Learning from Evolving Data Joao Gama Ernestina Menasalvas University of Porto, PORTUGAL Technical Univ of Madrid SPAIN Athena Vakali Myra Spiliopoulou University of Magdeburg, GERMANY University of Thessaloniki, GREECE and the Chairs of the HaCDAIS Workshop Mykola Pechenizkiy, Eindhoven Univ. of Technology, NETHERLANDS Indrė Žliobaitė, Eindhoven Univ. of Technology, NETHERLANDS ECML-PKDD Conference O UNIVERSIDADE DO PORTO Tutorial Barcelona, Sept. 24th, 2010



- data (as in Customer Relationship Management or Social Tagging present appropriate adaptive learning methods.
- vas, Spiliopoulou, Vakali Barcelona, 24th Sept. 2010

8





ж $(\mathbf{0})$ Tutorial Presenters (or HaCDAIS Workshop Co-chairs?) Mykola Pechenizkiy Assistant Professor at the Department of Computer Science, Eindhoven University of Technology, the Netherlands. He has broad research interests in data mining and its application to various (adaptive) information systems serving industry, commerse, medicine and education. He has been organizing several workshops and conferences in these areas. http://www.win.tue.nl/~mpechen/ Indrė Žliobaitė Postdoctoral Researcher at the Department of Computer Science, Eindhoven University of Technology, the Netherlands. She received her PhD from Vilnius University, Lithuania. Her main research interests include data mining under concept drift and context-aware prediction. She publishes in intelligent data analysis and pattern recognition venues.



Ivas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010

(c) Gama, N





Presentation Agenda (tentative)	
☑ Block 1: Introduction	
Block 2: Supervised learning on streams (loao Gama,Mykola Pechenizkiy,Indre Zliobalte)	
 Block 3: Unsupervised learning on streams (Myra Spiliopoulou) 	
Tiny Break	
Block 4: Mining evolving social data (Athena Vakali)	
Block 5: Mining under resource Constraints (Ernestina Menasalvas)	
Block 6: Conclusions and Outlook	
(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010	(8)

EC recommendation



















Stream Classification	
Processing each example:	
Small constant time	
Fixed amount of main memory	
Single scan of the data	
Without (or reduced) revisit old records.	
Processing examples at the speed they arrive	
Decision Models at anytime	
Ability to detect and react to concept drift	
Ideally, produce a model equivalent to the one that would be obtai by a batch data-mining algorithm	ned
(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010	(19)



(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010

VENTREDER RESIDENCE VALUE - REFERENCE 2010 (21)

Hoeffding bound	
Suppose we have made <i>n</i> independent observations of a random variable <i>r</i> whose range is R.	
The Hoeffding bound relates the mean in the sample with th mean in the population:	ne
With probability 1-δ	
• The true mean of <i>r</i> is in the range $\overline{r} \pm \varepsilon$ where $\varepsilon = \sqrt{R^2 \log(1/\delta)/(2N)}$	
 Independent of the probability distribution generating the examples. 	
(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010	(22)



Entropy: Sufficient Statistics	
Each leaf stores sufficient statistics to evaluate the splittin criterion	g
 For each attribute 	
If Nominal	
- Counter for each observed value per class	
If Continuous	
- Binary tree with counters of observed values	
- Discretization: e.g. 10 bins over the range of the	variable
- Univariate Quadratic Discriminant (UFFT)	
(1 Cama Manasahasi Shilionoulou Vakali Barcelona 24th Sant 2010	(24)









Properties of VFD	like Algorithms	;					
Low variance models:							
 Stable decisions with statistical 	l support.						
 No need for pruning; 							
 Decisions with statistical supp 	ort;						
Low overfiting:							
Examples are processed only once.							
 Decisions are taken using different set of examples 							
Convergence: VFDT becomes asymptotically close to that of							
a batch learner. The expected dis	agreement is δ/p , where p is						
the probability that an example fa	all into a leaf.						
Decision Processes Decide to expand	Batch Learning Choose the best split from all data in a given node	Hoeffding Trees Accumulate data till there is statistical evidence in favor to a particular soliting test					
Pruning Drift Detection	Mandatory Assumes data is stationary	No need Smooth adaptation to the most recent concepts					
(c) Gama, Menasalvas, Spiliopoulou, Vak	ali - Barcelona, 24th Sept. 2010		(29)				

















Snapshot	S
Store the m micro-clust in such a w	icro-clusters at different moments in time. These stored ter states are referred to as <i>snapshots</i> . Snapshots are stored ay that:
 Maint horizo 	ain sufficient amount of information about different time ons.
 Avoid 	the storage of an unnecessarily large number of time horizons
 This is 	achieved with the use of a geometric time frame.
Tume no. Sampshots (by clock time) 1 09 67 65 2 70 66 62 3 68 60 52 4 56 40 24 5 48 16	 Snapshots are classified into different frame numbers. The frame number of a particular class of snapshots defines the level of granularity in time at which the snapshots are maintained.
6 64.32	 Snapshots of frame number i are stored at clock times which are divisible by 2ⁱ but not by 2ⁱ⁺¹.
	» Each frame as limited capacity



Cl	assification on Demand
Со	nsider an example in which:
	 the current clock time is t_c, and
	- a horizon of length h in order to find the micro-clusters in the time period $(t_c\text{-}h, t_c).$
Fin	ding the micro-clusters for the time period of interest
	 Find the stored snapshot which occurs just before the time t_e-h.
	 For each micro-cluster in the current set S(t_c), we find the list of ids in each micro-cluster.
	- For each id in the list of ids , we find the corresponding micro-clusters in S(th), and subtract the CF vectors for the corresponding micro-clusters.
	 The resulting set of micro-clusters correspond to the time horizon (t-h, t).

