

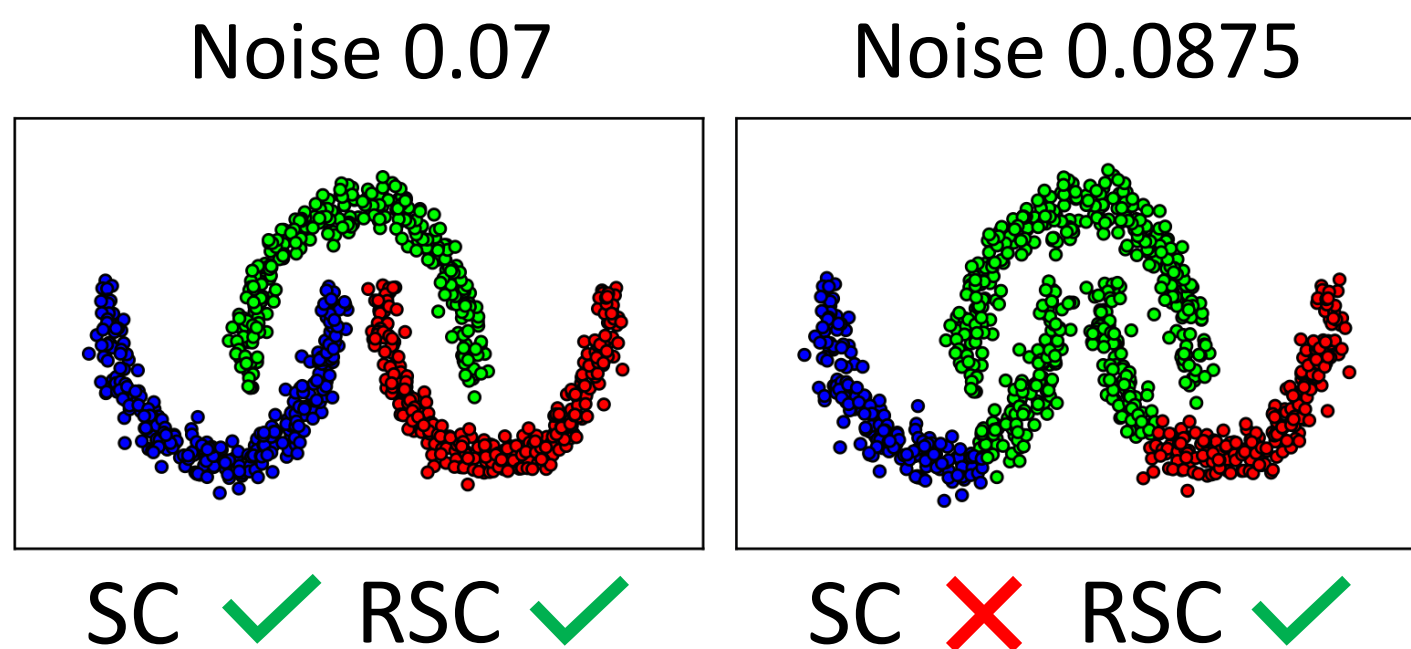
Robust Spectral Clustering for Noisy Data

