

Unsupervised Feature Selection with Adaptive Structure Learning

Liang Du and Yi-Dong Shen

Institute of Software, Chinese Academy of Sciences

10-13 August 2015, SIGKDD, Sydney Australian

Outline

- 1 Introduction and Related Work
- 2 Our Method
 - Proposed Formulation
 - Algorithm
- 3 Experimental Results

Feature Selection

- Data with high dimensionality are often encountered in many real applications.
 - Features are often correlated or redundant, or sometimes noisy.
- Such high dimensionality presents great challenges to learning methods.
 - The curse of dimensionality, computation and storage cost.
- Feature selection techniques can be used to effectively keep few informative features.
 - Supervised methods select those features to respect the the label information.
 - Unsupervised methods select those features to preserve the underlying structure of data.

Unsupervised Filter Feature Selection Methods

The basic procedure

- 1 It first employs *all the input features* to characterize the underlying structure of data,
 - e.g., the pairwise similarity, the graph Laplacian.
- 2 It then selects features to preserve such structures based on certain evaluation criteria.

Typical methods

- Laplacian Score [NIPS, 2006], SPEC [ICML, 2007], EVSC [ICML, 2011]

Unsupervised Embedded Feature Selection Methods (I)

The basic procedure

- ① It first employs *all the input features* to characterize the underlying structure of data.
- ② It then *jointly* selects features to preserve such structure.

Typical methods

- TraceRatio [IJCAI, 2008], UDFS [IJCAI, 2011]

Unsupervised Embedded Feature Selection Methods (II)

The basic procedure

- 1 It first employs *all the input features* to characterize the underlying structure of data.
- 2 It then flats the cluster structure via graph embedding or other clustering module.
 - an *intermediate cluster analysis sub-step* is involved.
- 3 It selects those features that are best aligned to the embedding via sparse spectral regression.

Typical methods

- MCFS [KDD, 2010], MRSF [IJCAI, 2011], FSSL [IJCAI, 2011], SPFS [TKDE, 2013], GLSPFS [TNNLS, 2014]

Unsupervised Embedded Feature Selection Methods (III)

The basic procedure

- ① It first employs *all the input features* to characterize the underlying structure of data.
- ② It then flats the cluster structure via graph embedding or other clustering module.
- ③ It selects those features that are best aligned to the embedding via sparse spectral regression.
- ④ The selected features are used to iteratively improve the intermediate cluster analysis sub-step.
 - the *intermediate cluster analysis* is updated with selected features.

Typical methods

- JELSR [IJCAI, 2011; TC, 2014], NDFS[AAAI, 2012], RUFS [IJCAI, 2013], CGSSL [TKDE, 2014], RSFS [ICDM, 2014]

Unsupervised Embedded Feature Selection Methods (IV)

The basic procedure

- ① It first employs *all the input features* to characterize the underlying structure of data.
- ② It then flats the cluster structure via graph embedding or other clustering module.
- ③ It selects those features that are best aligned to the embedding via sparse spectral regression.
- ④ The selected features are used to iteratively re-capture the underlying structure of data.
 - the *underlying structure of data* is re-captured with selected features.

Typical method

- LLCFS [PAMI, 2011]
 - It actually optimizes two different objectives for structure learning and feature selection. Its theoretic convergence can not be guaranteed and empirically performance is poor.

Limitation of existing unsupervised FS algorithms

- For the problem of unsupervised feature selection, we have to face the *chicken-and-egg dilemma* between **structure characterization** and **feature learning**.
 - on the one hand, one need the true structures of data to identify the informative features.
 - on the other hand, one need the informative features to accurately estimate the true structures of data.



- Most existing unsupervised feature selection methods failed to accurately estimate the structure of data **only with the informative features**.

Outline

- 1 Introduction and Related Work
- 2 Our Method
 - Proposed Formulation
 - Algorithm
- 3 Experimental Results

Basic Idea

- **Structure characterization**, we extract both the global and local structure of data.
 - extract the global structure via sparse representation.
 - extract the local structure with probabilistic neighborhood.
- **Unsupervised Feature Learning**, we use these structures to guide the search of relevant features.
 - flat the structures of data via graph embedding.
 - estimate the informative features via sparse spectral regression.
- Perform **Structure characterization** and **Unsupervised Feature Learning** in a unified framework.
 - The structures are adaptively refined according to the results of feature selection.
 - Better structure characterization often leads to select better features.