(42)











Novelty Detection	
Classification problems where the full set of class la unknown.	bels is
Automatic identification of unforeseen phenomena embedde large amount of normal data.	ed in a
Novelty is a relative concept with regard to our current know	ledge:
 It must be defined in the context of a representation of our cu knowledge. 	rrent
 Specially useful when novel concepts represent abnormal or u conditions 	nexpected
Expensive to obtain abnormal examples	
Probably impossible to simulate all possible abnormal condition	ons
(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010	(48)





 In a stable state, the contribution of each unit is likely to remain constant. Changes in the participation of decision units may indicate a conceptual change

(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010





OLLINDA	
Cluster-based novelty detection, Spinoza, Carvalho, Gama, SAC 08.	
Initial Phase: Supervised, batch mode	
 Start by modeling the normal condition. 	
 Learns a partial model about what is known. 	
 Based on a set of classified examples. 	
Second Phase: Process stream of unlabelled examples	
 For each incoming example: 	
If it is explained by the current model: classify the example and discard	
 If it is not explained: Store in a short-term memory 	
Time to Time	
- Find clusters in the examples stored in the Short Term Memory	
 Clusters far away from existing ones: Novel concept. 	
 Clusters closed to existing ones: Extend known concepts. 	
(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010	(54)











Motivation

View CD research from an application perspective What is the match between the mainstream CD research assumptions and properties of the applications? Identify promising future research directions from the application perspective

We will talk about

Why changes appear in different applications?

What are the properties of CD application tasks?

How the application tasks can be categorized in terms of these basic properties?

(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010

What is Concept Drift?

The closed world assumption in data mining

learn a model from examples described by a finite set of features

In reality some important properties are not observed

hidden variables that influence the concept

Hidden variables may change over time

- concepts learned at one time can become inaccurate
- possible changes in the characteristic properties of the concept

Concept Drift

- changes in the *hidden context* that can induce more or less radical changes in the *target concept*
- Virtual concept drift changes due to population drift
- Gama, Menasalvas, Spiliopoulou, Vakali Barcelona, 24th Sept. 2010

Desired Properties of a System Handling Concept Drift

Adapting to concept drift asap

must have assumptions of what and how may change

Being robust to noise and distinguishing it from concept drift • e.g. occasionally wrong selection or rating of an item, clicking a link, connection failure (mobile computing)

Elasticity

discouraging brittleness

(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010

Being capable to recognize and react to reoccurring contexts

such as seasonal differences









Types of apps Monitoring/ Personal assistance/ Management Ubiquitous									
Industries	control	personalization	and planning	applications					
Security, Police	Fraud detection, insider trading detection, adversary actions detection		Crime volume prediction	Authentica- tion, Intrusion detection					
Finance, Banking, Telecom, Credit Scoring, Insurance, Direct Marketing, Retail, Advertising, e-Commerce	Monitoring & management of customer segments, bankruptcy prediction	Product or service recommendation, including complimentary	Demand prediction, response rate prediction, budget planning	Location based services, related ads, mobile apps					
Education (higher, professional, child- ren, e-Learning) Entertainment, Media	Gaming the system, Drop out prediction	Music, VOD, movie, learning object recommendation, adaptive news access, personalized	Player- centered game design, learner- centered	Virtual reality simulations					

















Anubiouc Rea	SIST	anc	e P	redi	ction	(Tsyn	nbal e	t al.,	2008	3)
predict the sens	itivit	vo	far	hatho	aen to	an ant	ihiotic	hased	l on da	ita
about the antibi	otic	the		later	l natho	den an	nd the	damo	aranhi	c and
about the antibi	of +1	une ao a	- 150	natec	patilo	gen, a	iu the	uemo	grapin	c anu
ciinical reatures	01 11	le h	Jalle	ent.				pathogen	antibiotic	sensi
0	ate	SEX	age	ISINEW	days_total	days_ICU	main_dept	data	data	tivity
22.1	2002	m	25		171	81	å			3
22.1	2002	m	25		171	81	ő			3
22.1	2002	m	25	- i	171	81	ä			3
22.1	2002	4	61		261	52	3			3
28.1	2002	÷	61	0	261	52	3			3
28.1	2002	÷	61	õ	261	52	3			3
28.1	2002	÷	61	ŏ	261	52	3			1 i
28.1	2002	÷	61	ő	261	52	3			i i
28.1	2002	m	25	1	171	81	9			à
28.1	2002	m	25	- i	171	81	9			3
30.1	2002	m	25	- i	171	81	9			3
8.2	2002	m	30	0	209	209	9			3
8.2	2002	m	30	ō	209	209	9			i i
8.2	2002	m	30	ò	209	209	9			1
11.2	2002	f	0	ō	18	0	2			1
11.2	2002	f	ó	ō	18	ò	2			1
11.2	2002	f	0	0	18	0	2			1
new	data									2
new	data									?
	data									2













