

Dissimilarity-based sparse Subset Selection

paper report by
Mohsen Nabian

goal is

- in this paper, given pairwise dissimilarities d_{ij} between the elements of "source set" and "target set". we consider the problem of finding a subset of source set, called "representative" that can efficiently describe the target set.

Formulations:

→ source set

$$X = \{x_1, \dots, x_M\}$$

$$Y = \{y_1, \dots, y_N\}$$

↘ target set

$$\text{or } X = \begin{bmatrix} x_1 & x_2 & \dots & x_M \end{bmatrix}$$

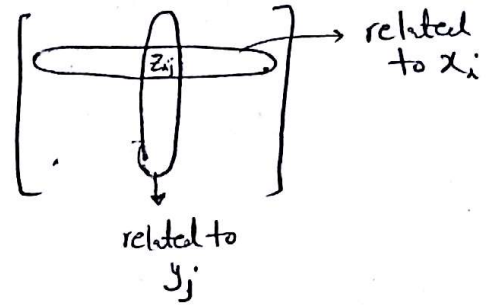
$$Y = \begin{bmatrix} y_1 & y_2 & \dots & y_N \end{bmatrix}$$

We are given pairwise dissimilarities: $\{d_{ij}\}_{i=1, \dots, M}^{j=1, \dots, N}$ between each x_i to each y_j

- d_{ij} shows how well x_i represents y_j .
- Smaller d_{ij} , the better x_i represent y_j .

$$\begin{matrix} \text{is given} \\ \downarrow \\ D \triangleq \begin{bmatrix} d_1^T \\ \vdots \\ d_M^T \end{bmatrix} = \begin{bmatrix} d_{11} & d_{12} & \dots & d_{1N} \\ \vdots & \vdots & & \vdots \\ d_{M1} & d_{M2} & \dots & d_{MN} \end{bmatrix} \in \mathbb{R}^{M \times N} \end{matrix}$$

we have to find $Z \in \mathbb{R}^{M \times N}$ matrix that meets the following objectives:



$$- z_{ij} \in \{0, 1\}$$

$$- \sum_{i=1}^M z_{ij} = 1 \quad \forall j \text{ means every element in } y, \text{ should be presented by only 1 } x_i$$

$\Rightarrow d_{ij} z_{ij}$ is the cost encoding y_j by x_i .
So if z_{ij} is not 0, d_{ij} should be small.

So we can formulate the problem as follows:

$$\left\{ \begin{array}{l} \min \quad \underbrace{\lambda \sum_{i=1}^M I(\|z_i\|_p)}_{\text{number of representative}} + \underbrace{\sum_{j=1}^N \sum_{i=1}^M d_{ij} z_{ij}}_{\text{total cost of encoding } Y} \quad (*) \\ \text{s.t.} \quad \sum_{i=1}^M z_{ij} = 1, \forall j; \quad z_{ij} \in \{0, 1\} \quad \forall i, j \end{array} \right.$$

we want $\sum_{i=1}^M I(\|z_i\|_p)$ to be minimized so that we have less points involved on representing Y target set.
we want more all zero rows in Z .

Problem (*) is NP-hard and non Convex :

Convex form :
$$(**) \begin{cases} \min_{\{z_{ij}\}} \lambda \sum \|z_i\|_p + \sum_{j=1}^N \sum_{i=1}^M d_{ij} z_{ij} \\ \text{s.t.} \sum_{i=1}^M z_{ij} = 1, \forall j; z_{ij} \geq 0 \forall i,j \end{cases}$$

