Applications of Dominant Set

Sebastiano Vascon, PhD

DAIS 09/05/2017







Recap on the Dominant Set technique

- Graph-based clustering technique
- A DS is subset of highly coherent nodes in a graph (high internal similarity and high external dissimilarity).
- Maximal clique in edge weighted graph
- Pros:
 - No need for k
 - Provide a quality value for each cluster (cohesiveness)
 - Provide a membership value for each element in a cluster
 - Undirected and directed graph
- Cons:
 - Require $O(n^2)$ to store the similarity matrix (does not scale for big data)

Recap on the Dominant Set technique

- Given an edge-weighted graph G=(V,E,w) with no self loop
- A DS is found optimizing the following problem (1):

 $\max x' A x$
s.t. $x \in \Delta^n$

where A is the affinity (similarity) matrix of G and x is a probability distribution over V (usually set as a uniform distribution).

- Solution to (1) can be found with dynamical systems like:
 - Replicator Dynamics [1]
 - Exponential Rep Dynamics [1]
 - Infection Immunization [2]

Recap on the Dominant Set technique

A dataset is modeled as a weighted graph $G = (V, E, \omega)$ with no self loop. The set of nodes V are the dataset's items and the edges are weighted by $\omega: V \times V \to \mathbb{R}_+$ that quantifies the pairwise similarity of the items. G is thus represented by an $n \times n$ adjacency matrix $A = (a_{ij})$



http://www.github.com/xwasco/DominantSetLibrary



Pattern Recognition



Human Behavior



Brain Connectomics



Nano science





Pattern Recognition



Human Behavior



Brain Connectomics



Nano science



Problem: Understanding the activity of Gephyrine and vGAT proteins.

Gephyrine and vGAT are two proteins that takes parts into the synapse activation.

Gephyrine is a post-synaptic protein that sustain the grid of GABA receptors that receive the chemical stimuli in a synapse.

Analyze the morphological changes of this grid during the synapses activation is of crucial importance for the Nanophysicists (e.g. discovering disease). These changes is reflected into the morphology and number of clusters of Gephyrine.



Finding an alignment with the v-GAT pre-synaptic protein clusters is important to understand when and where an accumulation of Gephyrine occurs.

Dataset: set of molecules position (x,y) for each channel (Gephyrine and vGAT)



(x,y) locations of each molecule Gephyrine

vGAT

8

Aim:

- 1. Extract clusters of Gephyrine and vGAT based on the single molecules detection
- 2. Find associations between clusters of the two channel

Solution:

- 1. Create a graph-based representation of the points for each channel G(V,E,w) in which $w_{ij} = \begin{cases} e^{-\frac{||i-j||}{2\sigma^2}} & \text{if } i \neq j \\ 0 & \text{otherwise} \end{cases}$ and extract the clusters using the DS
- 2. Apply a chain of post processing filtering to merge the smaller clusters and remove the meaningless ones.
- 3. Find clusters associations between the two channels providing statistics

Pipeline:



Pipeline:



Pipeline:



Remove clusters having a cohesiveness ($x^T A x$) values lower than a certain threshold θ .

This remove clusters with few and spread points.



DS find circular and compact clusters ... it is ok but ?

We merge clusters having the centroid (mean points) closer to a certain threshold or if their convex hull overlap for a certain %





Cluster Association between channel

 $1NN \le d_{RG}$



DS find circular and compact clusters ... it is ok but ?

We merge clusters having the centroid (mean points) closer to a certain threshold or if their convex hull overlap for a certain %





Cluster Association between channel







Pipeline:



Pipeline:

- Evaluate pairwise distances between green and red clusters centroid
- 2. For each green cluster assign the 1-NN red cluster

Cluster Association between channel

 $1NN \le d_{RG}$





Cluster statistics for Gephyrine's clusters:

- Number of points
- Convex Hull area
- Variance
- Distance of the closest vGAT's cluster

Cluster statistics for vGAT's clusters:

- Number of points
- Convex Hull area
- Variance
- Number of associated Gephyrine's cluster

Validation

- Nanophysicists annotate a set of images
- Completeness/Correctness



Pattern Recognition



Human Behavior



Brain Connectomics



Nano science





Pattern Recognition



Human Behavior



Brain Connectomics



Nano science



 k-NN classifier: Assign the class based on classes of the k nearest sample in the feature space.



- Problems of k-NN classifiers:
 - Sensitive to noise and outliers
 - Slow if the number of elements is high
- Solution:
 - Reducing the space of search by using prototypes
 - Create/select prototypes such that the noise and outliers are minimized.