References (Block 1) 2/3	References (Block 2) 1/3
Handling concept drift	Unsupervised learning – underpinnings (3) and applications (5,6):
Gama J., Castillo G. 2006. Learning with Local Drift Detection. ADMA 2006: 42-55	L AlSumait, D. Barbara, C. Domeniconi. On-lineLDA: Adaptive topic models for mining text streams with applications to topic detection and tracking. In ICDW108: Proc. of IEEE Int. Conf. on Data Mining. Pice. Dec. 2008. IEEE
Klinkenberg R. and Renz. I. 1998. Adaptive information filtering: Learning in the presence of concept drifts. In Learning for Text Categorization, pages 33–40. AAAI Press.	H. Cao, E. Chen, J. Yang, and H. Xiong. Enhancing recommender systems under volatile user interest clifts. In CIKM'09: Proc. of the 18th ACM Conf. on Information and Knowledge Management (CIKM'09), pages 1257–1266. 2009. ACM.
Kuncheva LI. 2008. Classifier ensembles for detecting concept change in streaming data: Overview and perspectives, Proc. 2nd Workshop SUEMA 2008 (ECAI 2008), Patras, Greece, 5–10	Aysegul Cayci, João Bartolo Gomes, Andrea Zanda, Ernestina Menasalvas, and Santiago Elbe. Situation-Aware Data Mining
Kuncheva LJ. 2004. Classifier ensembles for changing environments, Proceedings 5th Int. Workshop on Multiple Classifier Systems, MCS2004, Lecture Notes in Computer Science, Vol 3077, 1–15.	Jewice in conquictus a minimients, baccoming or mata citic 2005 J. Sun, D. Tao, and C. Faloutsos, Beyond streams and graphs; dynamic tensor analysis. In KDD '06: Proc. of 12th ACM SICKDD Int Crinf on Knowledne Discoversion and Data Minimo pranes 747–833. New York: NY, UKA 2006. ACM
Minku L L; White A P; Yao X. 2010. The Impact of Diversity on On-line Ensemble Learning in the Presence of Concept Drift., IEEE Transactions on Knowledge and Data Engineering, IEEE, v. 22, n. 5, p. 730–742.	Cohr, A. Hinneburg, R. Schult, M. Spilopoulou, Topic evolution in a stream of documents. In SDM'09: Proc. of SIAM Data Mining Cohr, nages 738–785. Renn, NV, USA. Anr. (Nex 2000)
Pechenizkiy M., Bakker J., Žliobaitė I., Ivannikov A. & Kärkkäinen T. 2009. Online Mass Flow Prediction in CFB Boilers with Explicit Detection of Sudden Concept Driff, SIGKDD Explorations 11(2).	Z. Lu, D. Agarwal, and J. Dhillon. A spatio-temporal approach to collaborative filtering. In RecSys '09: Proc. of 3rd ACM Conf. on Recommender Systems: nanes 12: 20. New York: Nov. 2009. ACM
Tsymbal A., Pechenizkiy M., Cunningham P. and Puuronen S. 2008. Dynamic Integration of Classifiers for Handling Concept Drift, Information Fusion, Special Issue on Applications of Ensemble Methods, 9(1), pp. 56-68.	Z. F. Siddigui, M. Spillopoulou. Combining multiple interrelated streams for incremental clustering. In SSDBM'09: Proc. of 21st Int. Conf. on scientific and Statistical Database Management. New Orleans. USA. Lune 2009
Widmer G. and Kubat M. 1996. Learning in the presence of concept drift and hidden contexts. Machine Learning, 23:69–101.	Ruiz, Carlos: Menasakas, Emestina: Snilionoulou, Myra, C-DenStream: Using Domain Knowledge for Clustering on Data
Žliobaitė I., Bakker J. and Pechenizkiy M. 2009. Towards Context Aware Food Sales Prediction, In Proceedings of IEEE International Conference on Data Mining (ICDM/09) Workshons. IEEE Computer Society, pp. 94-99.	Streams. Discovery Sicence '09. Porto oct. 2009
Žliobaitė I. 20140daptive Training Set Formation. PhD Tesis: Vilnius University.	Maria, Valencia, Lauth Ina; Ernestina Menasalvas. Ernerging user intentions: matching user queries with topic evolution in news text streams. International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems (JUFKS) Vol. 17pp. 59– 80. 2009
Žliobaitė I. 2009. Learning under Concept Drift: an overview. Technical Report: Vilnius University.	María Valencia, Ernestina Menasalvas, and Santiago Elbe. Approach for Discovering and Tracking Local Web Search Trends. LA- WEB '09. Mayico new 2009.
Žilobaitė I. and Pechenizkiy M. 2010. Handling Concept Drift in Adaptive Information Systems. Technical Report: Eindhoven University of Technology.	
(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010 (85)	(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010 (86)































asalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010





















MONIC (Spiliopoulou et al., 2006)

- 1. Partition the time axis into timepoints $\boldsymbol{t}_1,...,\boldsymbol{t}_m$
- 2. Cluster (adaptively?) the data at each snapshot.
- 3. Compare the sets of clusters (need not be of equal strength)

Matching model:

How to track a given cluster?

• When is a new cluster a mutation of an old one?

Transition model:

- Is an old cluster associated with exactly one new cluster?
- How did the cluster change with respect to other clusters?
- What kind of internal changes did the cluster experience?
 Gama, Menasalvas, Spiliopoulou, Vakali Barcelona, 24th Sept. 2010









Presentation Outline	
☑Block 1: Introduction	
☑Block 2: Supervised learning on streams	
Block 3: Unsupervised learning on streams	(Myra Spiliopoulou)
✓ Adapting clusters	
 Probabilistic models 	
 Learning on complex data 	
Block 4: Mining evolving social data	
Block 5: Mining under resource constraints	
Block 6: Conclusions and Outlook	
(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010	(120)















Online LDA (AlSumait et al., 2008)	
Evaluation on topic evolution	
 Interesting findings on the evolution of NIPS topics 	
Table 2. Examples of topics estimated by OLDA from NIPS corpus and its evolution over 13 years	
Topic 12: Reinforcement Learning 88: state learning system states time cycles recurrence failure weight algorithm 89: node system state rule learning nodes tim match transition 90: state learning rule system node algorithm change rules constoller dynamic	
91: learning state reinforcement system world time adaptive planning controller 92: state learning action task exploration tasks sutton elemental 02: lowing other sufficiency control line octing take actions take head.	
95: tearing state reinforcement control time action tasks of pointal based of the tearing state optimal control dynamic policy action time adaptive 95: tearing state optimal action policy control reinforcement grid dynamic 96: tearing state action policy reinforcement al policy time action time.	
97: Idearing state action reinforcement time policy optimal algorithm dynamic 97: Idearing state action reinforcement time policy optimal algorithm dynamic 98: state learning policy reinforcement optimal time action step control	
99: state learning policy action reinforcement optimal rl time 2000: state learning policy action reward time reinforcement belief	
(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010	(128)