Adaptive Global Structure Learning

- Global structure learning via sparse representation

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

We use the sparse reconstruction coefficients \mathbf{S} to extract the global structure of data.

- Adaptive global structure Learning with selected features

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

- The selected features \mathbf{W} should preserve the global structure captured by \mathbf{S} .
- The global structure \mathbf{S} can be refined with the selected informative features \mathbf{W} .

Adaptive Local Structure Learning

- Local structure learning via probabilistic neighborhood

$$\min_{\mathbf{P}} \sum_{i,j} (\|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \mathbf{P}_{ij} + \mu \mathbf{P}_{ij}^2), \text{ s.t. } \mathbf{P} \mathbf{1}_n = \mathbf{1}_n, \mathbf{P} \geq 0 \quad (3)$$

We use the sparse reconstruction coefficients \mathbf{P} to extract the local structure of data.

- Adaptive local structure learning with selected features

$$\min_{\mathbf{P}, \mathbf{W}} \sum_{i,j}^n (\|\mathbf{W}^T \mathbf{x}_i - \mathbf{W}^T \mathbf{x}_j\|_2^2 \mathbf{P}_{ij} + \mu \mathbf{P}_{ij}^2) + \gamma \|\mathbf{W}\|_{21} \quad (4)$$

$$\text{s.t. } \mathbf{P} \mathbf{1}_n = \mathbf{1}_n, \mathbf{P} \geq 0, \mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W} = \mathbf{I}$$

- The selected features \mathbf{W} should preserve the local structure captured by \mathbf{P} .
- The local structure \mathbf{P} can be refined with the selected informative features \mathbf{W} .

Unsupervised Feature Selection with Adaptive Structure Learning (FSASL)

- The formulation of FSASL

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

- The benefits of FSASL
 - Given \mathbf{S} and \mathbf{P} , FSASL selects those features to well respect both the global and local structure of data.
 - Given \mathbf{W} , FSASL estimates the global and local structure of data in a transformed space, i.e. $\mathbf{W}^T \mathbf{X}$, where the adverse effect of noisy features is largely alleviated by sparse regularization.

The Optimization Algorithm

We derive an alternative iterative algorithm to solve the problem.

- ① Given \mathbf{W} and \mathbf{P} , optimize it w.r.t \mathbf{S} .
- ② Given \mathbf{W} and \mathbf{S} , optimize it w.r.t \mathbf{P} .
- ③ Given \mathbf{S} and \mathbf{P} , optimize it w.r.t \mathbf{W} .

We repeat the following steps until it converges.

- Given \mathbf{W} and \mathbf{P} , the optimal value of \mathbf{S} can be obtained by solving the following LASSO problem.

$$\min_{\mathbf{s}_i} \quad ||\mathbf{x}'_i - \mathbf{X}' \mathbf{s}_i||^2 + \alpha |\mathbf{s}_i|, \quad \text{s.t.} \quad \mathbf{S}_{ii} = 0 \quad (5)$$

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

The Optimization Algorithm

- Given \mathbf{W} and \mathbf{S} , we have to solve

$$\begin{aligned} \min_{\mathbf{P}_i^T} \quad & \sum_{j=1}^n \|\mathbf{x}'_i - \mathbf{x}'_j\|^2 \mathbf{P}_{ij} + \mu \|\mathbf{P}_{ij}\|^2, \\ \text{s.t.} \quad & \mathbf{1}_n^T \mathbf{p}_i = 1, \mathbf{P}_{ij} \geq 0 \end{aligned} \quad (6)$$

which can be reformulated as the euclidean projection of a vector onto the probability simplex, where the optimal value of \mathbf{P} can be efficiently obtained without iterations.

- It should be pointed out that the regularization term μ can be empirically determined by the neighborhood size k , which is more intuitive and easy to tune.

