

# Unsupervised Feature Selection with Adaptive Structure Learning

Reporter: Songling Liu



Data Mining Lab  
Big Data Research Center

# Overview

- Current feature selection methods
- Motivation
- Method
- Optimization Algorithm
- Experiments
- Conclusion

# Current feature selection methods

Method type	Structure characterization	Intermediate Analysis	Feature search	Typical algorithm
Filter	Structure learning with All features		Ranking criteria	MaxVar, LapScore, SPEC, EVSC
Embedded Type I	Structure learning with All features		A learning model	TraceRatio, UDFS
Embedded Type II	Structure learning with All features	Clustering	A learning model	MCFS, MRFS, SPFS, FSSL, GLSPFS
Embedded Type III	Structure learning with All features	Clustering	A learning model	JELSR, NDFS, RUFs, CGSSL

# Current feature selection methods

Filter:

Addressing this issue by selecting the top ranked features based on some scores computed independently for each feature.

feature 1	feature 2	feature 3

If in variance:  $1 > 2 > 3$  and we choose the largest two

# Current feature selection methods

## Embedded Method:

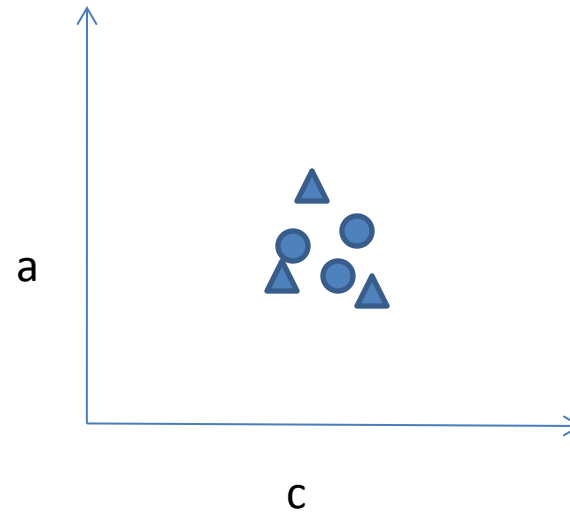
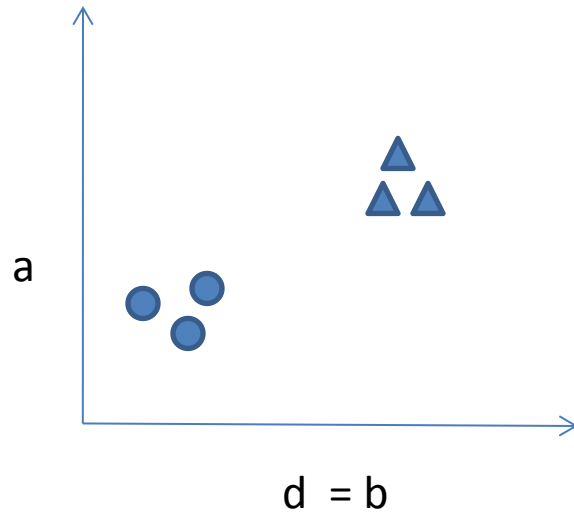
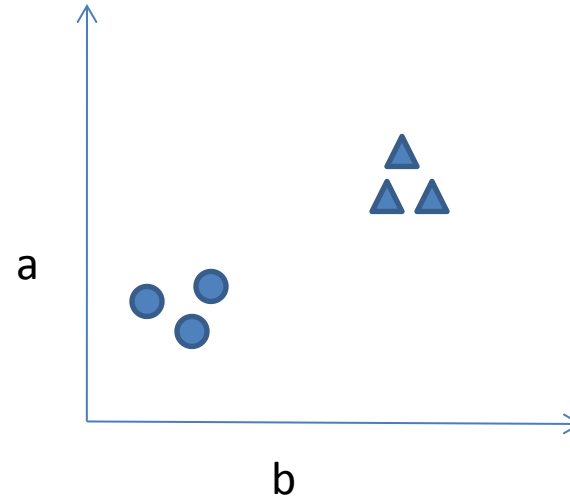
Using all features to estimate the underlying structures of data and select features which can preserve data structures.

Common drawback--Using all features which could be redundant and include noisy features.



# Current feature selection methods

Suppose a dataset with 4 features a,b,c and d.



