

Fast Clustering with Co-clustering via Discrete Non-negative Matrix Factorization for Image Identification – ICASSP 2020

Feiping Nie, Shenfei Pei, Rong Wang, Xuelong Li

Center for OPTical IMagery Analysis and Learning (OPTIMAL)
Northwestern Polytechnical University

shenfeipei@gmail.com

April 17, 2020

Fast Clustering with Co-clustering via Discrete Non-negative Matrix Factorization for Image Identification – ICASSP 2020

Feiping Nie, Shenfei Pei, Rong Wang, Xuelong Li

Center for **OPT**ical **IM**agery **A**nalysis and **L**earning (**OPTIMAL**)
Northwestern Polytechnical University

shenfeipei@gmail.com

April 17, 2020



Outline

- 1 Large-scale Datasets
- 2 Introduction to Clustering
 - Definition of clustering
 - Related work
- 3 Motivation
- 4 Our Model
- 5 Optimization
 - Contributions
- 6 Experiments

- 1 Large-scale Datasets
- 2 Introduction to Clustering
 - Definition of clustering
 - Related work
- 3 Motivation
- 4 Our Model
- 5 Optimization
 - Contributions
- 6 Experiments

Large-scale Datasets

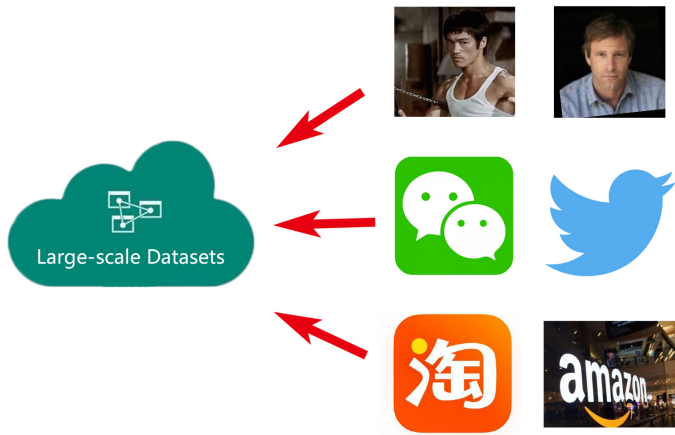


Figure 1: Large-scale dataset can be found everywhere in our lives

Large-scale Datasets

Why Clustering?

- Clustering has no requirement on data.
- Collecting unlabeled data is easy.

Supervised learning
(Labeled data)

**Mislabeling
time consuming
expensive**

Unsupervised learning
(Unlabeled data)

None

Requirements of data for supervised
and unsupervised learning tasks

Figure 2: Supervised Learning and Unsupervised learning

- 1 Large-scale Datasets
- 2 Introduction to Clustering
 - Definition of clustering
 - Related work
- 3 Motivation
- 4 Our Model
- 5 Optimization
 - Contributions
- 6 Experiments

Introduction to Clustering

Definition of clustering

Clustering is the task of grouping a set of objects in such a way that objects in the same group are more similar (in some sense) to each other than to those in other groups.

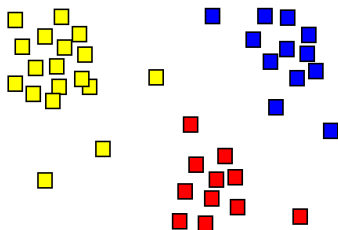


Figure 3: Schematic diagram of clustering

Introduction to Clustering

Related work

- The family of spectral clustering methods gains the most popularity.
- Despite its good performance, the time and space complexity of SC are $O(n^3)$ and $O(n^2)$, respectively.
- Much effort has been devoted for accelerating the spectral clustering algorithm, in recent years.

Approximate
eigenvalue
decomposition

(a)

Sampling
based methods

(b)

Sparse coding

(c)

Figure 4: Related work

- 1 Large-scale Datasets
- 2 Introduction to Clustering
 - Definition of clustering
 - Related work
- 3 Motivation**
- 4 Our Model
- 5 Optimization
 - Contributions
- 6 Experiments

Motivation

Key observation

- Bipartite spectral graph partition (BSGP) is the most famous co-clustering algorithm because of its remarkable performance.
- The similarity between the sample and the anchor can be treated as another description of sample.
- Recent studies have shown that using anchor graph to construct similar matrix can still yield promising results.