The Optimization Algorithm

- Given \mathbf{S} and \mathbf{P} , we have

$$\begin{aligned} \min_{\mathbf{W}} \quad & \text{Tr}(\mathbf{W}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{W}) + \gamma \|\mathbf{W}\|_{21} \\ \text{s.t.} \quad & \mathbf{W}^T \mathbf{X} \mathbf{X}^T \mathbf{W} = \mathbf{I} \end{aligned} \quad (7)$$

where $\mathbf{L}_S = (\mathbf{I} - \mathbf{S})(\mathbf{I} - \mathbf{S})^T$, $\mathbf{L}_P = \mathbf{D}_P - (\mathbf{P} + \mathbf{P}^T)/2$ and let $\mathbf{L} = \mathbf{L}_S + \beta \mathbf{L}_P$.

- Instead of solving a generalized eigen-problem, we solve it with the following two-steps:
 - 1 Solve the eigen-problem $\mathbf{L}\mathbf{Y} = \Lambda\mathbf{Y}$ to get \mathbf{Y} corresponding to the c smallest eigenvalues;
 - 2 Find \mathbf{W} which satisfies $\mathbf{X}^T \mathbf{W} = \mathbf{Y}$ by solving the following optimization problem:

$$\min_{\mathbf{W}} \quad \|\mathbf{Y} - \mathbf{X}^T \mathbf{W}\|^2 + \gamma \|\mathbf{W}\|_{21} \quad (8)$$

The Optimization Algorithm

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. (5);

For each i , update the i -th row of \mathbf{P} using the algorithm for the euclidean projection of a vector onto the probability simplex;

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

Compute \mathbf{W} according to Eq. (8);

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.

Outline

- 1 Introduction and Related Work
- 2 Our Method
 - Proposed Formulation
 - Algorithm
- 3 Experimental Results

Datasets

Table: Summary of the benchmark data sets and the number of selected features

Data Sets	sample	feature	class	selected features
MFEA	2000	240	10	[5, 10, ..., 50]
USPS49	1673	256	2	[5, 10, ..., 50]
UMIST	575	644	20	[5, 10, ..., 50]
JAFFE	213	676	10	[5, 10, ..., 50]
AR	840	768	120	[5, 10, ..., 50]
COIL	1440	1024	20	[5, 10, ..., 50]
LUNG	203	3312	5	[10, 20, ..., 150]
TOX	171	5748	4	[10, 20, ..., 150]

Compared Algorithms

- LapScore [NIPS, 2006]
- MCFS [KDD, 2010]
- LLCFS [PAMI, 2011]
- UDFS [IJCAI, 2011]
- NDFS [AAAI, 2012]
- SPFS [TKDE, 2013]
- RUFS [IJCAI, 2013]
- JELSR [TC, 2014]
- GLSPFS [TNNLS, 2014]
- FSASL

code: <https://github.com/csliangdu/FSASL>

Experiment Protocol

- With the selected features, we evaluate the performance in terms of k-means clustering measured by Accuracy (ACC) and Normalized Mutual Information (NMI).
- For the compared methods, we tune the parameters in relative-large ranges and record the best result according to the grid-search strategy.
- We report the clustering results aggregated from different number of selected features with significance test.

Experimental Results

Table: Aggregated clustering results measured by Accuracy (%) of the compared methods.

