

Gilles Bisson<sup>1,2</sup>, Fawad Hussain<sup>2</sup>

<sup>1</sup> Centre National de la Recherche Scientifique, France

<sup>2</sup> Laboratoire TIMC-IMAG, CNRS/UJF 5525

Université de Grenoble, France

## Co-clustering problem

- Simultaneously cluster rows (instances) and columns (features)
  - Allows to deal with sparse data
  - Improves clustering result
- Old Idea (Hartigan [1], etc) but has drawn recent attention.
  - Text Mining: Document Classification
  - Bio-Informatics : DNA chip analysis
- Widely used techniques :
  - Matrix Decomposition (BVD) [2]
  - Information theoretic (ITCC) [3]
  - Graph based (RSN) [4]

## Proposed method: $\chi$ -Sim

M	m <sup>1</sup>	m <sup>2</sup>	...	m <sup>c</sup>
m <sub>1</sub>	1	1	...	0
m <sub>2</sub>	0	0	...	1
...				
m <sub>r</sub>	0	3	...	0

- Classical similarity as a sum of the similarities between features

$$Sim(\mathbf{m}_i, \mathbf{m}_j) = F_s(m_{i1}, m_{j1}) + F_s(m_{i2}, m_{j2}) + \dots + F_s(m_{ic}, m_{jc})$$

- Extension to all pairs

$$Sim(\mathbf{m}_i, \mathbf{m}_j) = \frac{F_s(m_{i1}, m_{j1}) \cdot sc_{11}}{F_s(m_{i1}, m_{j1}) \cdot sc_{11} + F_s(m_{i1}, m_{j2}) \cdot sc_{12} + \dots + F_s(m_{i1}, m_{jc}) \cdot sc_{1c}} + \frac{F_s(m_{i2}, m_{j1}) \cdot sc_{21}}{F_s(m_{i2}, m_{j1}) \cdot sc_{21} + F_s(m_{i2}, m_{j2}) \cdot sc_{22} + \dots + F_s(m_{i2}, m_{jc}) \cdot sc_{2c}} + \dots + \frac{F_s(m_{ic}, m_{j1}) \cdot sc_{c1}}{F_s(m_{ic}, m_{j1}) \cdot sc_{c1} + F_s(m_{ic}, m_{j2}) \cdot sc_{c2} + \dots + F_s(m_{ic}, m_{jc}) \cdot sc_{cc}}$$

## Normalization

- Normalize by the maximum possible value which is the product of length of document vector or term vector

$$\forall i, j \in 1..r, sr_{ij} = Sim(\mathbf{m}_i, \mathbf{m}_j) / (|\mathbf{m}_i| |\mathbf{m}_j|)$$

- Column similarity is given by

$$\forall i, j \in 1..c, sc_{ij} = Sim(\mathbf{m}_i^t, \mathbf{m}_j^t) / (|\mathbf{m}_i^t| |\mathbf{m}_j^t|)$$

where  $|\mathbf{m}_i|$  = length of document  $i$   
and  $|\mathbf{m}_i^t|$  = number of documents featuring word  $i$

- Normalization insures that the similarity value is in [0,1]

## How to implement efficiently

- When functions  $F_s$  is defined as a product

$$F_s(m_{ij}, m_{kl}) = m_{ij} \cdot m_{kl}$$

- At any given iteration  $t$ ,

$$SR_t = M^t \cdot SC_{(t-1)} \cdot M \times NR \quad \forall i, j \in [1, r], NR_{ij} = 1 / (|\mathbf{m}_i| |\mathbf{m}_j|) \quad \text{Eq (1)}$$

$$SC_t = M \cdot SR_{(t-1)} \cdot M^t \times NR \quad \forall i, j \in [1, r], NC_{ij} = 1 / (|\mathbf{m}_i^t| |\mathbf{m}_j^t|) \quad \text{Eq (2)}$$

" $\times$ " here represents the Hadamard matrix product

- Complexity  $O(N^3)$  like any other classical similarity measure (cosine, Euclidean, etc)

## $\chi$ -Sim Algorithm

### Function $\chi$ -Sim

Input : data matrix  $M$ , no. of iterations  $N_{iter}$

Output : two similarity matrices  $SR$  and  $SC$  expressing the co-similarity between rows and columns of  $M$