Length I=1

Independent PLSA
 Adaptive PLSA

16 32 64 Learning & Recalibration Steps

128

(130

410

8

35

340



			100 100 100						
Adaptive PLSA (Gohr et al, 2009)									
Тор	Topic threads in SIGIR (2000-2007)								
Thread	1 2000	2001	2002	2003	2004	2005	2006	2007	
Evaluation	evalu search queri docu- ment relev approach re- sult method effect retriev	queri docu- ment relev evalu search retriev lan- guag result translat summari	queri retriev search term document result relev techniqu in- form translat	retriev queri system ef- fect evalu search collect word inform techniqu	retriev queri relev docu- ment collect inform sys- tem feedback evalu term	queri retriev term inform feedback per- form system document evalu relev	queri retriev document relev result evalu term feedback improv precis	queri retriev relev system rank evalu perform col- lect measur effect	
Presentation later empha on multimed	inform user retriev model document system base SiS roach process studi la	model re- triev inform user languag system doc- ument estim feedback problem	model in- form languag retriev user document queri col- lect method paramet	model inform retriev imag user languag video perform annot show	model retriev inform lan- guag imag approach document framework distribut propos	model retriev imag inform languag sp- proach annot gener propos show	model retriev inform docu- ment languag approach probabl base space framework	model re- triev inform languag base term imag propos approach show	
3 Supervise learning	perform queri select clas- difi method diassif text vector system databas	text perform classif similar improv stori approach combin fea- tur distribut	text index classifi result classifi data corpu ap- proach time invert	text clas- sif classifi method fea- tur result structur cat- egor term approach	featur classif text method learn algo- rithm per- form achiev automat train	method text classif featur detect extract base automat show data	text question document featur classif approach method task sentenc summar	featur learn text classif answer ap- proach paper task classifi base	
Web	web qualiti page system document techniqu question re- triev answer paper	web page topic answer link automat text question task gener	web extract data sys- tem item text analysi search page separ	ir web re- search search system ques- tion answer inform data task	web search page system user task document answer evalu result	search web user page engin result algorithm inform collect rank	search web user inform page result network rank algorithm structur	search web user queri inform page engin result relev log	
5 Clustering	document cluster word method retriev algo- rithm inform perform sys- tem improv	document method base threshold in- dex segment score term distribut sigorithm	document cluster topic algorithm set method score evalu model similar	document topic cluster method base algorithm set approach similar pro- pos	document cluster filter method al- gorithm user data semant index differ	method doc- ument topic algorithm optim semant similar data index result	cluster doc- ument data method sim- ilar semant propos algo- rithm index rate	document index cluster method algo- rithm result propos op- tim perform problem	31)





















One topic thread towards tensor-based stream clustering
Evolutionary clustering (with temporal smoothness)
D. Chakrabarti, R. Kumar, A. Tomkins, "Evolutionary clustering", 12th ACM SIGKDD Int. Conf (KDD'06), Philadelphia, PA, Aug. 2006
Incremental spectral clustering (with temporal smoothness)
 Y. Chi, X. Song, D. Zhou, K. Hino, B. Tseng, "Evolutionary Spectral Clustering by incorporating Temporal Smoothness", 13th ACM SIGKDD Int. Conf. (KDD'07), San Jose, CA, Aug. 2007.
Incremental learning with generative models and temporal smoothness
 YR. Lin, Y. Chi, S. Zhu, H. Sundaram, B. Tseng, "FacetNet: A Framework for Analyzing Communities and their Buolutions in Dynamic Networks", World Wide Web Int. Conf. (WWW'08), Beijing, China, Apr. 2008.
Incremental tensor-based clustering
 J. Sun, D. Tao, C. Faloutsos. "Beyond streams and graphs: dynamic tensor analysis", 12th ACM SIGKDD Int. Conf. (KDD'06), Philadelphia, Aug. 2006
 YR. Lin, J. Sun, P. Castro, R. Konuru, H. Sundaram, A. Kelliher, "MetaFac: Community Discovery via Relational Hypergraph Factorization", ACM SIGKDD Int. Conf. (KDD'09), Paris, France, June-July 2009
(r) Gama Menasakas Snilionnulnu Vakali – Barrelona 24th Sent 2010 (142)





Modeling temporal smoothness (Chi et al., '07)
Snapshot cost of clustering ξ at t': $\ \Sigma \ _{\Sigma}$ $\ \Sigma \ _{2}$
$CS(\xi^{*}, t^{*}) := SSE(\xi^{*}, t^{*}) = \sum_{j=1}^{K} \sum_{x \in C_{j,t^{*}}} \left\ \frac{Z_{y \in C_{j,t^{*}}}}{ C_{j,t} } - x \right\ $
Two functions for temporal cost of clustering ξ' at t' towards clustering ξ at t (t:=t'-1):
2) <u>Preserving Cluster Membership (PCM):</u>
$CT_{PCM}(\xi^{*},\xi^{*}) = -\sum_{X \in \mathbb{P}^{*}} \sum_{Y \in \mathbb{P}^{*}} \frac{ X \cap Y ^{2}}{ X \cdot Y }$ originally assuming that $ \xi = \xi^{*} = K$, and relaxing this later
on.
NOT ORIGINAL NOTATION ! (c) Cam, Marsalvas, Spilopoulou, Vakali - Earcelona, 24th Sept. 2010 (145)





























References (Block 3) (1/3)

(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010

C. Aggarwal, J. Han, J. Wang, P. Yu. "A Framework for Clustering Evolving Data Streams", 29th Int. Conf. on Very Large Data Bases (VLDB'03), Berlin, Germany, 2003.

F. Cao, M. Ester, W. Qian, A. Zhou. "Density-Based Clustering Over an Evolving Data Stream with Noise", SIAM Int. Conf. on Data Mining (SDM'06), 2006.

S. Guha, A. Meyerson, N. Mishra, R. Motwani, L. O' Callaghan. *Clustering data streams: Theory and practice*. IEEE Trans. of Knowledge and Data Eng., 15(3):515–528, 2003.