Equivalent Matrix form
$$\begin{cases} \min_Z \lambda \|Z\|_{1,p} + \text{tr}(D^T Z) \\ \text{s.t.} \mathbf{1}^T Z = \mathbf{1}^T, Z \geq 0 \end{cases}$$

where
$$\|Z\|_{1,p} \triangleq \sum_{i=1}^M \|z_i\|_p$$

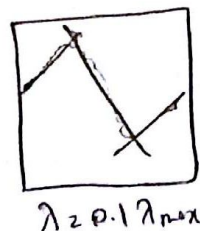
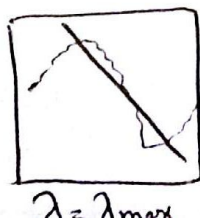
λ : $\begin{cases} \text{if } \lambda \rightarrow 0 \Rightarrow \text{each } y_j \text{ chooses the nearest } x_i \\ \text{if } \lambda \rightarrow \infty \Rightarrow \text{only one representative will be selected.} \end{cases}$

as example :

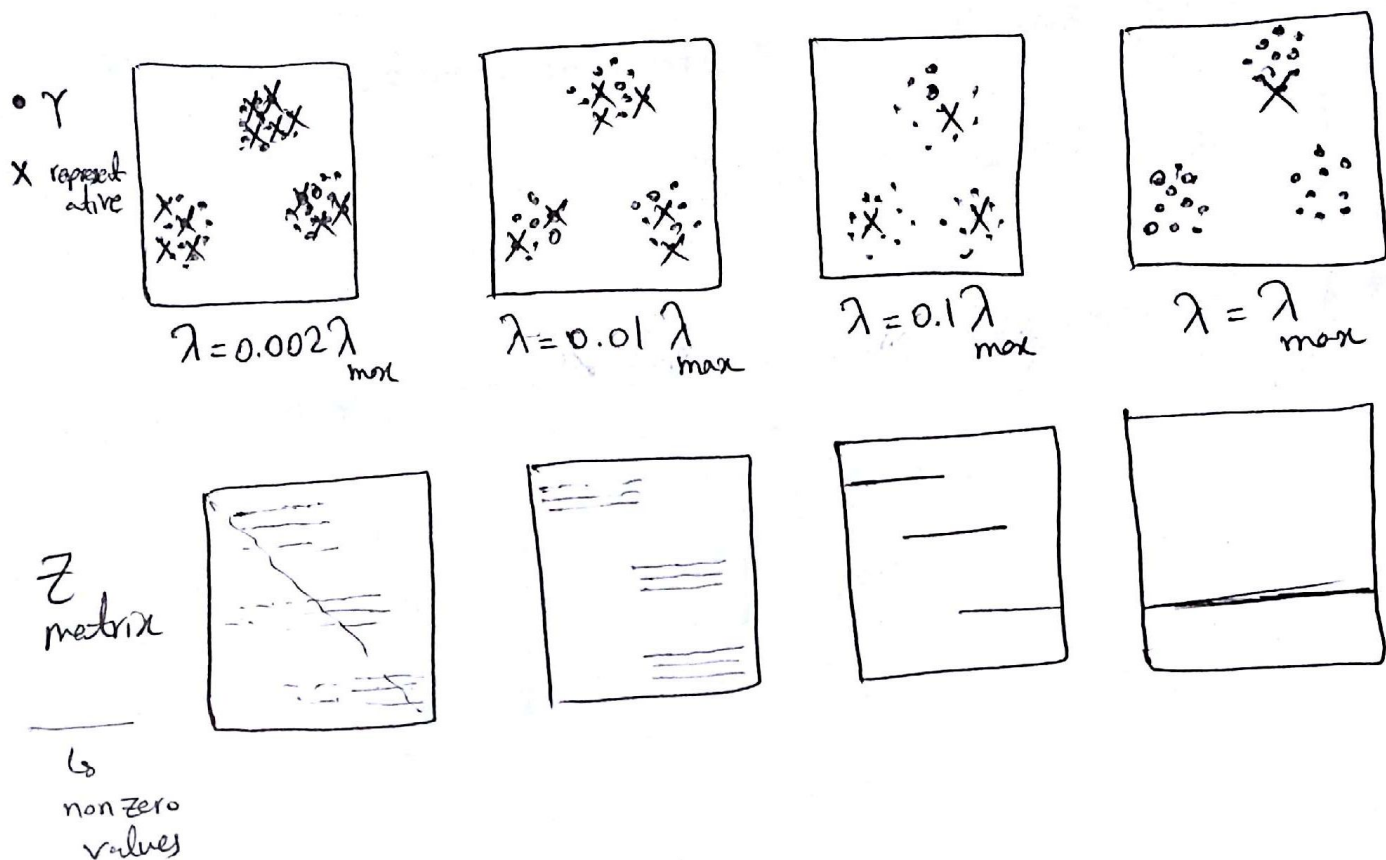
we have a noisy data $\{y_j\}_{j=1}^N$ and $X = \{\theta_i\}_{i=1}^N$ where each $\theta_i = (a_i, b_i)$ represent a model.

