



# Selective Sampling-based Scalable Sparse Subspace Clustering

Shin Matsushima



Maria Brbić



paper

codes



## OVERVIEW

### Sparse Subspace Clustering (SSC)

- ✓ high performance clustering for high dimensional data
- ☹ quadratic complexity w.r.t. number of data points

### Selective Sampling-based Scalable Sparse Subspace Clustering (S<sup>5</sup>C)

- ✓ Theoretical Guarantee → Theoretical Scalability
- ✓ Low computational cost → Experimental Scalability
- ✓ Good clustering performance

## KEY IDEA

Represent data point as linear combination among a small number of subsamples  $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$

$$\min_c \left\| \begin{pmatrix} c_1 \mathbf{x}_1 + c_2 \mathbf{x}_2 + c_3 \mathbf{x}_3 \\ + c_5 \mathbf{x}_5 + c_6 \mathbf{x}_6 + c_7 \mathbf{x}_7 + c_8 \mathbf{x}_8 \end{pmatrix} - \mathbf{x}_4 \right\|^2$$

→ Perform selective sampling for subsamples so that all data points are represented well

## S<sup>5</sup>C ALGORITHM

1. Randomly sample a datapoint  $\mathbf{x}_i$
2. Solve  $\min_c \left\| \begin{pmatrix} c_2 \mathbf{x}_2 + c_5 \mathbf{x}_5 + c_7 \mathbf{x}_7 \end{pmatrix} - \mathbf{x}_i \right\|^2$  among current collection of subsamples  $\{\mathbf{x}_2, \mathbf{x}_5, \mathbf{x}_7\}$
3. Find selective sample  $\mathbf{x}_{i'}$  that helps representation of  $\mathbf{x}_i$
4. Add  $\mathbf{x}_{i'}$  to current collection of subsamples  $\{\mathbf{x}_2, \mathbf{x}_5, \mathbf{x}_7\}$
5. Repeat 1.-4. T times starting from  $\{\}$
6. Solve  $\min_c \left\| \begin{pmatrix} c_2 \mathbf{x}_2 + c_5 \mathbf{x}_5 + c_7 \mathbf{x}_7 \end{pmatrix} - \mathbf{x}_i \right\|^2$  w.r.t. final subsamples for all data points

→ O(N) algorithm!

Iteratively select samples which seems to help better representations!

1. Initialize N-by-L matrix V as a random orthogonal matrix
2.  $V \leftarrow LV$
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	SSSC	S <sup>5</sup> C
Theoretical Scalability	X	X	X	✓
Experimental Scalability	X	✓	✓	✓

## THEORETICAL GUARANTEES

S<sup>5</sup>C Algorithm performs perfect clustering in O(N) running time in high probability under some generative model

Required number of subsample is O(1) w.r.t. N

SDP property

Semi Random model



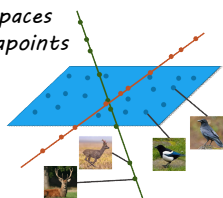
Data points are uniformly randomly generated from unit ball in each subspace

1.  $c_i$  for i-th data is nonzero only when i and i' shares the same subspace
2. c is not all zeros

$$-0.1 \mathbf{x}_1 + 0.3 \mathbf{x}_2 + 0.8 \mathbf{x}_3 \approx \mathbf{x}_4$$

## 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:

1. Representation learning
- 1.5. Derive an affinity graph
2. Spectral clustering

Application:

Clustering images  
Clustering documents...

$$-0.1 \mathbf{x}_1 + 0.3 \mathbf{x}_2 + 0.8 \mathbf{x}_3 \approx \mathbf{x}_4$$

Representation of



## SPARSE SUBSPACE CLUSTERING

Challenge: Quadratic complexity in the number of data points!

$$\min_c \left\| \begin{pmatrix} c_1 \mathbf{x}_1 + c_2 \mathbf{x}_2 + c_3 \mathbf{x}_3 \\ + c_5 \mathbf{x}_5 + c_6 \mathbf{x}_6 + c_7 \mathbf{x}_7 + c_8 \mathbf{x}_8 \end{pmatrix} - \mathbf{x}_4 \right\|^2 + \text{L1 regularizer on } c$$

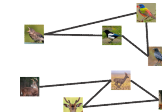
→ represent data point as linear combination among all data points

Solve

$$-0.1 \mathbf{x}_1 + 0.3 \mathbf{x}_2 + 0.8 \mathbf{x}_3 \approx \mathbf{x}_4$$

→ only those data points in the same subspace will remain

SDP property



→ perfect clustering is theoretically proven

## CLUSTERING PERFORMANCE

	Nystrom	AKK	SSC	SSC-OMP	SSC-ORGEN	SSSC	S <sup>5</sup> C
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

## EXPERIMENTAL SCALABILITY

