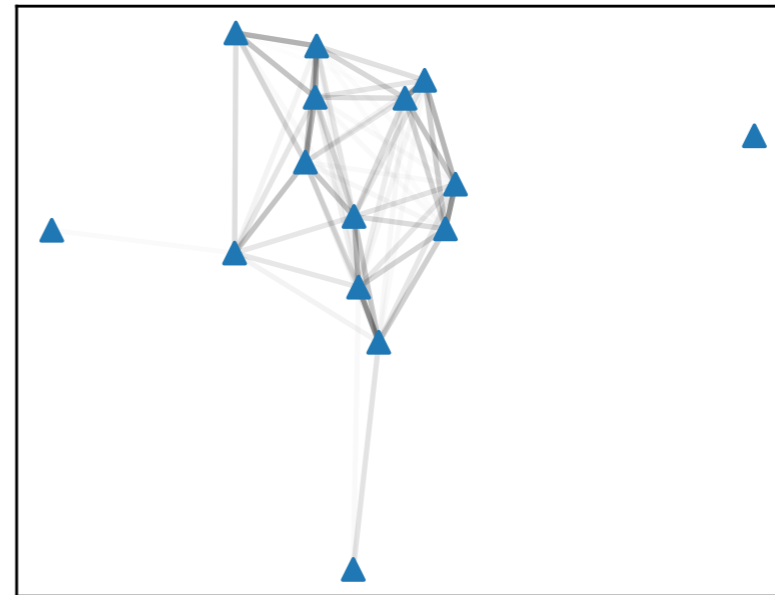
$$\min_{\mathbf{P} \in \mathbb{R}_+^{n \times n}} \langle \mathbf{P}, \mathbf{C} \rangle - \sigma \sum_i H(\mathbf{P}_{i:}) \quad \text{s.t.} \quad \mathbf{P}\mathbf{1} = \boldsymbol{\alpha}, \quad \mathbf{P}^\top \mathbf{1} = \boldsymbol{\beta}.$$

Symmetric Entropic OT

Symmetric OT Plan



Symmetric Entropic OT Plan



Sym. OT OT plan

Pairwise cost among points (SYMMETRIC)

$$\min_{\mathbf{P} \in \mathbb{R}_+^{n \times n}} \langle \mathbf{P}, \mathbf{C} \rangle \quad \text{s.t.} \quad \mathbf{P}\mathbf{1} = \mathbf{1}, \quad \mathbf{P}^\top \mathbf{1} = \mathbf{1}. \quad \text{Marginals}$$

$H(\mathbf{p}) = -\langle \mathbf{p}, \log \mathbf{p} - \mathbf{1} \rangle$ is the Shannon entropy.

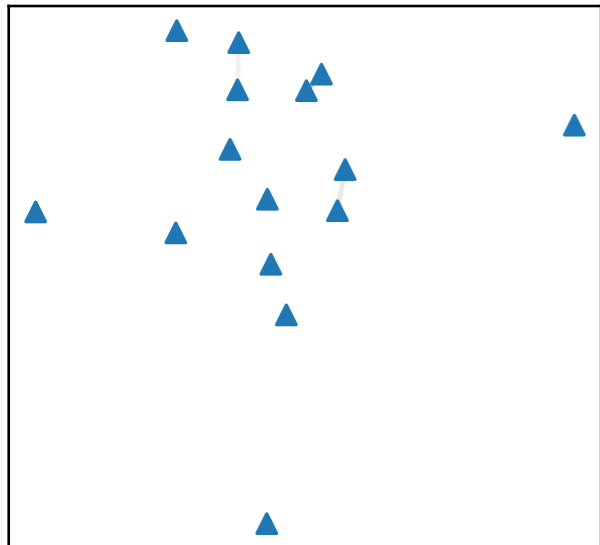
Sym. Entropic OT

Entropic regularizer

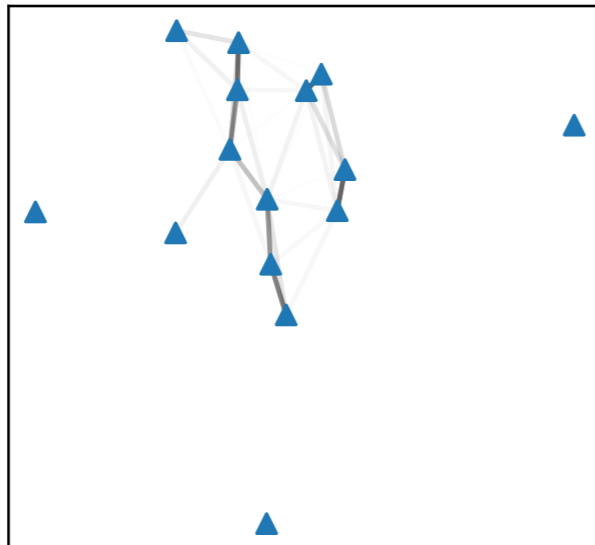
$$\min_{\mathbf{P} \in \mathbb{R}_+^{n \times n}} \langle \mathbf{P}, \mathbf{C} \rangle - \sigma \sum_i H(\mathbf{P}_{i:}) \quad \text{s.t.} \quad \mathbf{P}\mathbf{1} = \mathbf{1}, \quad \mathbf{P}^\top \mathbf{1} = \mathbf{1}.$$

Symmetric Entropic OT

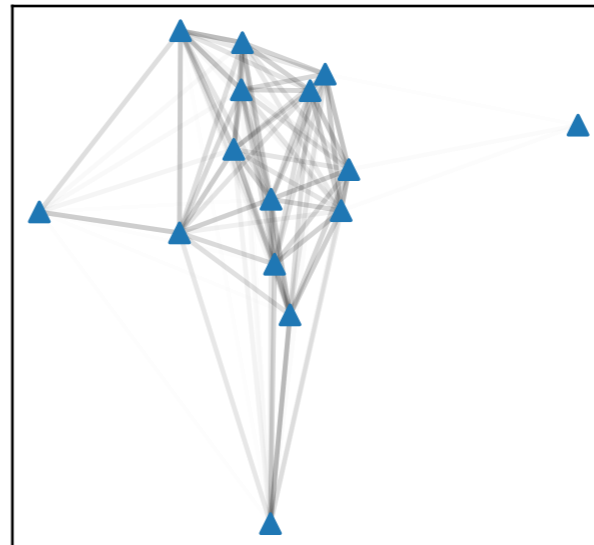
Sym. Entropic OT Plan, $\sigma=0.1$



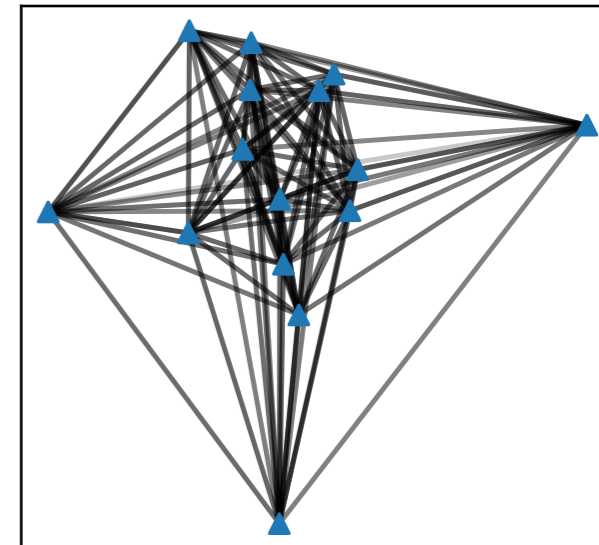
Sym. Entropic OT Plan, $\sigma=1.0$



Sym. Entropic OT Plan, $\sigma=10.0$



Sym. Entropic OT Plan, $\sigma=100.0$



Sym. Entropic OT

Entropic regularizer

$$\min_{\mathbf{P} \in \mathbb{R}_+^{n \times n}} \langle \mathbf{P}, \mathbf{C} \rangle - \sigma \sum_i H(\mathbf{P}_{i:}) \quad \text{s.t.} \quad \mathbf{P}\mathbf{1} = \mathbf{1}, \mathbf{P}^\top \mathbf{1} = \mathbf{1}.$$

Interpolates between:

- identity \mathbf{I}_n for $\sigma \rightarrow 0$.
- uniform $\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top$ for $\sigma \rightarrow \infty$.

Constrained Formulation

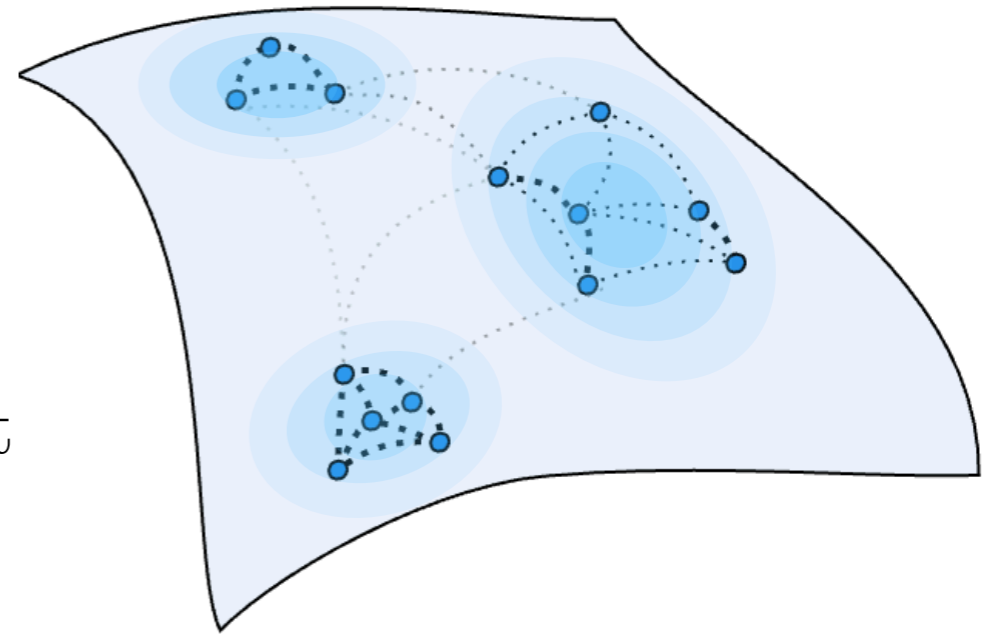
$$\begin{aligned} & \min_{\mathbf{P} \in \mathbb{R}_+^{n \times n}, \mathbf{P}\mathbf{1}=\mathbf{1}, \mathbf{P}^\top \mathbf{1}=\mathbf{1}} \langle \mathbf{P}, \mathbf{C} \rangle \\ & \text{s.t.} \quad \sum_i H(\mathbf{P}_{i:}) \geq \eta. \end{aligned}$$

Entropy OT plan

Entropic Affinity

Data has **varying noise levels**.

We can control the entropy in each point with **adaptive bandwidths**.



Definition [Hinton, Roweis 2002]

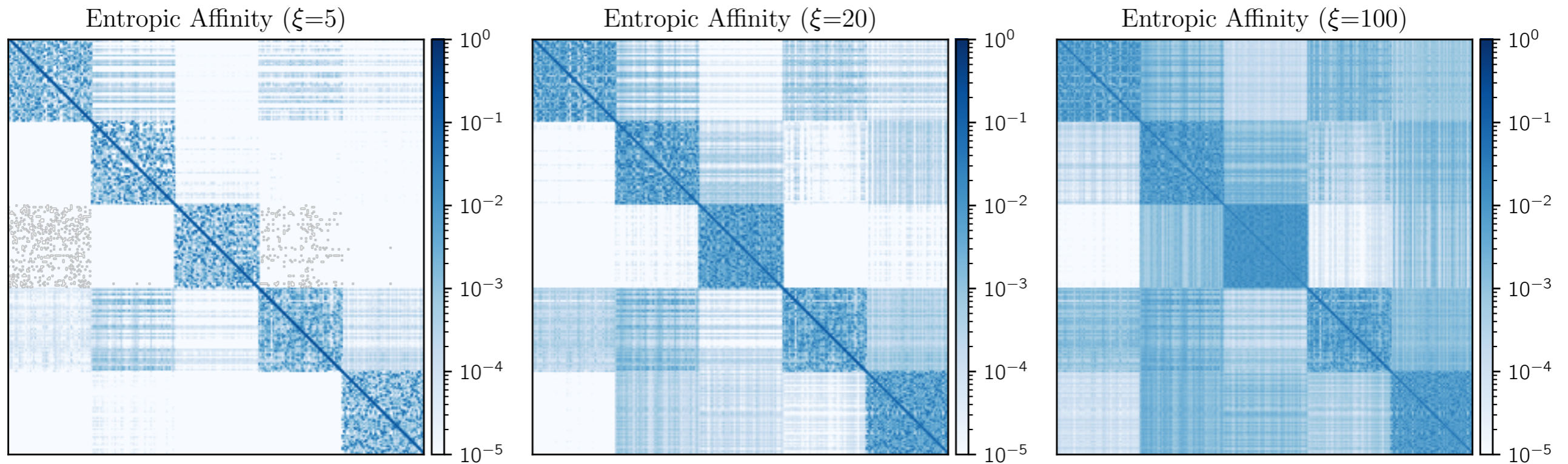
$$\forall i, \forall j, P_{ij}^e = \frac{\exp(-C_{ij}/\varepsilon_i^*)}{\sum_{\ell} \exp(-C_{i\ell}/\varepsilon_i^*)}$$

$$\text{with } \varepsilon_i^* \in \mathbb{R}_+^* \text{ s.t. } H(\mathbf{P}_{i:\cdot}^e) = \log \xi + 1.$$

$H(\mathbf{p}) = -\langle \mathbf{p}, \log \mathbf{p} - \mathbf{1} \rangle$ is the Shannon entropy.