- Given a dataset the DS are used to extract the cluster and the centroid.
- The k-NN classification is performed on the prototypes and not on the entire set



Given a data set the DS are used to extract the cluster and the centre of the centre of

- Given a datase and the cen
- The k-NN claprototypes a

DS are used to extract the cluster

performed on the entire set

- Given a dataset the DS are used to extract the cluster and the centroid.
- The k-NN classification is performed on the prototypes and not on the entire set

- Given a dataset the DS are used to extract the cluster and the centroid.
- The k-NN classification is performed on the prototypes and not on the entire set

Name	#ex	#att	Name	#ex	#att	Name	#ex	#att
appendicitis	106	7	haberman	306	3	sonar	208	60
$\operatorname{australian}$	690	14	heart	270	13	spambase	4597	57
banana	5300	2	mammographic	961	5	spectfheart	267	44
bands	539	19	monk-2	432	6	titanic	2201	3
bupa	345	6	pima	768	8	wisconsin	699	9

- 15 binary classification datasets from UCI
- 25 different prototype methods
- 1 common benchmark [1]

28

- Accuracy, Compression rate and Exec. Time
- Evaluation of 1-NN and 3-NN performances

Rank	1-NN	Classifier	Rank	3-NN (Classifier	Ran	k PS exec	time (sec)
#	A×K×R	Name	#	A×K×R	Name	#	Time	Name
1	$0,\!428$	'CHC'	1	$0,\!386$	'GGA'	14	$10,\!073$	'DROP3'
2	0,393	'SSMA'	2	$0,\!382$	'RMHC'	15	$52,\!625$	\mathbf{MAX}
3	0,393	'GGA'	3	$0,\!344$	'SSMA'	16	$52,\!625$	MAXCO
4	0,369	'RMHC'	4	$0,\!334$	\mathbf{AVG}	17	$52,\!625$	\mathbf{AVG}
5	$0,\!349$	'RNN'	5	0,328	'CHC'	18	$52,\!625$	WAVG
6	$0,\!340$	\mathbf{AVG}	6	$0,\!322$	\mathbf{MAX}	19	80,508	'RNG'
7	0,332	WAVG	7	0,306	MAXCO	20	$127,\!942$	'Recons.'
8	0,329	MAXCO	8	0,284	'RNN'	21	372,360	'CHC'
9	0,328	\mathbf{MAX}	9	0,268	WAVG	22	391,666	'RMHC'
10	0,296	'CCIS'	10	0,244	'CCIS'	23	$513,\!965$	'SSMA'
11	$0,\!243$	'MCNN'	11	0,239	'DROP3'	24	874,246	'RNN'
15	$0,\!238$	'CPruner'	12	$0,\!224$	'HMNEI'	25	$1525,\!734$	'GGA'
Table	2: 1-NN	classifier	Table	e 3: 3-NN	classifier	Tabl	e 4: Time	needed for
rankir	ıg.		ranki	ng.		the p	rototypes	extraction.

- Method strengthens:
 - <u>Compression rate is around 90%</u>
 - Good balance between accuracy, compression rate and exec time.
 - <u>Time is an order of magnitude faster than the best competitors.</u>
 - Method weakness:
 - Does not scale due to the quadratic requirement of the DS
 - Future work:
 - Extend the approach to handle multiple classes
- Publications:
 - S Vascon, M Cristani, M Pelillo, V Murino Using Dominant Sets for k-NN Prototype Selection - International Conference on Image Analysis and Processing (ICIAP) 2013

Pattern Recognition

Human Behavior

Brain Connectomics

Nano science

Pattern Recognition

Human Behavior

Brain Connectomics

Nano science

Brain Connectomics: White matter multi-subject clustering

- White matter (WM) is a component of the central nervous system and consists mostly of cells that transmit signals from one region of the cerebrum to another.
- The studies of WM fibers organization is important in the diagnosis of diseases like Alzheimer or Multiple Sclerosis.
- The problem:
 - Simplify the complexity
 - Find common structure across subjects

- Why ?
 - Neuroscientist needs an higher level of abstraction (manual investigation is prone to human error)
 - Automatic tool for white matter investigation and brain parcellation
 - Data-driven atlas of the brain avoiding neuroscientists bias

