EXPLORATION OF HIGH-DIMENSIONAL DATA

EXPLAINABLE STRUCTURING DISCOVERY OF RELEVANT CASES

Joris Falip, Frédéric Blanchard, Michel Herbin joris.falip@univ-reims.fr www.datacrunch.sh

March 20 2019

CReSTIC, University of Reims Champagne-Ardenne, France

Help domain experts to visualize, explore and manipulate data

Approach

- Leverage analogical reasoning
- Structure data according to their similarity
- Visually interpretable results

OVERVIEW



Definition

"Classification decisions are based on the similarity of stimuli to the stored exemplars." - R. M. Nosofsky

- Rely on past experiences
- Allows rapid decision making
- Can be applied to complex objects

- 1. Explainable algorithm and interpretable results
- 2. Efficient with high-dimensionality

1. Explainable algorithm and interpretable results

2. Efficient with high-dimensionality

Curse of dimensionality

"Having more dimensions usually means data sets tend to be sparse and all distances between data points tend to become harder to distinguish." - N. Tomasev et al.

1. Explainable algorithm and interpretable results

2. Efficient with high-dimensionality

Curse of dimensionality

"Having more dimensions usually means data sets tend to be sparse and all distances between data points tend to become harder to distinguish." - N. Tomasev et al.

No a priori knowledge \rightarrow no dimensionality reduction

Hubness

"High dimensionality causes some data points to be the nearest neighbor of many others points, thus becoming 'hubs'." - totally made up definition

Hubs exhibit interesting properties and can be used to improve traditional data mining approaches

- Compute on each dimension then aggregate
- Use ranks to avoid sparsity and outliers' exclusion

Steps

- 1. Ranking individuals on each dimension
- 2. Aggregate the results into a score : the Degree of Representativeness
- 3. Link individuals to their neighbor with the highest DoR

Dataset			
	Feat. 1	Feat. 2	
Ind1	67	173	
Ind2	53	161	

Dataset			
	Feat. 1	Feat. 2	
Ind1	67	173	
Ind2	53	161	



	Feature	e 2 - rankir	ng	
(Feature 1	- ranking		
	Ind1	Ind2		
				7

Dataset			
	Feat. 1	Feat. 2	
Ind1	67	173	
Ind2	53	161	



	Feature	e 2 - rankir	ng	
	Feature 1	- ranking		
	Ind1	Ind2		
Ind1	1st	4th		
man	130	401		
				F
				L
<u> </u>				

Dataset			
	Feat. 1	Feat. 2	
Ind1	67	173	
Ind2	53	161	



	Feature	2 - rankir	ng	
	Feature 1	- ranking		
	Ind1	Ind2		╞
Ind1	1st	4th		_
Ind2	2nd	1st		
maz	2110	101		

Dataset			
	Feat. 1	Feat. 2	
Ind1	67	173	
Ind2	53	161	



	Feature	e 2 - rankir	Ig	
	Feature 1	- ranking		
	Ind1	Ind2		
Ind1	1st	4th		F
Ind2	2nd	1et		
1102	2110	150		

Dataset			
	Feat. 1	Feat. 2	
Ind1	67	173	
Ind2	53	161	



	Feature	2 - rankir	ng	1
	Feature 1	- ranking		
	Ind1	Ind2		
Ind1	1st	4th		F
Ind2	2nd	1et		
11102	2110	131		
				r

=	

Feature 2 - scores						
	Feature 1 - scores					
	Ind1	Ind2		_		
Ind1						
Ind I						
Ind2				-		

Dataset			
	Feat. 1	Feat. 2	
Ind1	67	173	
Ind2	53	161	



Feature 2 - ranking					
Feature 1 - ranking					
	Ind1	Ind2		-	
Ind1	1st	4th		L	
Ind2	2nd	1st			
	2.110	101			

	Feature	e 2 - score	s			
	Feature 1	- scores				
	Ind1	Ind2		_		
Ind1	100			_		
Ind2		100		_		

Dataset			
	Feat. 1	Feat. 2	
Ind1	67	173	
Ind2	53	161	



Feature 2 - ranking					
	Feature 1	- ranking			
	Ind1	Ind2			
Ind1	1st	4th		⊢	
110	0.1	4.1			
ina2	2nd	ist			
				L	
				\sim	

=

Feature	e 2 - score	S		
Feature 1 - scores				
Ind1	Ind2			
100				
90	100			

Ind1 Ind2

Dataset			
	Feat. 1	Feat. 2	
Ind1	67	173	
Ind2	53	161	



	Feature 2 - ranking				
	Feature 1 - ranking				
	Ind1	Ind2			
Ind1	1st	4th		-	
Ind2	2nd	1et			
maz	2110	130			

$\overline{}$		7	
	\checkmark		

Feature 2 - scores						
	Feature 1	1 - scores				
	Ind1	Ind2				
Ind1	100	75		-		
Ind2	90	100		L		