P. Kranen, I. Assent, C. Baldauf, T. Seidl. "Self-Adaptive Anytime Stream Clustering", IEEE Int. Conf. on Data Mining (ICDM'09), Miami, Florida, Dec. 2009

M. Spiliopoulou, I. Ntoutsi, Y. Theodoridis, R. Schult. "MONIC – Modeling and Monitoring Cluster Transitions", 12th ACM SICKDD Int. Conf. on Knowledge Discovery and Data Mining (KDD'06), Philadelphia, PA, Aug. 2009



References (Block 3) (3/3)	References (Block 3) (further citations)
 Y. Chi, X. Song, D. Zhou, K. Hino, B. Tseng. "Evolutionary Spectral Clustering by Incorporating Temporal Smoothness", ACM SIGKDD Int. Conf. (KDD'07). San Jose, CA, Aug. 2007. YR. Lin, Y. Chi, S. Zhu, H. Sundaram, B. Tseng. "FacetNet: A Framework for Analyzing Communities and their Evolutions in Dynamic Networks", World Wide Web Int. Conf. (WWV'08), Beijing, China, Apr. 2008. YR. Lin, J. Sun, P. Castro, R. Konuru, H. Sundaram, A. Kelliher. "MetaFac: Community Discovery via Relational Hypergraph Factorization", ACM SIGKDD Int. Conf. (KDD'09), Paris, France, June–July 2009. Z. Siddiqui, M. Spiliopoulou. "Combining Multiple Interrelated Streams for Incremental Clustering", 21st Int. Conf. on Scientific and Statistical Database Management (SSDBM09), New Orleans, May 2009 Z. Siddiqui, Myra Spiliopoulou. "Clustering a Stream of Growing Objects", Discovery Science Int. Conf. (DS'09), Porto, Oct. 2009 	 G. Hulten, L. Spencer, P. Domingos. "Mining Time-Changing Data Streams", 7th ACM SIGKDD Int. Conf. (KDD '0), 97-106, ACM, New York, NY, USA, 2001 Mark A. Kroegel, On Propositionalization for Knowledge Discovery in Relational Databases. PhD thesis, University of Magdeburg, Germany. 2003. R. Nallagati et al. "Multiscale Topic Tomography", ACM SIGKDD Int. Conf. (KDD'07), San Jose, CA, Aug. 2007 Ning et al. "Incremental spectral clustering by efficiently updating the eigen-system" in Pattern Recognition 43(1), 113-127, 2010 Z. Siddiqui, M. Spillopoulou. "Tree Induction over Perennial Objects", In Proc. of 22d nltr. Conf. on Scientific and Statistical Database Management (SSDBM'10), LNCS 6187, 640-657, Heidelberg, June-July 2010 P. Symeonidis, A. Nanopoulos, Y. Manolopoulos, "Tag Recommendations based on Tensor Dimensionality Reduction", Int. Conf. on Recommender Systems (RecSys'08), Lausanne, Switzerland, Oct. 2008, D. Charkabazit, R. Kumar, A. Tomikins, "Evolutionary dustering", 12th ACM SIGKDD Int. Conf. (KDD'06), Philadelphia, PA, Aug. 2006 J. Sun, D. Tao, C. Falouzos: Beyond Streams and graphs: dynamic tensor analysis", 12th ACM SIGKDD Int. Conf. (KDD'06), Philadelphia, PA, Aug. 2006
(c) Gama, Menasalvas, Spiliopoulou, Vakali – Barcelona, 24th Sept. 2010 (163)	(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010 (164)





Motivation for	Mining Social Data	
 The availability of mining. 	f massive sizes of data gave new impe	tus to data
 e.g. more that than 25 billio 	an 400 million active Facebook users, sharii on pieces of content each month (Facebook Sta	ng on average more ttistics 2010]
 Mining social well Non-obvious res Collaboration form Wisdom of the 	b data can act as a barometer of the us ults may emerge. n and contribution of many individuals ation of <i>collective intelligence</i> <i>he crowd</i> : more accurate, unbiased source	sers' opinion.
 Social data minin recommender sy 	g results can be useful for applications stems, automatic event detectors, etc	s such as
 Various mining te clustering, statist 	echniques are/can be used: communit ical analysis, classification, association	y detection , rules mining,



Which information is richer?























Presentation Outline	
☑Block 1: Introduction	
☑Block 2: Supervised learning on streams	
☑Block 3: Unsupervised learning on streams	
Block 4: Mining evolving social data (Athena Vakali)	
✓ Structure and Models	
Community Detection in Evolving Social Graphs	
 Applications of Evolving Community Detection 	
Block 5: Mining under resource constraints	
Block 6: Conclusions and Outlook	
(r) Cama Manarahuri Shilionnulou Vakali Burcelona 24th Sant 2010 (1	82)





Community Detect	ion in Socia	al Graph	s	
Measures				
Quantitative, for assessin between nodes	g relations	Qualitat commu	ive, for evaluating nity structure	
co-occurrence		modulari	ity	
cosine similarity		local modularity		
tf- idf	cut-size			
betweeness centrality		node out	twardness	
structural similarity Methodologies			Extra challenges: ov of hierarchical com	erlapping munities
Graph partitioning	Spectral	algorithms	5	
Clustering	Methods	Methods based on statistical inference		
Divisive algorithms	Dynamic	algorithms	\$	



























































Stream-Group
Stream-Group: applied on dynamic weighted directed graphs [Duan09]
Uses:
 Random Walk with Restart to compute the graph's relevance matrix R
• r_{ij} expresses the probability that random walker will stay at i when starting from j
 relevance scores between within community nodes are usually higher than the ones between different communities
 an extension of <i>modularity</i> to evaluate the goodness of a partition
Procedure
1.For each new graph arrival, identify its community structure
2.Compute similarity between new structure and current segment's structure
3.If similarity over threshold, the community structure of the current segment is updated incrementally with the new data
4.else, a new segment is initiated

(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010

(216)

UNIVERSIDADI	. (1)	A	8
DO PORTO	NUTENICS	VX//	ARETOTEADO EAMETIETRICO BETEARONIDIE

Stream-Group

Stream-Group: applied on dynamic weighted directed graphs [Duan09] Uses:

- Random Walk with Restart to compute the graph's relevance matrix R
- r_{ij} expresses the probability that random walker will stay at *i* when starting from *j* relevance scores between within community nodes are usually higher than the
- ones between different communities

an extension of *modularity* to evaluate the good 2.Merge groups while goodness indicator increases