$\xi \in \llbracket 1, n \rrbracket$ is the **perplexity** parameter.

Entropic Affinity



Definition [Hinton, Roweis 2002]

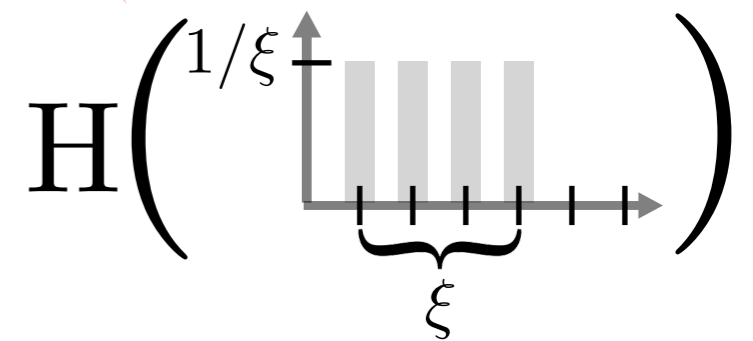
$$\forall i, \forall j, P_{ij}^e = \frac{\exp(-C_{ij}/\varepsilon_i^*)}{\sum_{\ell} \exp(-C_{i\ell}/\varepsilon_i^*)}$$

ξ effective neighbors

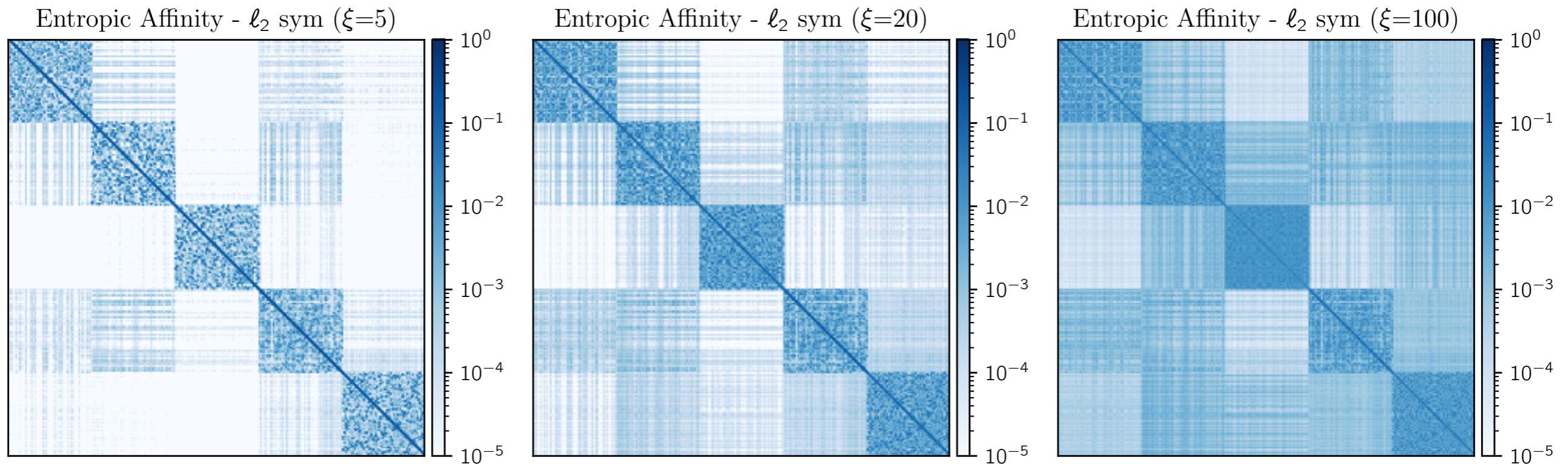
with $\varepsilon_i^* \in \mathbb{R}_+^*$ s.t. $H(\mathbf{P}_{i:\cdot}^e) = \log \xi + 1$.

$H(\mathbf{p}) = -\langle \mathbf{p}, \log \mathbf{p} - \mathbf{1} \rangle$ is the Shannon entropy.

$\xi \in \llbracket 1, n \rrbracket$ is the **perplexity** parameter.



Entropic Affinity



Definition [Hinton, Roweis 2002]

$$\forall i, \forall j, P_{ij}^e = \frac{\exp(-C_{ij}/\varepsilon_i^*)}{\sum_{\ell} \exp(-C_{i\ell}/\varepsilon_i^*)}$$

with $\varepsilon_i^* \in \mathbb{R}_+^*$ s.t. $H(\mathbf{P}_{i:}^e) = \log \xi + 1$.

\mathbf{P}^e is not symmetric.

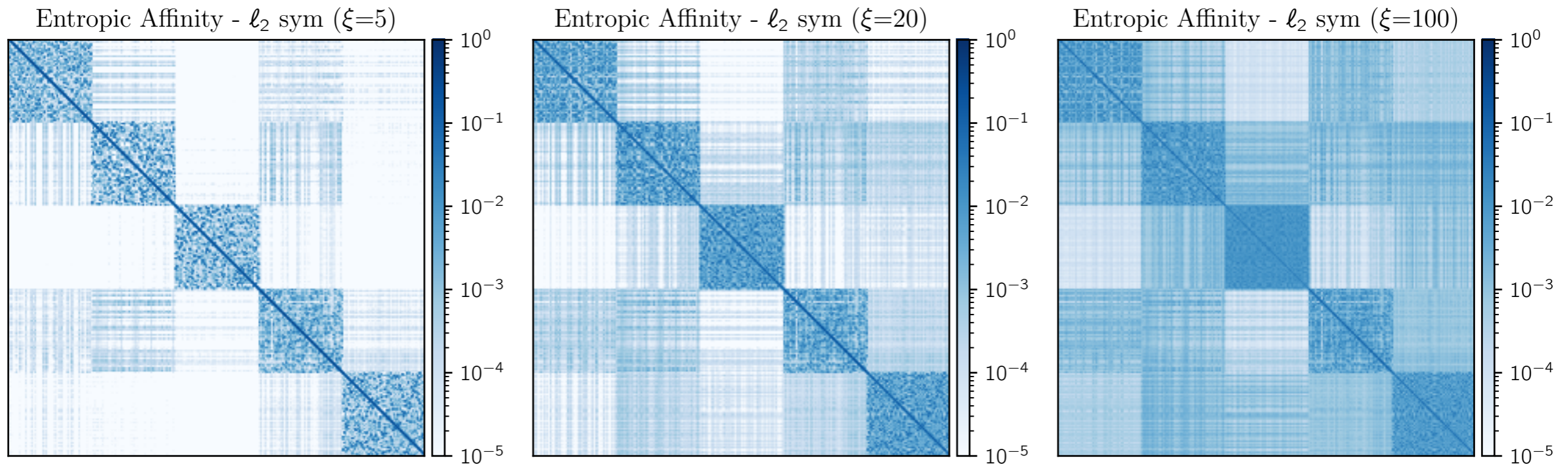
→ $\overline{\mathbf{P}}^e = \frac{1}{2}(\mathbf{P}^e + \mathbf{P}^{e\top})$ is used in practice. [Van der Maaten and Hinton, 2008]

t-SNE algorithm

↓

$$\overline{\mathbf{P}}^e = \text{Proj}_{\mathcal{S}}^{\ell_2}(\mathbf{P}^e)$$

Entropic Affinity



Definition [Hinton, Roweis 2002]

$$\forall i, \forall j, P_{ij}^e = \frac{\exp(-C_{ij}/\varepsilon_i^*)}{\sum_{\ell} \exp(-C_{i\ell}/\varepsilon_i^*)}$$

with $\varepsilon_i^* \in \mathbb{R}_+^*$ s.t. $H(\mathbf{P}_{i:}^e) = \log \xi + 1$.

\mathbf{P}^e is not symmetric.

→ $\overline{\mathbf{P}}^e = \frac{1}{2}(\mathbf{P}^e + \mathbf{P}^{e\top})$ is used in practice. [Van der Maaten and Hinton, 2008]

t-SNE algorithm

↓

$$\overline{\mathbf{P}}^e = \text{Proj}_{\mathcal{S}}^{\ell_2}(\mathbf{P}^e)$$

Breaks the construction of entropic affinities.

Affinity Panorama*

Gibbs kernel

$$\mathbf{K} = \exp(-\mathbf{C}/\sigma)$$

Doubly-Sto

$$\mathbf{P}^{ds} = \text{Proj}_{\mathcal{D}_S}^{\text{KL}}(\mathbf{K})$$

Entropic

$$\mathbf{P}^e = \text{Proj}_{\mathcal{H}_\xi}^{\text{KL}}(\mathbf{K})$$

Entropic (ℓ_2 Sym)

$$\overline{\mathbf{P}}^e = \text{Proj}_S^{\ell_2}(\mathbf{P}^e)$$

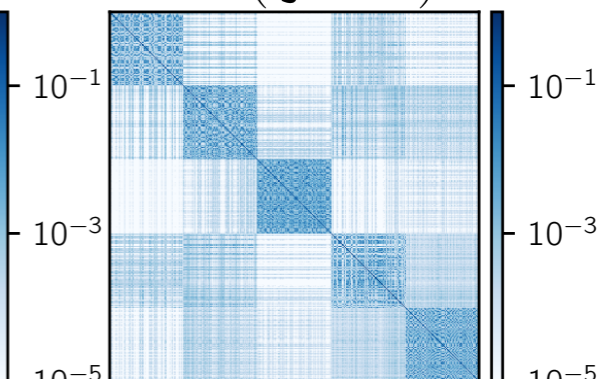
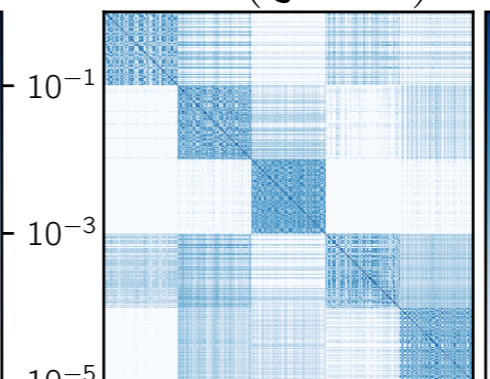
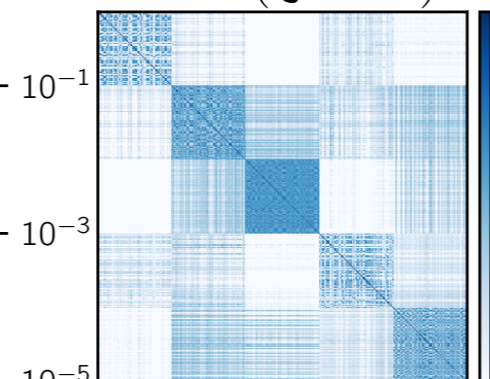
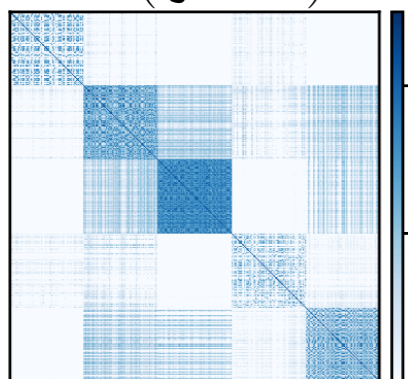
$\mathbf{K} (\bar{\xi}=20)**$

$\mathbf{P}^{ds} (\bar{\xi}=20)$

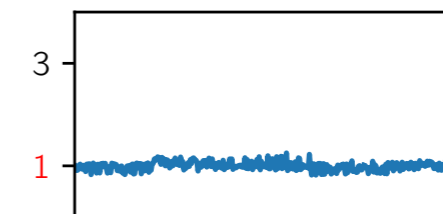
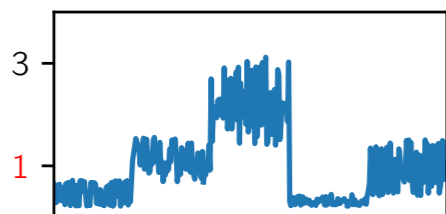
$\mathbf{P}^e (\xi=20)$