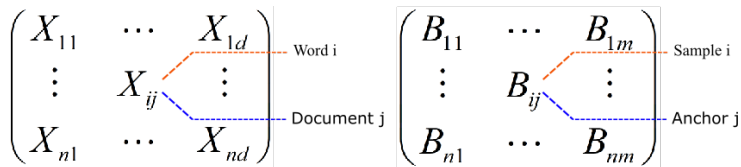


Figure 5: From BSGP to our model

Motivation

Our model

According to BSGP, it is not difficult to get our model (replace the data matrix X with the similarity matrix between samples and anchors B).

$$\min_{Y \in \Phi^{(n+m) \times c}} \sum_{k=1}^c \frac{y_k^T L y_k}{y_k^T D y_k}$$

where $L = D - W$, D is a diagonal matrix, $D_{ii} = \sum_{j=1}^{n+m} W_{ij}$,

$$W = \begin{pmatrix} 0 & B \\ B^T & 0 \end{pmatrix}$$

- 1 Large-scale Datasets
- 2 Introduction to Clustering
 - Definition of clustering
 - Related work
- 3 Motivation
- 4 Our Model**
- 5 Optimization
 - Contributions
- 6 Experiments

Our Model

Derivation

$$\min_{Y \in \Phi^{(n+m) \times c}} \sum_{k=1}^c \frac{y_k^T L y_k}{y_k^T D y_k} \quad (1)$$

$$\min_{Y \in \Phi^{(n+m) \times c}} \text{Tr}(Y^T L Y (Y^T D Y)^{-1}) \quad (2)$$

$$\max_{Y \in \Phi^{(n+m) \times c}} \text{Tr}(Y^T W Y (Y^T D Y)^{-1}) \quad (3)$$

Taking $Y^T = [P^T \quad Q^T]$ into Eq. (5.3), we have

$$Y^T W Y = P^T B Q + Q^T B^T P \quad (4)$$

$$Y^T D Y = P^T D^{(1)} P + Q^T D^{(2)} Q \quad (5)$$

where $D^{(1)}$ and $D^{(2)}$ are both diagonal matrices, $D_{ii}^{(1)} = \sum_{j=1}^m B_{ij}$,

$$D_{jj}^{(2)} = \sum_{i=1}^n B_{ij}$$

Our Model

With those notations in Eq. (4) and Eq. (5), the problem in Eq. (3) can be rewritten as follows:

$$\max_{P \in \Phi^{n \times c}, Q \in \Phi^{m \times c}} \text{Tr}(P^T B Q (P^T D^{(1)} P + Q^T D^{(2)} Q)^{-1}) \quad (6)$$

Relaxation

We relax the problem in eq. (6) into the following form by adding two terms $\text{Tr}(T^{-1} P^T P T^{-1} Q^T Q)$ and $\text{Tr}(B^T B)$, where T represents $P^T D^{(1)} P + Q^T D^{(2)} Q$.

$$\min_{P, Q} \|B - P(P^T D^{(1)} P + Q^T D^{(2)} Q)^{-1} Q^T\|_F^2 \quad (7)$$

Our Model

Relaxation

$$\min_{P,Q} \|B - P(P^T D^{(1)} P + Q^T D^{(2)} Q)^{-1} Q^T\|_F^2$$



$$\min_{P,Q,S} \|B - PSQ^T\|_F^2, \quad (8)$$

$$s.t. \quad P \in \Phi^{n \times c}, Q \in \Phi^{m \times c}, S \in R^{c \times c}$$

The optimization problem in Eq. (8) can be solved using standard techniques (alternating minimization).

