# Selective Sampling-based Scalable Sparse Subspace Clustering

**Shin Matsushima** — 東京大学 THE UNIVERSITY OF TOKYO

**Maria Brbić** — Stanford University

paper  code  poster
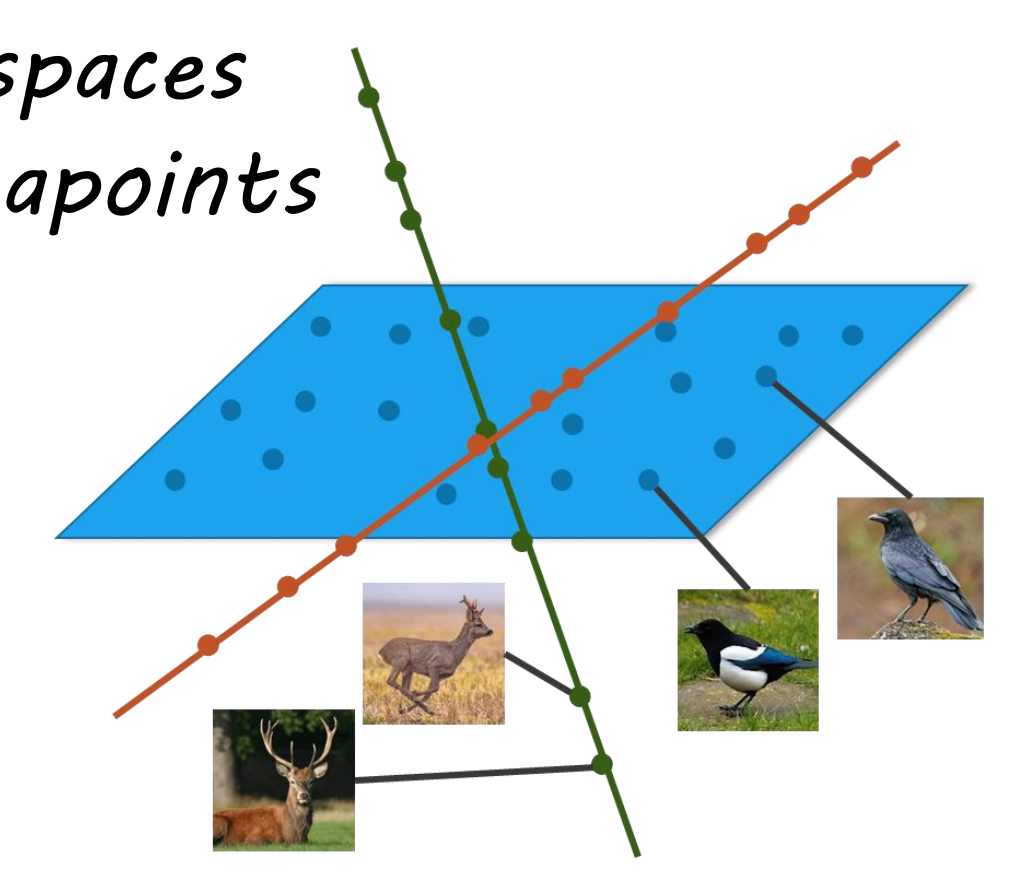
## OVERVIEW

**Sparse Subspace Clustering (SSC)**
- ✓ High performance clustering for high dimensional data with strong theoretical guarantees
- ☹ Quadratic complexity w.r.t. number of data points

▼

**Selective Sampling-based Scalable Sparse Subspace Clustering (S⁵C)**
- ✓ Theoretical Guarantee → Theoretical Scalability
- ✓ Low computational cost → Experimental Scalability
- ✓ Good clustering performance

## SUBSPACE CLUSTERING

$L$ subspaces
$N$ datapoints

**Assumption:** high-dimensional data points lie in the union of low-dimensional subspaces.

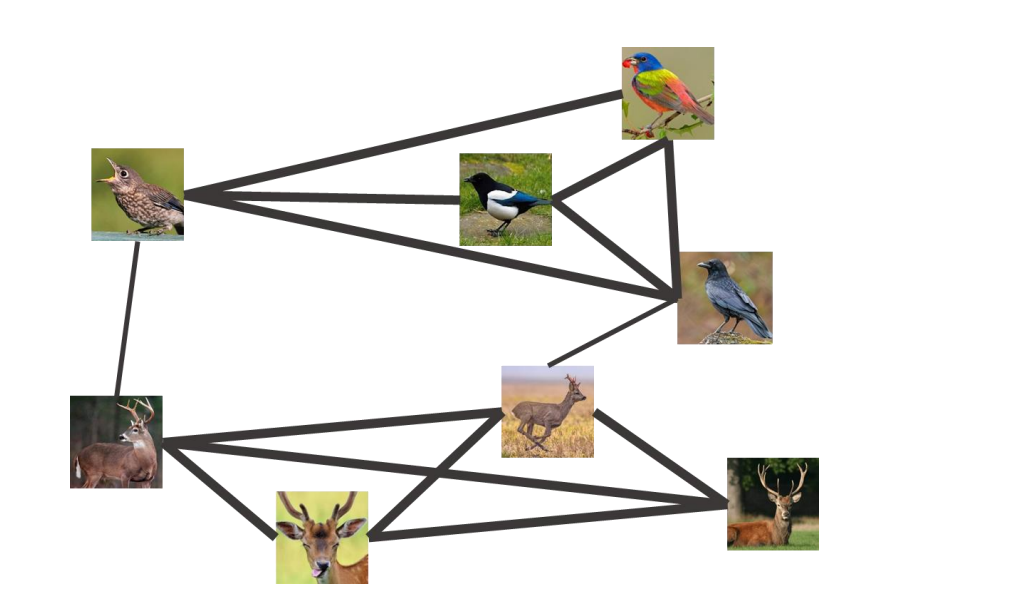**Goal:** identify subspaces and assign data points to the subspaces