Modeling Sparse Corruptions Improves Latent Embeddings

Aleksandar Bojchevski, Yves Matkovic, Stephan Günnemann

MOTIVATION

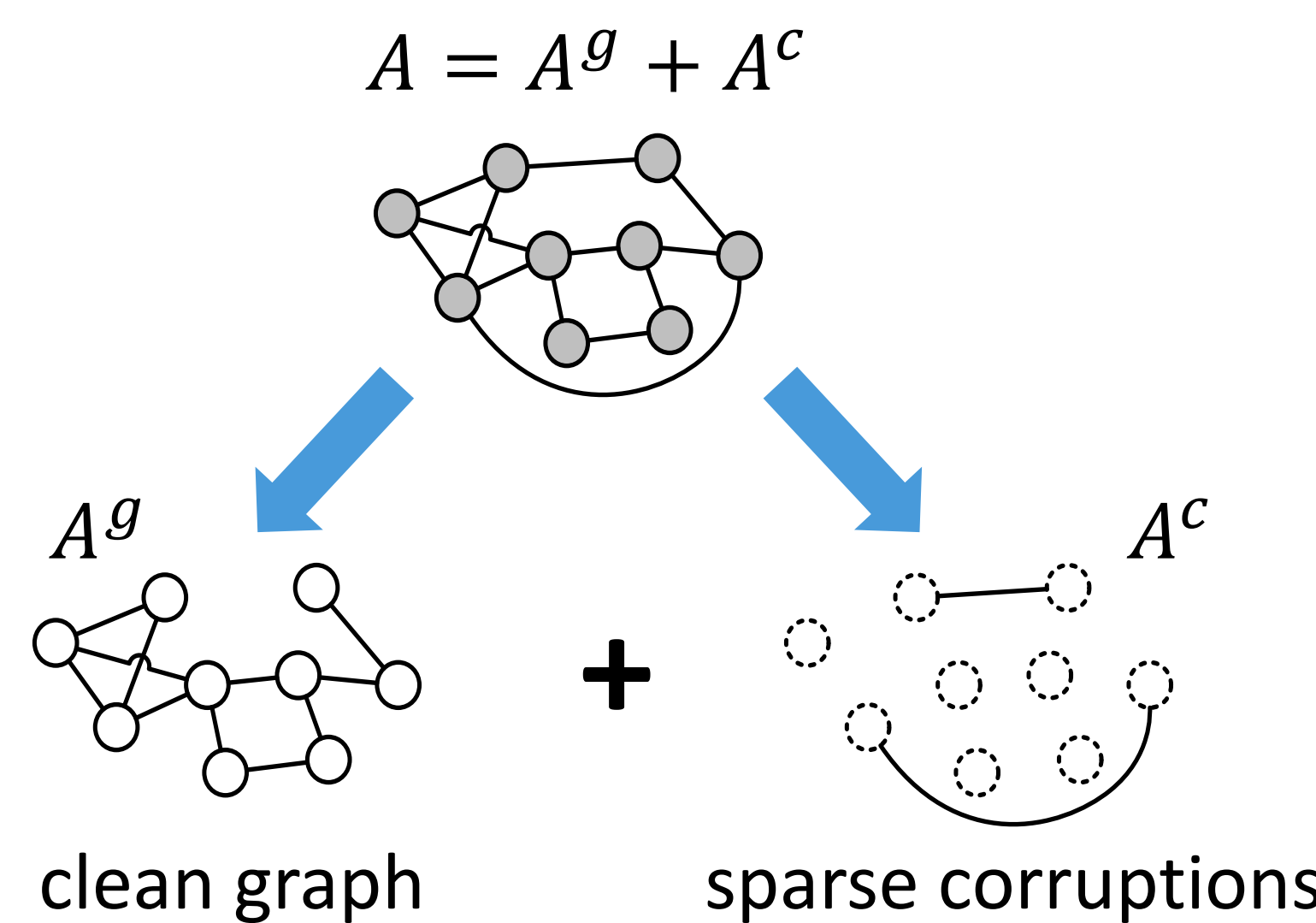
- Spectral clustering (SC) **widely** used, but highly **sensitive** to noisy data
- Noise distorts the embedding space and obfuscates the clustering structure



- We propose a robust version: **RSC**

PROBLEM FORMULATION

- Core idea: **Latent Decomposition**



- Jointly learn decomposition and embedding
- Decomposition steered by the underlying clustering

$$A^*, H^* = \underset{\substack{H \in \mathbb{R}^{n \times d}, \\ A^g \in (\mathbb{R}_{\geq 0})^{n \times n}}}{\operatorname{argmin}} \operatorname{Tr}(H^T \cdot L(A^g) \cdot H) \quad (1)$$

subject to: $H^T \cdot D(A^g) \cdot H = I$ and

$$A = A^g + A^c, \|A^c\|_0 \leq 2\theta, \|a_i^c\|_0 \leq \omega_i$$

- Robust formulation for all SC versions
- Result \rightarrow improved embedding

ALGORITHMIC SOLUTION

Alternating Optimization

- Update H , Given $A^g/A^c \rightarrow$ **Easy**
 - Trace minimization problem
 - Solution for H are the k first generalized eigenvectors of $L(A^g)$
- Update A^g/A^c , Given $H \rightarrow$ **(NP) Hard**
 - Express eigenvalues of A_{new}^g in closed form
 - A_{new}^g that minimizes (1) equivalent to maximizing:

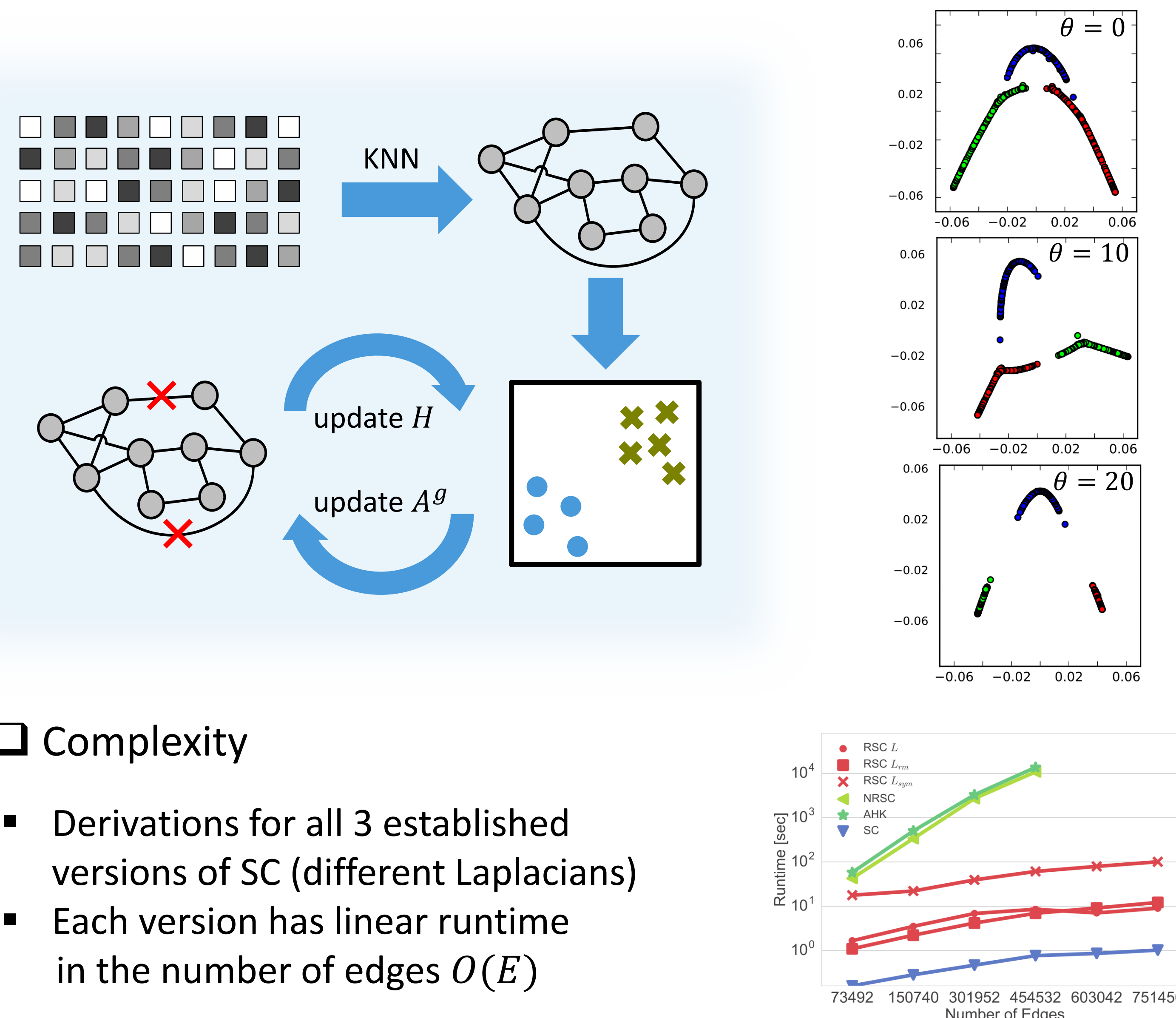
$$f([a_{uv}^c]_{u,v \in E}) = \sum_{u,v \in E} a_{uv}^c \left(\frac{\|h_u - h_v\|_2^2}{\text{nodes far away in the embedding space}} - \frac{\|\sqrt{\lambda} \circ h_u\|_2 - \|\sqrt{\lambda} \circ h_v\|_2}{\text{prefers edges close to critical region}} \right)$$

subject to $\|\cdot\|_0$ constraints

- Observation: Above problem is equivalent to **Multidimensional Knapsack** problem
- Greedy approximation scheme
- Best possible approximation ratio of $1/\sqrt{N+1}$