$\overline{\mathbf{P}}^e (\xi=20)$

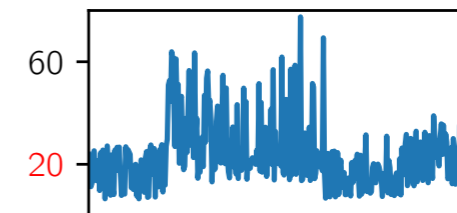
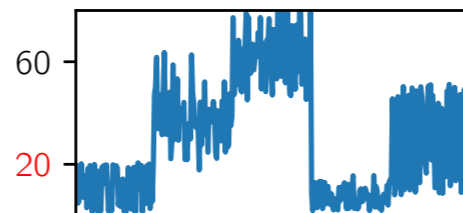
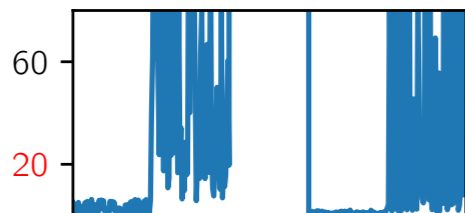
Affinity



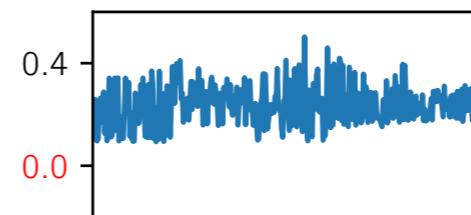
ℓ_1 norm



Perplexity



Symmetry



Sample

Sample

Sample

Sample

* On 5 classes of the COIL Dataset [Nene et al., 1996]

** $\bar{\xi}$ is average perplexity \rightarrow same global entropy as with $\xi = \bar{\xi}$.

Affinity Panorama*

Gibbs kernel

$$\mathbf{K} = \exp(-\mathbf{C}/\sigma)$$

Doubly-Sto

$$\mathbf{P}^{ds} = \text{Proj}_{\mathcal{D}_S}^{\text{KL}}(\mathbf{K})$$

Entropic

$$\mathbf{P}^e = \text{Proj}_{\mathcal{H}_\xi}^{\text{KL}}(\mathbf{K})$$

Entropic (ℓ_2 Sym)

$$\overline{\mathbf{P}}^e = \text{Proj}_{\mathcal{S}}^{\ell_2}(\mathbf{P}^e)$$

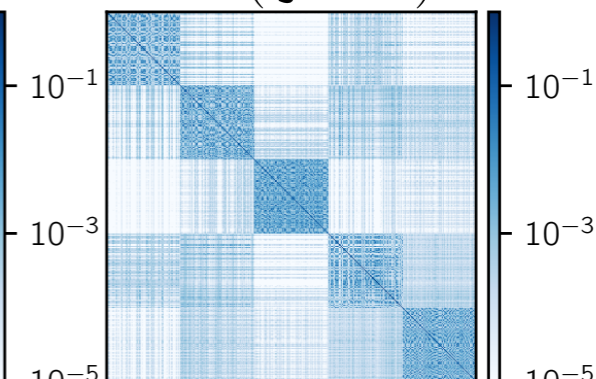
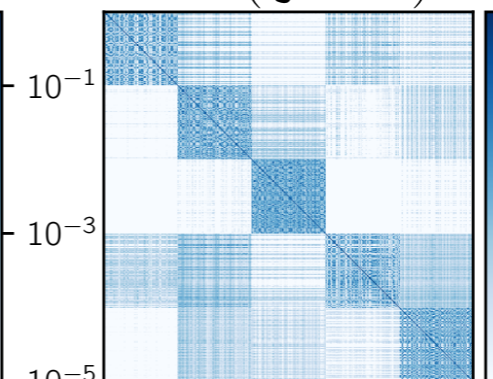
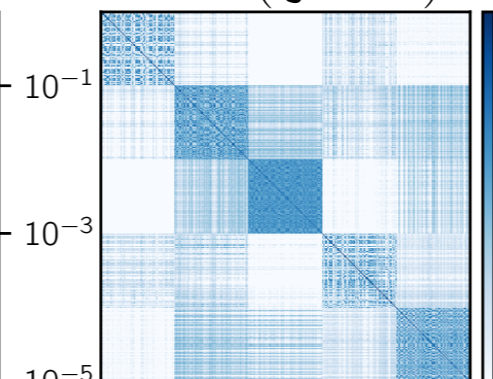
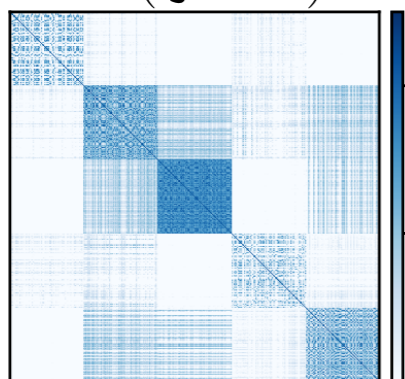
\mathbf{K} ($\bar{\xi}=20$)**

\mathbf{P}^{ds} ($\bar{\xi}=20$)

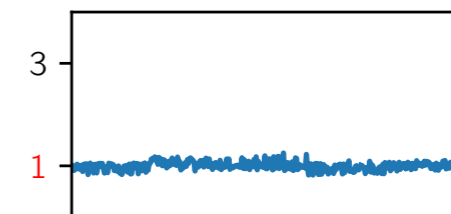
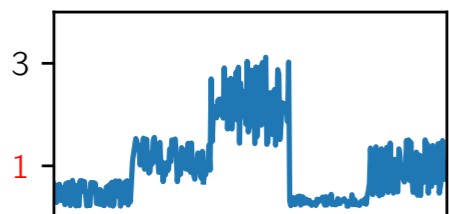
\mathbf{P}^e ($\xi=20$)

$\overline{\mathbf{P}}^e$ ($\xi=20$)

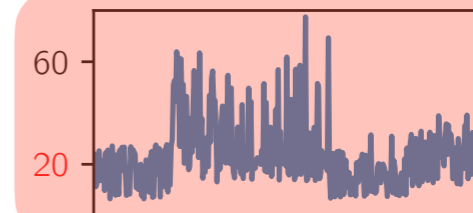
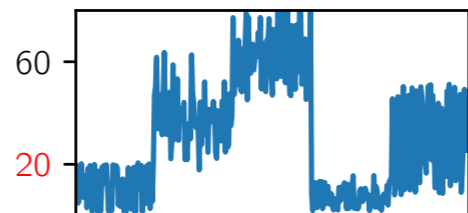
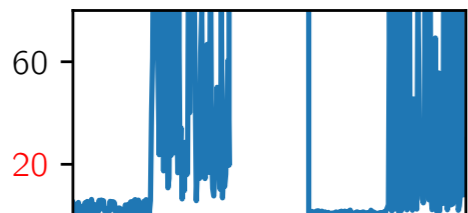
Affinity



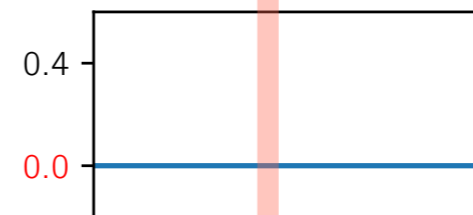
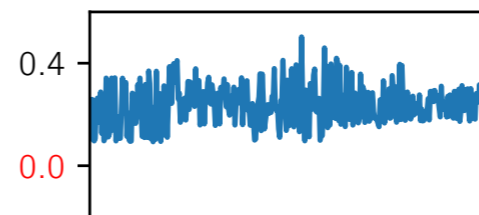
ℓ_1 norm



Perplexity



Symmetry



* On 5 classes of the COIL Dataset [Nene et al., 1996]

** $\bar{\xi}$ is average perplexity \rightarrow same global entropy as with $\xi = \bar{\xi}$.

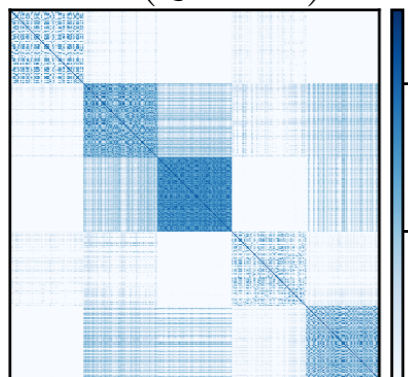
Entropies not controlled.

Affinity Panorama

Gibbs kernel

$$\mathbf{K} = \exp(-\mathbf{C}/\sigma)$$

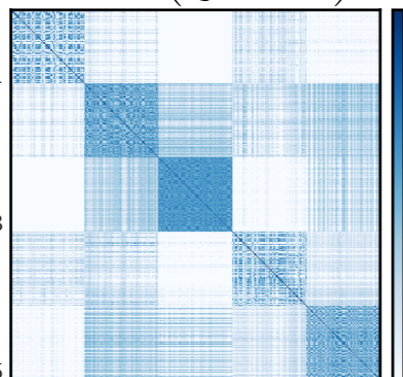
\mathbf{K} ($\bar{\xi}=20$)



Doubly-Sto

$$\mathbf{P}^{ds} = \text{Proj}_{\mathcal{D}_S}^{\text{KL}}(\mathbf{K})$$

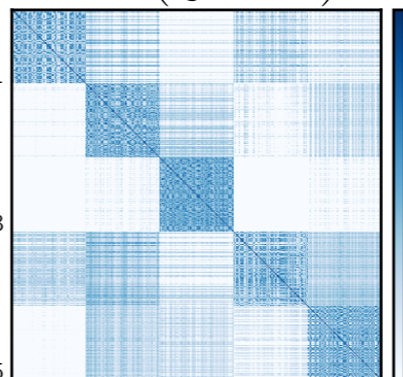
\mathbf{P}^{ds} ($\bar{\xi}=20$)



Entropic

$$\mathbf{P}^e = \text{Proj}_{\mathcal{H}_\xi}^{\text{KL}}(\mathbf{K})$$

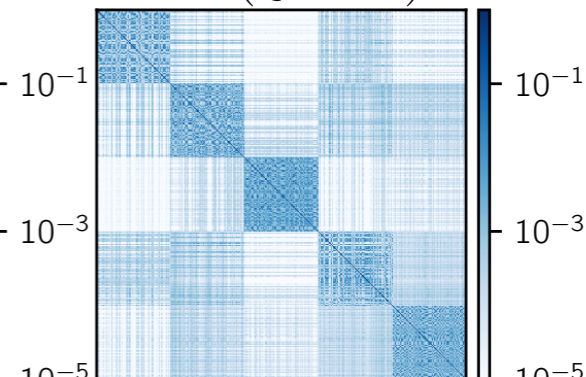
\mathbf{P}^e ($\xi=20$)



Entropic (ℓ_2 Sym)

$$\bar{\mathbf{P}}^e = \text{Proj}_{\mathcal{S}}^{\ell_2}(\mathbf{P}^e)$$

$\bar{\mathbf{P}}^e$ ($\xi=20$)

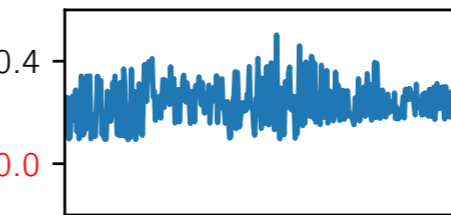
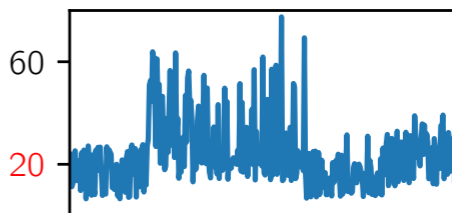
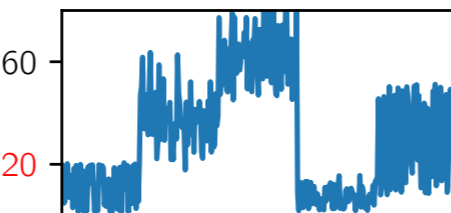
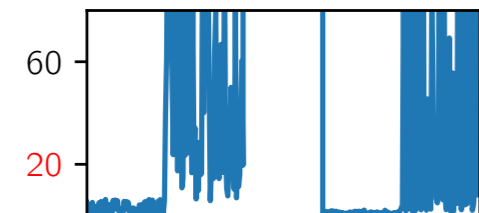
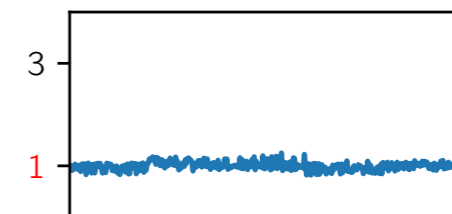
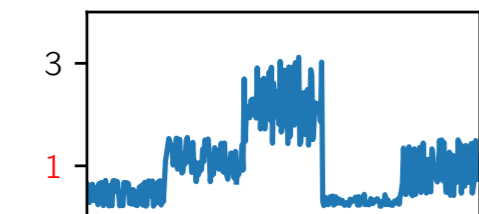


Affinity

ℓ_1 norm

Perplexity

Symmetry



Sample

Sample

Sample

Sample

Can we control ℓ_1 norm, entropy and symmetry ?