1.For each new graph arrival, identify its community structure

- 2.Compute similarity between new structure and current segment's structure
- 3.If similarity over threshold, the community structure of the current segment is
- updated incrementally with the new data
- 4.else, a new segment is initiated
- (c) Gama, Menasalvas, Spiliopoulou, Vakali Barcelona, 24th Sept. 2010



4.else, a new segment is initiated

(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010



(217















DenGraph	
Incremental density-based clustering method [Falkowski08]	
•Let $min\{num_interactions_{i \rightarrow j}num_interactions_{i \rightarrow j}\} = k_{ij}$	
 Node-distance matrix, where: dist(i,i) = 1 and dist(i,j) = 0 if k_{ij} = dist(i,j) = 1 / k_{ij} 	0, or else
•Clusters are built by connecting adjacent (μ, ε) -neighborhoods	[Ester96]
•A ($\mu\varepsilon$)-neighborhood is built around a core node and includes nodes within a radius of ε	at least μ
Border nodes are included in a neighborhood but are not cores	
Noise nodes are not included in any neighborhood	
(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010	(228)











Social network analysis		
Source	Community detection method	
synthetic datasets	Chi07, Duan09, Gorke10, Kim09, Lin08, Tang08, Yang09	
benchmark datasets (e.g. dolphin's network, IEEE VAST dataset)	Chi07, Yang09	
article online archives (e.g. DBLP, arXiv)	Gorke10, Kim09, Lin08, Palla07, Tang08, Wang08, Yang09	
mobile phone call networks	Palla07, Sun07, Wang08, Yang09	
imdb actor collaboration dataset	Wang08	
Enron e-mail dataset	Duan09, Falkowski08, Sun07, Tang08, Wang08, Yang09	
blog data	Chi07, Lin07, Lin08	
Flickr	Chakrabarti06	
mobile device proximity networks	Sun07	
company call networks	Yang09	
football match schedules	Kim09	
department e-mail dataset	Gorke10	
Digg	Lin09	
(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona	, 24th Sept. 2010 (234)	





$(\mathbf{0})$

Trend detection

- Social data fluctuate in their structure and frequency as they evolve. At each timeperiod there are some topics, images, tags, etc, that are most popular amongst users (trends). Data mining can be used for detecting trends in evolving social data.
- Trends can be identified globally or even locally (within communities) and they usually indicate what interests users the most at a given time

Trend identification in Twitter

- Lwit • Twitter : popular microblogging website where users are allowed to post short messages (up to 140 chars) and "follow" the posts of others
- Rich source of rapidly evolving social data which are also public.
- Suitable for trend detection
- · Recently, there have been many attempts to statistically analyze Twitter data · Evolving Twitter data to identify trending keywords for different weekdays [Java07]
- Microblogging as a form of electronic word-of-mouth for sharing consumer opinions concerning brands. Sentiment identification performed in Twitter posts to identify trending sentiments about brands [Jansen09].

alvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010



Trend detection

- Social data fluctuate in their structure and frequency as they evolve. At each timeperiod there are some topics, images, tags, etc, that are most popular amongst users (trends). Data mining can be used for detecting trends in evolving social data.
- Trends can be identified globally or even locally (within communities) and they usually indicate what interests users the most at a given time

Trend identification in the blogosphere

٩

- Traditional approach: calculation of the frequency of appearance of each term in all blogs for each time-step
- However, in [Chi06] different emphasis is put on every blog depending on its overall amount of contribution to the trend Data represented as a combination of information capturing: (i) temporal changes in them
- and (ii) characteristics of individual bloggers

(238

- Method based on SVD: produces scalar eigen-trends capturing overall trends, and authority scores representing the contribution of each blog to trends
- Method based on HOSVD: applied on keyword-specific blog-citation 6 esemblanc to HITS graphs to produce structural eigen-trends, hub scores and ithority scores capturing structural changes

c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010

(f)

Clustering of users exploiting the dimension of time

Social data can be analyzed for automatic synthesis of user profiles Case study: Social Tagging Systems (STS)

Clustering of users in STS according to topics of interest and the time cality of tagging activity [Koutsonikola09]

e.g. a user who tagged a set of photos depicting sports is probably interested in sports. However, if his tagging activity took place during the Olympic games, maybe he is simply interested in the Olympics and is not a regular sports fan.

Segmentation of time in frames

A user is related to a given tag if he has assigned at least one semantically close tag. The topic-distance between two users is calculated based on the similarity of their relations to all involved tags.

•Time-similarity between a user and a tag is calculated with the *cosine coefficient*. The time-distance between two users is calculated considering their similarity over all timeframes.

•Topic-based clustering with K-means, then refinement with time criterion (c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010

🔔 🚯 🛛 8 References (Block 4) (1/3)

[Au Yeung09] Au Yeung, C.M., Gibbins, N., and Shadbolt, N. 2009. Conte of 20th ACM Conference on Hypertext and Hypermedia, pp. 251–260.

[Chakrabarti06] Chakrabarti, D., Kumar, R., and Tomkins, A. 2006. Evolutionary clustering. In Proceedings of the 12th ACM SIGKDD international Conference on Knowledge Discovery and Data Mining . KDD '06. ACM, New York, NY, 554–560.

[Chi06] Chi, Y., Tseng, B. L., and Tatemura, J. 2006. Eigen-trend: trend analysis in the blogosphere based on singular value decompositions. In Proceedinas of the 15th ACM international Conference on information and Knowledge Management. CIKM '06. ACM. ew York, 68–77

(Ch07) Chi, Y., Zhu, S., Song, X., Tatemura, J., and Tseng, 8. L. 2007. Structural and temporal analysis of the biogosphere through community factorization. In Proceedings of the 13th ACM SGROD international Conference on Knowledge Discovery and Data Mining KDD 07. ACM. New York, NY, 163–127.

[Duan09] Duan, D., Li, Y., Jin, Y., and Lu, Z. 2009. Community mining on weighted directed graphs. In Proceeding of the 1st ACM international Workshop on Complex Networks Meet information & Knowledge Management . CNIKM '09. ACM, New York, NY, 11–18.