θ_i are obtained as follows: for each y_j and its $k=4$ nearest neighbors, we find θ_i by minimizing $|a^T y - b|$ for the $k+1$ points.

by changing λ we see:



in another experiment: $X = Y = \underline{3}$ gaussian representative



Dealing with outliers:

outlier in { Source set \rightarrow if X_i is outlier, it can not represent any Y_j , therefore automatically removed not selected
 Target set

\hookrightarrow we define e_j corresponding to each target Y_j and set: $\sum z_{ij} + e_j = 1$

optimization robust to outlier: (***)

$$\begin{cases} \min_{\{z_{ij}\}, \{e_j\}} & \lambda \sum_{i=1}^M \|z_i\|_p + \sum_{j=1}^N \sum_{i=1}^M d_{ij} z_{ij} + \sum_{j=1}^N \omega_j e_j \\ \text{s.t.} & \left(\sum_{i=1}^M z_{ij} \right) + e_j = 1, \forall j; \quad z_{ij} \geq 0 \\ & e_j \geq 0 \quad \forall j \end{cases}$$

Since $\left(\sum_{i=1}^M z_{ij}\right) + e_j = 1$,

if y_j is outlier, minimization will put $e_j = 1$ and $\left(\sum_{i=1}^M z_{ij}\right) = 0$ therefore, y_j is not represented at all.

opt (***) in Matrix form: (equivalent to ***)

$$\min_{z, e} \lambda \|z\|_{1,p} + \text{tr} \left(\begin{bmatrix} D \\ w^T \end{bmatrix}^T \begin{bmatrix} z \\ e^T \end{bmatrix} \right)$$

$$\text{s.t.} \quad 1^T \begin{bmatrix} z \\ e^T \end{bmatrix} = 1^T, \quad \begin{bmatrix} z \\ e^T \end{bmatrix} \geq 0$$

also w_j in optimization (***) can be chosen as

$$\underline{w_j = \beta e^{-\frac{\min_i d_{ij}}{\tau}}} \quad \leftarrow \text{suggestion}$$

another choice: $w_j = w$ for all j .

Clustering via Representative:

- ① Optimal solution z^* indicates elements of x that are representative to y .
- ② They also provides membership information.
~~ie~~ ie. for $y_j \rightarrow [z_{1j}^*, \dots, z_{mj}^*]$ is probability vector for presenting y_j by elements of X . \Rightarrow This is like soft assignment for clustering.

we can also do hard assignment by taking minimum d_{ij} for those selected $\{x_{e1}, \dots, x_{ek}\}$ for presenting y_j . i.e. we take the closest among selected x_i s to y_j (hard assignment)

DSS Implementation

author has introduced an algorithm to efficiently solve minimization (***) using ADMM.

This algorithm is provided in page 7 of the paper.

The computational complexity is of $O(MN)$ very fast

61

in table 1, the average computational time of the proposed ADMM algorithm is compared to the CVX (SeDuMi Solver) which showed ~~that~~ that proposed ADMM outperforms significantly in all trials with $\lambda = 0.01 \cdot \lambda_{\max}$.
 data size $N \times N$, experiment on:
$$\begin{cases} N = \{30, 50, 100, 200, 500, 1000, 2000\} \\ P = \{2, \infty\} \end{cases}$$

Regularization parameter effect

regularization puts a trade off between the number of representative and encoding cost.