Brain Connectomics: multi-subject clustering

- Other methods:
 - Hierarchical clustering [1]
 - Spectral clustering [2]
 - Stochastic processes [3]
- Problems:
 - Need an a-priori number of cluster
 - Need an a-priori level of the hierarchy
- Solution, a three steps algorithm:
 - 1. Reduce the complexity through brain abstraction
 - 2. Project the subject to a common space (landmark space)
 - 3. Performing a cross-subject clustering identifying the commonalities

[1] Guevara et al. Automatic fiber bundle segmentation in massive tractography datasets using a multisubject bundle atlas. NI2012
 [2] O'Donnell et al. Automatic tractography segmentation using a high-dimensional white matter atlas. IEEET.Med.Img 2007
 [3] Wang et al. Tractography segmentation using a hierarchical dirichlet processes mixture model.Neuroimage 2011

Brain Connectomics: fiber bundle extraction

Brain Connectomics: fiber bundle ext

100000 fibers

Tractography

Brain Connectom fiber bundle extra 200 bundles Intra-Subject Input Data Left Hemisphere Fibers Sluster Selection Tractography Introid Selected for Cross-subject

Brain Connectomics: multi-subject clustering

Brain Connectomics: Landmark Space multi-subject clustering Landmark Encoding Landmark **Atlas Registration** Landmark extraction Landmark each Diffusion from Cortical Regions bace #N Mouse Bra Ν

Brain Connectomics: multi-subject clustering

Close, T. G et al. A software tool to generate simulated white matter structures for the assessment of fibretracking algorithms. Neuroimage 2009

Brain connectomics: Conclusions

- Method strengthens:
 - Automatically extract the bundles (No prior on the number of clusters)
 - No need to register the tractography
 - Designed to cope with large set of subject (thanks to the first reduction)
 - Good performances if compared with both supervised and unsupervised clustering techniques.
- Method weakness:
 - The method due to the quadratic complexity cannot scale to bigger dataset.
- Future Work:
 - Applied to human data-sets to build an atlas of WM bundles for clinical applications
- Publications:
 - L Dodero, S Vascon, L Giancardo, A Gozzi, D Sona, V.Murino. Automatic white matter fiber clustering using dominant sets - Pattern Recognition in Neuroimaging (PRNI), 2013
 - S Vascon, L Dodero, V Murino, A Bifone, A Gozzi, D.Sona Automated multi-subject fiber clustering of mouse brain using dominant sets. Fr.NeuroInf 2015

Pattern Recognition

Human Behavior

Brain Connectomics

Nano science

Pattern Recognition

Human Behavior

Brain Connectomics

Nano science

The problem

Tasks:

- *Security & Surveillance*: who were with ?
- Scene understanding: are there any groups of people ?
- Behavior analysis: how we join and leave a group ?
- *Social robotics*: how to interact with humans ?
- •

Problems:

- Unreliability of the detectors
- Grouping in a low density space (few persons per scene)
- Respecting sociological constraints
- Respecting biological constraints

Constraints

a) Circular arrangement

b) Vis-a-vis arrangement

c) L-arrangement

d) Side-by-side arrangement

Sociological constraints:

F-Formation: "whenever two or more individuals in close proximity orient their bodies in such a way that each of them has an easy, direct and equal access to every other participant's transactional segment" [1]

Biological Constraints:

Human field of view is the range [120°- 190°] [2]

[1] Ciolek, T.M., Kendon, A.: Environment and the Spatial Arrangement of Conversational Encounters. Sociological Inquiry 50 (1980)
 [2] I.P. Howard and B.J. Rogers. Binocular Vision and Stereopsis. Oxford psychology series. Oxford University Press, (1995).

[1] Cristani et al: Social interaction discovery by statistical analysis of F-formations. In: Proc. Of BMVC, BMVA Press (2011)

[2] Hung, H., Krose, B.: Detecting F-formations as dominant sets. In: ICMI. (2011)

[3] Setti, F., Lanz, O., Ferrario, R., Murino, V., Cristani, M.: Multi-Scale F-Formation Discovery for Group Detection. In: ICIP. (2013

State of the art

F-Formation detection algorithms:

- Hough voting [1]
 - Samples vote for an o-space
 - O-space with the majority of votes is taken.
- Graph Based[2,3]
 - A scene is represented as a weighted graph G.
 - An F-F is found partitioning the graph (graph-cut, max clique)
- Multi-Scale [4]
 - Based on [2] Hough Voting schema but for different F-F sizes.