Data Sets	AllFea	LapScore	MCFS	LLCFS	UDFS	NDFS	SPFS	RUFS	JELSR	GLSPFS	FSASL
MFEA	68.73	51.78	51.04	60.38	48.94	67.13	68.20	64.58	67.01	61.00	69.94
		± 5.51	± 8.13	± 8.58	± 3.32	± 7.53	± 9.43	± 7.99	± 8.37	± 8.70	± 7.19
USPS49	77.70	0.00	0.00	0.00	0.00	0.01	0.22	0.00	0.01	0.00	1.00
		69.21	53.74	94.96	94.05	68.12	83.43	85.86	95.16	94.75	95.95
UMIST	42.40	± 8.95	± 3.50	± 1.44	± 1.13	± 8.18	± 6.66	± 2.58	± 0.55	± 0.61	± 0.48
		0.00	0.00	0.03	0.00	0.00	0.00	0.00	0.00	0.00	1.00
JAFPE	71.57	36.73	44.46	47.31	48.04	52.80	46.72	50.87	53.52	50.53	54.92
		± 1.18	± 3.26	± 0.83	± 1.92	± 2.26	± 1.70	± 1.95	± 1.54	± 0.59	± 1.89
AR	30.26	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.01	0.00	1.00
		67.62	73.56	64.79	75.48	74.98	73.93	75.75	77.77	75.46	79.29
COIL	59.17	± 8.49	± 4.83	± 4.08	± 1.63	± 2.15	± 2.85	± 2.53	± 1.87	± 1.61	± 2.24
		0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00
LUNG	72.46	25.29	29.05	34.22	30.87	32.34	31.06	34.84	34.19	34.12	36.11
		± 2.89	± 1.19	± 2.70	± 0.35	± 1.52	± 2.14	± 1.90	± 2.52	± 1.60	± 0.75
TOX	43.65	0.00	0.00	0.05	0.00	0.00	0.00	0.04	0.02	0.00	1.00
		45.60	51.50	50.84	31.40	44.22	56.94	59.20	59.53	57.96	60.93
Average	58.24	± 6.16	± 5.38	± 3.76	± 16.89	± 6.33	± 3.43	± 3.28	± 4.01	± 2.27	± 2.50
		0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.03	0.00	1.00
TOX	43.65	58.97	70.42	71.58	65.46	75.52	73.49	77.35	77.86	77.83	81.93
		± 5.24	± 3.41	± 5.85	± 3.88	± 1.57	± 3.43	± 2.62	± 3.12	± 2.70	± 1.63
Average	58.24	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00
		40.25	43.10	39.28	47.14	38.28	39.93	47.67	43.96	47.38	49.17
Average	58.24	± 0.65	± 1.86	± 0.49	± 0.75	± 1.64	± 1.13	± 0.83	± 1.56	± 1.93	± 0.67
		0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00
Average	58.24	49.43	52.11	57.92	55.17	56.67	59.21	62.02	63.63	62.38	66.03

Experimental Results

Table: 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	49.19	64.97	64.92	63.98	64.51	59.26	66.70
		± 4.77 0.00	± 9.14 0.00	± 9.76 0.00	± 3.83 0.00	± 7.54 0.03	± 8.27 0.11	± 7.22 0.00	± 9.07 0.06	± 7.59 0.00	± 6.71 1.00
USPS49	23.51	15.88	4.60	72.03	68.12	12.27	38.10	41.73	72.28	70.43	75.88
		± 17.98 0.00	± 2.57 0.00	± 5.56 0.03	± 4.46 0.00	± 9.62 0.00	± 16.66 0.00	± 7.23 0.00	± 2.24 0.00	± 2.57 0.00	± 2.28 1.00
UMIST	64.15	55.57	63.46	63.42	65.19	71.19	64.90	68.19	71.33	69.16	72.39
		± 2.32 0.00	± 4.93 0.00	± 1.42 0.00	± 2.96 0.00	± 2.77 0.01	± 3.06 0.00	± 2.61 0.00	± 2.06 0.00	± 0.97 0.00	± 2.39 1.00
JAFPE	81.52	77.28	79.04	66.97	84.25	82.53	80.01	82.00	85.23	83.20	86.42
		± 8.98 0.00	± 5.88 0.00	± 3.47 0.00	± 1.74 0.00	± 3.49 0.00	± 3.06 0.00	± 3.56 0.00	± 3.31 0.00	± 3.17 0.00	± 3.34 1.00
AR	65.48	63.59	66.41	69.01	67.49	67.89	66.94	69.54	69.02	69.44	70.78
		± 2.36 0.00	± 0.85 0.00	± 1.45 0.01	± 0.27 0.00	± 0.89 0.00	± 1.11 0.00	± 1.10 0.01	± 1.32 0.00	± 0.84 0.00	± 0.63 1.00
COIL	75.58	62.21	66.19	64.04	44.27	56.29	69.91	70.54	71.37	69.89	72.93
		± 4.98 0.00	± 6.78 0.00	± 4.34 0.00	± 12.61 0.00	± 6.91 0.00	± 4.38 0.00	± 4.48 0.00	± 4.97 0.00	± 4.00 0.00	± 4.44 1.00
LUNG	60.37	50.14	55.68	60.12	54.88	60.57	61.75	65.47	63.54	63.50	66.78
		± 4.13 0.00	± 2.31 0.00	± 4.65 0.00	± 4.21 0.00	± 1.54 0.00	± 3.32 0.00	± 1.87 0.00	± 2.94 0.00	± 2.99 0.00	± 1.72 1.00
TOX	15.87	10.92	16.53	9.68	22.16	9.07	10.13	23.58	17.46	23.49	25.79
		± 0.68 0.00	± 2.68 0.00	± 0.75 0.00	± 1.36 0.00	± 1.87 0.00	± 1.03 0.00	± 1.60 0.00	± 3.36 0.00	± 2.77 0.00	± 1.62 1.00
Average	57.10	48.67	50.83	57.26	56.94	53.10	57.08	60.63	64.34	63.55	67.21