- 1 Large-scale Datasets
- 2 Introduction to Clustering
 - Definition of clustering
 - Related work
- 3 Motivation
- 4 Our Model
- 5 Optimization**
 - Contributions
- 6 Experiments

Optimization

Update S

$$\min_{S \in \mathbb{R}^{c \times c}} \|B - PSQ^T\|_F^2. \quad (9)$$

Let $J(S)$ denote the objective function in Eq. (9). The derivative of $J(S)$ with respect to S is as follows:

$$\frac{\partial J(S)}{\partial S} = -2P^T BQ + 2P^T PSQ^T Q. \quad (10)$$

By setting the derivative of the objective function with respect to S to zero, we have

$$S_{ij} = \frac{(P^T BQ)_{ij}}{(P^T P)_{ii} (Q^T Q)_{jj}}. \quad (11)$$

Update P

$$\min_{P \in \Phi^{n \times c}} \|B - P(SQ^T)\|_F^2. \quad (12)$$

The solution can be determined by

$$P_{ij} = \begin{cases} 1 & j = \arg \min_k \|B_i - (SQ^T)_k\|_2^2, \\ 0 & \text{otherwise.} \end{cases} \quad (13)$$

Update Q

$$\min_{Q \in \Phi^{m \times c}} \|B - (PS)Q^T\|_F^2, \quad (14)$$

The solution is determined by

$$Q_{ij} = \begin{cases} 1 & j = \arg \min_k \|B^i - (PS)^k\|_2^2, \\ 0 & \text{otherwise,} \end{cases} \quad (15)$$

where $B^i((PS)^k)$ denote i -th (k -th) column of $B(PS)$.

Algorithm

Algorithm 1: Algorithm to solve the problem in Eq. (8)

Data: Data matrix $X \in \mathbb{R}^{n \times d}$, the number of anchors and nearest neighbors;

Result: Indicator matrices P and Q

Construct B according to [21] and initialize P and Q in a random way;

while *not converge* **do**

 Compute S by Eq. (11) ;

 Compute P by Eq. (13) ;

 Compute Q by Eq. (15) ;

end

Contributions

Pros

- The limitation of collaborative clustering can only be applied to specific scene can be broken by introducing anchor-based strategy.
- The final clustering result can be obtained directly without any post-processing
- The time and space complexity of FCDMF are both linear with respect to the number of samples.

- 1 Large-scale Datasets
- 2 Introduction to Clustering
 - Definition of clustering
 - Related work
- 3 Motivation
- 4 Our Model
- 5 Optimization
 - Contributions
- 6 Experiments**

Experiments

Comparison methods

- Traditional Spectral Clustering (SC) (NIPS-2002)
- Scalable Spectral Clustering with cosine similarity (SSC) (ICPR-2018)
- Improved Anchor-based Graph Clustering based on multiplicative update optimization (AGC-I) (RS-2019)
- Fast Spectral Clustering with anchor graph for large hyperspectral images (FSC) (GRSL-2017)
- Large scale Spectral Clustering via landmark-based sparse representation (LSC) (TC-2015)
- Fast Clustering with co-clustering via Discrete non-negative Matrix Factorization (FCDMF) (Our method)

Experiments

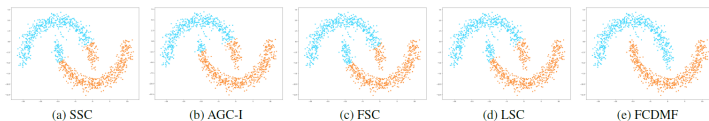


Figure 1: Results on the Two-moon synthetic data.

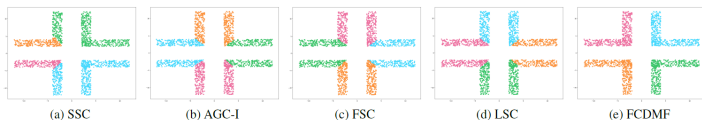


Figure 2: Results on the Four-corner synthetic data.

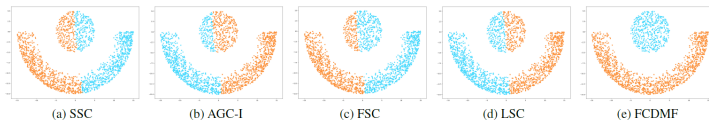


Figure 3: Results on the Crescent-fullmoon synthetic data.