Symmetric Entropic Affinity

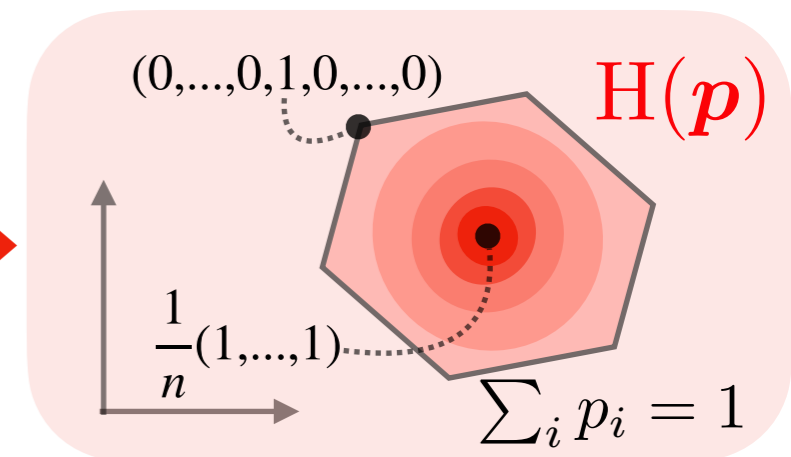
$$\mathcal{H}_\xi := \{\mathbf{P} \in \mathbb{R}_+^{n \times n} \text{ s.t. } \mathbf{P}\mathbf{1} = \mathbf{1} \text{ and } \forall i, H(\mathbf{P}_{i:}) \geq \log \xi + 1\}$$

Entropic Affinity as OT

$$\mathbf{P}^e = \arg \min_{\mathbf{P} \in \mathcal{H}_\xi} \langle \mathbf{P}, \mathbf{C} \rangle.$$

The constraints in \mathcal{H}_ξ are saturated at the optimum.

Symmetric matrices $\mathcal{S} = \{\mathbf{P} \text{ s.t. } \mathbf{P} = \mathbf{P}^\top\}$.



Definition

$$\mathbf{P}^{\text{se}} := \arg \min_{\mathbf{P} \in \mathcal{H}_\xi \cap \mathcal{S}} \langle \mathbf{P}, \mathbf{C} \rangle.$$

OURS

Symmetric Entropic Affinity

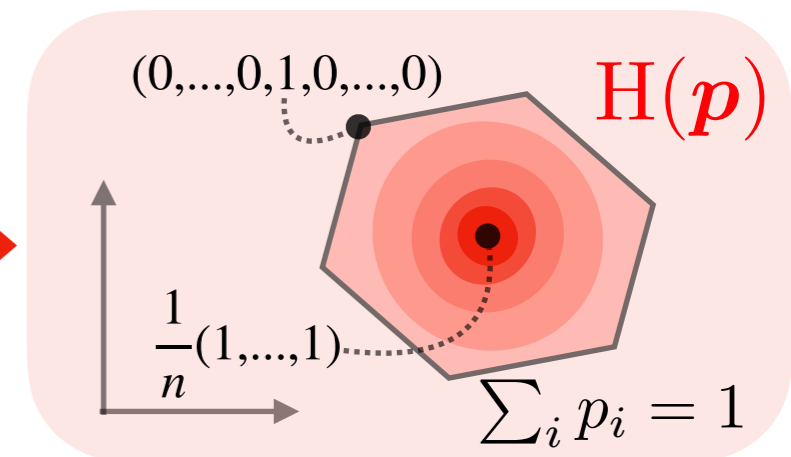
$$\mathcal{H}_\xi := \{\mathbf{P} \in \mathbb{R}_+^{n \times n} \text{ s.t. } \mathbf{P}\mathbf{1} = \mathbf{1} \text{ and } \forall i, H(\mathbf{P}_{i:}) \geq \log \xi + 1\}$$

Entropic Affinity as OT

$$\mathbf{P}^e = \arg \min_{\mathbf{P} \in \mathcal{H}_\xi} \langle \mathbf{P}, \mathbf{C} \rangle.$$

The constraints in \mathcal{H}_ξ are saturated at the optimum.

Symmetric matrices $\mathcal{S} = \{\mathbf{P} \text{ s.t. } \mathbf{P} = \mathbf{P}^\top\}$.



Definition

$$\mathbf{P}^{\text{se}} := \arg \min_{\mathbf{P} \in \mathcal{H}_\xi \cap \mathcal{S}} \langle \mathbf{P}, \mathbf{C} \rangle.$$

Enforce Symmetry

OURS

Symmetric Entropic Affinity

$$\mathcal{H}_\xi := \{\mathbf{P} \in \mathbb{R}_+^{n \times n} \text{ s.t. } \mathbf{P}\mathbf{1} = \mathbf{1} \text{ and } \forall i, H(\mathbf{P}_{i:}) \geq \log \xi + 1\}$$

Definition

OURS

$$\mathbf{P}^{\text{se}} := \arg \min_{\mathbf{P} \in \mathcal{H}_\xi \cap \mathcal{S}} \langle \mathbf{P}, \mathbf{C} \rangle.$$

Enforce Symmetry

Property

For at least $n - 1$ indices $i \in \llbracket n \rrbracket$, it holds $H(\mathbf{P}_{i:}^{\text{se}}) = \log \xi + 1$.

In practice, we have n saturated entropies.

Dual Ascent

$$\mathbf{P}^{\text{se}} = \exp((\boldsymbol{\lambda}^* \oplus \boldsymbol{\lambda}^* - 2\mathbf{C}) \oslash (\boldsymbol{\gamma}^* \oplus \boldsymbol{\gamma}^*))$$

where $\boldsymbol{\lambda}^*$ and $\boldsymbol{\gamma}^*$ are computed using dual ascent.

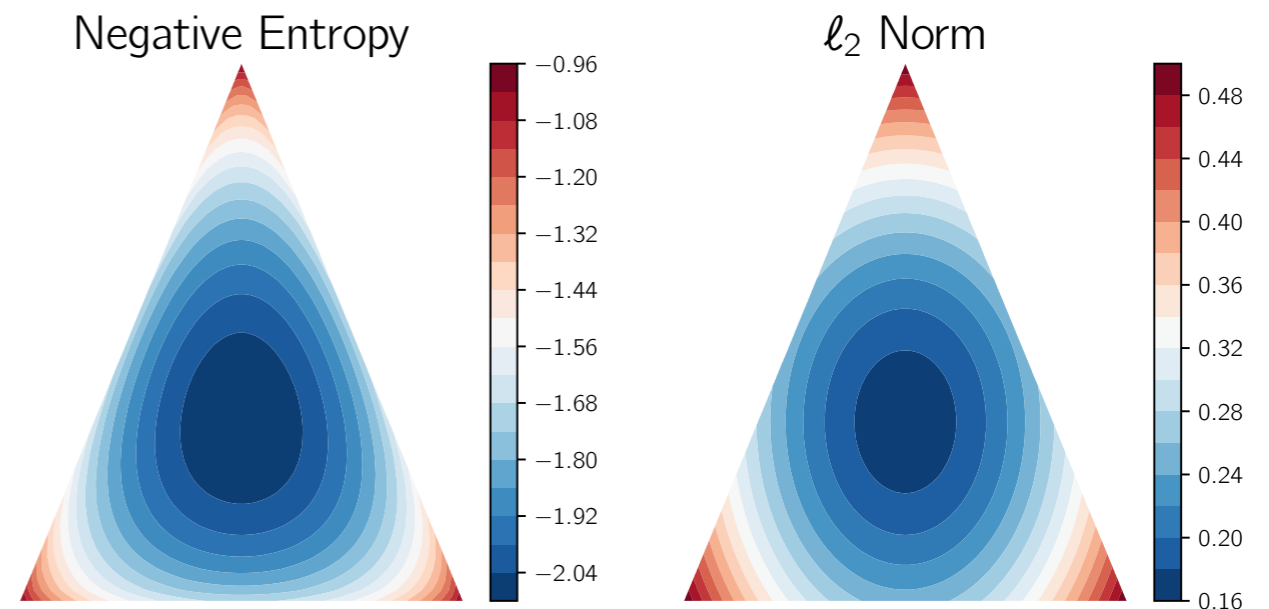
Formulation via Projection

For any set \mathcal{E} :

$$\text{Proj}_{\mathcal{E}}^{\text{KL}}(\mathbf{K}) = \arg \min_{\mathbf{P} \in \mathcal{E}} \langle \mathbf{P}, \log(\mathbf{P} \oslash \mathbf{K}) \rangle$$

$$\text{Proj}_{\mathcal{E}}^{\ell_2}(\mathbf{K}) = \arg \min_{\mathbf{P} \in \mathcal{E}} \|\mathbf{P} - \mathbf{K}\|_2$$

Gibbs kernel : $\mathbf{K}_{\sigma} := \exp(-\mathbf{C}/\sigma)$.



Entropic Affinity as Projection

$$\mathbf{P}^e = \text{Proj}_{\mathcal{H}_{\xi}}^{\text{KL}}(\mathbf{K}_{\sigma}) \text{ for any } 0 < \sigma \leq \min_i \varepsilon_i^* .$$

$\overline{\mathbf{P}}^e = \text{Proj}_{\mathcal{S}}^{\ell_2}(\mathbf{P}^e)$ is a mixture of KL and ℓ_2 projections.

Symmetric Entropic Affinity as Projection

$$\mathbf{P}^{\text{se}} = \text{Proj}_{\mathcal{H}_{\xi} \cap \mathcal{S}}^{\text{KL}}(\mathbf{K}_{\sigma}) \text{ for any } 0 < \sigma \leq \min_i \gamma_i^* .$$

Affinity Panorama*

t-SNE

OURS

Gibbs kernel

Doubly-Sto

Entropic

Entropic (ℓ_2 Sym)

Sym-Entropic

$$\mathbf{K} = \exp(-\mathbf{C}/\sigma)$$

$$\mathbf{P}^{ds} = \text{Proj}_{\mathcal{DS}}^{\text{KL}}(\mathbf{K})$$

$$\mathbf{P}^e = \text{Proj}_{\mathcal{H}_\xi}^{\text{KL}}(\mathbf{K})$$

$$\overline{\mathbf{P}}^e = \text{Proj}_{\mathcal{S}}^{\ell_2}(\mathbf{P}^e)$$

$$\mathbf{P}^{se} = \text{Proj}_{\mathcal{H}_\xi \cap \mathcal{S}}^{\text{KL}}(\mathbf{K})$$