[Ester96] Ester, M., Kriegel, H.-P., Sander, J., and Xiu, X. 1996. A density-based algorithm for discovering clusters in large spatial databases with noise. In Proceedings of the 2nd international Conference on Knowledge Discovery and Data Mining. KDD '96. AAAI Press, 226-231. [Falkowski06] Falkowski, T., Bartelheimer, J., and Spillopoulou., M. 2006. Community Dynamics Mining. In Proceedings of the 14th European Conference on Information Systems. ECIS '06.

[Falkowski08] Falkowski, T., Barth, A., and Spiliopoulou M. 2008. Studying Community Dynamics with an Incremental Graph Mining Algorithm. In Proceedings of the 14th Americas Conference on Information Systems. AMCIS.

[Fortunato07] Fortunato, S. & Castellano, C. 2007, Community structure in graphs, arXiv:0712.2716v1.

[Girvan 02] Girvan, M. & Newman, M. E. J. 2002. Community structure in social and biological networks. In Proceedings of the National Academy of Sciences of the United States of America, 99(12):7821–7826.

vas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010 (c) Gama, Menasa (240



(243

(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010

(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010

End of Block 4



(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010

(242

Presentation Outline		
Block 1: Introduction		
☑Block 2: Supervised learning of the second sec	on streams	
☑Block 3: Unsupervised learning on streams		
☑Block 4: Mining evolving social data		
Block 5: Mining under resource	constraints	
IntroductionApproaches	(Ernestina Menasalvas)	
Block 6: Conclusions and Outlo	ok	
(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Se	pt. 2010 (245)	

Thank you!

Questions?

Introduction and Motivation
Challenge:
•a stream is theoretically infinite so cannot be materialized.
Data have to be processed in a single pass using little
memory.
Based on this restriction one can identify two divergent objectives:
 the analysis should produce comprehensive and exact results and detect changes in the data as soon as possible.
The single pass demand together with resource limitations allow only to perform the analysis on an approximation of the stream or a window (finite subset of the stream).
(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010 (246)

Introduction and Motivation

Data streams are mostly generated or sent to resource constrained computing environments:

Data generated on-board astronomical spacecrafts

oulou, Vakali - Barcelona, 24th Sept. 2010

- Data generated in sensor networks. The additional constraint that sensor nodes consume their energy rapidly with data transmission
- Data received in resource-constrained environment represents a different category of applications:
 - Personal Digital Assistants PDAs: users might request sheer amounts of data of interest to be streamed to their mobile devices. Storing and retrieving these huge amounts of data are also infeasible in such an environment



Applications

•monitoring physiological data streams obtained from wearable sensing devices. Such monitoring can be either for:

- applications for pervasive healthcare management,
- Applications for seniors,
- emergency response personnel,
- soldiers in the battlefield
- or athletes

•Onboard analysis of data streams: data generation would exceed the bandwidth to transfer these streams of data to ground stations for analysis. They necessitate the need for onboard analysis of data streams.

c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010



Applications
•Vehicle Data Stream Mining [Kargupta]
Vehicle Health Monitoring and Maintenance:
Detecting unusual behavior for a subsystem
Fuel Consumption Analysis:
Is the vehicle burning fuel efficiently? Identify influencing factors
Detect influence of driver behavior on gas mileage
Driver Behavior Monitoring:
Route monitoring: Fixed and variable routes
Direct Cost Issues: e.g. Idling, braking habits, Safety Issues
Vehicle location related services
 Vehicular network security and privacy management

(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010



 Most of these algorithms are not designed with regard to adaptation to resource availability

(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010

(252)

Issues (Gaber)	Presentation Outline
• the input,	☑Block 1: Introduction
output, and	ØBlock 2: Supervised learning on streams
 processing settings of an algorithm 	☑Block 3: Unsupervised learning on streams
 could be changed according to measurements of resource availability in a time frame: 	ØBlock 4: Mining evolving social data
Find:	Block 5: Mining under resource constraints
 resource consumption patterns and 	Introduction
algorithm settings.	Approaches (Ernestina Menasalvas)
•Challenge:	Block 6: Conclusions and Outlook
Limited computational resources	
Limited bandwidth	
Change of the user's context (c) Gama, Menashus, Spillopoulou, Vakali – Barcelona, 24th Sept. 2010 (253)	(c) Gama, Menasalvas, Spiliopoulou, Vakali – Barcelona, 24th Sept. 2010 (254)







Data Processing model

Landmark model

- mines all frequent itemsets over the entire history of stream data from a specific time point called landmark to the present.
- not suitable for applications where people are interested only in the most recent information of the data streams, such as in the stock monitoring systems

Damped model, (Time-Fading model), mines frequent itemsets in stream data in which each transaction has a weight and this weight decreases with age.

- Older transactions contribute less weight toward itemset frequencies. Different weights for new and old transactions.
- suitable for applications in which old data has an effect on the mining results, but the effect decreases as time goes on.

Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010



sliding windows. Only part of the data streams within the sliding window are stored and processed at the time when the data flows in.

- The size of the sliding window may be decided according to applications and system resources.
- The mining result of the sliding window method totally depends on recently generated transactions in the range of the window;
- all the transactions in the window need to be maintained in order to remove their effects on the current mining results when they are out of range of the sliding window.

 (\mathbf{O}) $(\mathbf{0})$ Memory management Memory management [Nan Jiang and Le Gruenwald] Classical association rule mining algorithms on static data collect An efficient and compact data structure is needed to store, the count information for all itemsets and discard the nonupdate and retrieve the collected information: frequent itemsets and their count information after multiple bounded memory size and scans of the database. when we mine association rules in stream data: huge amounts of data streams coming continuously. 1. there is not enough memory space to store all the itemsets and their counts when a huge amount of data comes continuously. Failure in developing such a data structure will largely decrease the efficiency of the mining algorithm 2. the counts of the itemsets are changing with time when new The data structure needs to be incrementally maintained. stream data arrives. Therefore, we need to collect and store the least information possible • it is not possible to rescan the entire input due to the huge Some methods uses sizes of itemsets such as 3 or 2 to generate amount of data and requirement of rapid online querying speed. only this itemsets. asalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010 salvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010

One Pass Algorithm to Generate Association Rules [Nan Jiang and Le Gruenwald]

Association rules can be found in two steps:

asalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010

- 1.finding large itemsets (support is greater than user specified support) for a given threshold support and
- 2. generate desired association rules for a given confidence



Association generation and maintenance [Nan Jiang and Le Gruenwald]

Mining association rules involves a lot of memory and CPU costs.