# Data structure estimating VS Feature selection







# Motivation



- A unified learning framework which performs structure learning and feature selection **simultaneously**.
- Select good features from well-estimated data structure.
- Estimate data structure with good features.

# Problem I

How to preserve pairwise global structure?

PCA? MaxVar?

However, such dense similarity becomes less discriminative for high dimension data, especially when there are many unfavorable features in the original high dimensional space.



# Method

## Adaptive Global Structure Learning

using the sparse reconstruction coefficients to extract the global structure of data.

$$\min_{\mathbf{S}} \sum_{i=1}^n (\|\mathbf{x}_i - \mathbf{X}\mathbf{s}_i\|^2 + \alpha\|\mathbf{s}_i\|_1) \quad \text{s.t.} \quad \mathbf{S}_{ii} = 0$$

# Method

the selected features should preserve such global and sparse reconstruction structure

$$\begin{aligned} \min_{\mathbf{S}, \mathbf{W}} \quad & \sum_{i=1}^n \|\mathbf{W}^T \mathbf{x}_i - \mathbf{W}^T \mathbf{X} \mathbf{s}_i\|^2 + \alpha \|\mathbf{S}\|_1 + \gamma \|\mathbf{W}\|_{21} \\ \text{s.t.} \quad & \mathbf{S}_{ii} = 0, \mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W} = \mathbf{I} \end{aligned}$$

Compared with the last equation:

- 1) The global structure captured by  $\mathbf{S}$  can be used to **guide the search of relevant features**;
- 2) the global structure can also be better estimated.

# Problem II

How to preserve pairwise local structure?

LLE? Graph Laplacian?

They would be inevitably affected by the redundant and noisy features. Moreover, the iterative updating of discrete neighborhood relationship using the result of feature selection still suffers from the lack of theoretical guarantee of its convergence.



# Method

## Adaptive Local Structure Learning

Learn a euclidean distance induced probabilistic neighborhood matrix

$$\min_{\mathbf{P}} \sum_{i,j} (\|x_i - x_j\|_2^2 \mathbf{P}_{ij} + \mu \mathbf{P}_{ij}^2), \text{ s.t. } \mathbf{P} \mathbf{1}_n = \mathbf{1}_n, \mathbf{P} \geq 0$$

Using matrix  $\mathbf{P}$ , graph laplacian can be characterized as

$$\mathbf{L}_{\mathbf{P}} = \mathbf{D}_{\mathbf{P}} - (\mathbf{P} + \mathbf{P}^T)/2$$

Where  $\mathbf{D}_{\mathbf{P}}$ 's element is  $\sum_j (\mathbf{P}_{ij} + \mathbf{P}_{ji})/2$



# Method

Using the sparse matrix to select **informative features**

$$\begin{aligned} \min_{\mathbf{P}, \mathbf{W}} \quad & \sum_{i,j}^n (\|\mathbf{W}^T \mathbf{x}_i - \mathbf{W}^T \mathbf{x}_j\|_2^2 \mathbf{P}_{ij} + \mu \mathbf{P}_{ij}^2) + \gamma \|\mathbf{W}\|_{21} \\ \text{s.t.} \quad & \mathbf{P} \mathbf{1}_n = \mathbf{1}_n, \mathbf{P} \geq \mathbf{0}, \mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W} = \mathbf{I} \end{aligned}$$

we can learn a better probabilistic neighborhood graph for local structure characterization

# Method

## Unsupervised Feature Selection with Adaptive Structure Learning

$$\begin{aligned} \min_{\mathbf{W}, \mathbf{S}, \mathbf{P}} \quad & \left( \|\mathbf{W}^T \mathbf{X} - \mathbf{W}^T \mathbf{X} \mathbf{S}\|^2 + \alpha \|\mathbf{S}\|_1 \right) \\ & + \beta \sum_{i,j}^n \left( \|\mathbf{W}^T \mathbf{x}_i - \mathbf{W}^T \mathbf{x}_j\|^2 \mathbf{P}_{ij} + \mu \mathbf{P}_{ij}^2 \right) + \gamma \|\mathbf{W}\|_{21} \\ \text{s.t.} \quad & \mathbf{S}_{ii} = 0, \mathbf{P} \mathbf{1}_n = \mathbf{1}_n, \mathbf{P} \geq \mathbf{0}, \mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W} = \mathbf{I} \end{aligned}$$