(c) Graph Clustering

 Dominant sets using replicator dynamic

Modularity cut

(d) Calculate Affinity

timate of focus orien

Socially motivated

roximity

Group Detection: Our method

- **1.** Probabilistic model of Frustum of Visual Attention
- 2. Quantify interactions in a pairwise matrix using Information-Theoretic measures
- 3. Game-theoretic clustering for finding groups

Our method - 1 Frustum

- A person in a scene is described by his/her position (x,y) and the head orientation ϑ
- The frustum of visual attention is defined by an aperture (160°) and by a length *l* [1].

[1] Vinciarelli et al. Social Signal Processing: Survey of an emerging domain.IJCV 2009.

Our method - 2 Quantify Pairwise Interaction

Given two histogram P and Q their distance is:

Kullback-Leibler divergence (A-Sym)Jensen-Shannon divergence (Sym)
$$KL(P || Q) = \sum_{i=1}^{n} \left(\log(p_i) \frac{p_i}{q_i} \right)$$
 $JS(P,Q) = \frac{KL(P|| M) + KL(Q|| M)}{2}$
 $M = \frac{1}{2}(P+Q)$

• A measure of affinity is obtained through a Gaussian Kernel $a_{P,Q} = exp\left\{-\frac{d(P,Q)}{\sigma}\right\}$ where P,Q are the frustum of two persons, d(...) could be either

KL or JS and σ act as normalization term.

Our method - 2 Quantify Pairwise Interaction

Frame + Frustum

Payoff matrix

Experiments

Dataset	#Sequences	s #Frames	Consecutive	e Automated
		imes seq.	Frames	Tracking
CoffeeBreak	2	45,74	Y	Y
CocktailParty	/ 1	320	Y	Y
GDet	5	132,115,79,17,60	Ν	Y
PosterData	82	1	Ν	Ν
Synth	10	10	Ν	Ν

Evaluation criteria:

57

As in [1] a group is correctly detected if at least $\left\lceil \frac{2}{3} |G| \right\rceil$ of its members matches the ground truth.

Metrics: Precision, Recall, F1-Score (averaged over the frames)

^[1] Setti, F., Hung, H., Cristani, M.: Group Detection in Still Images by F-formation Modeling: a Comparative Study. In: WIAMIS. (2013)

Results

	Coffe	eBreak	(S1+S2)
Method \star	Prec	Rec	F1
IRPM	0.60	0,41	0,49
HFF	0,82	0,83	0,82
DS	0,68	0,65	0,66
MULTISCALE	0,82	0,77	0,80
R-GTCG	0,86	0,88	0,87
	σ	=0.2 , <i>l</i> =	=145

PosterData			Gdet			
Prec	Rec	F1	Prec	Rec	F1	
-	-	-	-	-	-	
0,93	0,96	0,94	0,67	0,57	0,62	
0,93	0,92	0,92	-	-	-	
-	-	-	-	-	-	
0,92	0,96	0,94	0,76	0,76	0,76	
σ =0.25 , <i>l</i> =115			σ=	=0.7 <i>l</i> =1	80	

	Cocktail Party		
Method \star	Prec	Rec	F1
IRPM	-	-	-
HFF	0,59	0,74	0,66
MULTISCALE	0,69	0,74	0,71
R-GTCG	0,87	0,82	0,84
	σ=0.6 , <i>l</i> =170		

	Synth				
Prec	Rec	F1			
0,71	0,54	0,61			
0,73	0,83	0,78			
0,86	0,94	0,90			
1,00	1,00	1,00			
$\sigma = 0.1$, $l = 75$					

Conclusions

- Method strengthens:
 - Based on sociological and biological constraints
 - No assumption on the size or shape of the F-F
 - Designed to cope with very different realistic scenario
 - Work on top of a tracker or person detection algorithms (15-20 fps)
 - State of the art in all publicly available datasets.
 - Comparable performances on non dedicated datasets.
- Method weakness:
 - Pairwise Affinity matrix does not scale on thousands of detections per frame (unlikely situation)
- Future work:
 - Group tracking
- Publications:
 - A Game-Theoretic Probabilistic Approach for Detecting Conversational Groups.
 S Vascon, Z Eyasu, M Cristani, H Hung, M Pelillo, V Murino. Asian Conference in Computer Vision 2014
 - Detecting conversational groups in images and sequences: A robust game-theoretic approach.
 S Vascon, EZ Mequanint, M Cristani, H Hung, M Pelillo, V Murino Computer Vision and Image Understanding 2015
 - Group detection and tracking with sociological features. S.Vascon, L.Bazzani. Book chapter submitted