This is especially a problem in data streams since the processing time is limited to one online scan.

Approaches of the static world: frequent updating

data stream environment, stream data are added continuously, and therefore, if we update association rules too frequently, the cost of computation will increase drastically.

Some methods assume little concept drifting, that is to say the change of data distribution is relatively small.

(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010





WINNEW Winney Winney

(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010



(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010

vas. Spilio

COFI-tree mining M. [El-Hajj and O.R. Zaiane]

FP-trees used in the mining process can all fit in memory. COFI algorithm—as an alternative to the FP-growth algorithm COFI consists of three main phrases:

- the construction of an FP-tree representing the original database,
- The construction of a COFI-tree (Co-Occurrence Frequent Item tree) for each frequent item, and
- the mining of frequent patterns from each COFI-tree.

ulou, Vakali - Barcelona, 24th Sept. 2010

In the first phrase, a global FP-tree is constructed in the same way as in the FP-growth algorithm. Thus 2 database scan are required—one scan for finding the frequency of each item and another scan for building the FP-tree

COFI tree the COFI-tree contains: 1. the item, 2. its frequency count and

3. its participation counter. This counter is initialize to 0, and is incremented every time the node is raversed/participated. At the end of the mining process for the COFI-tree of x, the value of this counter is equal to its frequency count. Note that, the COFI algorithm requires at most two trees (i.e., the global FPtree and the COFI-tree for a specific item) to co-exist at any time during the mining process, whereas FP-growth usually keeps more than two trees.

ona, 24th Sept. 2010









































INIVERSIDADE		(\mathbf{A})	8	
O PORTO	POUTÉCNICK	1.1.1.1	ARIENOTEADO EAMEDIETRIRO	
	Systems I faire	SAP -	TODOLOGIA	

AOG based algorithms

- LightWeight Clustering (LWC): the threshold is used to specify the minimum distance between the cluster center and the data element/record;
- LightWeight Classification (LWClass): In addition of using the threshold in specifying the distance, the class label is checked. If the class label of the stored items and the new item that are similar (within the accepted distance) is the same, the weight of the stored item is increased along with the weighted average of the other attributes, otherwise the weight is decreased and the new item is ignored;
- LightWeight Frequent patterns (LWF): the threshold is used to determine the number of counters for the heavy hitters.



























RA-cluster algorithm RA-Cluster:combines resource-awareness, adaptation and real- time all in a holistic approach. Process:	Repeat Repeat Get next DS record DSRec Find ShortDist which is the shortest distance between DSRec and micro-cluster centers If ShortDist < Radiusthreshold Assign DS record to that micro-cluster Update micro-cluster statistics E create new micro-cluster	Gaber]
 starts with using an initial threshold to run the algorithm and after a fixed time frame the resource consumption patterns of the CPU, memory, and battery given that we run in a resource constrained environment is assessed. 	End Until (END-OF-TIME-FRAME) Calculate NoFMem, NoFCPU, NoFBatt IF NoFMem < RTMem Reclaim outlier memory Increase Radiusthreshold (discourage micro-luster creation) ElseIf if available memory increases Decrease Radiusthreshold (encourage micro-cluster creation)	
According to the above assessment, the algorithm settings are changed to cope with the data rate.	If WorDU < RTCPU Decrease randomization factor (less processing per unit) ElseIf unused CPU power increases Increase randomization factor (more processing per unit)	
RA-Cluster is an incremental online micro-clustering algorithm that has all the required parameters to enable resource- awareness	End If NoFBatt < RTBatt Decrease sampling rate (slower consumption pattern) Elself remaining battery life increases Increase sampling rate (faster consumption pattern) End Until (END-OF-STREAM)	
(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010 (309)	(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010	(310)





















References (Block 5) (4/5)

Gaber, M, M., Krishnaswamy, S., and Zaslavsky, A., Resource–Aware Mining of Data Streams, Journal of Universal Computer Science, Vol. 11, No. 8 (2005), pp. 1440–1453, ISSN 0948–695x, Special Issue on Knowledge Discovery in Data Streams, Jesus S. Aguilar–Ruiz and Joao Gama (Eds.), Verlag der Technischen Universität Graz, Know–Center Graz, Austria, August 2005.

Hamed Chok, Le Gruenwald: Spatio-temporal association rule mining framework for real-time sensor network applications. CIKM 2009: 1761-1764

Nan Jiang, Le Gruenwald: CFI-Stream: mining closed frequent itemsets in data streams. KDD 2006: $592{-}597$

Wei-Guang Teng, Ming-Syan Chen, and Philip S. Yu; Resource-Aware Mining with Variable Granularities in Data Streams; SIAM Int'l Conf. on Data Mining; 2004

(c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010







 $(\mathbf{0})$ (\mathbf{O}) Outlook II - Old challenges, not yet solved Outlook II - Old challenges, not yet solved • Dealing with time: Multi-horizon and multi-granularity analysis Supervised learning for evolving data Scalability for large data volumes Dealing with emerging and evolving concepts Decrease complexity Delayed labeling Increase parallelism Label availability Robustness, esp. if the learners are complex Cost-benefit trade-off of the model update Learning for online or realtime applications Change description Visualization of WHAT ? The streams, the objects, the models Visualization Predicting re–occurring contexts Visualization for WHOM? Multi-label prediction The expert, the data owner, the decision-maker, the casual observer Reliability and uncertainty Evaluation: Measures & Benchmarks (c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010 (c) Gama, Menasalvas, Spiliopoulou, Vakali - Barcelona, 24th Sept. 2010