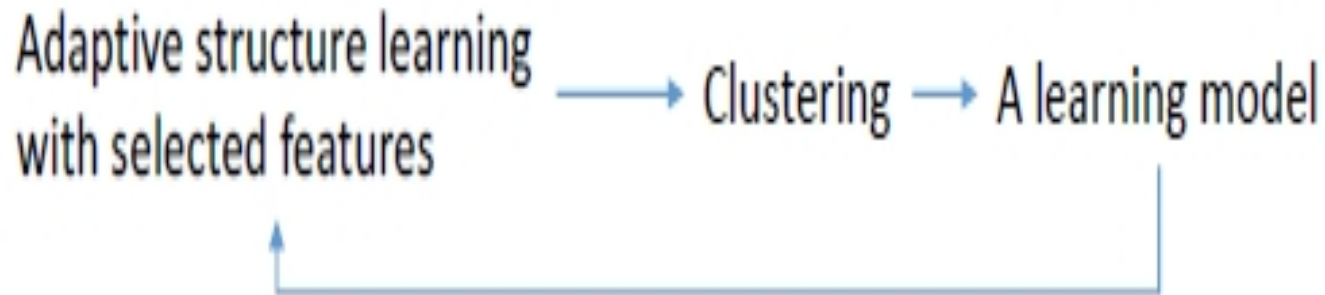
# Method

## When S and P fixed

It selects those **features** to well respect both the global and local structure of data;

## When W fixed

It learns the **global and local structure** of data in a transformed space



# Optimization Algorithm

First, when  $W$  and  $P$  are fixed,

$$\min_{\mathbf{s}_i} \|\mathbf{x}'_i - \mathbf{X}'\mathbf{s}_i\|^2 + \alpha|\mathbf{s}_i|, \quad \text{s.t.} \quad \mathbf{S}_{ii} = 0$$

Where  $\mathbf{X}' = \mathbf{W}^T\mathbf{X}$ .

We get sparsity representation of instances, which means global structure

# Optimization Algorithm

Next, when  $W$  and  $S$  are fixed,

$$\begin{aligned} \min_{\mathbf{P}_i^T} \quad & \sum_{j=1}^n \|\mathbf{x}'_i - \mathbf{x}'_j\|^2 \mathbf{P}_{ij} + \mu \|\mathbf{P}_{ij}\|^2, \\ \text{s.t.} \quad & \mathbf{1}_n^T \mathbf{P}_i = 1, \mathbf{P}_{ij} \geq 0 \end{aligned}$$

With transformation

$$\min_{\mathbf{P}_i^T} \quad \frac{1}{2} \|\mathbf{P}_i^T - \mathbf{a}_i^T\|^2, \quad \text{s.t.} \quad \mathbf{P}_i^T \mathbf{1}_n = 1, 0 \leq \mathbf{P}_{ij}^T \leq 1$$

Where  $\mathbf{A}_{ij} = -\frac{1}{2\mu} \|\mathbf{x}'_i - \mathbf{x}'_j\|^2,$

# Optimization Algorithm

Next, when  $S$  and  $P$  are fixed,

$$\begin{aligned} \min_{\mathbf{W}} \quad & \|\mathbf{W}^T \mathbf{X} - \mathbf{W}^T \mathbf{X} \mathbf{S}\|^2 + \beta \sum_{i,j}^n \|\mathbf{W}^T \mathbf{x}_i - \mathbf{W}^T \mathbf{x}_j\|^2 \mathbf{P}_{ij} + \gamma \|\mathbf{W}\|_{21} \\ \text{s.t.} \quad & \mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W} = \mathbf{I} \end{aligned} \quad (9)$$

Using  $\mathbf{L}_S = (\mathbf{I} - \mathbf{S})(\mathbf{I} - \mathbf{S})^T$ ,  $\mathbf{L}_P = \mathbf{D}_P - (\mathbf{P} + \mathbf{P}^T)/2$  and let  $\mathbf{L} = \mathbf{L}_S + \beta \mathbf{L}_P$ , the above problem can be rewritten as

$$\begin{aligned} \min_{\mathbf{W}} \quad & \text{Tr}(\mathbf{W}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{W}) + \gamma \|\mathbf{W}\|_{21} \quad (10) \\ \text{s.t.} \quad & \mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W} = \mathbf{I} \end{aligned}$$

# Optimization Algorithm

---

**Algorithm 2** The optimization algorithm of FSASL

---

**Input:** The data matrix  $\mathbf{X} \in \mathbb{R}^{d \times n}$ , the regularization parameters  $\alpha, \beta, \gamma, \mu$ , the dimension of the transformed data  $c$ .