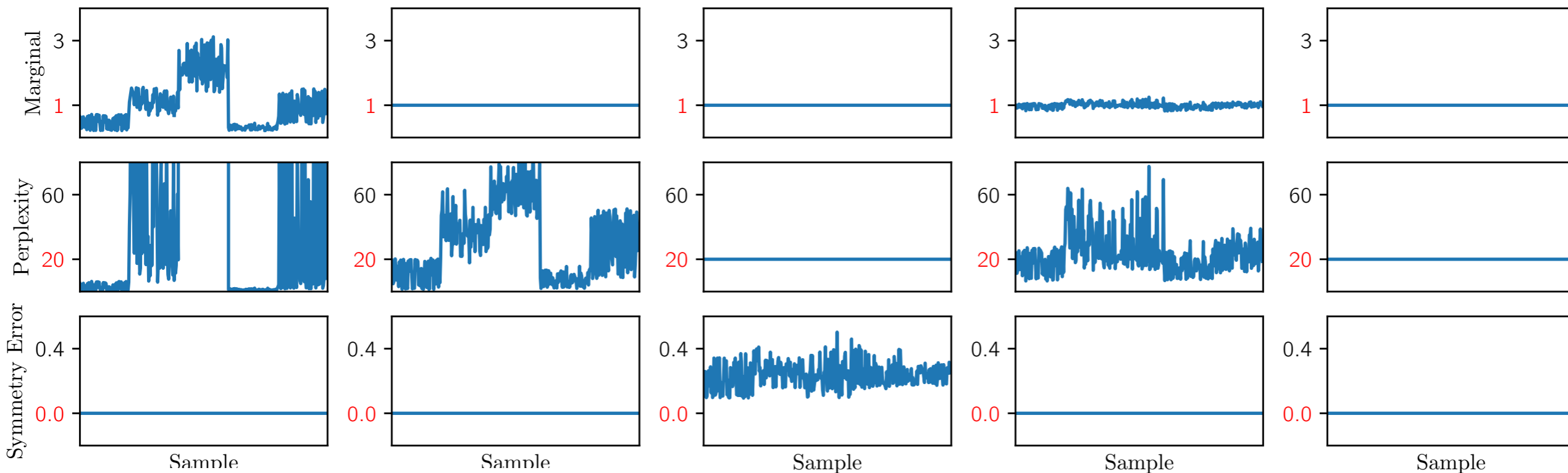
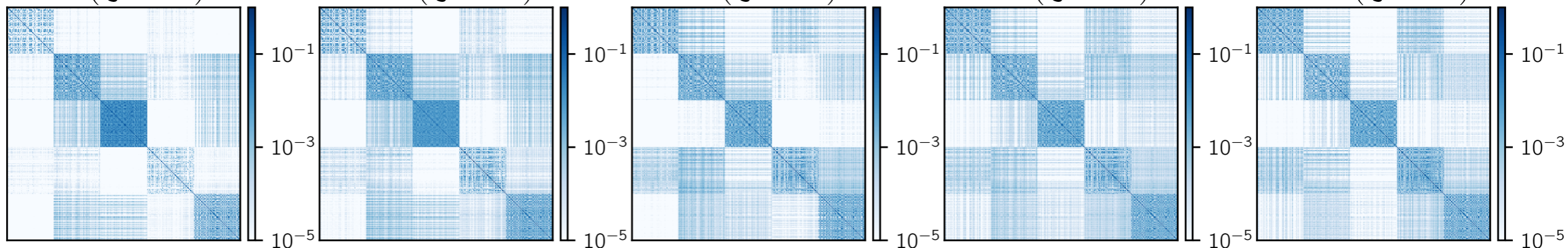
\mathbf{K} ($\bar{\xi}=20$)**

\mathbf{P}^{ds} ($\bar{\xi}=20$)

\mathbf{P}^e ($\xi=20$)

$\overline{\mathbf{P}}^e$ ($\xi=20$)

\mathbf{P}^{se} ($\xi=20$)



* On 5 classes of the COIL Dataset [Nene et al., 1996]

** $\bar{\xi}$ is average perplexity \rightarrow same global entropy as with $\xi = \bar{\xi}$.

Affinity Panorama*

OURS

Gibbs kernel

$$\mathbf{K} = \exp(-\mathbf{C}/\sigma)$$

Doubly-Sto

$$\mathbf{P}^{ds} = \text{Proj}_{\mathcal{DS}}^{\text{KL}}(\mathbf{K})$$

Entropic

$$\mathbf{P}^e = \text{Proj}_{\mathcal{H}_\xi}^{\text{KL}}(\mathbf{K})$$

Entropic (ℓ_2 Sym) Sym-Entropic

$$\bar{\mathbf{P}}^e = \text{Proj}_{\mathcal{S}}^{\ell_2}(\mathbf{P}^e)$$

$$\mathbf{P}^{se} = \text{Proj}_{\mathcal{H}_\xi \cap \mathcal{S}}^{\text{KL}}(\mathbf{K})$$

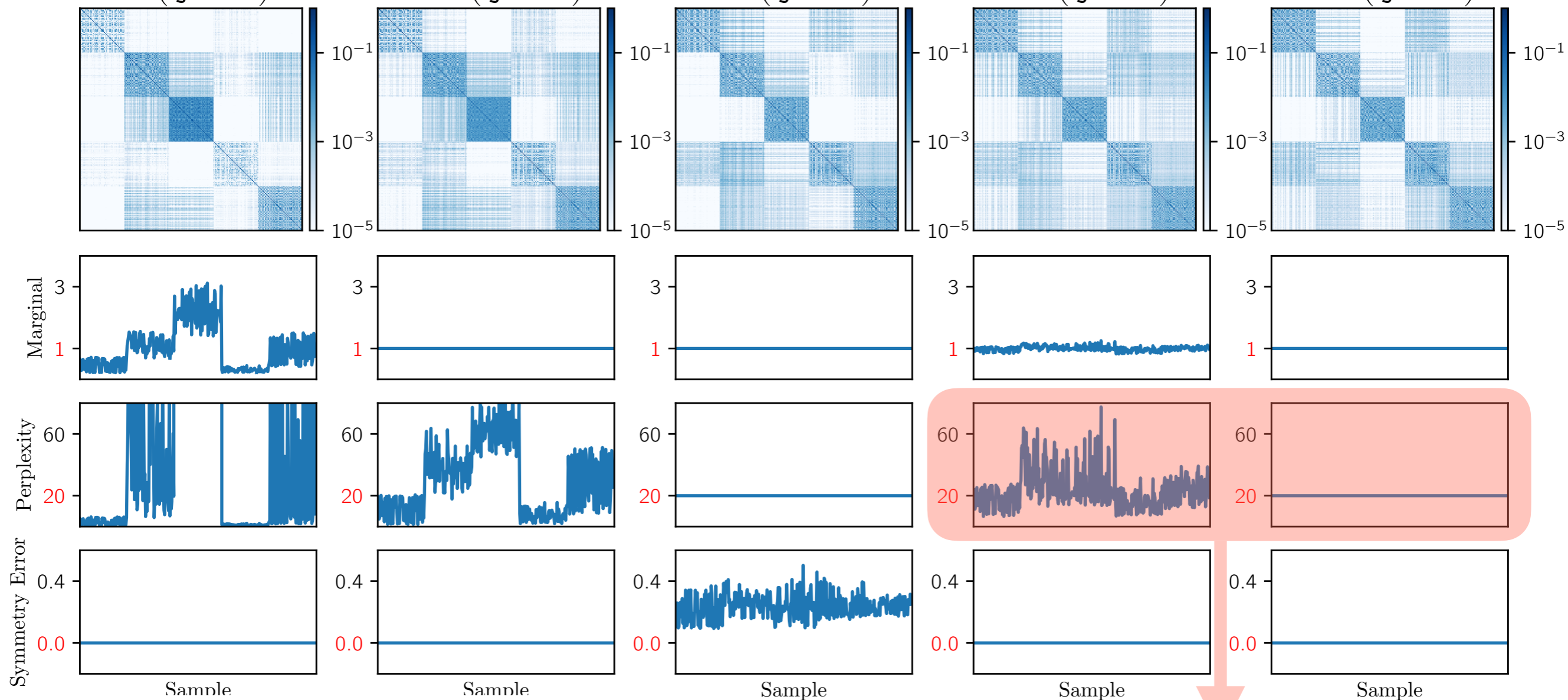
\mathbf{K} ($\bar{\xi}=20$)**

\mathbf{P}^{ds} ($\bar{\xi}=20$)

\mathbf{P}^e ($\xi=20$)

$\bar{\mathbf{P}}^e$ ($\xi=20$)

\mathbf{P}^{se} ($\xi=20$)



* On 5 classes of the COIL Dataset [Nene et al., 1996]

** $\bar{\xi}$ is average perplexity \rightarrow same global entropy as with $\xi = \bar{\xi}$.

Effective control over entropies.

Affinity Panorama*

OURS

Gibbs kernel

$$\mathbf{K} = \exp(-\mathbf{C}/\sigma)$$

Doubly-Sto

$$\mathbf{P}^{ds} = \text{Proj}_{\mathcal{DS}}^{\text{KL}}(\mathbf{K})$$

Entropic

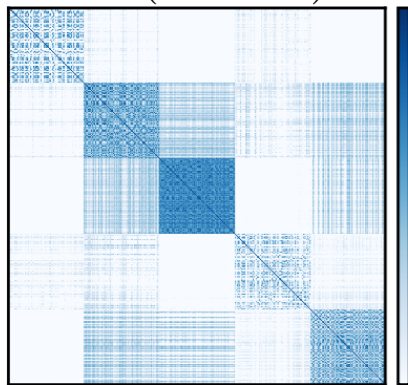
$$\mathbf{P}^e = \text{Proj}_{\mathcal{H}_\xi}^{\text{KL}}(\mathbf{K})$$

Entropic (ℓ_2 Sym) Sym-Entropic

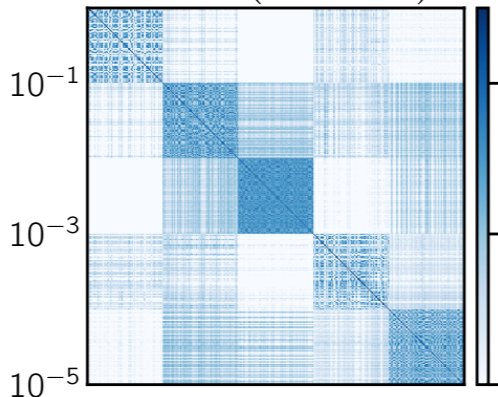
$$\bar{\mathbf{P}}^e = \text{Proj}_{\mathcal{S}}^{\ell_2}(\mathbf{P}^e)$$

$$\mathbf{P}^{se} = \text{Proj}_{\mathcal{H}_\xi \cap \mathcal{S}}^{\text{KL}}(\mathbf{K})$$