Theorem 1: if
$$\begin{cases} \lambda_{\max, 2} \triangleq \max_{i \neq l} \frac{\sqrt{N}}{2} \frac{\|d_i - d_l^*\|_2^2}{\gamma^T(d_i - d_l^*)} \\ \lambda_{\max, \infty} \triangleq \max_{i \neq l} \frac{\|d_i - d_l^*\|_1}{2} \end{cases} \quad l^* = \arg \min_i \gamma^T d_i$$

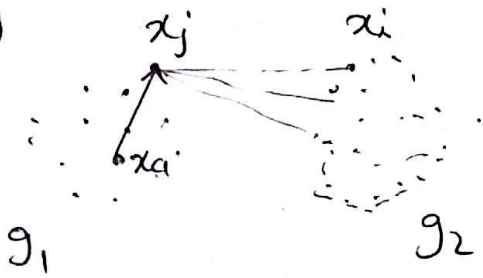
Then entire γ will be presented only by one single x^*

Clustering Guarantee

under 2 conditions, clustering is guaranteed.

Figure 7 shows the conditions.

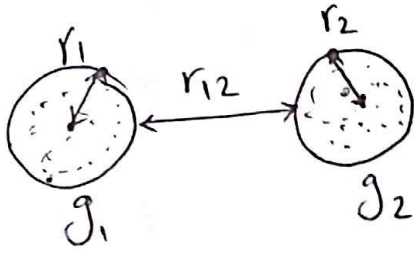
①



distance x_j to x_{ci} is lower than distance x_j to any point in g_2

Same condition for g_2 against g_1

②



$$r_{12} < r_1 + r_2$$

Experiments :

1) classification using representative:

$K = 15$ class 4485 images 210^{or} 410 images per class

80% training 20% test

we find representative of the training data in each class and use them as reduced data to perform NN classification on the test data.

after selecting η fraction of training samples in each class using each algorithm.

$$err(\eta) = accuracy(1) - accuracy(\eta)$$

8/

Table 2 showed the classification performance using different methods:

Algorithm	Rare	K-medoids	AP	DS3
$\eta = 0.05$			*	*
$\eta = 0.10$				*
$\eta = 0.20$				*
$\eta = 0.35$				*

(*) → best performance

also

Confusion matrix
 $\eta = 0.05$

Confusion matrix
 $\eta = 0.35$

Confusion matrix
 $\eta = 1$

② Modeling segmentation of dynamic data:

- modeling and segmentation of human activities in motion capture data
- data set = time series of different subjects each performing several activities.

- used 14 most informative joints

Seq number	7	2	-	-	-	-	14
# frames	865	2115	-	-	-	-	1204
# activities	4	8	-	-	-	-	4

SC error/.

SBiC error/.

Kmedoid error/.

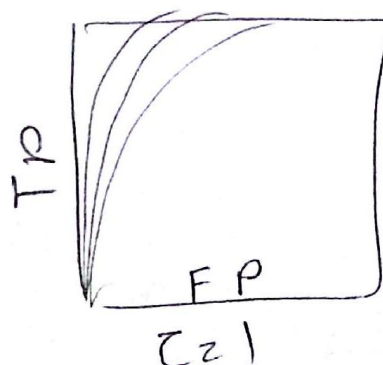
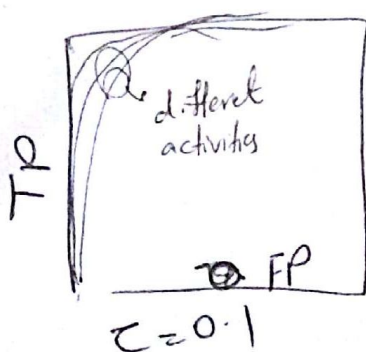
AP error

DSS error/.

↳ outperformed others in almost all ~~experiments~~ sequences.

outlier

The ~~author~~ author also showed algorithms robustness to outlier by plotting True positive rate vs ~~to~~ false positive for some outliers in prev. experiment for $\tau = 0.1$ & $\tau = 1.0$



⇒ Algorithm is very robust to outlier