**repeat**

For each  $i$ , update the  $i$ -th column of  $\mathbf{S}$  by solving the problem in Eq. (6);

For each  $i$ , update the  $i$ -th row of  $\mathbf{P}$  using Algorithm 1;

Compute the overall graph Laplacian  $\mathbf{L} = \mathbf{L}_S + \beta \mathbf{L}_P$ ;

Compute  $\mathbf{W}$  by Eq. (12) or Eq. (14);

**until** Converges

**Output:** Sort all the  $d$  features according to  $\|\mathbf{w}_i\|_2$  ( $i = 1, \dots, d$ ) in descending order and select the top  $m$  ranked features.

---

# Experiments

Data Sets	AllFea	LapScore	MCFS	LLCFS	UDFS	NDFS	SPFS	RUFS	JELSR	GLSPFS	FSASL
MFEA	68.73	51.78	51.04	60.38	64.94	67.13	<b>68.20</b>	64.58	67.01	61.00	<b>69.94</b>
		$\pm 5.51$	$\pm 8.13$	$\pm 8.58$	$\pm 3.32$	$\pm 7.53$	$\pm 9.43$	$\pm 7.99$	$\pm 8.37$	$\pm 8.70$	$\pm 7.19$
		0.00	0.00	0.00	0.00	0.01	<b>0.22</b>	0.00	0.01	0.00	<b>1.00</b>
USPS49	77.70	69.21	84.77	94.96	94.05	68.12	83.43	85.86	95.16	94.75	<b>95.95</b>
		$\pm 8.95$	$\pm 1.59$	$\pm 1.44$	$\pm 1.13$	$\pm 8.18$	$\pm 6.66$	$\pm 2.58$	$\pm 0.55$	$\pm 0.61$	$\pm 0.48$
		0.00	0.00	0.03	0.00	0.00	0.00	0.00	0.00	0.00	<b>1.00</b>
UMIST	42.40	36.73	44.46	47.31	48.04	52.80	46.72	50.87	53.52	50.53	<b>54.92</b>
		$\pm 1.18$	$\pm 3.26$	$\pm 0.83$	$\pm 1.92$	$\pm 2.26$	$\pm 1.70$	$\pm 1.95$	$\pm 1.54$	$\pm 0.59$	$\pm 1.89$
		0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.01	0.00	<b>1.00</b>
JAFFE	71.57	67.62	73.56	64.79	75.48	74.98	73.93	75.75	77.77	75.46	<b>79.29</b>
		$\pm 8.49$	$\pm 4.83$	$\pm 4.08$	$\pm 1.63$	$\pm 2.15$	$\pm 2.85$	$\pm 2.53$	$\pm 1.87$	$\pm 1.61$	$\pm 2.24$
		0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	<b>1.00</b>
AR	30.26	25.29	29.05	<b>34.22</b>	30.87	32.34	31.06	34.84	34.19	34.12	<b>36.11</b>
		$\pm 2.89$	$\pm 1.19$	$\pm 2.70$	$\pm 0.35$	$\pm 1.52$	$\pm 2.14$	$\pm 1.90$	$\pm 2.52$	$\pm 1.60$	$\pm 0.75$
		0.00	0.00	<b>0.05</b>	0.00	0.00	0.00	0.04	0.02	0.00	<b>1.00</b>
COIL	59.17	45.60	51.50	50.84	48.40	52.22	56.94	59.20	59.53	57.96	<b>60.93</b>
		$\pm 6.16$	$\pm 5.38$	$\pm 3.76$	$\pm 16.89$	$\pm 6.33$	$\pm 3.43$	$\pm 3.28$	$\pm 4.01$	$\pm 2.27$	$\pm 2.50$
		0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.03	0.00	<b>1.00</b>
LUNG	72.46	58.97	70.42	71.58	65.46	75.52	73.49	77.35	77.86	77.83	<b>81.93</b>
		$\pm 5.24$	$\pm 3.41$	$\pm 5.85$	$\pm 3.88$	$\pm 1.57$	$\pm 3.43$	$\pm 2.62$	$\pm 3.12$	$\pm 2.70$	$\pm 1.63$
		0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	<b>1.00</b>
TOX	43.65	40.25	43.10	39.28	47.14	38.28	39.93	49.17	43.96	47.38	<b>50.12</b>
		$\pm 0.65$	$\pm 1.86$	$\pm 0.49$	$\pm 0.75$	$\pm 1.64$	$\pm 1.13$	$\pm 0.83$	$\pm 1.56$	$\pm 1.93$	$\pm 0.67$
		0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	<b>1.00</b>
Average	58.24	49.43	55.98	57.92	59.29	56.67	59.21	62.2	63.62	62.38	<b>66.15</b>

# Experiments

Table 2: Aggregated clustering results measured by Normalized Mutual Information (%) of the compared methods.