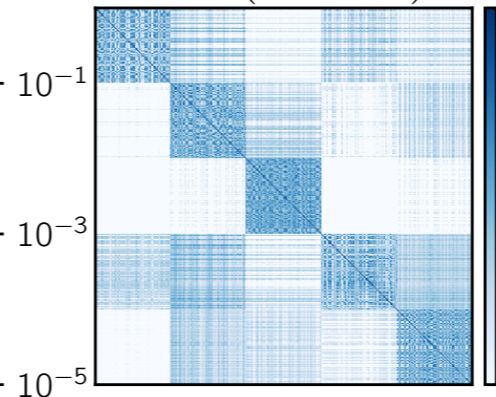
$\mathbf{K} (\bar{\xi}=20)^{**}$



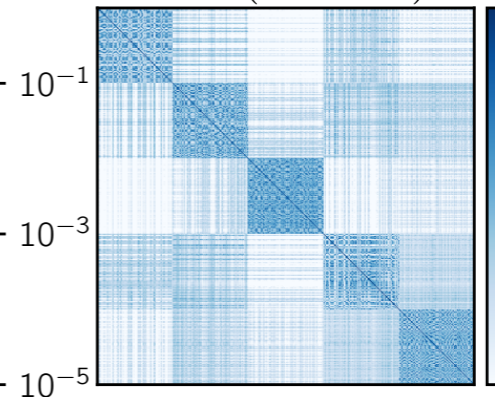
$\mathbf{P}^{ds} (\bar{\xi}=20)$



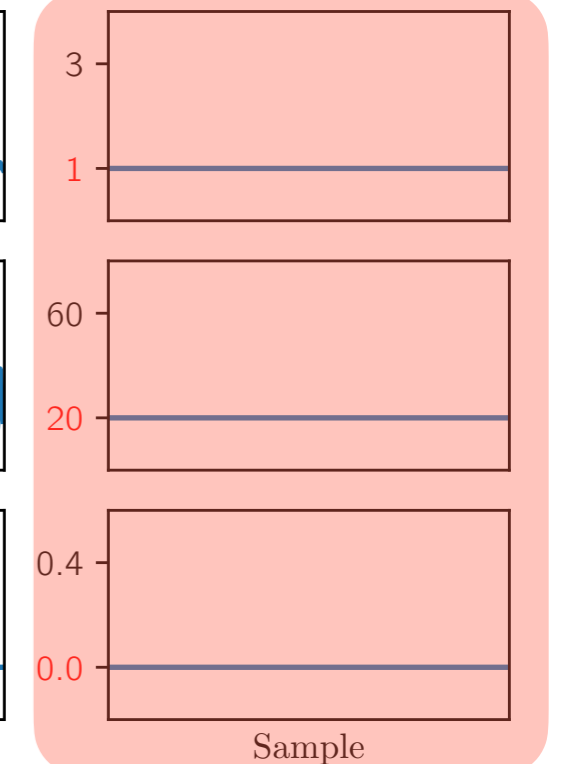
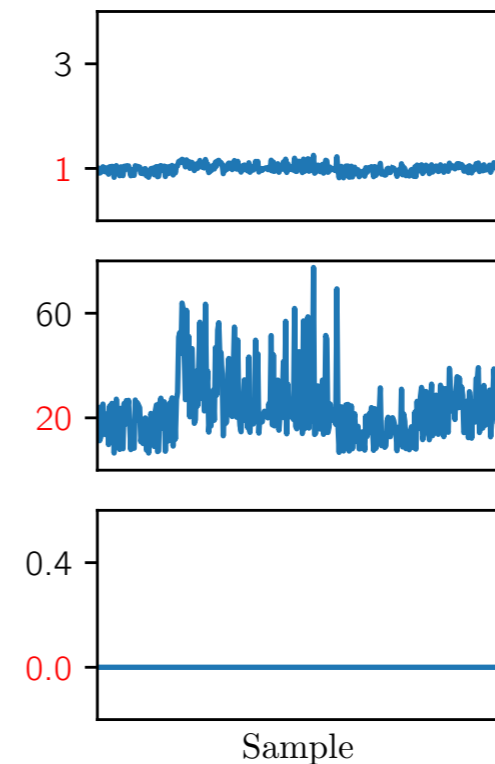
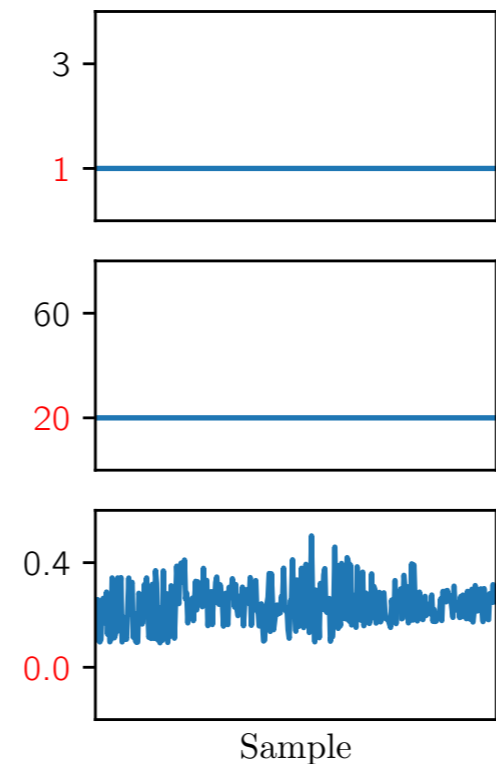
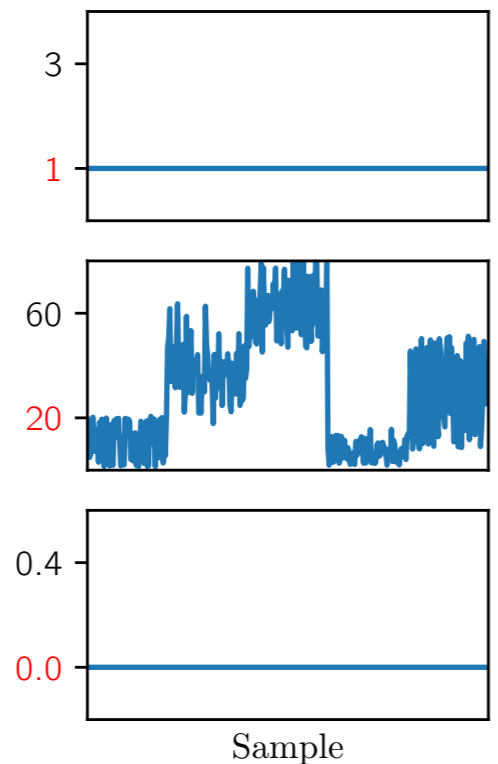
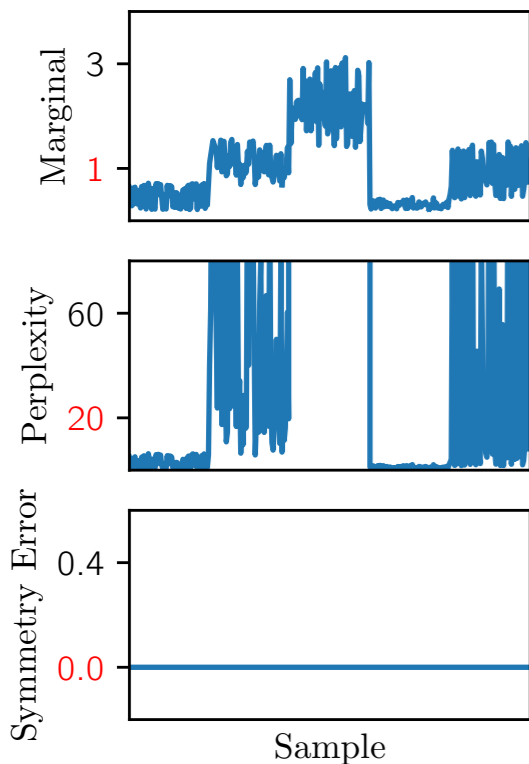
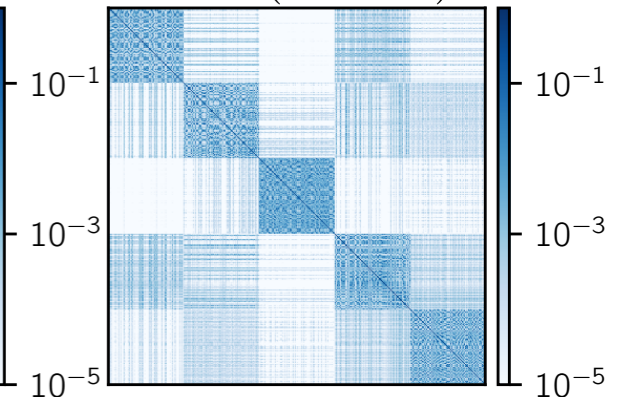
$\mathbf{P}^e (\xi=20)$



$\bar{\mathbf{P}}^e (\xi=20)$



$\mathbf{P}^{se} (\xi=20)$



➔ **Controls ℓ_1 norm, entropy and symmetry at the same time.**

Visualization on toy example

Doubly-Stochastic

$$\min_{\mathbf{P} \geq 0, \mathbf{P}\mathbf{1}=\mathbf{1}, \mathbf{P} \in \mathcal{S}} \langle \mathbf{P}, \mathbf{C} \rangle$$

$$\sum_i H(\mathbf{P}_{i:}) \geq n(\log \xi + 1)$$

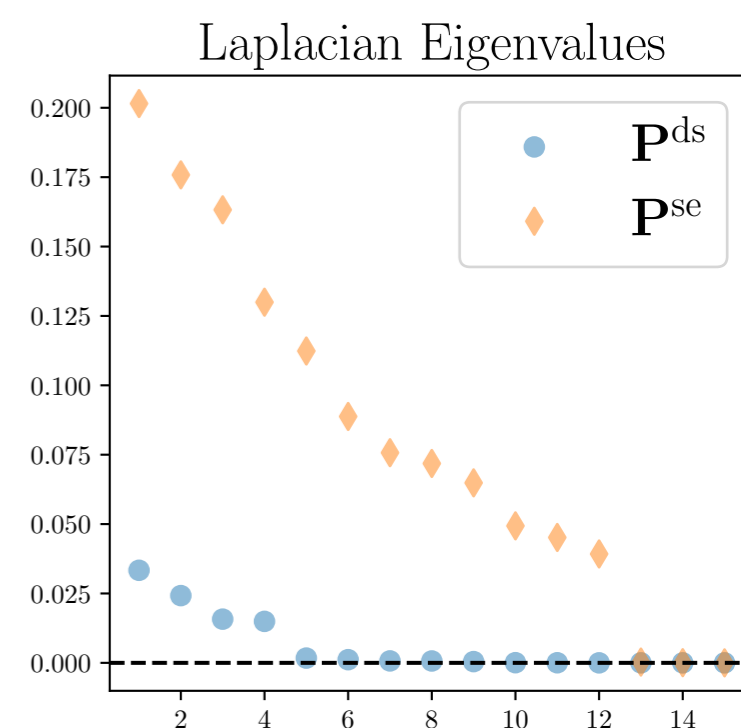
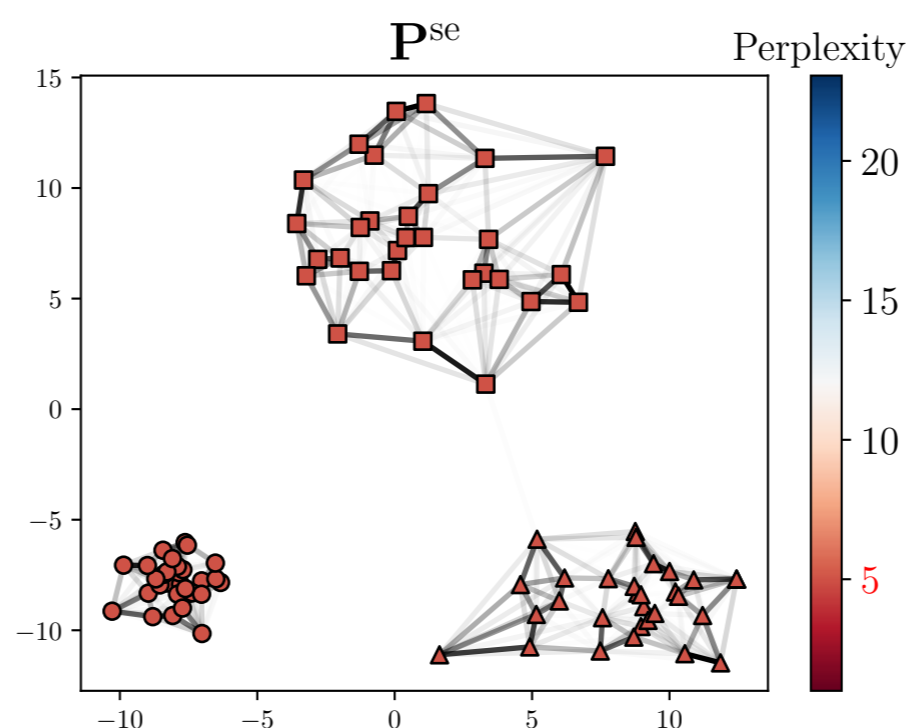
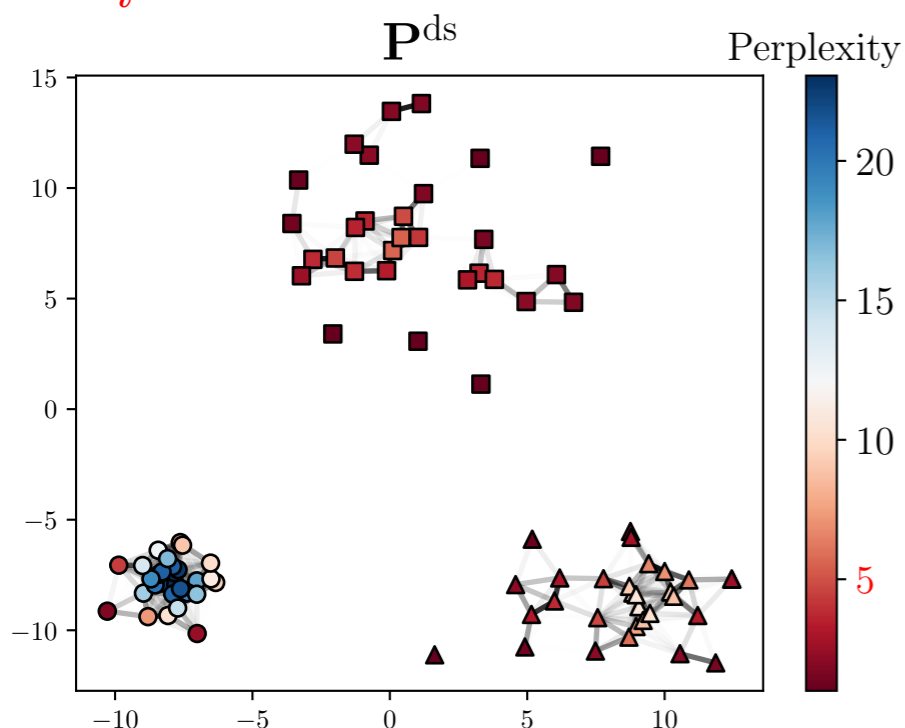
Symmetric-Entropic (OURS)

$$\min_{\mathbf{P} \geq 0, \mathbf{P}\mathbf{1}=\mathbf{1}, \mathbf{P} \in \mathcal{S}} \langle \mathbf{P}, \mathbf{C} \rangle$$

$$\forall i, H(\mathbf{P}_{i:}) \geq \log \xi + 1$$

Symmetric OT

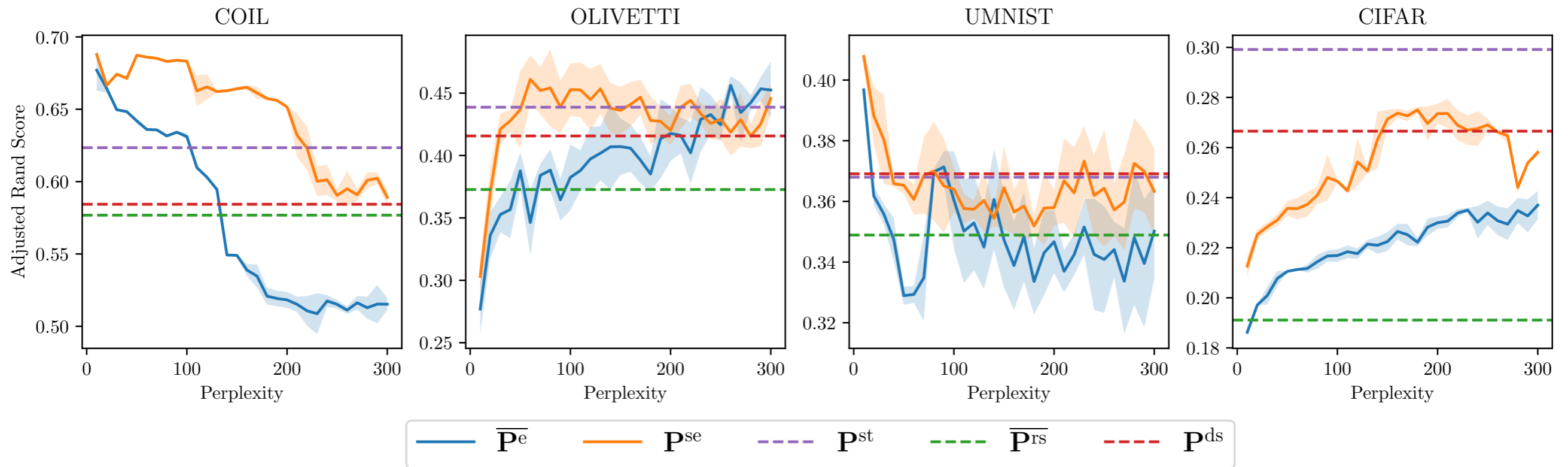
Global vs Pointwise



- Perplexity is set to $\xi = 5$ (average for \mathbf{P}^{ds}).
- \mathbf{P}^{se} can adapt to the varying noise levels.
- \mathbf{P}^{ds} retrieves many unwanted clusters.