**Algorithm:**

$-0.1 \cdot + 0.3 \cdot + 0.8 \cdot \doteq$

*representation of*

1. Representation learning
2. Derive an affinity graph
3. Spectral clustering

**Applications:**
Clustering images
Clustering documents...

## KEY IDEA

Represent data point as a linear combination of a **small number of subsamples** {🐦,🐦,🐦}

$$\min_{\mathbf{c}} \left\| \left( c_1 \cdot + c_2 \cdot + c_3 \cdot + c_5 \cdot + c_6 \cdot + c_7 \cdot + c_8 \cdot \right) - \cdot \right\|^2$$

+ L1 regularizer on $\mathbf{c}$

To generate subsamples, perform **selective sampling** so that all data points are represented well

## S⁵C ALGORITHM

1. Randomly sample a data point 🐦
2. Solve $\min_{\mathbf{c}} \left\| \left( c_2 \cdot + c_5 \cdot + c_7 \cdot \right) - \cdot \right\|^2$ among current collection of subsamples {🐦,🐦,🐦}
3. Find **selective sample** 🐦 that helps representation of 🐦 best *According to the absolute value of the gradient*
4. Add 🐦 to current collection of subsamples {🐦,🐦,🐦}
5. Repeat 1.-4. **T** times starting from { }
6. Solve $\min_{\mathbf{c}} \left\| \left( c_2 \cdot + c_3 \cdot + c_5 \cdot + c_7 \cdot \right) - \cdot \right\|^2$ w.r.t. final subsamples for all data points

→ O(NT) algorithm!

**Iteratively select a subsample which seems to improve the current representation!**

1. Initialize N-by-L matrix **V** as a random orthogonal matrix
2. $\mathbf{V} \leftarrow (2\mathbf{I} - \mathbf{L})\mathbf{V}$
3. Orthogonalize **V** by QR decomposition
4. Repeat 2.-3. until convergence
5. Apply K-means for **V**

**Graph Laplacian L has O(N) nonzero elements → O(N) algorithm!**

## RELATED WORK

| | SSC | SSC-OMP | SSC-ORGEN | SSSC | **S⁵C** |
|---|---|---|---|---|---|
| Theoretical Scalability | ✗ | ✗ | ✗ | ✗ | ✓ |
| Experimental Scalability | ✗ | ✓ | ✓ | ✓ | ✓ |

## THEORETICAL SCALABILITY

**S⁵C Algorithm performs *perfect clustering* in *O(N) running time* in high probability under *some generative model***

*Required number of subsamples ($\doteq$T) is O(1) w.r.t. N*

**Semi Random model**

Data points are uniformly randomly generated from unit ball in each subspace

*Subspace detection property (SDP)*

1. $c_{i'}$ for i-th data is nonzero only when i and i' share the same subspace
2. **c** is not all zeros

$-0.1 \cdot + 0.3 \cdot + 0.8 \cdot \doteq$

## CLUSTERING PERFORMANCE

Clustering error (%)

| | Nystrom | AKK | SSC | SSC-OMP | SSC-ORGEN | SSSC | **S⁵C** |
|---|---|---|---|---|---|---|---|
| Yale B | 76.8 | 85.7 | 33.8 | 35.9 | 37.4 | 59.6 | 39.3 |
| Hopkins 155 | 21.8 | 20.6 | 4.1 | 23.0 | 20.5 | 21.1. | 14.6 |
| COIL-100 | 54.5 | 53.1 | 42.5 | 57.9 | 89.7 | 67.8 | 45.9 |
| Letter-rec | 73.3 | 71.7 | / | 95.2 | 68.6 | 68.4 | 67.7 |
| CIFAR-10 | 76.6 | 75.6 | / | | 82.4 | 82.4 | 75.1 |
| MNIST | 45.7 | 44.6 | / | / | 28.7 | 48.7 | 40.4 |
| Devanagari | 73.5 | 72.8 | / | / | 58.6 | 84.9 | 67.2 |

## SPARSE SUBSPACE CLUSTERING

**Solve for all data points → O(N²) algorithm! → not scalable** ☹

$$\min_{\mathbf{c}} \left\| \left( c_1 \cdot + c_2 \cdot + c_3 \cdot + c_5 \cdot + c_6 \cdot + c_7 \cdot + c_8 \cdot \right) - \cdot \right\|^2$$

+ L1 regularizer on **c**

*Represent data point as a linear combination of all data points*

⇒ $-0.1 \cdot +0.3 \cdot +0.8 \cdot \doteq$

*Only some data points in the same subspace remain*

⇒ *Find eigenvectors corresponding to L smallest eigenvalues of graph Laplacian* **L**

⇒ *Perfect clustering is theoretically proven (SDP property)*

Elhamifar and Vidal. TPAMI (2013)

## EXPERIMENTAL SCALABILITY

quadratic   linear

Time (s)

- S⁵C
- SSSC
- ORGEN
- Nyström
- AKK
- SSC
- OMP

Number of Datapoints