#### GUIDING USER GROUPINGS

#### LEARNING AND COMBINING CLASSIFICATION FOR ITEMSET STRUCTURING

#### Mathias Verbeke, Ilija Subašić and Bettina Berendt



MUSE September 24, 2012, Bristol





Delicious C. comiskdd		BibSonomy - search :: Reliter even	search	EN DE
		THE BLUE SOCIAL BOORMARK AND PUBLICATION SHARING SYSTEM.	Sector Se	: password sign in
Sign in to search your links		home groups popular +		sign in
Links	48 Results - view all	BOOKMARKS Ø	PUBLICATIONS .	@ C18 & & . ??
ECML PKDD 2009 13 saves http://www.ecmlpkidd2009.net/ conference 2009 datamining pkdd machinelearning research machine-learning ecml knowledge-discovery conferences	¢ 80 0	R Tutorial: Tables R Tutorial: Tables R Tutorial: http://www.cyclamo.org/tubrialR/tables.htmil/manipu an hour and 53 minutes ago by provisen R tutorials	A wearable sensor for unobtrusive, long-term <sup>2</sup> M.Z. Poh, N.C. Swenson, and R.W. Picard Bornedical Engineerin 8 hours and 18 minutes ago by yish education epistemic learning neuroscience ped	news  Release 2.0.28 on Sep 5, 2012  Feature of the Week: BioSonomy TeXlipse Plugin on Aug 29, 2012  Feature of the week: Entry Type Pictograms on Aug 6, 2012
ECML PKDD Discovery Challenge 2009 10 saves http://www.kde.cs.uni-kassel.da/ws/dc00/ folksonomy research tags competition conference machinelearning challenge kdd recommendation dataset	ф В <sub>0</sub>	Quick-R: Boxplots Boxplots and variations - http://statmenods.net/graphs/boxplot.html an hour and 53 minutes ago by preveen R plots	Characterizing bird migration phenology usin <sup>1</sup> E Routsen, A Linden, T Ergon, N Janzén, JO VR, J Knepe, JE R 10 hours ago by peter raiph artival_times birds climate_charge phenology	<ul> <li>busy tags         <ul> <li>(alpha   freq) (cloud   list)</li> <li>2012 analysis at blog took tooks</li> <li>computer conference data database</li> <li>dbip design education evaluation</li> </ul> </li> </ul>
ECML PKDD 2011 8 saves http://www.ecmlpkdd2011.org/ conference 2011 machinelearning learning pkdd data ECML machine conferences mining	÷ •	Guick-R: Correlations 1 Person, Spearman, Kandal, Polyserial, Polychoric Correlations an hour and 54 minutes ago by proveen R correlations	Regression quantiles * * * * * * * * * * * * * * * * * * *	jova knowledge language learning linux management woteing MyOWN network news online ontology ostimization privacy processing programming r recommender science search simulation social software systems technology teory tools
Wikis, Biogs, Bookmarking Tools - Mining the Web 2.0 - Workshop at ECML PKDD 200         8 saves http://www.kds.cs.uri-kassel.de/ws/wbbtmine2008/         conference       2008         web2.0       network         workshop       mining         data       analysis         conferences       ECML/PKDD02 Tutorial on Text Mining and Internet Content filtering         7 saves http://www.esi.uem.es/-jmgomez/tutorials/ecmlpkdd02/         textmining       flickr:saves1	÷ • •	toread beland You areth signed in Sign in Help	BibSonomy	user visualization web
resources 1 Home The Tour Sign Up Explore - Upload		Search Search III		
Explore / Tags / bristolSort by: Most recent - Most interestingDristol clusters Dristol clustery goodnessImage: Image: I	Enal Franchus Fr Franchus Fr Franchus Fr Franchus Fr	Sidesher   Image: Sidesher	Applications Applications Applications Applications Aboutools.pdf Aboutools.pdf About Xcode.app Applications Appli	
Price Gauranteed at Priceline ® Priceline.com/San-Diego	From Rich Andrews	n Rich Andrews	Maciltopic » 🕞 Developer » 🕞 Fatras » 🗁 Vicode Index Templates	k k (
Affordable Hotels & Great Service in Bristol. Best Web Ratest www.ChelceHotels.com/Bristol			10 items, 25 G8 available	di.