Dataset			
	Feat. 1	Feat. 2	
Ind1	67	173	
Ind2	53	161	



	Feature	e 2 - rankir	Ig	
	Feature 1	- ranking		
	Ind1	Ind2		
Ind1	1st	4th		
Ind2	and	1 ot		
muz	Znu	151		
				-
$\left(- \right)$)	

\prec	7	
Feature	e 2 - score	s
Feature 1	- scores	
Ind1	Ind2	

	Feature 1	- scores	
	Ind1	Ind2	 -
Ind1	100	75	 _
Ind2	90	100	
			r

Dataset			
	Feat. 1	Feat. 2	
Ind1	67	173	
Ind2	53	161	



	Feature	e 2 - rankir	ng	
	Feature 1	- ranking		
	Ind1	Ind2		┝
1				
Ind1	1st	4th		
Ind2	2nd	1st		L





	Feature	e 2 - score	s	1
	Feature ?	1 - scores		
	Ind1	Ind2		\vdash
Ind1	100	75		⊢
Ind2	90	100		-
				7
				'

Dataset				
	Feat. 1	Feat. 2		
Ind1	67	173		
Ind2	53	161		



	Feature	e 2 - rankir	ng	1
	Feature 1	- ranking		
	Ind1	Ind2		
Ind1	1st	4th		-
Ind2	2nd	1st		
				/





	Feature	e 2 - score	s	
	Feature	I - scores		
	Ind1	Ind2		
Ind1	100	75		
la dO	00	100		
Indz	90	100		

Dataset				
	Feat. 1	Feat. 2		
Ind1	67	173		
Ind2	53	161		



	Feature	2 - rankir	ng	1
	Feature 1	- ranking		
	Ind1	Ind2		\vdash
Ind1	1st	4th		-
Ind2	2nd	1st		
				/
			/	



Ind1

Ind2



	Feature	e 2 - score	es	1
	Feature 1	- scores		
	Ind1	Ind2		
Ind1	100	75		
Ind2	90	100		

Dataset					
	Feat. 1 Feat. 2				
Ind1	67	173			
Ind2	53	161			



Feature 2 - ranking						
	Feature 1 - ranking					
	Ind1	Ind2		┝		
1						
Ind1	1st	4th				
Ind2	2nd	1st		L		
			/			





Easture 2 searce						
Feature 1 - scores						
	Ind1	Ind2				
Ind1	100	75				
Ind2	90	100				

	Feat. 1	Feat. 2	Feat. 3	Sum
Ind1	18	48	30	96
Ind2	12	17	5	34
Ind3	49	36	80	165
Ind4	59	79	30	168

	Feat. 1	Feat. 2	Feat. 3	Sum
Ind1	18	48	30	96
Ind2	12	17	5	34
Ind3	49	36	80	165
Ind4	59	79	30	168

	Feat. 1	Feat. 2	Feat. 3	Sum
Ind1	18	48	30	96
Ind2		17	5	22
Ind3	49	36	80	165
Ind4	59	79	30	168

	Feat. 1	Feat. 2	Feat. 3	Sum
Ind1	18	48	30	96
Ind2			5	5
Ind3	49	36	80	165
Ind4	59	79	30	168

	Feat. 1	Feat. 2	Feat. 3	Sum
Ind1	18	48	30	96
Ind2			5	34
Ind3	49	36		85
Ind4	59	79	30	168





EVALUATION





	exemplar structuring		nearest neighbor			
dataset	number	size	diameter	number	size	diameter
normal	18	16.7	5.2	22	13.6	2.8
residential	64	5.8	2.7	105	3.5	2.1
communities	256	8.7	3.2	400	5.5	2.5

PROTOTYPE



CONCLUSION

In a nutshell

- Structuring of high-dimensional data
- Interpretable results for domain experts
- Graph structure suitable for exploration

CONCLUSION

In a nutshell

- Structuring of high-dimensional data
- Interpretable results for domain experts
- Graph structure suitable for exploration

Future work

- Automatically find best tradeoff for neighborhood size K
- Use all values of K to determine the DoR score

- "Rapid Decision Making on the Fire Ground" 2010 - G. Klein et al.
- "On the Surprising Behavior of Distance Metrics in High Dimensional Space"
 2001 - C. C. Aggarwal et al.
- "The Role of Hubness in Clustering High-Dimensional Data"
 - 2014 N. Tomasev et al.

EXPLORATION OF HIGH-DIMENSIONAL DATA

EXPLAINABLE STRUCTURING DISCOVERY OF RELEVANT CASES

Joris Falip, Frédéric Blanchard, Michel Herbin joris.falip@univ-reims.fr www.datacrunch.sh

March 20 2019

CReSTIC, University of Reims Champagne-Ardenne, France