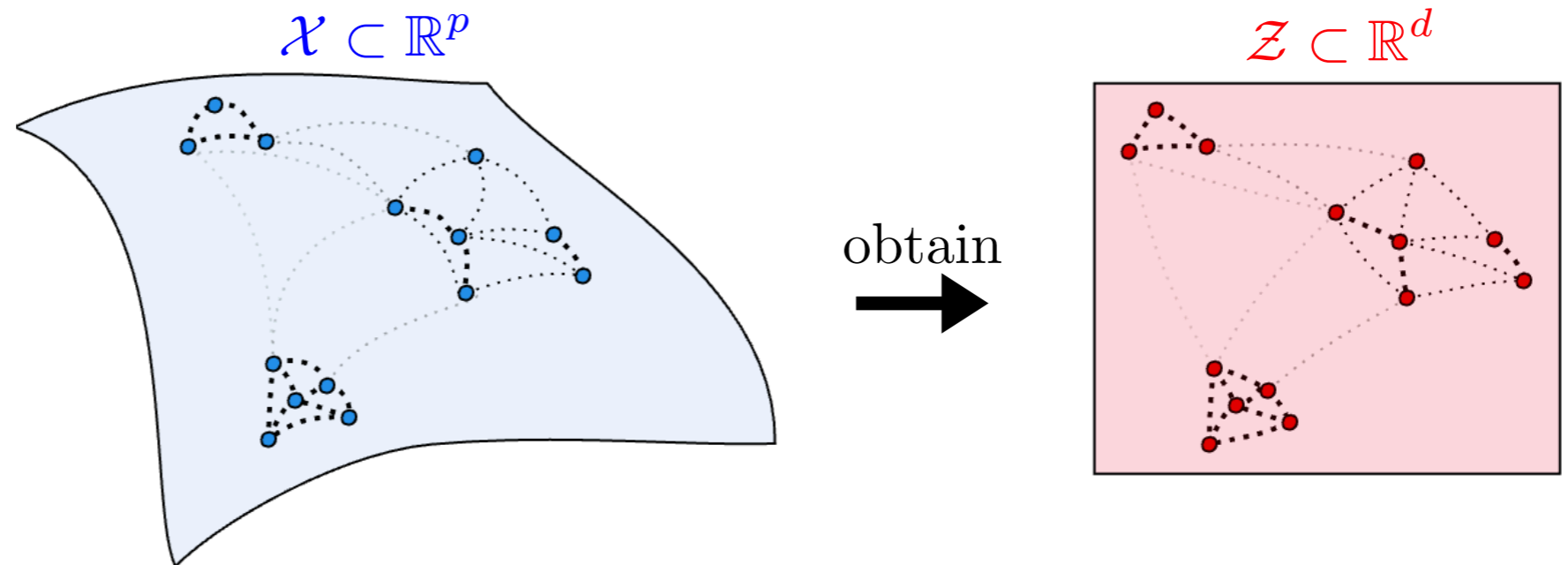
null eigenvalues
=
clusters

Spectral Clustering Results

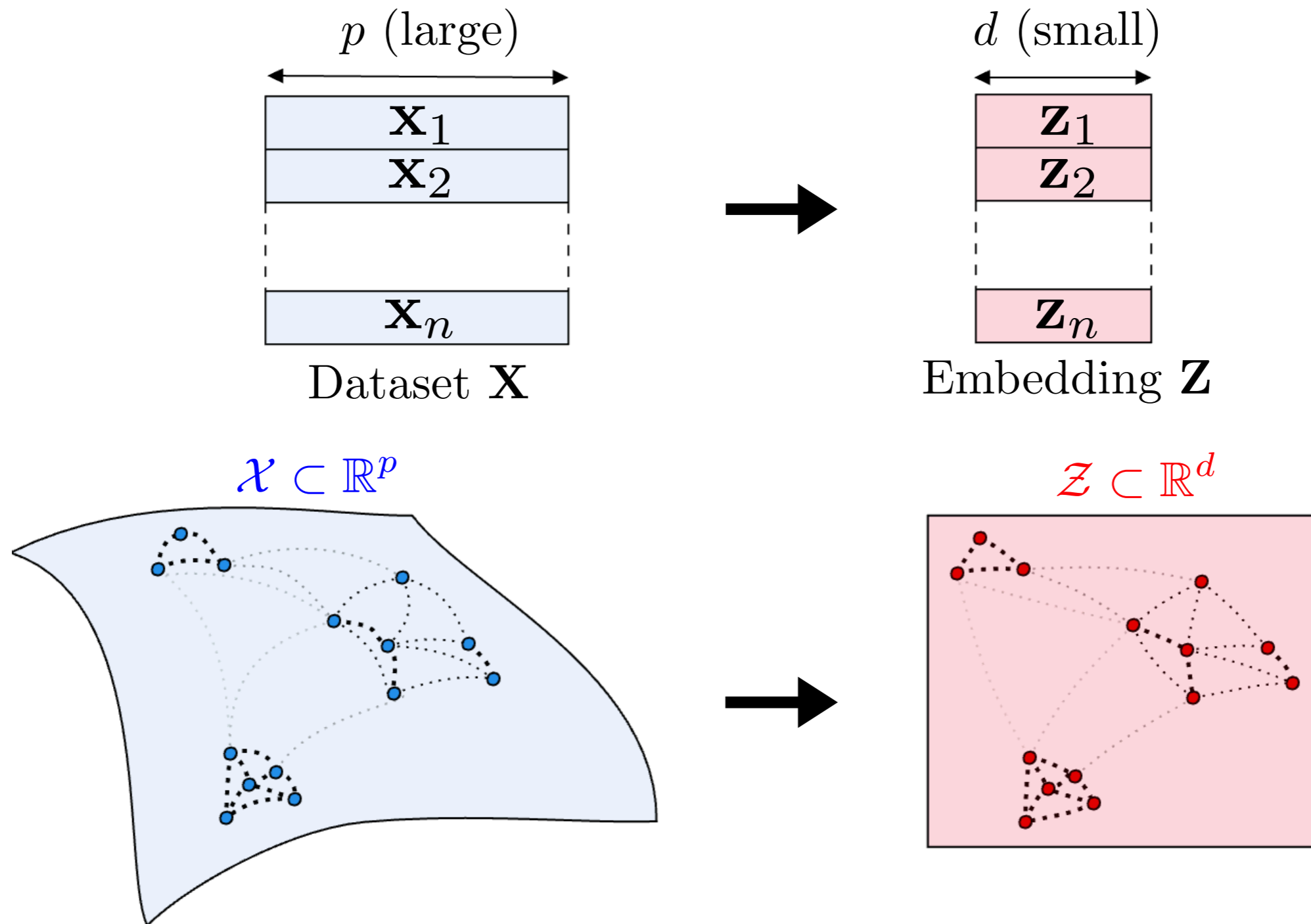


- $\overline{\mathbf{P}}^{rs}$ is the ℓ_2 symmetrized row-stochastic Gaussian kernel with constant bandwidth.
- \mathbf{P}^{st} is the self-tuning affinity. [\[Zelnik-Manor et al., 2004\]](#)
- \mathbf{P}^{se} consistently outperforms other affinities (except \mathbf{P}^{st} on CIFAR).

Part II: Application to Dimensionality Reduction

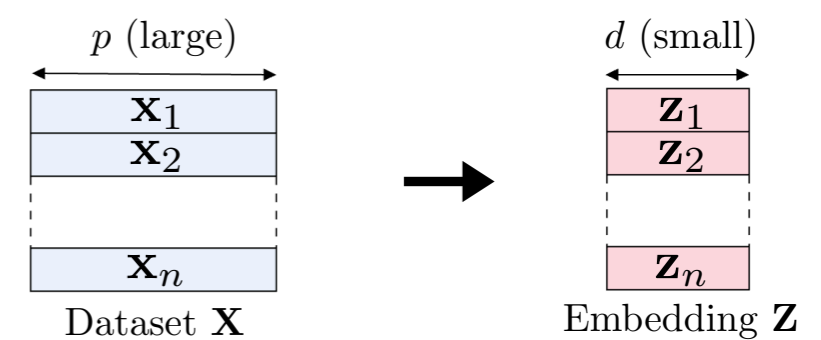


Dimensionality Reduction



The goal is to capture the geometry in \mathcal{X} and reproduce it in \mathcal{Z} .

SNE & SNEkhorn



Cost matrix between embeddings: $[\mathbf{C}_Z]_{ij} = \|\mathbf{Z}_{i:} - \mathbf{Z}_{j:}\|_2^2$.

Stochastic Neighbor Embedding (SNE)

$$\min_{\mathbf{Z} \in \mathbb{R}^{n \times d}} \text{KL}(\overline{\mathbf{P}}^e | \tilde{\mathbf{Q}}_Z)$$

where $[\tilde{\mathbf{Q}}_Z]_{ij} = \exp(-[\mathbf{C}_Z]_{ij}) / \sum_{\ell, t} \exp(-[\mathbf{C}_Z]_{\ell t})$.

SNAREseq Single Cell data

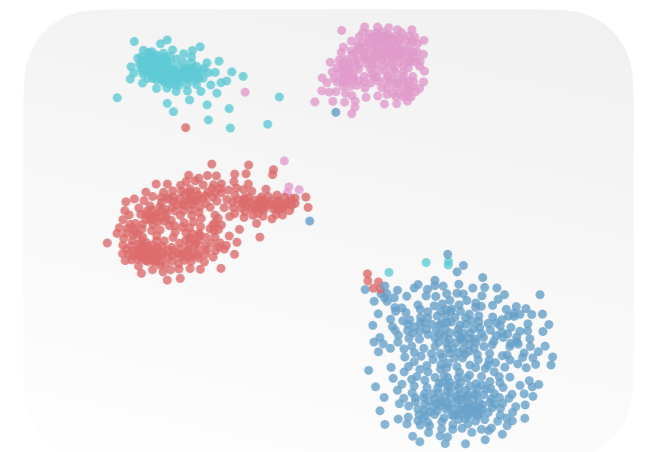


SNEkhorn

OURS

$$\min_{\mathbf{Z} \in \mathbb{R}^{n \times d}} \text{KL}(\mathbf{P}^{\text{se}} | \mathbf{Q}_Z^{\text{ds}})$$

where $\mathbf{Q}_Z^{\text{ds}} = \exp(\mathbf{f}_Z \oplus \mathbf{f}_Z - \mathbf{C}_Z)$ is the DS affinity.