Initialize  $SR$  and  $SC$  with the identity matrix

for  $iter = 1$  to  $N_{iter}$

    Update- $SR_{iter}(M, SC_{iter-1})$  // Equation 1

    Update- $SC_{iter}(M, SR_{iter})$  // Equation 2

end for

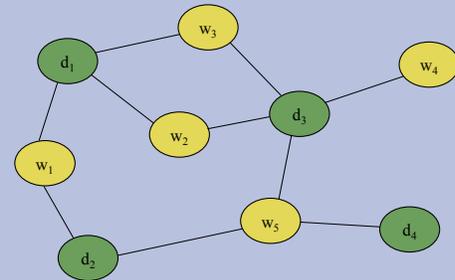
## Experimental Results

- Corpus
  - NewsGroup (M2, M5, M10, NG1, NG2, NG3) and SMART (Classic3)
- Methodology
  - Generate matrices  $SR$  and  $SC$
  - Generate clusters using Hierarchical clustering (ward's linkage)
  - Evaluate result using Precision value (Pr) and Normalized Mutual Information (NMI)

	Runtime		$\chi$ -Sim	Cosine	LSA	ITCC	BVD	RSN
M2	12.01s	Best Pr	0.96	0.72	0.94	0.92	0.95	-
		Pr <sub>(avg)</sub>	0.93 ± 0.01	0.62 ± 0.04	0.78 ± 0.15	0.71 ± 0.15		
		NMI <sub>(avg)</sub>	0.65 ± 0.05	0.15 ± 0.05	0.37 ± 0.23	0.21 ± 0.21		
M5	10.93s	Best Pr	0.98	0.76	0.95	0.56	0.93	-
		Pr <sub>(avg)</sub>	0.94 ± 0.04	0.6 ± 0.08	0.82 ± 0.09	0.48 ± 0.06		
		NMI <sub>(avg)</sub>	0.85 ± 0.07	0.55 ± 0.11	0.65 ± 0.14	0.27 ± 0.09		
M10	13.53s	Best Pr	0.80	0.55	0.67	0.33	0.67	-
		Pr <sub>(avg)</sub>	0.75 ± 0.04	0.45 ± 0.07	0.54 ± 0.09	0.28 ± 0.04		
		NMI <sub>(avg)</sub>	0.65 ± 0.05	0.43 ± 0.07	0.45 ± 0.08	0.16 ± 0.04		
Classic3	164.38s	Pr	0.99	0.89	0.98	0.97	-	-
		NMI	0.96	0.63	0.93	0.90		
		NMI	0.96	0.63	0.93	0.90		
NG1	6.60s	Best Pr	1	0.97	0.98	0.83	-	-
		Pr <sub>(avg)</sub>	1 ± 0	0.90 ± 0.11	0.95 ± 0.02	0.67 ± 0.12		
		NMI <sub>(avg)</sub>	1 ± 0	0.62 ± 0.20	0.75 ± 0.08	0.13 ± 0.13	0.64 ± 0.16	
NG2	14.56s	Best Pr	0.94	0.72	0.87	0.72	-	-
		Pr <sub>(avg)</sub>	0.93 ± 0.01	0.60 ± 0.08	0.82 ± 0.03	0.57 ± 0.08		
		NMI <sub>(avg)</sub>	0.81 ± 0.04	0.51 ± 0.07	0.64 ± 0.03	0.35 ± 0.09	0.75 ± 0.07	
NG3	20.20s	Best Pr	0.91	0.62	0.80	0.65	-	-
		Pr <sub>(avg)</sub>	0.85 ± 0.06	0.60 ± 0.04	0.72 ± 0.05	0.55 ± 0.06		
		NMI <sub>(avg)</sub>	0.77 ± 0.02	0.55 ± 0.03	0.59 ± 0.03	0.48 ± 0.05	0.70 ± 0.04	

## Graphical Meaning of an iteration

- Matrix  $M$  can be seen as a bipartite graph
- 3-4 iterations are enough [5]
- First Iteration: direct occurrences
- Second Iteration: Second order occurrences
- ...



## Conclusion

- New approach to use co-similarity for co-clustering
- Allows to use various similarity based clustering methods
- No need to define number of word clusters
- Good results even on smaller number of examples

## References

- Hartigan J.A. Direct clustering of a data matrix. Journal of American Statistical Association, 67(337):123-129, 1972
- Long B., Zhongfei Z. and Yu P. S. Co-clustering by Block Value Decomposition. In SIGKDD, august 21-24, 2005.
- Dhillon I. S., Mallela S. and Modha D. S. Information Theoretic Co-clustering. Proceedings of The 9th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Aug 24-27, 2003, pp 89-98
- Long Bo, Wu Xiaoyun, Zhang Zhongfei and Yu Phillip: Unsupervised Learning on K-partite Graphs, Proceedings of SIGKDD August 20-23, 2006.
- Kontostathis A., Pottenger W.M. A framework for understanding LSI performance. Information Processing and Management. Volume 42, number 1, 2006, pp 56-73.