Experiments

Table 2: The average accuracy (\pm standard deviation) of several fast spectral clustering methods (The best result on each data set is highlighted in bold).

	SC	SSC	AGCI	FSC	LSC	FCDMF
BinAlpha	0.449(\pm 0.015)	0.288(\pm 0.013)	0.416(\pm 0.012)	0.405(\pm 0.013)	0.407(\pm 0.015)	0.421 (\pm 0.012)
FACE-94	0.904(\pm 0.012)	0.719(\pm 0.014)	0.746(\pm 0.017)	0.749(\pm 0.016)	0.778(\pm 0.017)	0.914 (\pm 0.011)
FACE-95	0.532(\pm 0.012)	0.235(\pm 0.010)	0.396(\pm 0.013)	0.387(\pm 0.012)	0.305(\pm 0.013)	0.487 (\pm 0.010)
FEI	0.509(\pm 0.020)	0.050(\pm 0.007)	0.379(\pm 0.015)	0.406(\pm 0.017)	0.388(\pm 0.017)	0.499 (\pm 0.014)
FERET	0.303(\pm 0.007)	0.163(\pm 0.008)	0.214(\pm 0.004)	0.196(\pm 0.009)	0.187(\pm 0.008)	0.252 (\pm 0.006)
FingerPrint	0.569(\pm 0.028)	0.172(\pm 0.009)	0.334(\pm 0.022)	0.306(\pm 0.012)	0.383(\pm 0.022)	0.525 (\pm 0.023)
Grimace	0.929(\pm 0.012)	0.130(\pm 0.010)	0.905(\pm 0.026)	0.923(\pm 0.022)	0.820(\pm 0.029)	0.967 (\pm 0.002)
GTDB	0.548(\pm 0.015)	0.199(\pm 0.011)	0.345(\pm 0.015)	0.336(\pm 0.015)	0.327(\pm 0.015)	0.472 (\pm 0.015)
IMM	0.593(\pm 0.020)	0.109(\pm 0.009)	0.319(\pm 0.015)	0.451(\pm 0.027)	0.407(\pm 0.023)	0.547 (\pm 0.018)
JAFFE	0.825(\pm 0.023)	0.640(\pm 0.034)	0.888 (\pm 0.023)	0.832(\pm 0.012)	0.705(\pm 0.043)	0.845(\pm 0.023)
JAFFE2	0.178(\pm 0.000)	0.174(\pm 0.002)	0.174(\pm 0.004)	0.175(\pm 0.008)	0.187(\pm 0.003)	0.197 (\pm 0.000)
MPEG-7	0.565(\pm 0.011)	0.477 (\pm 0.017)	0.194(\pm 0.010)	0.181(\pm 0.007)	0.181(\pm 0.008)	0.429(\pm 0.011)
ORL	0.588(\pm 0.023)	0.514(\pm 0.022)	0.427(\pm 0.017)	0.497(\pm 0.021)	0.463(\pm 0.023)	0.538 (\pm 0.017)
PALM	0.786(\pm 0.016)	0.861 (\pm 0.023)	0.600(\pm 0.014)	0.704(\pm 0.017)	0.635(\pm 0.020)	0.747(\pm 0.012)
Pixraw10P	0.910(\pm 0.045)	0.141(\pm 0.003)	0.729(\pm 0.058)	0.632(\pm 0.019)	0.659(\pm 0.058)	0.930 (\pm 0.000)
UMIST	0.409(\pm 0.014)	0.409(\pm 0.021)	0.390(\pm 0.011)	0.387(\pm 0.011)	0.370(\pm 0.011)	0.449 (\pm 0.011)
YALE	0.511(\pm 0.023)	0.332(\pm 0.017)	0.393(\pm 0.023)	0.389(\pm 0.022)	0.397(\pm 0.024)	0.489 (\pm 0.021)

Experiments

Table 4: The average normalized mutual information (\pm standard deviation) of several spectral clustering methods (The best result on each data set is highlighted in bold).