- Complexity

- Derivations for all 3 established versions of SC (different Laplacians)
- Each version has linear runtime in the number of edges $O(E)$

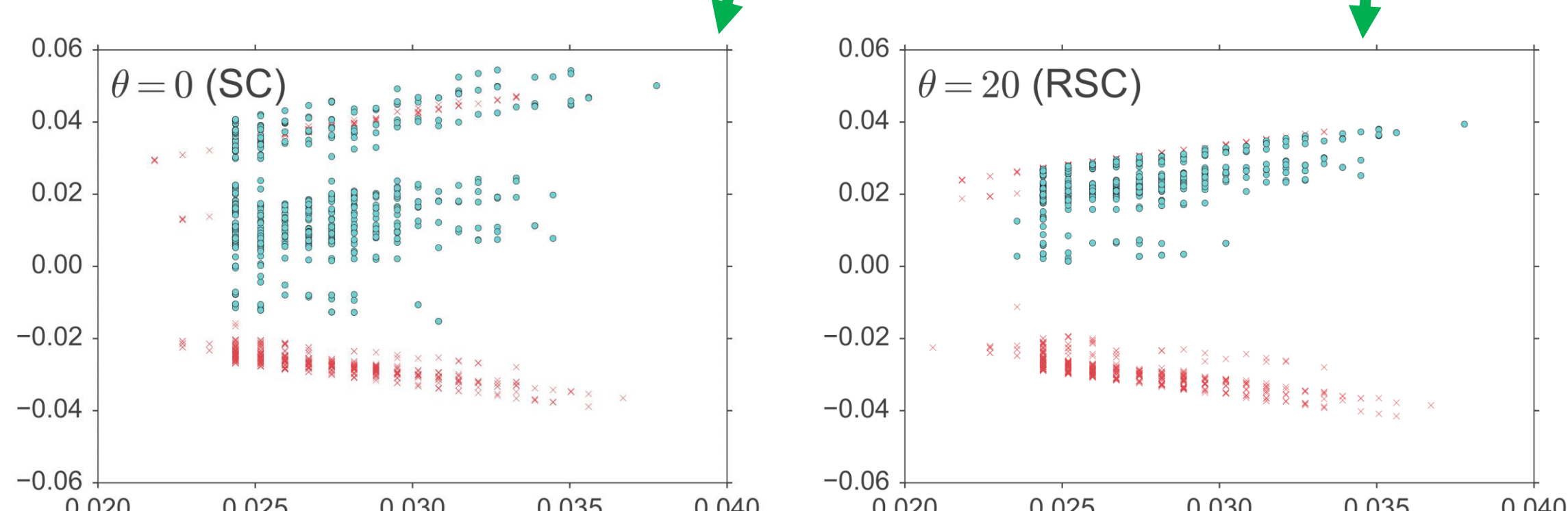


How can we evaluate the quality of the **clustering**?

- RSC improves the clustering as measured by the NMI

Dataset	AHK	NRSC	SC		RSC
Moons	0.53	0.99	0.47	+112.8 %	1.00
Banknote	0.53	0.47	0.46	+ 32.6 %	0.61
USPS	0.77	0.83	0.78	+ 8.9 %	0.85
MNIST	0.71	0.76	0.70	+ 11.4 %	0.78
Pendigits	0.94	0.94	0.93	+ 3.2 %	0.96

- RSC improves discrimination in the embedding space



RESULTS

How can we evaluate the quality of the **embeddings**?

- Global Separation

- Degree of separability between clusters
- Robust silhouette coefficient

$$P_{c,c'}(x) = \underset{x \% \text{ smallest}}{\operatorname{average}} [dist(h_i, h_j)]_{i \in C_c, j \in C_{c'}}$$

$$GS_c(x) = \frac{P_{c,c'}(x) - P_{c,c}(x)}{\max\{P_{c,c'}(x), P_{c,c}(x)\}}$$

$$c' = \underset{c \neq c'}{\operatorname{argmin}} P_{c,c'}(x)$$

- Local Purity

- Homogeneous local neighborhood

$$occ_x(c, i) = |\{j \in NN_x(i) \cup \{i\} \mid c_j = c\}|$$

$$pur_x(i) = \frac{1}{x+1} \max_{c \in C} occ_x(c, i)$$

$$PUR(x) = \frac{1}{N} \sum_{i=1}^N pur_x(i)$$