Data Sets	AllFea	LapScore	MCFS	LLCFS	UDFS	NDFS	SPFS	RUFS	JELSR	GLSPFS	FSASL
MFEA	70.33	53.74	54.72	52.77	54.19	64.97	<b>64.92</b>	63.98	<b>64.51</b>	59.26	<b>66.70</b>
		$\pm 4.77$	$\pm 9.14$	$\pm 9.76$	$\pm 3.83$	$\pm 7.54$	$\pm 8.27$	$\pm 7.22$	$\pm 9.07$	$\pm 7.59$	$\pm 6.71$
		0.00	0.00	0.00	0.00	0.03	<b>0.11</b>	0.00	<b>0.06</b>	0.00	<b>1.00</b>
USPS49	23.51	15.88	63.14	72.03	68.12	62.27	68.10	71.73	72.28	70.43	<b>75.88</b>
		$\pm 17.98$	$\pm 1.05$	$\pm 5.56$	$\pm 4.46$	$\pm 9.62$	$\pm 16.66$	$\pm 7.23$	$\pm 2.24$	$\pm 2.57$	$\pm 2.28$
		0.00	0.00	0.03	0.00	0.00	0.00	0.00	0.00	0.00	<b>1.00</b>
UMIST	64.15	55.57	63.46	63.42	65.19	71.19	64.90	68.19	71.33	69.16	<b>72.39</b>
		$\pm 2.32$	$\pm 4.93$	$\pm 1.42$	$\pm 2.96$	$\pm 2.77$	$\pm 3.06$	$\pm 2.61$	$\pm 2.06$	$\pm 0.97$	$\pm 2.39$
		0.00	0.00	0.00	0.00	0.01	0.00	0.00	0.00	0.00	<b>1.00</b>
JAFPE	81.52	77.28	79.04	66.97	84.25	82.53	80.01	82.00	85.23	83.20	<b>86.42</b>
		$\pm 8.98$	$\pm 5.88$	$\pm 3.47$	$\pm 1.74$	$\pm 3.49$	$\pm 3.06$	$\pm 3.56$	$\pm 3.31$	$\pm 3.17$	$\pm 3.34$
		0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	<b>1.00</b>
AR	65.48	63.59	66.41	69.01	67.49	67.89	66.94	69.54	69.02	69.44	<b>70.78</b>
		$\pm 2.36$	$\pm 0.85$	$\pm 1.45$	$\pm 0.27$	$\pm 0.89$	$\pm 1.11$	$\pm 1.10$	$\pm 1.32$	$\pm 0.84$	$\pm 0.63$
		0.00	0.00	0.01	0.00	0.00	0.00	0.01	0.00	0.00	<b>1.00</b>
COIL	75.58	62.21	66.19	64.04	44.27	56.29	69.91	70.54	71.37	69.89	<b>72.93</b>
		$\pm 4.98$	$\pm 6.78$	$\pm 4.34$	$\pm 12.61$	$\pm 6.91$	$\pm 4.38$	$\pm 4.48$	$\pm 4.97$	$\pm 4.00$	$\pm 4.44$
		0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	<b>1.00</b>
LUNG	60.37	50.14	55.68	60.12	54.88	60.57	61.75	65.47	63.54	63.50	<b>66.78</b>
		$\pm 4.13$	$\pm 2.31$	$\pm 4.65$	$\pm 4.21$	$\pm 1.54$	$\pm 3.32$	$\pm 1.87$	$\pm 2.94$	$\pm 2.99$	$\pm 1.72$
		0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	<b>1.00</b>
TOX	15.87	10.92	16.53	9.68	22.16	9.07	10.13	25.79	17.46	23.49	<b>27.37</b>
		$\pm 0.68$	$\pm 2.68$	$\pm 0.75$	$\pm 1.36$	$\pm 1.87$	$\pm 1.03$	$\pm 1.60$	$\pm 3.36$	$\pm 2.77$	$\pm 1.62$
		0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	<b>1.00</b>
Average	57.10	48.67	58.14	57.26	57.56	59.35	60.83	64.65	64.34	63.55	67.41