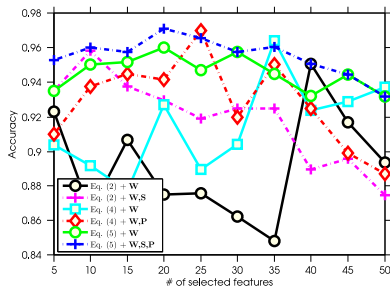
The Effect of Adaptive Structure Learning

Question: Does the adaptive structure learning lead to select more informative features?

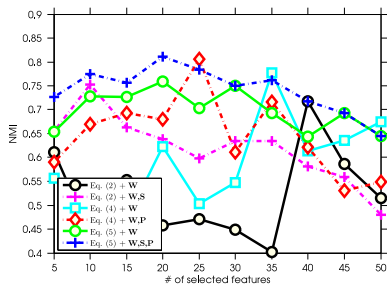
We design six different settings to empirically investigate the effect of adaptive structure learning.

- Global structure guided FS, with/without adaptive structure learning
- Local structure guided FS, with/without adaptive structure learning
- Global and local structures guided FS, with/without adaptive structure learning

The Effect of Adaptive Structure Learning



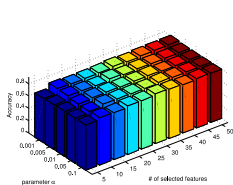
(a)



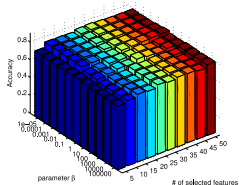
(b)

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

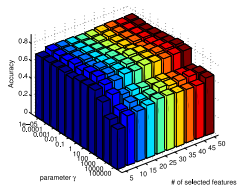
Parameter Sensitivity



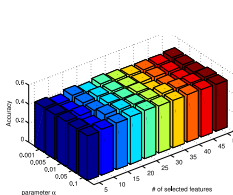
(a)



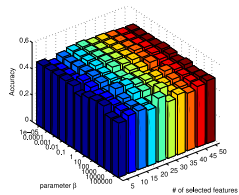
(b)



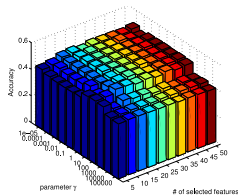
(c)



(d)



(e)



(f)

Figure: Accuracy of different parameters on JAFFE (a-c) and TOX (d-f).

Summary

- We investigate most of existing unsupervised embedded methods and further classify them into four closely related but different types. These analyses provide more insight into what should be further emphasized on the development of more essential unsupervised feature selection algorithm.
- We propose a novel unified learning framework, called unsupervised Feature Selection with Adaptive Structure Learning (FSASL in short), to fulfil the gap between two essential sub tasks, i.e. structure learning and feature learning. In this way, these two tasks can be mutually improved.
- Comprehensive experiments on benchmark data sets show that our method achieves statistically significant improvement over state-of-the-art feature selection methods.

Acknowledgement

We would like to thank Prof. Feiping Nie and Prof. Mingyu Fan for their helpful suggestions to improve this paper.

This work is supported in part by the China National 973 program 2014CB340301, the Natural Science Foundation of China (NSFC) grant 61379043, 61322211.

Thanks!

Q&A