Delicious C, ecmi-pkdd			Bibs	Sonomy .	search :: fulltert	earch search		Cold and the second		EN DE
Circuit a second second fields			THE BLUE SC	ICIAL BOOKMARK AND PUBLICA	TION SHARING SYSTEM.			X.	: password	sign in
sign in to search your links			home group	s popular <del>+</del>						sign in
Links		48 Results - view a	" 🔷	BOOKMARKS	•	1	PUBLICATIONS	٠	@ C18. 61	a ?
ECML PKDD 2009 13 saves http://www.ecm/pkdd2009.net/ conference 2009 datamining pkdd machinelearning machine-learning ecml knowledge-discovery conferences	research	© ,*		R Tutorial: Tables R Tutorial - http://www.cyclamo.or an hour and 53 minutes ago by R tutorials	phutorialR/tables.html#manips praveen		A wearable sensor for unobt M.Z. Poh, N.C. Swenson, and R.W. Pics 8 hours and 18 minutes ago by yish education epistemic learning neur	nusive, long-term <sup>2</sup> d Dometical Engineerin	news  Release 20.28 on Sep 5, 20  Feature of the Week: BibSon Plugin on Aug 29, 2012  Feature of the week: Entry Ty Aug 6, 2012	12 my TeXlipse pe Pictograms on
ECML PKDD Discovery Challenge 2009 10 saves http://www.kde.cs.uni-kassel.de/ws/dc09/ folksonomy research taos competition conference challenge kdd moormendation dataset	machineleaming	0 J		Quick-R: Boxplots Boxplots and variations - http://stat an hour and 53 minutes ago by	methods retignaphs/boxplot.html praveen	Gil	Characterizing bird migration E Knuber, A Linder, T Ergon, N Jonzé 10 hours ago by peter raigh	phenology usin <sup>1</sup>	E busy tags (alpha   frec) (cloud   list) 2012 analysis at blog enerce da educator	book books ta database evaluation
ECML PKDD 2011 8 saves http://www.ecmonosisticitiesgi conference 2011 machinelearning learning plots	lata ECML	Stru	ıctur	ing is	natu	ral			html inform pe languag http://www.secondine privacy op tig r rect science inch simulit is technolog	rector internet e learning ontology ontology symmender tion social g theory tools
Wikis, Blogs, Bookma         S saves http://www.kd         conference       20         data       analysis         ECML/PKDD02 Tutorial on Text Mining and Internet Content filtering         7 saves http://www.edi.uem.es/-jmgomec/tutorials/ecmlpkd02/         textmining         Internet         Sign Up         Explore         Upload	PKDD 200		course roland You aren't	Kommendar: Totalisusfai Internentari miche Waskange 2 heurs and 6 innenansis ago toy Demokratis Rachtepopularia signed in Sign in Heis Search	Rommey tul Obama	Bib	Sonor	arules in strong		
Explore / Tags / bristol Sort by Most recent - Most interesting Distol Clusters Distol Cluster	Image: Second system       From Earchu         Image: Second system       From Earchu         Image: Second system       Image: Second system         Image: Second system	Fram Farchu	Stdeek		V DEVICES MacBook P FLACES V SEARCH FOR O Today Vestenday Past Week All Images All Movies All Documents	Applications Developer Library lost+found System Users	Aboutools.pdf     About Xcode.app     Applications     Documentation     Examples     Extras     Headers     Library     Makefiles     Platforms     Private     SDKs     Tools     usr	G4BitConversion Core Image CoreAudio Dictionent Kit Java Kernelugging LegacyAPtSurvey Modem Scripts PreferencePanes Xcode Iplates	Carbonindex Cocoalpindex Cocoalavaindex Cocoalavaindex CPlusPlusindex IndexadMe.rtf Install_templates JavaDTpindex JavaWepindex	
Book Your Hotel with Priceline. Best Price Gauranteed at Priceline & Priceline.com/San-Clego Choice Hotels® Bristol Affordable Hotels & Great Service in Bristol. Best Web Rates!	From Skopele	Fram Rich Andrews	From Rich Andrews			🗟 Macilook + 🔛 De	ii veloper + 🔛 Extras + 🔛 Xcode I 10 items, 25 G8 availa	i idex Templates Ne		



#### CONTRIBUTIONS

I. A new DM approach that learns an intensional model of user groupings and uses this to group new items Identify structuring dynamics

2. New divergence measure

3. A study of grouping behaviour in a social bookmarking system























#### ... AND ITS DYNAMICS

#### • Goal: insight in structuring dynamics

Advision of the intellectual structuring process

- Two types of guides:
  - I. own prior structuring
  - 2. structuring of peers

### AT A GLANCE

Combination of 2 data mining tasks:

 Learn model of structuring (classification)
 = intension: set of conditions for an object to belong to a certain class
 (vs. extension: list of objects in class)

2. Use intension or extension to structure new items

A. based on own structuring

B. based on k peers

#### GROUPING GUIDANCE BASIC NOTATION

- U: the set of all users (used symbols: u, v, w)
- T: the set of all time points {0, 1, ..., tmax}, where tmax represents the time at which the last item arrives
- D: the set of all items (used symbol: d)  $D_u^t \subseteq D$ : the set of all  $d \in D$  already considered by  $u \in U$  at  $t \in T$  $d_u^t \in (D \setminus D_u^t)$ : the item assigned to the structure by user u at t

# GROUPINGS AND CLASSIFIERS

- G: (machine-induced) groupings for each user's items
- C: classifiers (i.e. intensions) learned for these groupings
  - OG: Observed Grouping
  - GS: Simulated Grouping, guided by self
  - Gn: Simulated Grouping, guided by *n* peers

# INITIAL CLASSIFIER LEARNING

Goal: determine intensional definitions for the user-generated groupings

Each group is regarded as a class for which a definition needs to be calculated

Definitions used to assign new items to these gorups

= selection of peer guides

= selection of peer guides

Requires divergence measure between groupings of non-identical item sets

= selection of peer guides

Requires divergence measure between groupings of non-identical item sets

... but existing measures require large overlap between sets

= selection of peer guides

Requires divergence measure between groupings of non-identical item sets

... but existing measures require large overlap between sets

Inter-guide measure of diversity:

$$egin{aligned} udiv(u,v) &= rac{1}{2}(\ 1/|G_u^t|\sum_{x\in G_u^t}min_{y\in G_v^t}gdiv(x,y)\ &+ 1/|G_v^t|\sum_{y\in G_v^t}min_{x\in G_u^t}gdiv(y,x) \ ) \end{aligned}$$

#### CLASSIFICATION

Selected classifiers are used to classify the item under consideration

Two cases:

- Self-guided classification
- Peer-guided classification











![](_page_32_Figure_1.jpeg)

![](_page_33_Figure_1.jpeg)

![](_page_34_Figure_1.jpeg)

#### DATASET

#### CiteULike dataset

sampled with p-core subgraphs to overcome sparsity

# users	377
# documents	,400
# tags	12,982
timeframe	01/2009 - 02/2010

## INITIAL GROUPING

Tagging as implicit structuring

- First 7 months to learn initial grouping
- Modularity clustering

## INITIAL GROUPING

Tagging as implicit structuring

- First 7 months to learn initial grouping
- Modularity clustering

Initial classifier learning

High dimensional input space (BoW of abstracts)
 Naive bayes

# SIMULATING GROUPINGS

- Groups represented by language models
- Jensen-Shannon divergence as inter-group divergence  $JS(\Theta_x, \Theta_y) = \frac{1}{2}KL(\Theta_x, \Theta_z) + \frac{1}{2}KL(\Theta_y, \Theta_z).$
- Normalized Mutual Information to compare groupigns

$$NMI(G,G') = H(G) + H(G') - \frac{H(G,G')}{\sqrt{H(G)H(G')}}$$

#### RESULTS

![](_page_39_Figure_1.jpeg)

#### SIMILARITY DISTRIBUTION

#### Similarity between final OG and GS grouping

![](_page_40_Figure_2.jpeg)

Similarity between final OG and G1 grouping

![](_page_40_Figure_4.jpeg)

Similarity between final OG and G5 grouping

![](_page_40_Figure_6.jpeg)

![](_page_40_Figure_7.jpeg)

![](_page_40_Figure_8.jpeg)

Similarity between final OG and G20 grouping

![](_page_40_Figure_10.jpeg)

# CONCLUSIONS

Investigate and simulate collaborative structuring

Learning and combining classifiers for itemset structuring

- New divergence measure
- Tested on social-bookmarking platform for literature management

#### CONCLUSIONS LIMITATIONS

- Observed groupings based on tag assignements
- Simple classifier

... but provides initial insights into grouping behaviour and behaviour of users in social bookmarking systems

### FUTURE WORK

• Applications: (tag) recommendation and social search

Adds new level to individual and social measures

- Regrouping based on peers
- Hybrid measure for itemset structuring

![](_page_44_Picture_0.jpeg)

![](_page_44_Picture_1.jpeg)