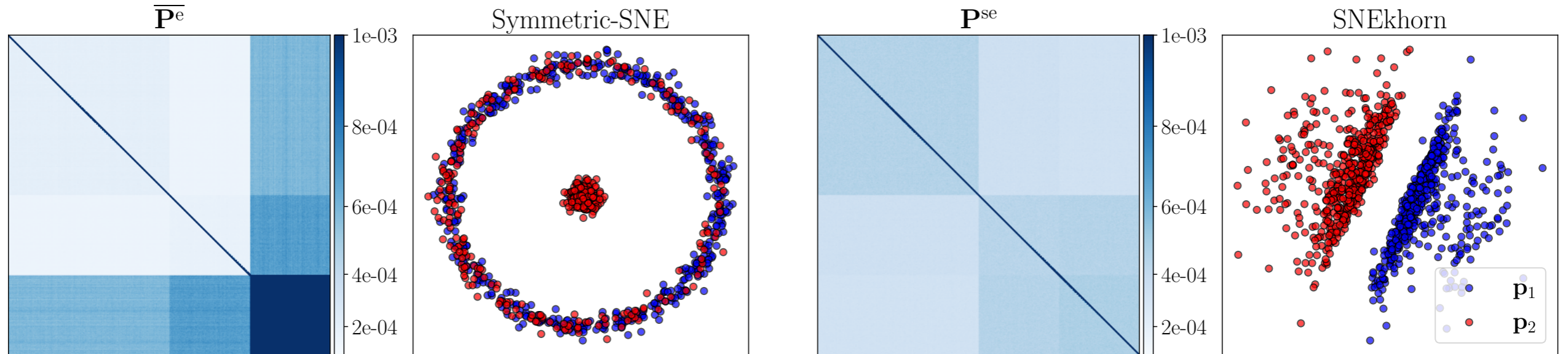


Extension to t-SNE / t-SNEkhorn with heavy-tailed kernels:

[Van der Maaten and Hinton, 2008]

$$[\mathbf{C}_Z]_{ij} = \left(\log(1 + \|\mathbf{Z}_{i:} - \mathbf{Z}_{j:}\|_2^2) \right)_{ij}$$

Toy example: varying noise levels



\mathbf{p}_1 and \mathbf{p}_2 taken in the 10^4 -dimensional probability simplex.

$$\mathbf{x}_i = \tilde{\mathbf{x}}_i / (\sum_j \tilde{x}_{ij}), \quad \tilde{\mathbf{x}}_i \sim \begin{cases} \mathcal{M}(1000, \mathbf{p}_1), & 1 \leq i \leq 500 \\ \mathcal{M}(1000, \mathbf{p}_2), & 501 \leq i \leq 750 \\ \mathcal{M}(2000, \mathbf{p}_2), & 751 \leq i \leq 1000. \end{cases}$$

SNE is misled by the batch effect unlike SNEkhorn.

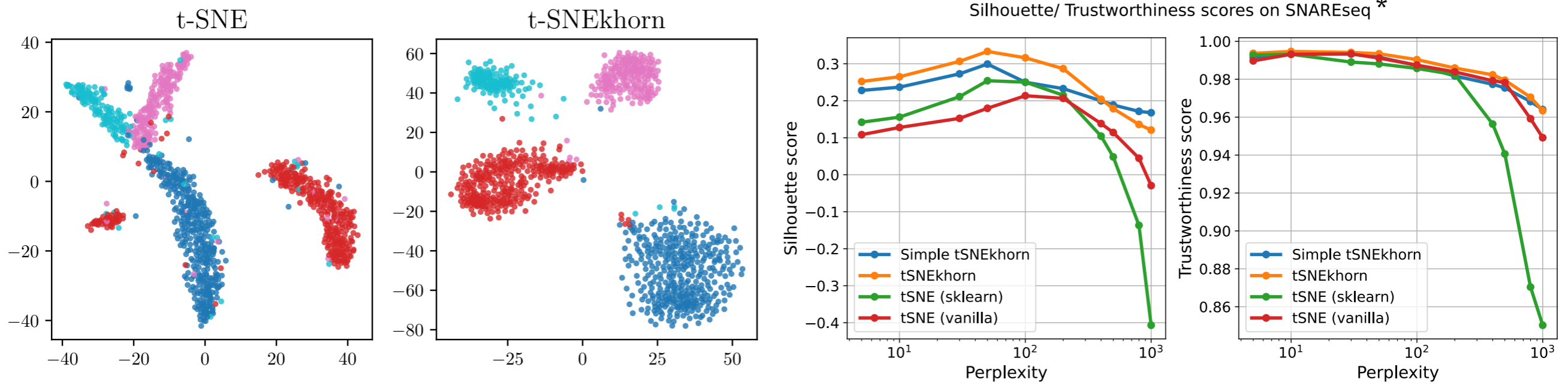
Dimension Reduction Results

	Silhouette ($\times 100$)			Trustworthiness ($\times 100$)		
	UMAP	t-SNE	t-SNEkhorn	UMAP	t-SNE	t-SNEkhorn
COIL	20.4 ± 3.3	30.7 ± 6.9	52.3 ± 1.1	99.6 ± 0.1	99.6 ± 0.1	99.9 ± 0.1
OLIVETTI	6.4 ± 4.2	4.5 ± 3.1	15.7 ± 2.2	96.5 ± 1.3	96.2 ± 0.6	98.0 ± 0.4
UMNIST	-1.4 ± 2.7	-0.2 ± 1.5	25.4 ± 4.9	93.0 ± 0.4	99.6 ± 0.2	99.8 ± 0.1
CIFAR	13.6 ± 2.4	18.3 ± 0.8	31.5 ± 1.3	90.2 ± 0.8	90.1 ± 0.4	92.4 ± 0.3
Liver (14520)	49.7 ± 1.3	50.9 ± 0.7	61.1 ± 0.3	89.2 ± 0.7	90.4 ± 0.4	92.3 ± 0.3
Breast (70947)	28.6 ± 0.8	29.0 ± 0.2	31.2 ± 0.2	90.9 ± 0.5	91.3 ± 0.3	93.2 ± 0.4
Leukemia (28497)	22.3 ± 0.7	20.6 ± 0.7	26.2 ± 2.3	90.4 ± 1.1	92.3 ± 0.8	94.3 ± 0.5
Colorectal (44076)	67.6 ± 2.2	69.5 ± 0.5	74.8 ± 0.4	93.2 ± 0.7	93.7 ± 0.5	94.3 ± 0.6
Liver (76427)	39.4 ± 4.3	38.3 ± 0.9	51.2 ± 2.5	85.9 ± 0.4	89.4 ± 1.0	92.0 ± 1.0
Breast (45827)	35.4 ± 3.3	39.5 ± 1.9	44.4 ± 0.5	93.2 ± 0.4	94.3 ± 0.2	94.7 ± 0.3
Colorectal (21510)	38.0 ± 1.3	42.3 ± 0.6	35.1 ± 2.1	85.6 ± 0.7	88.3 ± 0.9	88.2 ± 0.7
Renal (53757)	44.4 ± 1.5	45.9 ± 0.3	47.8 ± 0.1	93.9 ± 0.2	94.6 ± 0.2	94.0 ± 0.2
Prostate (6919)	5.4 ± 2.7	8.1 ± 0.2	9.1 ± 0.1	77.6 ± 1.8	80.6 ± 0.2	73.1 ± 0.5
Throat (42743)	26.7 ± 2.4	28.0 ± 0.3	32.3 ± 0.1	91.5 ± 1.3	88.6 ± 0.8	86.8 ± 1.0
scGEM	26.9 ± 3.7	33.0 ± 1.1	39.3 ± 0.7	95.0 ± 1.3	96.2 ± 0.6	96.8 ± 0.3
SNAREseq	6.8 ± 6.0	35.8 ± 5.2	67.9 ± 1.2	93.1 ± 2.8	99.1 ± 0.1	99.2 ± 0.1

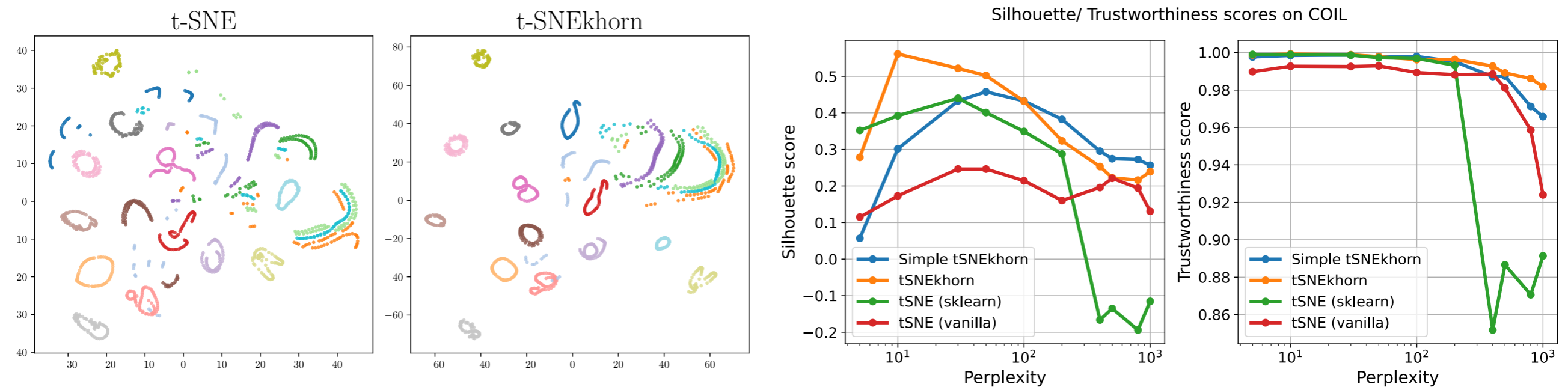
t-SNEkhorn outperforms **t-SNE** and **UMAP** on various real-world datasets.

Dimension Reduction Results

SNAREseq Single Cell data



COIL-20 Image data



* Simple t-SNEkhorn \rightarrow t-SNEkhorn with same affinity as t-SNE for the embeddings \mathbf{Z} (not doubly stochastic).

Conclusion

We provide a new **symmetric** affinity matrix, **controlling** for each row/column:

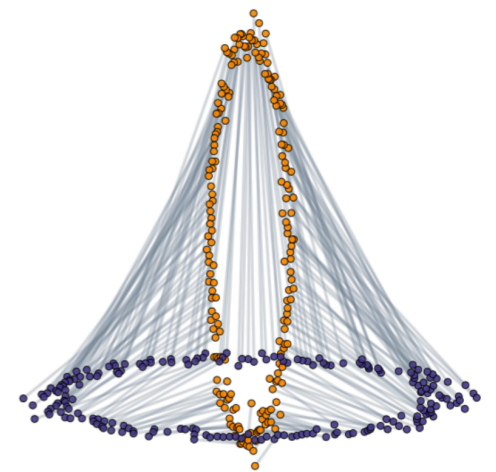
ℓ_1 norm & Shannon entropy.

We show its **robustness to heteroscedastic noise** (crucial for single cell data).

Based on this affinity, we propose a **new DR method : SNEkhorn**.

Python code available at :

<https://github.com/PythonOT/SNEkhorn>

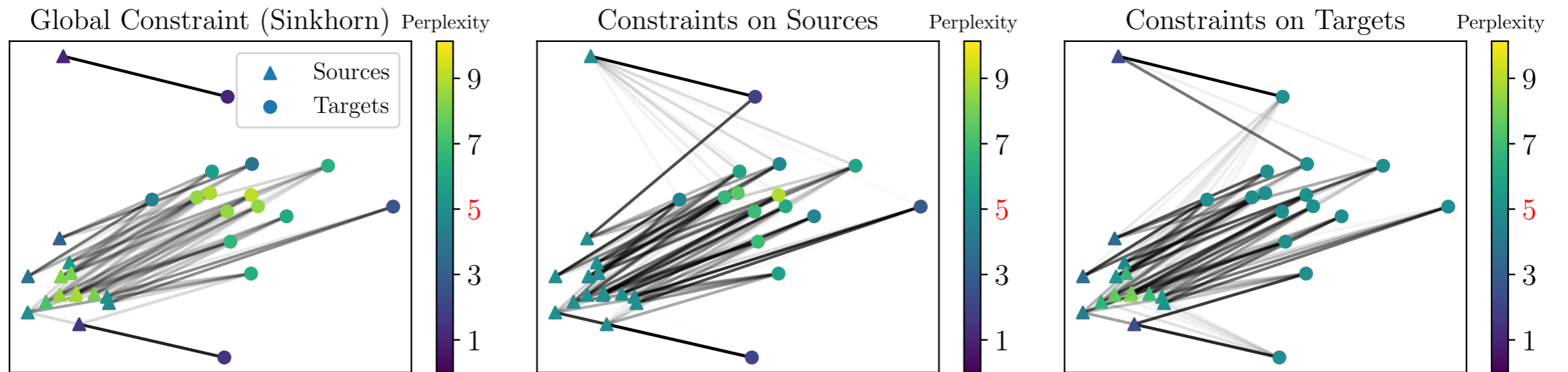


**SNEkhorn : Dimension Reduction with
Symmetric Entropic Affinities
NeurIPS 2023**



Part III: Future Works : OT with Adaptive Regularisation

OT with Adaptive Regularisation



Application to Domain Adaptation