# Experiments

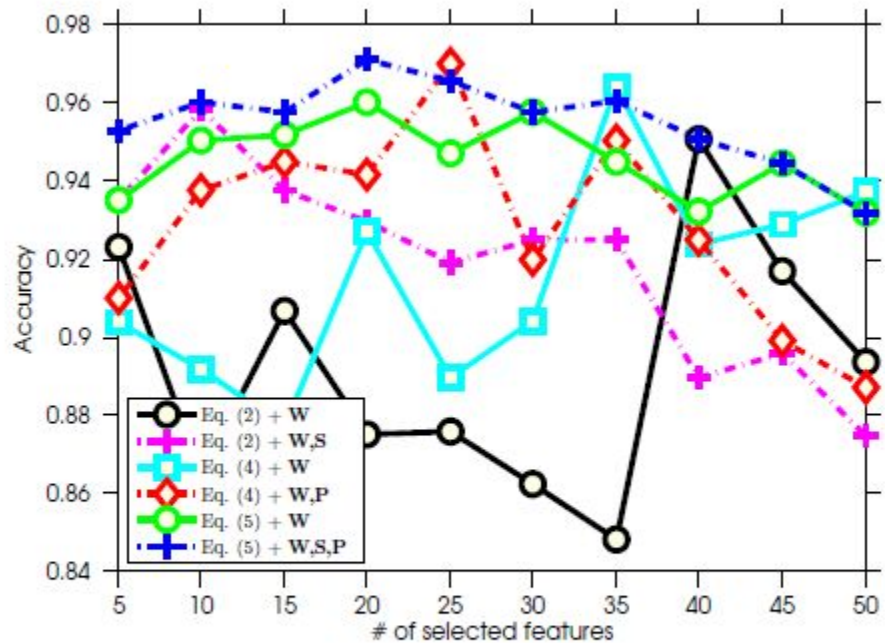


Figure 3: Clustering accuracy w.r.t. 6 different settings of FSASL on USPS200.

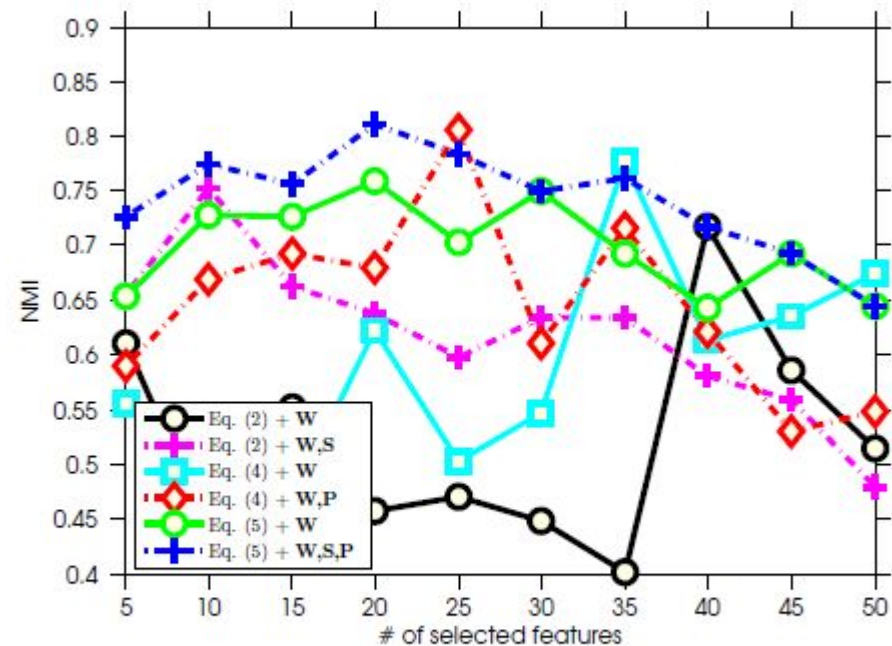


Figure 4: Clustering NMI w.r.t. 6 different settings of FSASL on USPS200.



# Q&A