	SC	SSC	AGCI	FSC	LSC	FCDMF
BinAlpha	0.592(\pm 0.007)	0.443(\pm 0.011)	0.578(\pm 0.006)	0.58 (\pm 0.006)	0.562(\pm 0.007)	0.570(\pm 0.006)
FACE-94	0.969(\pm 0.004)	0.881(\pm 0.005)	0.938(\pm 0.004)	0.939(\pm 0.005)	0.937(\pm 0.005)	0.973 (\pm 0.003)
FACE-95	0.743(\pm 0.005)	0.492(\pm 0.010)	0.676(\pm 0.007)	0.679(\pm 0.009)	0.618(\pm 0.010)	0.706 (\pm 0.005)
FEI	0.716(\pm 0.009)	0.103(\pm 0.020)	0.656(\pm 0.007)	0.677(\pm 0.010)	0.663(\pm 0.010)	0.696 (\pm 0.007)
FERET	0.693(\pm 0.003)	0.501(\pm 0.021)	0.641(\pm 0.003)	0.585(\pm 0.012)	0.559(\pm 0.013)	0.669 (\pm 0.003)
FingerPrint	0.697(\pm 0.017)	0.347(\pm 0.016)	0.577(\pm 0.022)	0.569(\pm 0.013)	0.595(\pm 0.021)	0.669 (\pm 0.014)
Grimace	0.954(\pm 0.004)	0.214(\pm 0.014)	0.964(\pm 0.008)	0.974 (\pm 0.007)	0.921(\pm 0.009)	0.972(\pm 0.002)
GTDB	0.712(\pm 0.008)	0.416(\pm 0.014)	0.626(\pm 0.008)	0.605(\pm 0.011)	0.59(\pm 0.011)	0.664 (\pm 0.007)
IMM	0.767(\pm 0.010)	0.253(\pm 0.028)	0.608(\pm 0.01)	0.728(\pm 0.018)	0.691(\pm 0.016)	0.744 (\pm 0.009)
JAFFE	0.861(\pm 0.010)	0.650(\pm 0.027)	0.877 (\pm 0.005)	0.868(\pm 0.009)	0.799(\pm 0.016)	0.822(\pm 0.018)
JAFFE2	0.013(\pm 0.000)	0.052(\pm 0.003)	0.015(\pm 0.002)	0.013(\pm 0.005)	0.035(\pm 0.003)	0.087 (\pm 0.001)
MPEG-7	0.738(\pm 0.004)	0.672 (\pm 0.009)	0.485(\pm 0.018)	0.439(\pm 0.019)	0.431(\pm 0.019)	0.652(\pm 0.004)
ORL	0.770(\pm 0.011)	0.708(\pm 0.012)	0.680(\pm 0.008)	0.735(\pm 0.012)	0.693(\pm 0.014)	0.747 (\pm 0.008)
PALM	0.924(\pm 0.006)	0.958 (\pm 0.006)	0.844(\pm 0.004)	0.903(\pm 0.005)	0.871(\pm 0.007)	0.891(\pm 0.004)
Pixraw10P	0.928(\pm 0.017)	0.159(\pm 0.005)	0.862(\pm 0.026)	0.816(\pm 0.008)	0.808(\pm 0.026)	0.935 (\pm 0.000)
UMIST	0.649(\pm 0.008)	0.591(\pm 0.013)	0.638(\pm 0.011)	0.645 (\pm 0.012)	0.623(\pm 0.011)	0.642(\pm 0.011)
YALE	0.572(\pm 0.016)	0.406(\pm 0.015)	0.494(\pm 0.020)	0.486(\pm 0.020)	0.481(\pm 0.019)	0.533 (\pm 0.013)

Conclusions

Conclusions

- Our model relaxed the original objective function to a non-negative matrix factorization problem.
- In our model, the final clustering result can be obtained directly without any post-processing.
- An efficient optimization algorithm whose time and space complexity are both linear with respect to the number of samples
- Substantial performance

Thanks for Listening!
Questions?